



Outline of the lecture course

STATISTICS 2

Winter semester 2022/23

Preliminary version: February 1, 2023

If you find **errors in the outline**, please send a short note
by email to johannes@math.uni-heidelberg.de.

MATHEMATIKON, Im Neuenheimer Feld 205, 69120 Heidelberg
phone: +49 6221 54.14.190 – fax: +49 6221 54.14.101
email: johannes@math.uni-heidelberg.de
webpage: sip.math.uni-heidelberg.de

Table of contents

1	Asymptotic properties of M- and Z-estimators	1
§01	Introduction	1
§02	Consistency	7
§03	Asymptotic normality	10
2	Asymptotic properties of tests	15
§04	Contiguity	15
§05	Local asymptotic normality (LAN)	23
§06	Asymptotic relative efficiency	25
§07	Rank tests	27
§08	Asymptotic power of rank tests	30
3	Nonparametric estimation by projection	33
§09	Review	33
§10	Noisy version of the parameter	35
§11	Orthogonal projection	39
§12	Orthogonal projection estimator	42
§13	Minimax optimal estimation	50
§14	Data-driven estimation	58
4	Nonparametric density estimation	73
§15	Noisy density coefficients	73
§16	Projection density estimator	75
§17	Minimax optimal density estimation	80
§18	Data-driven density estimation	83
5	Nonparametric regression	95
§19	Noisy regression coefficients	95
§20	Projection regression estimator	98
§21	Minimax optimal regression	103
§22	Data-driven regression	107
A	Probability theory	121
§19	Fundamentals	121
§20	Convergence of random variables	123
§21	Conditional expectation	127

Chapter 1

Asymptotic properties of M- and Z-estimators

Asymptotic properties of M- and Z-estimators are presented generalising the minimum contrast approach introduced in the lecture [Statistik 1](#). For a more detailed exposition we refer to the text book [van der Vaart \[1998\]](#).

Overview

§01	Introduction	1
§01 01	Motivation / illustration	1
§01 02	Notation / definition	5
§02	Consistency	7
§03	Asymptotic normality	10
§03 01	Testing procedures	12

§01 Introduction

§01|01 Motivation / illustration

§01.01 **Example (Linear model)**. The dependence of the variation of a real random variable Y_1 (response) on the variation of a random vector $X_1 = (X_{1j})_{j \in \llbracket k \rrbracket}$ in \mathbb{R}^k (explanatory variable) is often described by a linear relationship $\mathbb{E}(Y_1|X_1) = \sum_{j \in \llbracket k \rrbracket} \gamma_j X_{1j} = X_1^t \gamma$ or equivalently $Y_1 = X_1^t \gamma + \varepsilon_1$ where ε_1 is a real random error satisfying $\mathbb{E}(\varepsilon_1|X_1) = 0$. We aim to infer on the unknown parameter of interest $\gamma \in \mathbb{R}^k$ from $n \in \mathbb{N}$ i.i.d. copies (Y_i, X_i) , $i \in \llbracket n \rrbracket$. Writing $Y := (Y_i)_{i \in \llbracket n \rrbracket}$ and $X^t = (X_1 \cdots X_n)$ we have $\mathbb{E}(Y|X) = X\gamma$. Any (measurable) choice

$$\hat{\gamma} \in \arg \inf_{\gamma \in \mathbb{R}^k} \hat{M}_n(\gamma) \quad \text{with} \quad \hat{M}_n(\gamma) := \frac{1}{n} \sum_{i \in \llbracket n \rrbracket} (Y_i - X_i^t \gamma)^2 = \frac{1}{n} \|Y - X\gamma\|^2 \quad (01.01)$$

is called a **Least Squares Estimator (LSE)**, where $\arg \inf$ denotes the subset of vectors in \mathbb{R}^k attaining the function's smallest value. If $X^t X = \sum_{i \in \llbracket n \rrbracket} X_i X_i^t$ is strictly positive definite (hence, invertible) then $\hat{\gamma} = (X^t X)^{-1} X^t Y = \left(\sum_{i \in \llbracket n \rrbracket} X_i X_i^t \right)^{-1} \sum_{i \in \llbracket n \rrbracket} Y_i X_i$ is the unique LSE. Under “usual“ conditions ([Example §20.14](#)) holds $\frac{1}{n} \sum_{i \in \llbracket n \rrbracket} X_i X_i^t \xrightarrow{\mathbb{P}} \mathbb{E}(X_1 X_1^t) =: \Omega$ (LLN). If in addition $\mathbb{E}(\varepsilon_i^2|X_i) = \sigma^2$, then $\frac{1}{\sqrt{n}} \sum_{i \in \llbracket n \rrbracket} \varepsilon_i X_i \xrightarrow{d} N_{(0, \sigma^2 \Omega)}$ (CLT). Applying Slutsky's lemma [§20.10](#) and the continuous mapping theorem [§20.09](#) holds $\sqrt{n}(\hat{\gamma} - \gamma) \xrightarrow{d} N_{(0, \sigma^2 \Omega^{-1})}$ for $\Omega > 0$. Further inference on $\hat{\gamma}$ (hypothesis testing, confidence intervals, etc.) is typically based on this asymptotic result. However, a linear relationship $\mathbb{E}(Y|X) = X\gamma$ is often too restrictive. \square

§01.02 **Example (Generalised linear model)**. Consider a real random variable Y_1 and a random vector X_1 in \mathbb{R}^k obeying $\mathbb{E}(Y_1|X_1) = g(X_1^t \gamma)$ for a known link function $g : \mathbb{R} \rightarrow \mathbb{R}$. We aim to infer on the unknown parameter of interest $\gamma \in \mathbb{R}^k$ from $n \in \mathbb{N}$ i.i.d. copies (Y_i, X_i) , $i \in \llbracket n \rrbracket$.

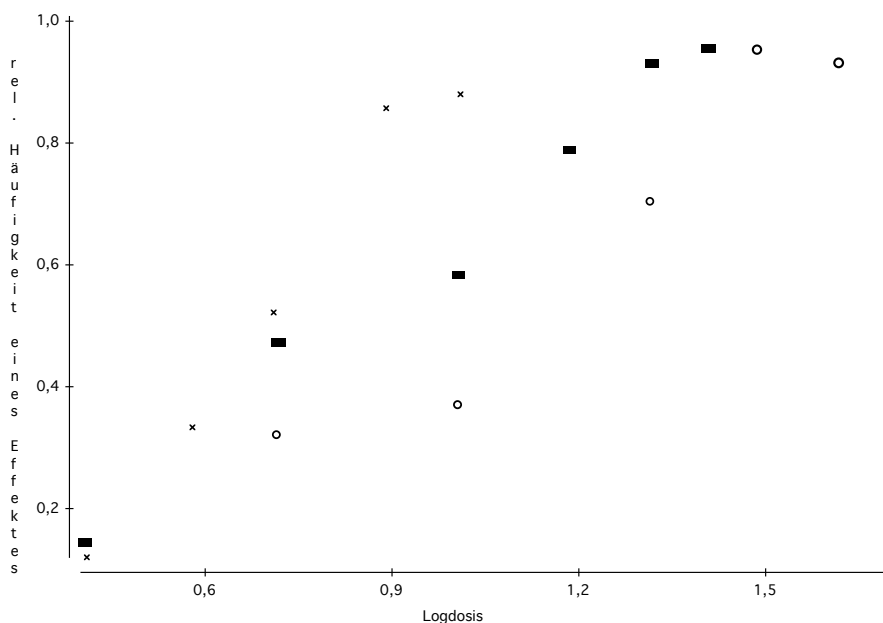
As an illustration let us consider the effect of three different drugs on the behaviour of certain animals. In a trial each drug is given in different dose to certain animals and the number of effected animals is counted. The Table 1.1 summarises the results. Let Y_{jk} denote the counts of an effect among n_{jk} animals applying a log-dose X_{jk} , $j \in \llbracket J_k \rrbracket$ of the drug $k \in \llbracket K \rrbracket$. Assuming an “independent and identical” behaviour of the n_{jk} animals it seems reasonable to model Y_{jk} as Binomial-distributed random variable, $Y_{jk} \sim \text{Bin}(n_{jk}, \pi_{jk})$ for short, with unknown percentage $\pi_{jk} \in (0, 1)$. It may be reasonable to assume that $n_{jk}\pi_{jk} = \mathbb{E}(Y_{jk} | X_{jk}) = g(\gamma_k + \gamma_0 X_{jk})$ where $(\gamma_k)_{k \in \llbracket K \rrbracket}$ is a drug specific factor and γ_0 is a common effect of the log-dose for all drugs. The model is called “probit” and “logit”, respectively, if g is the standard-normal distribution function and the logit-distribution function ($x \mapsto \frac{e^x}{1+e^x}$). As in [Example §01.01](#) inference on $\gamma = (\gamma_k)_{k \in \llbracket 0, K \rrbracket}$ is often based on a LSE, i.e., any (measurable) choice $\hat{\gamma} \in \arg \inf_{\gamma \in \mathbb{R}^{K+1}} \hat{M}_n(\gamma)$ with $\hat{M}_n(\gamma) := \frac{1}{K} \sum_{k \in \llbracket K \rrbracket} \frac{1}{J_k} \sum_{j \in \llbracket J_k \rrbracket} (Y_{jk} - g(\gamma_k + \gamma_0 X_{jk}))^2$.

Table 01 [§01]

drug	log-dose	effect	no effect	drug	log-dose	effect	no effect
1	1.01	44	6	2	1	18	30
1	0.89	42	7	2	0.71	16	33
1	0.71	24	22	3	1.4	48	2
1	0.58	16	32	3	1.31	43	3
1	0.41	6	44	3	1.18	38	10
2	1.7	48	0	3	1	27	19
2	1.61	47	3	3	0.71	22	24
2	1.48	47	2	3	0.4	7	40
2	1.31	34	14				

Number of animals exhibit an (no) effect in dependence of the drug’s log-dose.

Figure 01 [§01]

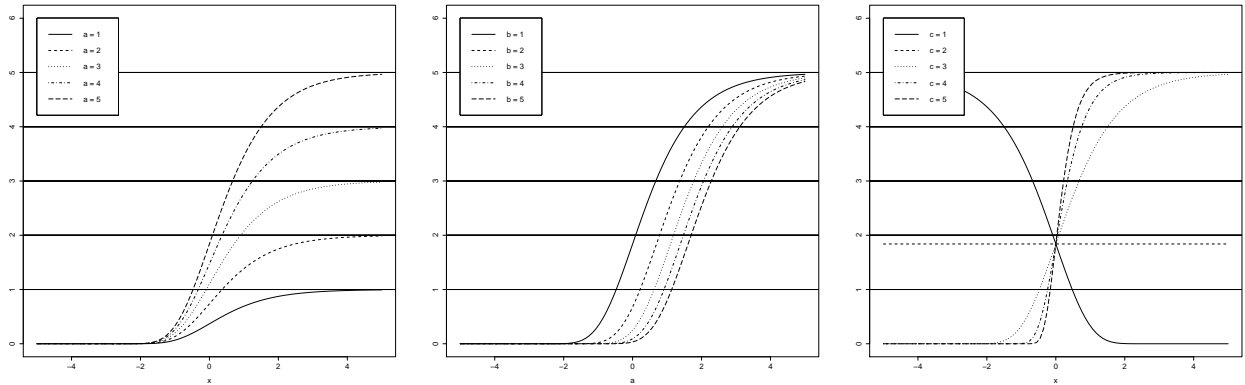


Relative frequency of the effects in dependence of the log-dose, drug 1: x; 2: o; 3: -. □

§01.03 **Example (Nonlinear regression).** Consider a real random variable Y_1 and a random vector X_1 in

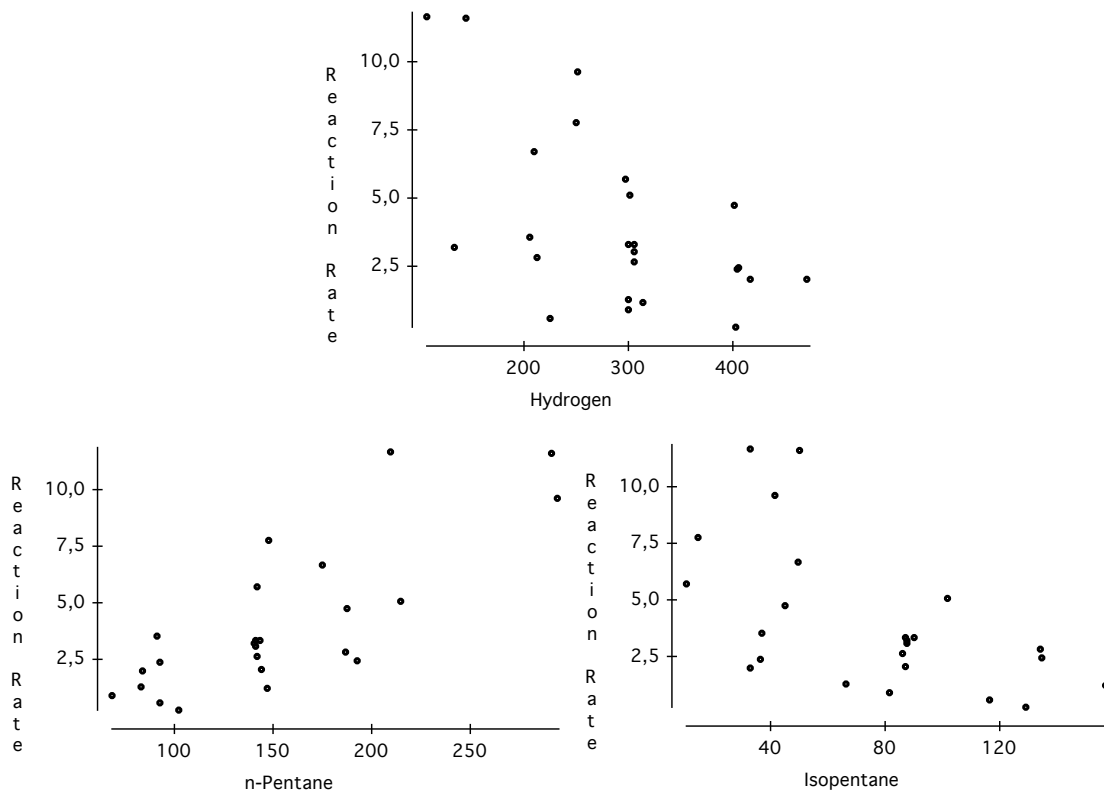
\mathbb{R}^k obeying $\mathbb{E}(Y_1|X_1) = g(X_1, \gamma)$ for a given link function $g : \mathbb{R}^k \times \mathbb{R}^p \rightarrow \mathbb{R}$. We aim to infer on the unknown parameter $\gamma \in \mathbb{R}^p$ from $n \in \mathbb{N}$ i.i.d. copies $(Y_i, X_i), i \in \llbracket n \rrbracket$. The next figure shows the widely used Gompertz function $g(x, (a, b, c)) = a \exp(-b \exp(x \log(c)))$.

Figure 02 [§01]



As an illustration consider the following data of a reaction rate of a catalytic isomerisation of n -pentane into an isopentane given the partial pressure of hydrogen, n -pentane, and isopentane (see Carr [1960]). Isomerisation is a chemical process where a complex chemical product is transformed into basic elements. The reaction rate depends on several factors as for example, the partial pressure and the concentration of a catalyser (hydrogen).

Figure 03 [§01]



Reaction rate in dependence of the partial hydrogen, n -pentane and isopentane pressure.

Table 02 [§01]

Reaction rate				Reaction rate			
Partial pressure				Partial pressure			
hydrogen	n-pentane	isopentane		hydrogen	n-pentane	isopentane	
3,541	205,8	90,9	37,1	5,686	297,3	142,2	10,5
2,397	404,8	92,9	36,3	1,193	314	146,7	157,1
6,694	209,7	174,9	49,4	2,648	305,7	142	86
4,722	401,6	187,2	44,9	3,303	300,1	143,7	90,2
0,593	224,9	92,7	116,3	3,054	305,4	141,1	87,4
0,268	402,6	102,2	128,9	3,302	305,2	141,5	87
2,797	212,7	186,9	134,4	1,271	300,1	83	66,4
2,451	406,2	192,6	134,9	11,648	106,6	209,6	33
3,196	133,3	140,8	87,6	2,002	417,2	83,9	32,9
2,021	470,9	144,2	86,9	9,604	251	294,4	41,5
0,896	300	68,3	81,7	7,754	250,3	148	14,7
5,084	301,6	214,6	101,7	11,59	145,1	291	50,2

Isomerisation reaction rate of an n -pentane into an isopentane.

A commonly used modelling for a reaction rate Y is the Hougen-Watson model where a special case is given by

$$\mathbb{E}(Y_i | (X_{i1}, X_{i2}, X_{i3})) = \frac{\gamma_1 \gamma_3 (X_{i2} - X_{i3} / 1.632)}{1 + \gamma_2 X_{i1} + \gamma_3 X_{i2} + \gamma_4 X_{i3}}, \quad i \in \llbracket n \rrbracket, \tag{01.02}$$

where X_{i1} , X_{i2} and X_{i3} is the partial pressure of hydrogen, isopentane and n -pentane, respectively, and $(\gamma_j)_{j \in \llbracket 4 \rrbracket}$ is the unknown parameter of interest. As in [Example §01.01](#) inference on γ is often based on a LSE, i.e., any (measurable) choice $\hat{\gamma} \in \arg \inf_{\gamma \in \mathbb{R}^4} \hat{M}_n(\gamma)$ with $\hat{M}_n(\gamma) := \frac{1}{n} \sum_{i \in \llbracket n \rrbracket} (Y_i - g(X_i, \gamma))^2$. □

§01.04 **Example (Quantile regression).** Consider a real random variable Y_1 and a random vector X_1 in \mathbb{R}^k obeying $Y_1 = X_1^t \gamma + \varepsilon_1$ with quantile condition $\mathbb{P}(\varepsilon_1 \leq 0 | X_1) = \alpha$ for a given probability $\alpha \in (0, 1)$ or equivalently $\mathbb{P}(Y_1 \leq X_1^t \gamma | X_1) = \alpha$ meaning that the conditional- α -quantile of Y_1 given X_1 equals $X_1^t \gamma$. Let q_α denote the α -quantile of $\mathbb{P}^Z \in \mathcal{W}(\mathcal{B})$, i.e., $\mathbb{P}(Z \leq q_\alpha) = \alpha$. Define $\tau_\alpha(z) := (1 - \alpha)z^- + \alpha z^+$ where $\tau_\alpha(z) = (1 - \alpha)|z|$ if $z \leq 0$ and $\tau_\alpha(z) = \alpha z$ otherwise. Under regularity conditions the function $q \mapsto \mathbb{E}(\tau_\alpha(Z - q))$ attains its minimum at the value $q = q_\alpha$. Roughly, the α -quantile satisfies $0 = \frac{\partial}{\partial q} \mathbb{E}(\tau_\alpha(Z - q)) \Big|_{q=q_\alpha}$, since

$$\begin{aligned} \frac{\partial}{\partial q} \mathbb{E}(\tau_\alpha(Z - q)) &= (1 - \alpha) \frac{\partial}{\partial q} \int_{-\infty}^q (q - z) f(z) dz + \alpha \frac{\partial}{\partial q} \int_q^\infty (z - q) f(z) dz \\ &= (1 - \alpha) \int_{-\infty}^q f(z) dz - \alpha \int_q^\infty f(z) dz \\ &= (1 - \alpha) \mathbb{P}(Z \leq q) - \alpha \mathbb{P}(Z > q) = \mathbb{P}(Z \leq q) - \alpha. \end{aligned}$$

Thereby, a reasonable estimator of γ is any (measurable) choice $\hat{\gamma} \in \arg \inf_{\gamma \in \mathbb{R}^k} \hat{M}_n(\gamma)$ with $\hat{M}_n(\gamma) = \frac{1}{n} \sum_{i \in \llbracket n \rrbracket} \tau_\alpha(Y_i - X_i^t \gamma)$. □

§01.05 **Example (Generalised Method of Moments).** Given a random vector Z_1 in \mathbb{R}^p and a function $h^J = (h_j)_{j \in \llbracket J \rrbracket} : \mathbb{R}^k \times \mathbb{R}^p \rightarrow \mathbb{R}^J$ let the unknown parameter of interest $\gamma \in \mathbb{R}^k$ satisfy $\mathbb{P}^{Z_1} h_j(\gamma) = \mathbb{E}(h_j(\gamma, Z_1)) = 0$ for all $j \in \llbracket J \rrbracket$, or $\mathbb{P}^{Z_1} h^J(\gamma) = \mathbb{E}(h^J(\gamma, Z_1)) = 0$ for short. Supposing an

i.i.d. sample $(Z_i)_{i \in [n]}$ any (measurable) choice $\hat{\gamma}$ satisfying $\hat{\mathbb{P}}_n h_j(\hat{\gamma}) = \frac{1}{n} \sum_{i \in [n]} h_j(\hat{\gamma}, Z_i) = 0$ for all $j \in [J]$, or $\hat{H}_n(\hat{\gamma}) = 0$ with $\hat{H}_n(\gamma) := \hat{\mathbb{P}}_n h^J(\gamma) = \frac{1}{n} \sum_{i \in [n]} h^J(\gamma, Z_i)$, $\gamma \in \mathbb{R}^k$, for short, is called **moment estimator**. In case a moment estimator does not exist, setting $\hat{M}_n(\gamma) := (\hat{\mathbb{P}}_n h^J(\gamma))^t W_n (\hat{\mathbb{P}}_n h^J(\gamma))$ for a given weighting matrix W_n one might consider any (measurable) choice $\hat{\gamma} \in \arg \inf_{\gamma \in \mathbb{R}^k} \hat{M}_n(\gamma)$ called a **Generalised Method of Moments (GMM) estimator**. \square

§01|02 Notation / definition

§01.06 **Reminder**. Denote by $\mathcal{W}(\mathcal{X})$ the set of all probability measures on a measurable space $(\mathcal{X}, \mathcal{X})$. For a non-empty index set Θ a family $\mathbb{P}_\Theta := (\mathbb{P}_\theta)_{\theta \in \Theta}$ of probability measures on \mathcal{X} is formally defined by the map $\Theta \rightarrow \mathcal{W}(\mathcal{X})$ with $\theta \mapsto \mathbb{P}_\theta$. Here and subsequently, for each $\theta \in \Theta$ denotes \mathbb{E}_θ the expectation with respect to \mathbb{P}_θ . For a random variable X taking its values in $(\mathcal{X}, \mathcal{X})$ we write shortly $X \odot \mathbb{P}_\theta$, if $X \sim \mathbb{P}_\theta$ for some $\theta \in \Theta$. If the random variables $(X_i)_{i \in [n]}$ form an *independent and identically distributed* (i.i.d.) sample of $X \sim \mathbb{P}$ with values in $(\mathcal{X}, \mathcal{X})$, then $\mathbb{P}^{\otimes n} = \otimes_{i \in [n]} \mathbb{P}$ denotes the joint product probability measure of the family $(X_i)_{i \in [n]}$ taking its values in the measurable product space $(\mathcal{X}^n, \mathcal{X}^{\otimes n})$. We write $(X_i)_{i \in [n]} \stackrel{i.i.d.}{\sim} \mathbb{P}$ or $(X_i)_{i \in [n]} \sim \mathbb{P}^{\otimes n}$ for short. We denote by $\mathbb{P}_\Theta^{\otimes n} := (\mathbb{P}_\theta^{\otimes n})_{\theta \in \Theta}$ a family of product probability measures on $\mathcal{X}^{\otimes n}$. Any random variable S on $(\mathcal{X}, \mathcal{X})$ taking values in a measurable space $(\mathcal{S}, \mathcal{S})$, i.e., \mathcal{X} - \mathcal{S} -measurable function $S : \mathcal{X} \rightarrow \mathcal{S}$, is called *observation* or *statistic*. We denote by $\mathbb{P}_\Theta^S := (\mathbb{P}_\theta^S)_{\theta \in \Theta}$ the family of probability measures on $(\mathcal{S}, \mathcal{S})$ induced by S . A map $\gamma : \Theta \rightarrow \Gamma$ and its value $\gamma(\theta)$ for each $\theta \in \Theta$ is called *parameter* and *parameter value of interest*, respectively. A *parameter of interest* $\gamma : \Theta \rightarrow \Gamma$ is called *identifiable*, if for any $\theta_1, \theta_2 \in \Theta$ from $\gamma(\theta_1) \neq \gamma(\theta_2)$ follows $\mathbb{P}_{\theta_1} \neq \mathbb{P}_{\theta_2}$. \square

§01.07 **Definition**. The triple $(\mathcal{X}, \mathcal{X}, \mathbb{P}_\Theta)$ is called a *statistical experiment* or *statistical model*. The non-empty set Θ and \mathcal{X} is called *parameter* and *sample space*, respectively. A statistical model $(\mathcal{X}, \mathcal{X}, \mathbb{P}_\Theta)$ is called *adequate* for a random variable X , if $X \odot \mathbb{P}_\Theta$. Given a family $\mathbb{P}_\Theta^{\otimes n}$ of product probability measures $(\mathcal{X}^n, \mathcal{X}^{\otimes n}, \mathbb{P}_\Theta^{\otimes n})$ is called a *statistical product experiment*. We denote by $(\mathcal{S}, \mathcal{S}, \mathbb{P}_\Theta^S)$ the statistical model induced by a $(\mathcal{S}, \mathcal{S})$ -valued statistic S on $(\mathcal{X}, \mathcal{X})$. A statistic $\hat{\gamma}$ on $(\mathcal{X}, \mathcal{X})$ with values in the measurable space (Γ, \mathcal{G}) is called *estimator* or *estimation function* for the identifiable parameter of interest γ . A statistical model $(\mathcal{X}, \mathcal{X}, \mathbb{P}_\Theta)$ (and the family \mathbb{P}_Θ) is called *dominated*, if a σ -finite measure μ on \mathcal{X} exists, $\mu \in \mathcal{M}_\sigma(\mathcal{X})$ for short, such that for each $\theta \in \Theta$ the probability measure \mathbb{P}_θ is absolutely continuous with respect to μ , i.e., $\mathbb{P}_\theta \ll \mu$. We write shortly $\mathbb{P}_\Theta \ll \mu$. Any version of the Radon-Nikodym densities

$$L(\theta, x) := \frac{d\mathbb{P}_\theta}{d\mu}(x) \quad x \in \mathcal{X}, \theta \in \Theta$$

considered as function of θ parametrised by x is called *likelihood* or *likelihood function* where typically it is understand as a random function $L : \Theta \rightarrow \overline{\mathcal{X}^+}$ with $\theta \mapsto L(\theta) := L(\theta, \bullet)$. Its logarithm $\ell := \log L$ (with convention $\log(0) := -\infty$) is called *log-likelihood* or *log-likelihood function*. The *likelihood* and *log-likelihood* in the corresponding dominated product experiment $(\mathcal{X}^n, \mathcal{X}^{\otimes n}, \mathbb{P}_\Theta^{\otimes n})$ are $\prod_{i \in [n]} L(\theta, x_i)$ and $\sum_{i \in [n]} \ell(\theta, x_i)$, $\theta \in \Theta$, $x^n \in \mathcal{X}^n$, respectively. \square

§01.08 **Reminder**. Let $(\mathcal{X}, \mathcal{X}, \mathbb{P}_\Theta)$ be dominated by $\mu \in \mathcal{M}_\sigma(\mathcal{X})$. If μ is finite, then $\mu \ll \mathbb{P}_\mu := \frac{1}{\mu(\mathcal{X})} \mu \in \mathcal{W}(\mathcal{X})$ and hence \mathbb{P}_Θ is also dominated by \mathbb{P}_μ . If μ is not finite, then there exists a countable and measurable partition $\{\mathcal{X}_m, m \in \mathbb{N}\}$ of \mathcal{X} with $\mu(\mathcal{X}_m) \in \mathbb{R}_0^+$ for all $m \in \mathbb{N}$. For each $m \in \mathbb{N}$ define $\mathbb{P}_\mu(\bullet | \mathcal{X}_m) \in \mathcal{W}(\mathcal{X})$ with $A \mapsto \mathbb{P}_\mu(A | \mathcal{X}_m) := \frac{\mu(A \cap \mathcal{X}_m)}{\mu(\mathcal{X}_m)}$. Then we have $\mu \ll \mathbb{P}_\mu := \sum_{m \in \mathbb{N}} 2^{-m} \mathbb{P}_\mu(\bullet | \mathcal{X}_m) \in \mathcal{W}(\mathcal{X})$, since $\mathbb{P}_\mu(A) = 0$ implies $\mu(A \cap \mathcal{X}_m) = 0$ for all $m \in \mathbb{N}$ and thus

$\mu(A) = 0$. Therewith, we have shown, that for each $\mu \in \mathcal{M}_\sigma(\mathcal{X})$ there is $\mathbb{P}_\mu \in \mathcal{W}(\mathcal{X})$ with $\mu \ll \mathbb{P}_\mu$ which automatically dominates \mathbb{P}_θ too. On the other hand, there is a probability measure $\mathbb{P}_0 = \sum_{i \in \mathbb{N}} c_i \mathbb{P}_i$ with $c_i \in \mathbb{R}^+$, $\theta_i \in \Theta$ for all $i \in \mathbb{N}$ and $\sum_{i \in \mathbb{N}} c_i = 1$, and thus $\mathbb{P}_0 \ll \mu$, such that $\mathbb{P}_\theta \ll \mathbb{P}_0$ for all $\theta \in \Theta$ (e.g. **Statistik 1, Satz §08.04**). We call any such probability measure \mathbb{P}_0 *privileged dominating measure*. Therefore, we eventually assume with out loss of generality that the dominating measure is indeed a probability measure. \square

§01.09 **Example (MLE)**. Let $(\mathcal{X}, \mathcal{X}, \mathbb{P}_\theta)$ be a statistical model dominated by $\mu \in \mathcal{M}_\sigma(\mathcal{X})$ with likelihood $L(\theta) = d\mathbb{P}_\theta/d\mu$ and log-likelihood $\ell(\theta) = \log L(\theta)$ for $\theta \in \Theta$ and let (Θ, \mathcal{T}) be a measurable space. Any statistic $\hat{\theta}$ on $(\mathcal{X}, \mathcal{X})$ with values in (Θ, \mathcal{T}) is called **Maximum-Likelihood-Estimator (MLE)** for θ , if $L(\hat{\theta}) = \sup_{\theta \in \Theta} L(\theta)$ μ -a.s. meaning $L(\hat{\theta}(x), x) = \sup_{\theta \in \Theta} L(\theta, x)$ for μ -a.e. $x \in \mathcal{X}$, or equivalently $\ell(\hat{\theta}) = \sup_{\theta \in \Theta} \ell(\theta)$ μ -a.s.. Considering a statistical product experiment $(\mathcal{X}^n, \mathcal{X}^{\otimes n}, \mathbb{P}_\theta^{\otimes n})$ dominated by $\mu^{\otimes n} \in \mathcal{M}_\sigma(\mathcal{X}^{\otimes n})$ and setting $\hat{M}_n(\theta) := -\hat{P}_n \ell(\theta)$, i.e. $\hat{M}_n(\theta, x^n) = -\frac{1}{n} \sum_{i \in [n]} \ell(\theta, x_i)$ for $x^n \in \mathcal{X}^n$, the MLE $\hat{\theta}$ is determined by $\hat{\theta} \in \arg \inf_{\theta \in \Theta} \hat{M}_n(\theta)$ μ -a.s.. However, in general it is not guaranteed that MLE is unique or even exits. The MLE depends on the version of the likelihood, but there exists often a canonical choice. Furthermore, $\gamma(\hat{\theta})$ is called MLE for a parameter of interest $\gamma : \Theta \rightarrow \Gamma$, if $\gamma(\hat{\theta})$ is a statistic on $(\mathcal{X}^n, \mathcal{X}^{\otimes n})$ with values in (Γ, \mathcal{G}) . \square

§01.10 **Remark**. In all the examples the estimator $\hat{\gamma}$ of the parameter of interest γ is determined by $\hat{\gamma} \in \arg \inf_{\gamma \in \Gamma} \hat{M}_n(\gamma)$ for some random function $\gamma \mapsto \hat{M}_n(\gamma) \in \overline{\mathcal{X}}$ of the data. Obviously, rather than minimising (or maximising) a criterion function we might search for a zero of the associated normal or estimating equations, that is, $\hat{\gamma}$ is determined as a zero of a random vector function $\gamma \mapsto \hat{H}_n(\gamma) \in \overline{\mathcal{X}^k}$. Note that estimator is defined \mathbb{P}_θ -a.s. only, meaning that one can change the estimator on a \mathbb{P}_θ -zero set N , i.e., $\mathbb{P}_\theta(N) = 0$ for all $\theta \in \Theta$. \square

§01.11 **Definition**. Let $(\mathcal{X}_n, \mathcal{X}_n, \mathbb{P}_\theta^n = (\mathbb{P}_\theta^n)_{\theta \in \Theta})$ for all $n \in \mathbb{N}$ be a statistical model over the same parameter space Θ and let $\gamma : \Theta \rightarrow \Gamma$ be a parameter of interest. We call a function $M : \Theta \times \Gamma \rightarrow \overline{\mathbb{R}}$ and $H : \Theta \times \Gamma \rightarrow \overline{\mathbb{R}^k}$ *criterion function*, if for all $\theta \in \Theta$ the function $M(\theta) : \gamma \mapsto M(\theta, \gamma)$, respectively $H(\theta) : \gamma \mapsto H(\theta, \gamma)$, has in $\gamma(\theta)$ an unique minimum, respectively an unique zero. A sequence $(\hat{M}_n)_{n \in \mathbb{N}}$ and $(\hat{H}_n)_{n \in \mathbb{N}}$ of functions $\hat{M}_n : \Gamma \times \mathcal{X}_n \rightarrow \overline{\mathbb{R}}$ and $\hat{H}_n : \Gamma \times \mathcal{X}_n \rightarrow \overline{\mathbb{R}^k}$ is called *random criterion function* or *criterion process*, if the following two conditions are satisfied:

(CP1) For all $\gamma \in \Gamma$ is $\hat{M}_n(\gamma) : x \mapsto \hat{M}_n(\gamma, x)$, respectively $\hat{H}_n(\gamma) : x \mapsto \hat{H}_n(\gamma, x)$, a statistic, that is, $\hat{M}_n(\gamma) \in \overline{\mathcal{X}_n}$, respectively $\hat{H}_n(\gamma) \in \overline{\mathcal{X}_n^k}$.

(CP2) For all $\gamma \in \Gamma$ and $\theta \in \Theta$ it holds $\hat{M}_n(\gamma) \xrightarrow{\mathbb{P}_\theta^n} M(\theta, \gamma)$, respectively $\hat{H}_n(\gamma) \xrightarrow{\mathbb{P}_\theta^n} H(\theta, \gamma)$.

Every (measurable) choice $\hat{\gamma}_n : \mathcal{X}_n \rightarrow \Gamma$ (if it exists) is called a *M-estimator*, respectively a *Z-estimator*, if it satisfies

$$\hat{M}_n(\hat{\gamma}_n) = \inf_{\gamma \in \Gamma} \hat{M}_n(\gamma) \quad \mathbb{P}_\theta^n\text{-a.s.}, \quad \text{respectively} \quad \hat{H}_n(\hat{\gamma}_n) = 0 \quad \mathbb{P}_\theta^n\text{-a.s.},$$

or more generally, if it is, respectively, a near minimum and near zero, that is, $\hat{M}_n(\hat{\gamma}_n) \leq \inf_{\gamma \in \Gamma} \hat{M}_n(\gamma) + o_{\mathbb{P}^n}(1)$ and $\hat{H}_n(\hat{\gamma}_n) = o_{\mathbb{P}^n}(1)$. \square

§01.12 **Remark**. There exists a measurable version of a minimum of an almost surely continuous function on a compact set (see Witting and Müller-Funk [1995], Satz 6.7). Note that in **Definition §01.11** the criterion process \hat{M}_n (respectively \hat{H}_n) is defined for each $n \in \mathbb{N}$ on a different measurable space. We write, however, shortly $\hat{M}_n(\gamma) \xrightarrow{\mathbb{P}_\theta^n} M(\theta, \gamma)$, if for each $\varepsilon \in \mathbb{R}_0^+$ holds $\mathbb{P}_\theta^n(|\hat{M}_n(\gamma) - M(\theta, \gamma)| > \varepsilon) \xrightarrow{n \rightarrow \infty} 0$. Let us briefly consider a sample $(X_i)_{i \in [n]} \odot \mathbb{P}_\theta^{\otimes n}$ of a

random variable $X \odot \mathbb{P}_\theta$. Keeping **Notation** §19.05 in mind $\mathbb{P}f$ and $\widehat{\mathbb{P}}_n f$ denotes the integral of $f \in \mathcal{L}_1(\mathcal{X}, \mathbb{P})$ with respect to \mathbb{P} and the empirical measure $\widehat{\mathbb{P}}_n(x^n) = \frac{1}{n} \sum_{i \in [n]} \delta_{x_i}$, $x^n \in \mathcal{X}^n$, respectively. Revisiting each of the **Examples** §01.01 to §01.04 there is a function $m : \Gamma \times \mathcal{X} \rightarrow \overline{\mathbb{R}}$ with $m(\gamma) \in \mathcal{L}_1(\mathcal{X}, \mathbb{P})$, $\gamma \in \Gamma$, such that the criterion process \widehat{M}_n and the associated criterion function M is for each $\gamma \in \Gamma$ given by $\widehat{M}_n(\gamma) = \widehat{\mathbb{P}}_n m(\gamma)$, i.e. $\widehat{M}_n(\gamma, x^n) = \frac{1}{n} \sum_{i \in [n]} m(\gamma, x_i)$, $x^n \in \mathcal{X}^n$, and $M(\theta, \gamma) = \mathbb{P}_\theta m(\gamma) = \int_{\mathcal{X}} m(\gamma, x) \mathbb{P}_\theta(dx)$, respectively. Analogously, a moment estimator as in **Example** §01.05 is a Z -estimator. By construction in each example is the condition **(CP1)** and with the help of the LLN (see **Remark** §20.06) also the condition **(CP2)** satisfied. Note that the GMM estimator in **Example** §01.05 is also a M -estimator with criterion process satisfying **(CP1)** and **(CP2)**. \square

§01.13 **Definition.** For two probability measure \mathbb{P}_0 and \mathbb{P}_1 on a measurable space $(\mathcal{X}, \mathcal{X})$ is the function

$$\text{KL}(\mathbb{P}_0|\mathbb{P}_1) = \begin{cases} \mathbb{P}_0 \left(\log \frac{d\mathbb{P}_0}{d\mathbb{P}_1} \right) = \int \log \left(\frac{d\mathbb{P}_0}{d\mathbb{P}_1} \right) d\mathbb{P}_0, & \text{if } \mathbb{P}_0 \ll \mathbb{P}_1, \\ +\infty, & \text{otherwise} \end{cases}$$

called *Kullback-Leibler-divergence* of \mathbb{P}_0 with respect to \mathbb{P}_1 . \square

§01.14 **Reminder.** The Kullback-Leibler-divergence satisfies $\text{KL}(\mathbb{P}_0|\mathbb{P}_1) \geq 0$ as well as $\text{KL}(\mathbb{P}_0|\mathbb{P}_1) = 0$ if and only if $\mathbb{P}_0 = \mathbb{P}_1$, but $\text{KL}(\bullet|\bullet)$ is not symmetric. Moreover, for product measures holds $\text{KL}(\mathbb{P}_{0,1} \otimes \mathbb{P}_{0,2}|\mathbb{P}_{1,1} \otimes \mathbb{P}_{1,2}) = \text{KL}(\mathbb{P}_{0,1}|\mathbb{P}_{1,1}) + \text{KL}(\mathbb{P}_{0,2}|\mathbb{P}_{1,2})$ (e.g. **Statistik 1, Lemma** §17.03). \square

§01.15 **Example (MLE, §01.09 continued).** Let $(\mathcal{X}^n, \mathcal{X}^{\otimes n}, \mathbb{P}_\theta^{\otimes n})$ be a statistical product experiment dominated by a privileged measure $\mathbb{P}_\theta \in \mathcal{W}(\mathcal{X})$ (see **Reminder** §01.08) with likelihood $L(\theta) = d\mathbb{P}_\theta/d\mathbb{P}_\theta$, log-likelihood $\ell = \log(L)$ and parameter of interest θ (i.e., $\gamma = \text{id}_\Theta$). Furthermore, for all $\theta, \theta_o \in \Theta$ let \mathbb{P}_θ and \mathbb{P}_{θ_o} be mutually dominated (i.e. $\mathbb{P}_\theta \ll \mathbb{P}_{\theta_o}$ and $\mathbb{P}_{\theta_o} \ll \mathbb{P}_\theta$, for short $\mathbb{P}_\theta \ll\!\!\!\ll \mathbb{P}_{\theta_o}$), which implies $\mathbb{P}_{\theta_o} \ll\!\!\!\ll \mathbb{P}_\theta$, and hence $-\text{KL}(\mathbb{P}_{\theta_o}|\mathbb{P}_\theta) = \text{KL}(\mathbb{P}_\theta|\mathbb{P}_{\theta_o})$. Then $\widehat{M}_n(\theta) := -\widehat{\mathbb{P}}_n \ell(\theta) \in \overline{\mathcal{X}^{\otimes n}}$ with

$$x^n \mapsto \widehat{M}_n(\theta, x^n) = -\frac{1}{n} \sum_{i \in [n]} \ell(\theta, x_i)$$

is a criterion process associated to the criterion function $M(\theta_o, \theta) := \text{KL}(\mathbb{P}_\theta|\mathbb{P}_{\theta_o}) - \text{KL}(\mathbb{P}_{\theta_o}|\mathbb{P}_\theta)$ assuming here and subsequently that the parameter θ is identifiable, that is, from $\mathbb{P}_{\theta_1} = \mathbb{P}_{\theta_2}$ follows $\theta_1 = \theta_2$. Identifiability is a natural condition since it is a necessary condition for the existence of a consistent estimator. However, if θ is identifiable then $\theta \mapsto M(\theta_o, \theta)$ attains its minimum $M(\theta_o, \theta_o) = -\text{KL}(\mathbb{P}_{\theta_o}|\mathbb{P}_{\theta_o})$ uniquely at θ_o (keeping **Reminder** §01.14 in mind). The corresponding M -estimator is thus just a MLE. \square

§02 Consistency

Here and subsequently, let (Γ, d) be a metric space endowed with its Borel- σ -algebra $\mathcal{G} := \mathcal{B}_\Gamma$, let $(\mathcal{X}_n, \mathcal{X}_n, \mathbb{P}_\theta^n = (\mathbb{P}_\theta^n)_{\theta \in \Theta})$ for all $n \in \mathbb{N}$ be a statistical model over the parameter space Θ and let $\gamma : \Theta \rightarrow \Gamma$ be an identifiable parameter of interest.

§02.01 **Reminder.** For each $n \in \mathbb{N}$ let $\widehat{\gamma}_n$ be an estimator of γ , i.e. a statistic on $(\mathcal{X}_n, \mathcal{X}_n)$ with values in (Γ, \mathcal{G}) . The sequence $(\widehat{\gamma}_n)_{n \in \mathbb{N}}$ of estimators is called *(weakly) consistent*, if for all $\varepsilon \in \mathbb{R}_0^+$ holds $\mathbb{P}_\theta^n(d(\widehat{\gamma}_n, \gamma(\theta)) > \varepsilon) = o(1)$ as $n \rightarrow \infty$ for all $\theta \in \Theta$. Note that the estimator $\widehat{\gamma}_n$ can be defined for each $n \in \mathbb{N}$ on a different measurable space. We write, however, shortly $d(\widehat{\gamma}_n, \gamma(\theta)) = o_{\mathbb{P}_\theta^n}(1)$ as $n \rightarrow \infty$. Moreover, saying „ $\widehat{\gamma}_n$ is consistent“ always means the sequence $(\widehat{\gamma}_n)_{n \in \mathbb{N}}$ is (weakly) consistent. \square

Consider an M-estimator $\hat{\gamma}_n$ for a random criterion function \hat{M}_n with associated criterion function M , that is, $\hat{M}_n(\gamma) \xrightarrow{\mathbb{P}^n} M(\theta, \gamma)$ holds point-wise for each $\gamma \in \Gamma$. For example, due to the LLN $\hat{M}_n(\gamma) = \hat{\mathbb{P}}_n m(\gamma) \xrightarrow{\mathbb{P}^{\otimes n}} \mathbb{E}_\theta m(\gamma) = M(\theta, \gamma)$ provided $m(\gamma) \in \mathcal{L}_1(\mathcal{X}, \mathbb{P}_\theta)$. The hope is that a minimising value of $\hat{M}_n(\gamma)$ then converges to the minimising value of $M(\theta, \gamma)$. However, in general point-wise convergence will not be sufficient.

§02.02 **Theorem.** Under the assumptions and notations of *Definition §01.11* any *M-estimator* $\hat{\gamma}_n$ of γ , i.e., $\hat{M}_n(\hat{\gamma}_n) \leq \hat{M}_n(\gamma(\theta)) + o_{\mathbb{P}^n}(1)$, is *consistent*, i.e., $d(\hat{\gamma}_n, \gamma(\theta)) = o_{\mathbb{P}^n}(1)$, if in addition the following two conditions are satisfied:

$$(CO1) \sup_{\gamma \in \Gamma} |\hat{M}_n(\gamma) - M(\theta, \gamma)| = o_{\mathbb{P}^n}(1) \quad (\text{uniform convergence in probability});$$

$$(CO2) \inf_{\gamma \in \Gamma: d(\gamma, \gamma(\theta)) \geq \varepsilon} M(\theta, \gamma) > M(\theta, \gamma(\theta)) \text{ for any } \varepsilon \in \mathbb{R}_{>0}^+ \quad (\text{identification}).$$

§02.03 **Proof of Theorem §02.02.** is given in the lecture. □

§02.04 **Corollary.** Under the assumptions and notations of *Definition §01.11* any *Z-estimator* $\hat{\gamma}_n$ of γ , i.e., $\hat{H}_n(\hat{\gamma}_n) = o_{\mathbb{P}^n}(1)$, is *consistent*, i.e., $d(\hat{\gamma}_n, \gamma(\theta)) = o_{\mathbb{P}^n}(1)$, if in addition the following two conditions are satisfied:

$$(CO1) \sup_{\gamma \in \Gamma} \|\hat{H}_n(\gamma) - H(\theta, \gamma)\| = o_{\mathbb{P}^n}(1) \quad (\text{uniform convergence in probability});$$

$$(CO2) \inf_{\gamma \in \Gamma: d(\gamma, \gamma(\theta)) \geq \varepsilon} \|H(\theta, \gamma)\| > 0 = \|H(\theta, \gamma(\theta))\| \text{ for any } \varepsilon \in \mathbb{R}_{>0}^+ \quad (\text{identification}).$$

§02.05 **Proof of Corollary §02.04.** is given in the lecture. □

§02.06 **Lemma.** If (i) Γ is compact, (ii) $M(\theta, \gamma) > M(\theta, \gamma(\theta))$ for all $\gamma \in \Gamma \setminus \{\gamma(\theta)\}$, and (iii) $\gamma \mapsto M(\theta, \gamma)$ is continuous, then (CO2) in *Theorem §02.02* holds.

§02.07 **Proof of Lemma §02.06.** is left as an exercise. □

§02.08 **Example (MLE, §01.15 continued).** Assuming in addition that the parameter space Θ is compact and that the criterion function $\theta \mapsto M(\theta, \theta) := \text{KL}(\mathbb{P}_\theta | \mathbb{P}_\theta) - \text{KL}(\mathbb{P}_\theta | \mathbb{P}_\theta)$ is continuous then employing *Lemma §02.06* the condition (CO2) of *Theorem §02.02* is satisfied. □

§02.09 **Lemma.** (CO1) in *Theorem §02.02* is satisfied, if the following conditions hold:

- (i) (Γ, d) is a compact metric space,
- (ii) $\gamma \mapsto M(\theta, \gamma)$ is continuous and $\hat{M}_n(\gamma) = M(\theta, \gamma) + o_{\mathbb{P}^n}(1)$ for all $\gamma \in \Gamma$, and
- (iii) $\lim_{\delta \downarrow 0} \limsup_{n \rightarrow \infty} \mathbb{P}_\theta^n \left(\sup_{\gamma_1, \gamma_2 \in \Gamma: d(\gamma_1, \gamma_2) \leq \delta} |\hat{M}_n(\gamma_1) - \hat{M}_n(\gamma_2)| \geq \varepsilon \right) = 0$ for all $\varepsilon \in \mathbb{R}_{>0}^+$.

§02.10 **Proof of Lemma §02.09.** is given in the lecture. □

§02.11 **Example.** Given $(\mathcal{X}^n, \mathcal{X}^{\otimes n}, \mathbb{P}_\theta^{\otimes n})$ and $\gamma : \Theta \rightarrow \Gamma$ for each $\gamma \in \Gamma$ let $m(\gamma) \in \mathcal{X}$ be a real function $x \mapsto m(\gamma, x)$ belonging to $\mathcal{L}_1(\mathcal{X}, \mathbb{P}_\theta)$. Consider $\hat{M}_n(\gamma) := \hat{\mathbb{P}}_n m(\gamma)$, i.e. $\hat{M}_n(\gamma, x^n) = \frac{1}{n} \sum_{i \in [n]} m(\gamma, x_i)$, $x^n \in \mathcal{X}^n$, and $M(\theta, \gamma) := \mathbb{E}_\theta m(\gamma)$ where due to the LLN §20.06 $\hat{M}_n(\gamma) = M(\theta, \gamma) + o_{\mathbb{P}^n}(1)$ for each $\gamma \in \Gamma$. Suppose in addition the following conditions:

- (i) (Γ, d) is a compact metric space,
- (ii) $\gamma \mapsto m(\gamma, x)$ is continuous for \mathbb{P}_θ -a.e. $x \in \mathcal{X}$,
- (iii) there is $H \in \mathcal{L}_1(\mathcal{X}, \mathbb{P}_\theta)$ with $\sup_{\gamma \in \Gamma} |m(\gamma, x)| \leq |H(x)|$ for \mathbb{P}_θ -a.e. $x \in \mathcal{X}$, or equivalently, $\sup_{\gamma \in \Gamma} |m(\gamma)|$ belongs to $\mathcal{L}_1(\mathcal{X}, \mathbb{P}_\theta)$.

Then, (I) $\gamma \mapsto \mathbb{P}_\theta \mathfrak{m}(\gamma) = M(\theta, \gamma)$ is continuous and (CO1) $\sup_{\gamma \in \Gamma} |\widehat{M}_n(\gamma) - M(\theta, \gamma)| = o_{\mathbb{P}^{\otimes n}}(1)$. Indeed, by dominated convergence (see §20.20) (ii) and (iii) imply together (I). Consider (CO1). Define the random variable $\Delta_\delta^n := \sup_{\gamma_1, \gamma_2 \in \Gamma: d(\gamma_1, \gamma_2) \leq \delta} |\widehat{M}_n(\gamma_1) - \widehat{M}_n(\gamma_2)| \in \overline{\mathcal{X}^{\otimes n}}$. We show below for all $\varepsilon, \eta \in \mathbb{R}_0^+$ exists $\delta \in \mathbb{R}_0^+$ with $\limsup_{n \rightarrow \infty} \mathbb{P}_\theta^{\otimes n}(\Delta_\delta^n \geq \varepsilon) \leq \eta$ which in turn by Lemma §02.09 implies the claim (CO1). Let $\varepsilon, \eta \in \mathbb{R}_0^+$. Keeping $\Delta_\delta^1 \in \overline{\mathcal{X}}$ with $x \mapsto \Delta_\delta^1(x) = \sup_{\gamma_1, \gamma_2 \in \Gamma: d(\gamma_1, \gamma_2) \leq \delta} |m(\gamma_1, x) - m(\gamma_2, x)|$ in mind and applying the elementary triangular inequality we have $\Delta_\delta^n \leq \widehat{\mathbb{P}}_n \Delta_\delta^1$ point-wise on \mathcal{X}^n . Moreover, due to (i) and (ii) for \mathbb{P}_θ -a.e. $x \in \mathcal{X}$ the function $\gamma \mapsto m(\gamma, x)$ is uniformly continuous on Γ , and thus $\lim_{\delta \rightarrow 0} \Delta_\delta^1(x) = 0$. Therewith, dominated convergence (see §20.20), which can be applied due to (iii), implies $\lim_{\delta \rightarrow 0} \mathbb{P}_\theta \Delta_\delta^1 = 0$. In particular there is $\delta \in \mathbb{R}_0^+$ such that $\mathbb{P}_\theta \Delta_\delta^1 \leq \eta \varepsilon$, which in turn implies $\mathbb{P}_\theta^{\otimes n} \Delta_\delta^n \leq \mathbb{P}_\theta^{\otimes n}(\widehat{\mathbb{P}}_n \Delta_\delta^1) = \mathbb{P}_\theta \Delta_\delta^1 \leq \eta \varepsilon$. Employing Markov's inequality §20.18 the last estimate implies the claim, that is, for all $\varepsilon, \eta \in \mathbb{R}_0^+$ exists $\delta \in \mathbb{R}_0^+$ with $\limsup_{n \rightarrow \infty} \mathbb{P}_\theta^{\otimes n}(\Delta_\delta^n \geq \varepsilon) \leq \eta$. If in addition to (i)-(iii) and, hence (I)

(iv) there is $\gamma(\theta) \in \Gamma$ with $M(\theta, \gamma) > M(\theta, \gamma(\theta))$ for all $\gamma \in \Gamma \setminus \{\gamma(\theta)\}$,

then applying Lemma §02.06 it holds (CO2) $\inf_{\gamma \in \Gamma: d(\gamma, \gamma(\theta)) \geq \varepsilon} M(\theta, \gamma) > M(\theta, \gamma(\theta))$. To summarise, with (CO1) and (CO2) the conditions of Theorem §02.02 are satisfied. Consequently, any *M-estimator* $\widehat{\gamma}_n$, i.e., $\widehat{M}_n(\widehat{\gamma}_n) \leq \inf_{\gamma \in \Gamma} \widehat{M}_n(\gamma) + o_{\mathbb{P}^{\otimes n}}(1)$, and thus $\widehat{M}_n(\widehat{\gamma}_n) \leq \widehat{M}_n(\gamma(\theta)) + o_{\mathbb{P}^{\otimes n}}(1)$, is a *consistent estimator of γ* , i.e., $d(\widehat{\gamma}_n, \gamma(\theta)) = o_{\mathbb{P}^{\otimes n}}(1)$. \square

§02.12 **Lemma.** (CO1) in Corollary §02.04 is satisfied, if the following conditions hold:

- (i) (Γ, d) is a compact metric space,
- (ii) $\gamma \mapsto H(\theta, \gamma)$ is continuous and $\|\widehat{H}_n(\gamma) - H(\theta, \gamma)\| = o_{\mathbb{P}^n}(1)$ for all $\gamma \in \Gamma$, and
- (iii) $\lim_{\delta \downarrow 0} \limsup_{n \rightarrow \infty} \mathbb{P}_\theta^n \left(\sup_{\gamma_1, \gamma_2 \in \Gamma: d(\gamma_1, \gamma_2) \leq \delta} \|\widehat{H}_n(\gamma_1) - \widehat{H}_n(\gamma_2)\| \geq \varepsilon \right) = 0$ for all $\varepsilon \in \mathbb{R}_0^+$.

§02.13 **Proof of Lemma §02.12.** is left as an exercise. \square

§02.14 **Example.** Given $(\mathcal{X}^n, \overline{\mathcal{X}^{\otimes n}}, \mathbb{P}_\theta^{\otimes n})$, $\gamma : \Theta \rightarrow \Gamma$ and $(X_i)_{i \in [n]} \sim \mathbb{P}_\theta^{\otimes n}$ for $\theta \in \Theta$, for each $\gamma \in \Gamma$ let $h(\gamma) \in \overline{\mathcal{X}^k}$ be a numerical function belonging to $\mathcal{L}_1^k(\mathbb{P}_\theta)$ for all $\gamma \in \Gamma$. Consider $\widehat{H}_n(\gamma) := \widehat{\mathbb{P}}_n h(\gamma)$, i.e. $\widehat{H}_n(\gamma, x^n) = \frac{1}{n} \sum_{i \in [n]} h(\gamma, x_i)$, $x^n \in \mathcal{X}^n$, and $H(\theta, \gamma) := \mathbb{P}_\theta h(\gamma)$ where due to the LLN §20.06 $\|\widehat{H}_n(\gamma) - H(\theta, \gamma)\| = o_{\mathbb{P}^{\otimes n}}(1)$ for each $\gamma \in \Gamma$. Suppose in addition the following conditions:

- (i) (Γ, d) is a compact metric space,
- (ii) $\gamma \mapsto h(\gamma, x)$ is continuous for \mathbb{P}_θ -a.e. $x \in \mathcal{X}$,
- (iii) $\sup_{\gamma \in \Gamma} \|h(\gamma)\|$ belongs to $\mathcal{L}_1(\mathbb{P}_\theta)$.

Then, arguing line by line as in Example §02.11 (I) $\gamma \mapsto \mathbb{P}_\theta h(\gamma) = H(\theta, \gamma)$ is continuous and (CO1) $\sup_{\gamma \in \Gamma} \|\widehat{H}_n(\gamma) - H(\theta, \gamma)\| = o_{\mathbb{P}^{\otimes n}}(1)$. If in addition to (i)-(iii) and hence (I)

(iv) there is $\gamma(\theta) \in \Gamma$ with $\|H(\theta, \gamma)\| > 0 = \|H(\theta, \gamma(\theta))\|$ for all $\gamma \in \Gamma \setminus \{\gamma(\theta)\}$,

then applying Lemma §02.06 it holds (CO2) $\inf_{\gamma \in \Gamma: d(\gamma, \gamma(\theta)) \geq \varepsilon} \|H(\theta, \gamma)\| > 0 = \|H(\theta, \gamma(\theta))\|$. To summarise, with (CO1) and (CO2) the conditions of Corollary §02.04 are satisfied. Consequently, any *Z-estimator* $\widehat{\gamma}_n$, i.e., $\widehat{H}_n(\widehat{\gamma}_n) = o_{\mathbb{P}^{\otimes n}}(1)$ is a *consistent estimator of γ* , i.e., $d(\widehat{\gamma}_n, \gamma(\theta)) = o_{\mathbb{P}^{\otimes n}}(1)$. \square

§02.15 **Remark.** The conditions (CO1) and (CO2) of Corollary §02.04 (respectively, (CO1) and (CO2) of Theorem §02.02) being sufficient to ensure consistency might be weakened in specific situations as we see next. \square

§02.16 **Proposition.** Let $\Gamma \subseteq \mathbb{R}$ and $\widehat{H}_n(\gamma) = H(\theta, \gamma) + o_{\mathbb{P}^n}(1)$ for all $\gamma \in \Gamma$ where H is a deterministic function. Assume in addition that either

(Ia) $\gamma \mapsto \widehat{H}_n(\gamma)$ is continuous and has exactly one zero $\widehat{\gamma}_n$, or

(Ib) $\gamma \mapsto \widehat{H}_n(\gamma)$ is non-decreasing with $\widehat{H}_n(\widehat{\gamma}_n) = o_{\mathbb{P}^n}(1)$,

and that (II) $H(\theta, \gamma(\theta) - \varepsilon) < 0 < H(\theta, \gamma(\theta) + \varepsilon)$ for every $\varepsilon \in \mathbb{R}_{>0}^+$. Then, $\widehat{\gamma}_n = \gamma(\theta) + o_{\mathbb{P}^n}(1)$.

§02.17 **Proof of Proposition §02.16.** is given in the lecture. \square

§02.18 **Example.** Consider $\mathbb{P} \in \mathcal{W}(\mathcal{B})$ and $h(\gamma, t) := \text{sign}(t - \gamma)$ with $\text{sign}(t) := \mathbb{1}_{\{t \geq 0\}} - \mathbb{1}_{\{t < 0\}}$ for all $\gamma, t \in \mathbb{R}$. The sample median $\widehat{\gamma}_n$ is a (near) zero of the map $\gamma \mapsto \widehat{H}_n(\gamma) := \widehat{\mathbb{P}}_n h(\gamma)$, i.e. $\widehat{H}_n(\gamma, x^n) = \frac{1}{n} \sum_{i \in [n]} h(\gamma, x_i)$, $x^n \in \mathbb{R}^n$. Considering $H(\gamma) = \mathbb{P}h(\gamma) = \mathbb{P}((\gamma, \infty)) - \mathbb{P}((-\infty, \gamma))$ we have obviously $\widehat{H}_n(\gamma) = H(\gamma) + o_{\mathbb{P}^n}(1)$ for each $\gamma \in \Gamma$. Keeping in mind that $\gamma \mapsto \widehat{H}_n(\gamma)$ is non-increasing from Proposition §02.16 follows consistency of the sample median $\widehat{\gamma}_n$, i.e., $\widehat{\gamma}_n = \gamma_o + o_{\mathbb{P}^n}(1)$, if for any $\varepsilon \in \mathbb{R}_{>0}^+$ in addition $H(\gamma_o - \varepsilon) > 0 > H(\gamma_o + \varepsilon)$ or equivalently $\mathbb{P}((-\infty, \gamma_o - \varepsilon)) < 1/2 < \mathbb{P}((-\infty, \gamma_o + \varepsilon))$. In other words, the sample median $\widehat{\gamma}_n$ is a consistent estimator of the population median, if it is unique. \square

§03 Asymptotic normality

Here and subsequently, for $k, n \in \mathbb{N}$ let $\Gamma \subseteq \mathbb{R}^k$ be endowed with its Borel- σ -algebra $\mathcal{G} := \mathcal{B}_\Gamma$, let $(\mathcal{X}^n, \mathcal{X}^{\otimes n}, \mathbb{P}_\theta^{\otimes n})$ be a statistical product experiment over the parameter space Θ and let $\gamma : \Theta \rightarrow \Gamma$ be an identifiable parameter of interest.

§03.01 **Heuristics.** Consider $\widehat{H}_n(\gamma) = \widehat{\mathbb{P}}_n h(\gamma)$, i.e. $\widehat{H}_n(\gamma, x^n) = \frac{1}{n} \sum_{i \in [n]} h(\gamma, x_i)$, $x^n \in \mathcal{X}^n$, and $H(\theta, \gamma) = \mathbb{P}_\theta h(\gamma)$ for $\gamma \in \Gamma$ and $\theta \in \Theta$. Let $\widehat{\gamma}_n$ be a zero of $\gamma \mapsto \widehat{H}_n(\gamma)$, i.e., $\widehat{\gamma}_n$ is a Z-estimator. Assume in addition that $\widehat{\gamma}_n = \gamma(\theta) + o_{\mathbb{P}^n}(1)$ where $\gamma(\theta)$ is a zero of $\gamma \mapsto H(\theta, \gamma)$. Heuristically, consider a *Taylor expansion* of a real-valued H around $\gamma(\theta) \in \Gamma \subseteq \mathbb{R}$, that is, $0 = \widehat{H}_n(\widehat{\gamma}_n) = \widehat{H}_n(\gamma(\theta)) + (\widehat{\gamma}_n - \gamma(\theta)) \widehat{H}_n(\gamma(\theta)) + \frac{1}{2} (\widehat{\gamma}_n - \gamma(\theta))^2 \ddot{H}_n(\widetilde{\gamma}_n)$ for some $\widetilde{\gamma}_n$ between $\gamma(\theta)$ and $\widehat{\gamma}_n$. Thus, rewriting the last identity $\sqrt{n}(\widehat{\gamma}_n - \gamma(\theta)) = -\sqrt{n} \widehat{H}_n(\gamma(\theta)) (\widehat{H}_n(\gamma(\theta)) + \frac{1}{2} (\widehat{\gamma}_n - \gamma(\theta)) \ddot{H}_n(\widetilde{\gamma}_n))^{-1}$. If $h(\gamma(\theta))$ belongs to $\mathcal{L}_2(\mathbb{P}_\theta)$, then due to the CLT it holds $-\sqrt{n}(\widehat{H}_n(\gamma(\theta)) - H(\theta, \gamma(\theta))) = -\sqrt{n}(\widehat{\mathbb{P}}_n h(\gamma(\theta)) - \mathbb{P}_\theta h(\gamma(\theta))) \xrightarrow{d} N_{(0, \mathbb{P}_\theta h^2(\gamma(\theta)))}$. If moreover $\dot{h}(\gamma(\theta)) \in \mathcal{L}_1(\mathbb{P}_\theta)$, then by the LLN $\widehat{H}_n(\gamma(\theta)) = \widehat{\mathbb{P}}_n h(\gamma(\theta)) = \mathbb{P}_\theta h(\gamma(\theta)) + o_{\mathbb{P}^n}(1)$. If in addition $\ddot{H}_n(\widetilde{\gamma}_n) = O_{\mathbb{P}^n}(1)$ then employing Slutsky's lemma §20.10 it follows $\sqrt{n}(\widehat{\gamma}_n - \gamma(\theta)) \xrightarrow{d} N_{(0, (\mathbb{P}_\theta h^2(\gamma(\theta)))^{-2} \mathbb{P}_\theta h^2(\gamma(\theta)))}$. In the sequel, γ is a vector and h vector-valued. Consequently, $\dot{h}(\gamma(\theta))$ is a matrix and we denote by $\|\dot{h}(\gamma(\theta))\|_F$ its *Frobenius norm*, where $\|M\|_F := (\sum_{j \in [J]} \sum_{k \in [K]} M_{jk}^2)^{1/2}$ for any matrix $M = (M_{jk}) \in \mathbb{R}^{(J, K)}$. \square

§03.02 **Theorem.** Under the assumptions and notations of Definition §01.11 with $\Gamma \subseteq \mathbb{R}^k$ let $\widehat{\gamma}_n$ be a consistent Z-estimator of γ , i.e. $\widehat{\gamma}_n = \gamma(\theta) + o_{\mathbb{P}^n}(1)$, with $\widehat{H}_n(\widehat{\gamma}_n) = o_{\mathbb{P}^n}(n^{-1/2})$. Assume the criterion process \widehat{H}_n is continuous differentiable in a neighbourhood U of $\gamma(\theta) \in \text{int}(\Gamma)$ with derivative $\dot{\widehat{H}}_n := \frac{\partial}{\partial \gamma} \widehat{H}_n \in \overline{\mathcal{X}}^{(k, k)}$ and satisfies the following two conditions:

(AN1) $\sqrt{n} \widehat{H}_n(\gamma(\theta)) \xrightarrow{d} N_{(0, \Omega_\theta)}$ under $\mathbb{P}_\theta^{\otimes n}$ for some positive semidefinite $\Omega_\theta \in \mathbb{R}^{(k, k)}$,

(AN2) $\sup_{\gamma \in U} \|\widehat{H}_n(\gamma) - \dot{H}(\theta, \gamma)\|_F = o_{\mathbb{P}^n}(1)$ for some continuous matrix-valued function $\gamma \mapsto \dot{H}(\theta, \gamma)$ with regular $\dot{H}(\theta, \gamma(\theta))$ having \dot{H}_θ^{-1} as inverse.

Then $\sqrt{n}(\widehat{\gamma}_n - \gamma(\theta)) + \sqrt{n} \dot{H}_\theta^{-1} \widehat{H}_n(\gamma(\theta)) = o_{\mathbb{P}^n}(1)$ and $\sqrt{n}(\widehat{\gamma}_n - \gamma(\theta)) \xrightarrow{d} N_{(0, \dot{H}_\theta^{-1} \Omega_\theta (\dot{H}_\theta^{-1})^t)}$.

§03.03 **Proof of Theorem §03.02.** is given in the lecture. \square

§03.04 **Corollary.** Under the assumptions and notations of *Definition §01.11* with $\Gamma \subseteq \mathbb{R}^k$ let $\hat{\gamma}_n$ be a consistent M-estimator of γ , i.e. $\hat{\gamma}_n = \gamma(\theta) + o_{\mathbb{P}^{\otimes n}}(1)$, with $\hat{M}_n(\hat{\gamma}_n) = \inf_{\gamma \in \Gamma} \hat{M}_n(\gamma)$. Assume the criterion process \hat{M}_n is twice continuously differentiable in a neighbourhood U of $\gamma(\theta) \in \text{int}(\Gamma)$ with derivatives $\dot{\hat{M}}_n := \frac{\partial}{\partial \gamma} \hat{M}_n \in \overline{\mathcal{X}}^k$ (score function) and $\ddot{\hat{M}}_n := \frac{\partial^2}{\partial^2 \gamma} \hat{M}_n \in \overline{\mathcal{X}}^{(k,k)}$ and satisfies in addition the following two conditions:

(AN1) $\sqrt{n} \dot{\hat{M}}_n(\gamma(\theta)) \xrightarrow{d} N_{(0, \Omega_\theta)}$ under $\mathbb{P}_\theta^{\otimes n}$ for some positive semidefinite $\Omega_\theta \geq 0$,

(AN2) $\sup_{\gamma \in U} \|\ddot{\hat{M}}_n(\gamma) - \ddot{M}(\theta, \gamma)\|_F = o_{\mathbb{P}^{\otimes n}}(1)$ for some continuous matrix-valued function $\gamma \mapsto \ddot{M}(\theta, \gamma)$ with regular $\ddot{M}(\theta, \gamma(\theta))$ having \ddot{M}_θ^{-1} as inverse.

Then $\sqrt{n}(\hat{\gamma}_n - \gamma(\theta)) \xrightarrow{d} N_{(0, \ddot{M}_\theta^{-1} \Omega_\theta \ddot{M}_\theta^{-1})}$.

§03.05 **Proof of Corollary §03.04.** is given in the lecture. \square

§03.06 **Example (§02.11 continued).** Given $(\mathcal{X}^n, \overline{\mathcal{X}}^{\otimes n}, \mathbb{P}_\theta^{\otimes n})$ and $\gamma : \Theta \rightarrow \Gamma$ for each $\gamma \in \Gamma$ let $m(\gamma) \in \mathcal{L}_1(\mathbb{P}_\theta)$ be a real function. Consider $\hat{M}_n(\gamma) = \hat{\mathbb{P}}_n m(\gamma)$ and $M(\theta, \gamma) = \mathbb{P}_\theta m(\gamma)$ where due to the LLN $\hat{M}_n(\gamma) = M(\theta, \gamma) + o_{\mathbb{P}^{\otimes n}}(1)$ for each $\gamma \in \Gamma$. Suppose in addition that

(i) Γ is compact,

(ii) $\gamma \mapsto m(\gamma, x)$ is twice continuously differentiable in a neighbourhood U of $\gamma(\theta) \in \text{int}(\Gamma)$ for \mathbb{P}_θ -a.e. $x \in \mathcal{X}$ with derivatives $\dot{m} := \frac{\partial}{\partial \gamma} m$ and $\ddot{m} := \frac{\partial^2}{\partial^2 \gamma} m$

(iii) $\dot{m}(\gamma(\theta)) \in \mathcal{L}_2(\mathbb{P}_\theta)$ with $\mathbb{P}_\theta \dot{m}(\gamma(\theta)) = 0$ and $\Omega_\theta := \mathbb{P}_\theta \dot{m}(\gamma(\theta)) \dot{m}(\gamma(\theta))^t \geq 0$,

(iv) $\sup_{\gamma \in U} \|\ddot{m}(\gamma)\|_F \in \mathcal{L}_1(\mathbb{P}_\theta)$ and $\ddot{M}_\theta := \mathbb{P}_\theta \ddot{m}(\gamma(\theta))$ is regular with inverse \ddot{M}_θ^{-1} .

hold true. If the M-estimator satisfies $\hat{\gamma}_n = \gamma(\theta) + o_{\mathbb{P}^{\otimes n}}(1)$ then $\sqrt{n}(\hat{\gamma}_n - \gamma(\theta)) \xrightarrow{d} N_{(0, \ddot{M}_\theta^{-1} \Omega_\theta \ddot{M}_\theta^{-1})}$ due to **Corollary §03.04** since the conditions (AN1)-(AN2) are satisfied. Indeed, following **Example §02.11**, (iv) implies the condition (AN2) and due to the CLT the condition (AN1) follows from (iii). However, estimators of \ddot{M}_θ and Ω_θ are necessary in order to use the asymptotic distribution to conduct inference. A typical approach to obtain these estimators is as follows. First replacing \mathbb{P}_θ by $\hat{\mathbb{P}}_n$, the quantity $\hat{\dot{M}}_n(\gamma) := \hat{\mathbb{P}}_n \dot{m}(\gamma)$ and $\hat{\Omega}_n(\gamma) = \hat{\mathbb{P}}_n \dot{m}(\gamma) \dot{m}(\gamma)^t$ is just an empirical counterpart of $\dot{M}_\gamma(\gamma) = \mathbb{P}_\theta \dot{m}(\gamma)$ and $\dot{M}_\theta(\gamma) = \mathbb{P}_\theta \dot{m}(\gamma) \dot{m}(\gamma)^t$, respectively. Secondly, replace γ by its estimator $\hat{\gamma}_n$ we obtain $\hat{\dot{M}}_n := \hat{\dot{M}}_n(\hat{\gamma}_n)$ and $\hat{\Omega}_n := \hat{\Omega}_n(\hat{\gamma}_n)$ as estimator of $\dot{M}_\theta = \dot{M}_\theta(\gamma(\theta))$ and $\Omega_\theta = \Omega_\theta(\gamma(\theta))$, respectively. If in addition to (i)-(iv) the following condition holds

(v) $\sup_{\gamma \in U} \|\dot{m}(\gamma)\|$ belongs to $\mathcal{L}_2(\mathbb{P}_\theta)$.

Then $\sup_{\gamma \in U} \|\hat{\dot{M}}_n(\gamma) - \dot{M}_\theta(\gamma)\|_F = o_{\mathbb{P}^{\otimes n}}(1)$ and $\sup_{\gamma \in U} \|\hat{\Omega}_n(\gamma) - \Omega_\theta(\gamma)\|_F = o_{\mathbb{P}^{\otimes n}}(1)$ following line by line the arguments in **Example §02.11**. From these uniform convergences and $\hat{\gamma}_n = \gamma(\theta) + o_{\mathbb{P}^{\otimes n}}(1)$ follows $\hat{\dot{M}}_n = \dot{M}_\theta + o_{\mathbb{P}^{\otimes n}}(1)$ and $\hat{\Omega}_n = \Omega_\theta + o_{\mathbb{P}^{\otimes n}}(1)$ which in turn implies $\hat{V}_n := \hat{\dot{M}}_n^{-1} \hat{\Omega}_n \hat{\dot{M}}_n^{-1} = \ddot{M}_\theta^{-1} \Omega_\theta \ddot{M}_\theta^{-1} + o_{\mathbb{P}^{\otimes n}}(1)$. Consequently, by applying Slutsky's lemma §20.10 we have $\sqrt{n} \hat{V}_n^{-1/2} (\hat{\gamma}_n - \gamma(\theta)) \xrightarrow{d} N_{(0, \text{Id})}$. \square

§03.07 **Example (MLE, §01.15 continued).** Let $(\mathcal{X}^n, \overline{\mathcal{X}}^{\otimes n}, \mathbb{P}_\theta^{\otimes n})$ with $\mathbb{P}_\theta \ll \mathbb{P}_\theta$ for all $\theta \in \Theta$, likelihood $L(\theta) = d\mathbb{P}_\theta/d\mathbb{P}$, log-likelihood $\ell = \log L$ and parameter of interest θ (i.e., $\gamma = \text{id}_\Theta$) as in **Example §01.15**. Consider the MLE $\hat{\theta}_n$ which maximises the (joint) log-likelihood $\theta \mapsto \hat{\mathbb{P}}_n \ell(\theta)$. Let the following conditions be satisfied:

(i) (Θ, d) is a compact metric space,

- (ii) the parameter θ is identifiable, i.e., $\theta_1 \neq \theta_2$ implies $\mathbb{P}_{\theta_1} \neq \mathbb{P}_{\theta_2}$
- (iii) the map $\theta \mapsto \ell(\theta, x)$ is continuous for \mathbb{P}_θ -a.e. $x \in \mathcal{X}$,
- (iv) $\sup_{\theta \in \Theta} |\ell(\theta)|$ belongs to $\mathcal{L}_1(\mathbb{P}_\theta)$.

Then combining the arguments in the **Examples §02.08 and §02.11** the conditions (CO1) and (CO2) of **Theorem §02.02** are satisfied, which in turn implies consistency of the MLE $\hat{\theta}_n = \theta + o_{\mathbb{P}^{\otimes n}}(1)$. In addition let the following conditions be fulfilled

- (v) for \mathbb{P}_θ -a.e. $x \in \mathcal{X}$ the map $\theta \mapsto \ell(\theta, x)$ is twice continuously differentiable in a neighbourhood U of $\theta \in \Theta$ with derivatives $\dot{\ell}_\theta := \frac{\partial}{\partial \theta} \ell$ and $\ddot{\ell}_\theta := \frac{\partial^2}{\partial \theta^2} \ell$,
- (vi) $\sup_{\theta \in U} \|\dot{\ell}_\theta\| \in \mathcal{L}_2(\mathbb{P})$ and $\sup_{\theta \in U} \|\ddot{\ell}_\theta\|_F \in \mathcal{L}_1(\mathbb{P})$,
- (vii) the Fisher-information matrix $\mathcal{J}_\theta := \mathbb{P}_\theta(\dot{\ell}_\theta \dot{\ell}_\theta^t)$ is strictly positive definite.

Then the conditions (AN1) and (AN2) of **Corollary §03.04**, and the identity $\mathcal{J}_\theta = -\mathbb{P}_\theta \ddot{\ell}_\theta$ are satisfied (for details see **Statistik 1 Satz §17.22**). Therewith, the MLE satisfies $\sqrt{n}(\hat{\theta}_n - \theta) = \sqrt{n} \mathcal{J}_\theta^{-1} \widehat{\mathbb{P}}_n \dot{\ell}_\theta + o_{\mathbb{P}^{\otimes n}}(1)$ and, consequently, $\sqrt{n}(\hat{\theta}_n - \theta_o) \xrightarrow{d} N_{(0, \mathcal{J}_\theta^{-1})}$. □

§03.08 Remark. The conditions (v) and (vi) in **Example §03.07** can be weakened replacing differentiability by Hellinger-differentiability. Keeping the *Hellinger-distance* $H(\mathbb{P}_\theta, \mathbb{P}_{\theta_o}) = \|L^{1/2}(\theta) - L^{1/2}(\theta_o)\|_{\mathcal{L}_2(\mathbb{P})}$ in mind, where $L^{1/2}(\theta) \in \mathcal{L}_2(\mathbb{P})$ using $\|L^{1/2}(\theta)\|_{\mathcal{L}_2(\mathbb{P})}^2 = \mathbb{P}_\theta(L(\theta)) = 1 < \infty$, the family \mathbb{P}_θ is called *Hellinger-differentiable with derivative* $\dot{\ell}_{\theta_o}$ in $\theta_o \in \text{int}(\Theta) \subseteq \mathbb{R}^k$, if $\dot{\ell}_{\theta_o} \in \mathcal{L}_2^k(\mathbb{P}_{\theta_o})$ and hence $\dot{\ell}_{\theta_o} L^{1/2}(\theta_o) \in \mathcal{L}_2^k(\mathbb{P})$ such that

$$\begin{aligned} \lim_{\theta \rightarrow \theta_o} \int_{\mathcal{X}} \left| \frac{L^{1/2}(\theta, x) - L^{1/2}(\theta_o, x) - \frac{1}{2} \langle \dot{\ell}_{\theta_o}(x), \theta - \theta_o \rangle L^{1/2}(\theta_o, x)}{\|\theta - \theta_o\|} \right|^2 \mathbb{P}_\theta(dx) \\ = \lim_{h \rightarrow 0} \frac{\|L^{1/2}(\theta_o + h) - L^{1/2}(\theta_o) - \frac{1}{2} \langle \dot{\ell}_{\theta_o}, h \rangle L^{1/2}(\theta_o)\|_{\mathcal{L}_2(\mathbb{P})}^2}{\|h\|^2} = 0 \end{aligned}$$

The map $x \mapsto \dot{\ell}_{\theta_o}(x)$ is also called *score function*. Keeping $\dot{\ell}_{\theta_o} \in \mathcal{L}_2^k(\mathbb{P}_{\theta_o})$ in mind the *Fisher-information* matrix $\mathcal{J}_{\theta_o} = \mathbb{P}_{\theta_o}(\dot{\ell}_{\theta_o} \dot{\ell}_{\theta_o}^t)$ is well-defined. Note that, the score function and the Fisher-information matrix are independent of the dominating measure \mathbb{P} . □

§03|01 Testing procedures

§03.09 Heuristics. Let $(\mathcal{X}_n, \mathcal{X}_n, \mathbb{P}_\theta^n)$ for all $n \in \mathbb{N}$ be a statistical model over the parameter space Θ and let $\gamma : \Theta \rightarrow \Gamma$ be an identifiable parameter of interest. Given a map $A : \Gamma \rightarrow \mathbb{R}^p$ we eventually test the hypothesis $H_0 : A(\gamma) = 0$ against the alternative $H_1 : A(\gamma) \neq 0$. Typical examples include $A(\gamma) = \gamma - \gamma_o$ for a given value γ_o , or more generally, linear hypothesis $A(\gamma) = M\gamma - a_o$ for a given value a_o and matrix M . It covers in particular testing the j -th coordinate of $\gamma = (\gamma^j)_{j \in [k]}$, i.e., $A(\gamma) = \gamma^j - \gamma_o^j$. Under regularity conditions it seems reasonable to assume an estimator $\hat{\gamma}_n$ of γ having under \mathbb{P}_θ^n the property $\sqrt{n}(A(\hat{\gamma}_n) - A(\gamma(\theta))) \xrightarrow{d} N_{(0, \Sigma_\theta)}$ with invertible asymptotic covariance matrix Σ_θ . If we have in addition an estimator $\widehat{\Sigma}_n = \Sigma_\theta + o_{\mathbb{P}^n}(1)$ at hand. Then under the hypothesis H_0 , i.e., for \mathbb{P}_θ^n with $A(\gamma(\theta)) = 0$, a *Wald test* exploits the property $\widehat{W}_n := nA(\hat{\gamma}_n)^t \widehat{\Sigma}_n^{-1} A(\hat{\gamma}_n) \xrightarrow{d} \chi_p^2$ where χ_p^2 is a Chi-square-distribution with p degrees of freedom. Precisely, a *Wald test* rejects the hypothesis $H_0 : A(\gamma) = 0$ if \widehat{W}_n exceeds the $1-\alpha$ -Quantile $\chi_{p, 1-\alpha}^2$ of a χ_p^2 -distribution. Obviously, the Wald test does exactly meets the asymptotic level α , i.e., $\lim_{n \rightarrow \infty} \mathbb{P}_\theta^n(\widehat{W}_n \geq \chi_{p, 1-\alpha}^2) = \mathbb{P}(W \geq \chi_{p, 1-\alpha}^2) = \alpha$ where $W \sim \chi_p^2$. However, the behaviour of the test statistic \widehat{W}_n under the alternative H_1 is still an open questions, which we intent to study in the next sections. □

§03.10 **Example** (§03.06 *continued*). Let $(\mathcal{X}^n, \mathcal{X}^{\otimes n}, \mathbb{P}_\theta^{\otimes n})$, $\gamma : \Theta \rightarrow \Gamma$ be an identifiable parameter of interest and let $m(\gamma) \in \mathcal{L}_1(\mathbb{P}_\theta)$ for all $\gamma \in \Gamma$. For each $\gamma \in \Gamma$ let $\hat{M}_n(\gamma) = \hat{\mathbb{P}}_n m(\gamma)$ and $M(\theta, \gamma) = \mathbb{E}_\theta m(\gamma)$. Under the conditions (i)-(v) in **Example** §03.06 an M-estimator $\hat{\gamma}_n \in \arg \inf_{\gamma \in \Gamma} \hat{M}_n(\gamma)$ satisfies $\sqrt{n}(\hat{\gamma}_n - \gamma(\theta)) \xrightarrow{d} N_{(0, \hat{M}_\theta^{-1} \Omega_\theta \hat{M}_\theta^{-1})}$ under $\mathbb{P}_\theta^{\otimes n}$. Moreover, we have eventually access to estimators $\hat{\hat{M}}_n = \ddot{M}_\theta + o_{\mathbb{P}^{\otimes n}}(1)$ and $\hat{\hat{\Omega}}_n = \Omega_\theta + o_{\mathbb{P}^{\otimes n}}(1)$. Let $A : \Gamma \rightarrow \mathbb{R}^p$ be continuously differentiable in a neighbourhood of $\gamma(\theta)$ then applying the delta method §20.16 we obtain $\sqrt{n}(A(\hat{\gamma}_n) - A(\gamma(\theta))) \xrightarrow{d} N_{(0, \Sigma_\theta)}$ under $\mathbb{P}_\theta^{\otimes n}$ with $\Sigma_\theta := \dot{A}_{\gamma(\theta)} \ddot{M}_\theta^{-1} \Omega_\theta \ddot{M}_\theta^{-1} \dot{A}_{\gamma(\theta)}^t$. From $\dot{A}_{\hat{\gamma}_n} = \dot{A}_{\gamma(\theta)} + o_{\mathbb{P}^{\otimes n}}(1)$ follows $\hat{\hat{\Sigma}}_n := \dot{A}_{\hat{\gamma}_n} \hat{\hat{M}}_n^{-1} \hat{\hat{\Omega}}_n \hat{\hat{M}}_n^{-1} \dot{A}_{\hat{\gamma}_n}^t = \Sigma_\theta + o_{\mathbb{P}^{\otimes n}}(1)$ and, thus $\sqrt{n} \hat{\hat{\Sigma}}_n^{-1/2} (A(\hat{\gamma}_n) - A(\gamma(\theta))) \xrightarrow{d} N_{(0, \text{Id}_p)}$ which under H_0 , i.e., for $\mathbb{P}_\theta^{\otimes n}$ with $A(\gamma(\theta)) = 0$, implies $\widehat{W}_n := n A(\hat{\gamma}_n)^t \hat{\hat{\Sigma}}_n^{-1} A(\hat{\gamma}_n) \xrightarrow{d} \chi_p^2$. \square

Chapter 2

Asymptotic properties of tests

Asymptotic properties of tests under local alternatives are presented complementing the Neyman-Pearson theory introduced in the lecture [Statistik 1](#). For a more detailed exposition we refer to the text books [Witting and Müller-Funk \[1995\]](#) and [van der Vaart \[1998\]](#).

Overview

§04	Contiguity	15
§04 01	Preliminaries: likelihood ratios and differentiable models	15
§04 02	Contiguity	19
§05	Local asymptotic normality (LAN)	23
§06	Asymptotic relative efficiency	25
§07	Rank tests	27
§08	Asymptotic power of rank tests	30

§04 Contiguity

§04|01 Preliminaries: likelihood ratios and differentiable models

§04.01 **Motivation.** Considering a statistical model $(\mathcal{X}_n, \mathcal{X}_n, \mathbb{P}_\theta^n)$, a parameter of interest $\gamma : \Theta \rightarrow \Gamma$, a partition $\{\mathcal{H}^0, \mathcal{H}^1\}$ of the parameter values of interests $\Gamma = \mathcal{H}^0 \uplus \mathcal{H}^1$ (i.e. $\Gamma = \mathcal{H}^0 \cup \mathcal{H}^1$, $\emptyset = \mathcal{H}^0 \cap \mathcal{H}^1$ and $\mathcal{H}^0 \neq \emptyset \neq \mathcal{H}^1$) we are interested in a (randomised) test $\varphi_n \in \mathcal{X}_n^+$ (i.e. $\varphi_n : \mathcal{X}_n \rightarrow [0, 1]$) of the hypothesis $H_0 : \mathcal{H}^0$ against the alternative $H_1 : \mathcal{H}^1$. Under regularity conditions we may have at hand an estimator $\hat{\gamma}_n$ of γ with known asymptotic distribution. Typically the estimator $\hat{\gamma}_n$ allows us to construct a test statistic T_n with known asymptotic distribution under H_0 , i.e. under \mathbb{P}_θ^n with $\gamma(\theta) \in \mathcal{H}^0$. Exploiting the asymptotic distribution an associated test $\varphi_n = \mathbb{1}_{\{T_n \notin C_\alpha\}}$ does eventually not exceed asymptotically a given level $\alpha \in (0, 1)$ under the hypothesis H_0 . However, we like to investigate also its power under the alternative H_1 , i.e. under a specific \mathbb{P}_θ^n with $\gamma(\theta) \in \mathcal{H}^1$. □

§04.02 **Reminder.** Let ν and μ be measures on $(\mathcal{X}, \mathcal{X})$.

- (a) For any positive numerical function $f \in \overline{\mathcal{X}}^+$ the map $B \mapsto f\mu(B) := \mu(\mathbb{1}_B f) = \int_B f d\mu$ defines a measure $f\mu$ on $(\mathcal{X}, \mathcal{X})$. Any $f = d\nu/d\mu \in \overline{\mathcal{X}}^+$ satisfying $\nu = f\mu$ is called *density* of ν with respect to μ , or *μ -density* for short.
- (b) We say ν is *dominated* by μ , symbolically $\nu \ll \mu$, if for each $B \in \mathcal{X}$ with $\mu(B) = 0$ follows $\nu(B) = 0$. The measures μ and ν are called *equivalent* or *mutually dominated*, symbolically $\mu \ll \nu$, if both $\nu \ll \mu$ and $\mu \ll \nu$.
- (c) We say ν and μ are *orthogonal* or *singular*, symbolically $\nu \perp \mu$, if there exists $\mathcal{X} = \mathcal{X}_\mu \uplus \mathcal{X}_\nu$ with $\mathcal{X}_\mu, \mathcal{X}_\nu \in \mathcal{X}$ and $\mu(\mathcal{X}_\nu) = 0 = \nu(\mathcal{X}_\mu)$. Evidently, we have $\nu \perp \mu$ if and only if there exists $N \in \mathcal{X}$ with $\mu(N) = 0$ such that $\nu = \mathbb{1}_N \nu$.

We note that $g \in \mathcal{L}_1(\mathbb{f}\mu)$ if and only if $gf \in \mathcal{L}_1(\mu)$. In this case holds $\mathbb{f}\mu(g) = \int g d(\mathbb{f}\mu) = \int (gf) d\mu = \mu(gf)$ (Klenke [2012], Satz 4.15, p. 93). Let additionally $\nu \in \mathcal{M}_\sigma(\mathcal{X})$ be a σ -finite measure on $(\mathcal{X}, \mathcal{X})$. If $\mathbb{f}_1\mu = \nu = \mathbb{f}_2\mu$ for $\mathbb{f}_1, \mathbb{f}_2 \in \overline{\mathcal{X}^+}$, then $\mathbb{f}_1 = \mathbb{f}_2$ μ -a.e.. In other words a density is unique up to μ -a.e. equivalence (Klenke [2012], Satz 7.29, p. 159). If in addition $\mu \in \mathcal{M}_\sigma(\mathcal{X})$, then by *Lebesgue's decomposition theorem* there exists $\nu^\nu, \nu^\perp \in \mathcal{M}_\sigma(\mathcal{X})$ such that $\nu = \nu^\nu + \nu^\perp$ with $\nu^\perp \perp \mu$ and $\nu^\nu = \mathbb{f}\mu$ where $\mathbb{f} \in \overline{\mathcal{X}^+}$ and $\mathbb{f} \in \mathbb{R}^+$ μ -a.e.. (Klenke [2012], Satz 7.33, p. 160) Furthermore, there is a *Radon-Nikodym-density* $\mathbb{f} \in \overline{\mathcal{X}^+}$ with $\nu = \mathbb{f}\mu$ and $\mathbb{f} \in \mathbb{R}^+$ μ -a.e. if and only if $\nu \ll \mu$ (Klenke [2012], Korollar 7.34, p. 161). If $\mathbb{f} \in \mathcal{X}^+$ is a Radon-Nikodym-density of ν with respect to μ , i.e. $\nu = \mathbb{f}\mu$, then the positive real function $\mathbb{f}\mathbb{1}_{\{\mathbb{f} \in \mathbb{R}^+\}} \in \mathcal{X}^+$ is it too. Consequently, without loss of generality we consider here and subsequently a positive real version of the Radon-Nikodym-density. Furthermore, given $\mathbb{f} = d\nu^\nu/d\mu \in \mathcal{X}^+$ let us define a numerical function $L := \mathbb{f}\mathbb{1}_{N^c} + \infty\mathbb{1}_N \in \overline{\mathcal{X}^+}$ with $\mu(N) = 0 = \nu^\perp(N^c)$ where $\{L = \infty\} = N$ and the Lebesgue decomposition writes $\nu = L\mu + \mathbb{1}_{\{L = \infty\}}\nu$, i.e. for all $A \in \mathcal{X}$ we have $\nu(A) = \mu(\mathbb{1}_A L) + \nu(A \cap \{L = \infty\})$. \square

§04.03 **Definition.** Let $\mathbb{P}_0, \mathbb{P}_1 \in \mathcal{W}(\mathcal{X})$ be probability measures on $(\mathcal{X}, \mathcal{X})$. Any positive numerical random variable $L \in \overline{\mathcal{X}^+}$ satisfying

$$\mathbb{P}_0(L = \infty) = 0 \text{ and } \mathbb{P}_1 = L\mathbb{P}_0 + \mathbb{1}_{\{L = \infty\}}\mathbb{P}_1 \quad \text{for all } B \in \mathcal{X} \quad (04.01)$$

is called a *likelihood ratio (LR)* of \mathbb{P}_1 with respect to \mathbb{P}_0 , symbolically $d\mathbb{P}_1/d\mathbb{P}_0 := L$. \square

Here and subsequently, let $\mathbb{P}_0, \mathbb{P}_1 \in \mathcal{W}(\mathcal{X})$ and $L := d\mathbb{P}_1/d\mathbb{P}_0$ be a likelihood ratio of \mathbb{P}_1 with respect to \mathbb{P}_0 . We first note that $\mathbb{P}_0(L) = \mathbb{P}_1(L < \infty) \in [0, 1]$ and $\mathbb{P}_0(L \in \mathbb{R}^+) = 1$ by definition, and also $\mathbb{P}_1(L = 0) = L\mathbb{P}_0(L = 0) + \mathbb{P}_1(\{L = 0\} \cap \{L = \infty\}) = 0$.

§04.04 **Property.**

- (i) $\mathbb{P}_0 \perp \mathbb{P}_1 \Leftrightarrow \exists B \in \mathcal{X} : \mathbb{P}_0(B) = 0$ (hence $L\mathbb{P}_0(B) = 0$) and $\mathbb{P}_1(B) = 1$ (hence $\mathbb{P}_1(B \cap \{L = \infty\}) = 1$) $\Leftrightarrow \mathbb{P}_1(L = \infty) = 1 \Leftrightarrow \mathbb{P}_0(L) = 0$;
- (ii) $\mathbb{P}_0 \not\perp \mathbb{P}_1 \Leftrightarrow \forall B \in \mathcal{X} : \mathbb{P}_0(B) = 0$ implies $\mathbb{P}_1(B) < 1$ (particularly for $B = \{L = \infty\}$) $\Leftrightarrow \mathbb{P}_1(L = \infty) < 1 \Leftrightarrow \mathbb{P}_0(L) > 0$;
- (iii) $\mathbb{P}_1 \ll \mathbb{P}_0 \Leftrightarrow \forall B \in \mathcal{X} : \mathbb{P}_0(B) = 0$ implies $\mathbb{P}_1(B) = 0$ (particularly for $B = \{L = \infty\}$) $\Leftrightarrow \mathbb{P}_1(L = \infty) = 0 \Leftrightarrow \mathbb{P}_0(L) = 1$.

§04.05 **Remark.** Note that both \mathbb{P}_0 and \mathbb{P}_1 are dominated by $\mathbb{P}_\mu := \frac{1}{2}(\mathbb{P}_0 + \mathbb{P}_1) \in \mathcal{W}(\mathcal{X})$. Let $\mathbb{f}_i \in \mathcal{X}^+$ denote a \mathbb{P}_μ -density of \mathbb{P}_i , $i \in \{0, 1\}$ (c.f. **Reminder** §04.02), then

$$L_* = \frac{\mathbb{f}_1}{\mathbb{f}_0} \mathbb{1}_{\{\mathbb{f}_0 \in \mathbb{R}_0^+\}} + \infty \mathbb{1}_{\{\mathbb{f}_0 = 0\} \cap \{\mathbb{f}_1 \in \mathbb{R}_0^+\}} \quad (04.02)$$

is a likelihood ratio of \mathbb{P}_1 with respect to \mathbb{P}_0 , i.e., $L_* = d\mathbb{P}_1/d\mathbb{P}_0$. Indeed, $L_* \in \overline{\mathcal{X}^+}$ satisfies $\mathbb{P}_0(L_* = \infty) \leq \mathbb{P}_0(\mathbb{f}_0 = 0) = 0$ and for all $B \in \mathcal{X}$

$$\begin{aligned} L_*\mathbb{P}_0(B) + \mathbb{P}_1(B \cap \{L_* = \infty\}) &= \mathbb{f}_0\mathbb{P}_\mu\left(\frac{\mathbb{f}_1}{\mathbb{f}_0} \mathbb{1}_{B \cap \{\mathbb{f}_0 \in \mathbb{R}_0^+\}}\right) + \mathbb{P}_1(B \cap \{\mathbb{f}_0 = 0\} \cap \{\mathbb{f}_1 \in \mathbb{R}_0^+\}) \\ &= \mathbb{f}_1\mathbb{P}_\mu(B \cap \{\mathbb{f}_0 \in \mathbb{R}_0^+\}) + \mathbb{P}_1(B \cap \{\mathbb{f}_0 = 0\}) = \mathbb{P}_1(B). \end{aligned}$$

Consequently, L_* is always a version of the likelihood ratio $d\mathbb{P}_1/d\mathbb{P}_0$. In general the likelihood ratio $d\mathbb{P}_1/d\mathbb{P}_0$ (and similar $d\mathbb{P}_0/d\mathbb{P}_1$) is uniquely determined by (04.01) up to $(\mathbb{P}_0 + \mathbb{P}_1)$ -a.e. equivalence (e.g. **WTheorie 1, Lemma** §03.15 or Witting [1985] Satz 1.110 a), p. 112). Moreover, the positive numerical random variable $L_*^{-1} = \frac{\mathbb{f}_0}{\mathbb{f}_1} \mathbb{1}_{\{\mathbb{f}_1 \in \mathbb{R}_0^+\}} + \infty \mathbb{1}_{\{\mathbb{f}_1 = 0\} \cap \{\mathbb{f}_0 \in \mathbb{R}_0^+\}}$ is a version of the likelihood ratio

$d\mathbb{P}_0/d\mathbb{P}_1$ switching the roles of \mathbb{P}_0 and \mathbb{P}_1 . Consequently, (iii) can equivalently be written as $\mathbb{P}_1 \ll \mathbb{P}_0 \Leftrightarrow \mathbb{P}_1(d\mathbb{P}_0/d\mathbb{P}_1 = 0) = \mathbb{P}_1(L_*^{-1} = 0) = \mathbb{P}_1(L_* = \infty) = 0$. However, given any version $L = d\mathbb{P}_1/d\mathbb{P}_0$ of the likelihood ratio the measure \mathbb{P}_1 can be written as a sum $\mathbb{P}_1 = \mathbb{P}_1^a + \mathbb{P}_1^\perp$ of two measures $\mathbb{P}_1^a, \mathbb{P}_1^\perp \in \mathcal{M}_\sigma(\mathcal{X})$ where $\mathbb{P}_1^a := L\mathbb{P}_0$ and $\mathbb{P}_1^\perp := \mathbf{1}_{\{L=\infty\}}\mathbb{P}_1$ with $\mathbb{P}_1^\perp(B) = \mathbb{P}_1(B \cap \{L = \infty\})$, $B \in \mathcal{X}$ is, respectively, the absolute continuous and singular part of \mathbb{P}_1 with respect to \mathbb{P}_0 (Lebesgue decomposition). \square

§04.06 **Property.** The two measures $\mathbb{P}_1^a := L\mathbb{P}_0$ and $\mathbb{P}_1^\perp := \mathbf{1}_{\{L=\infty\}}\mathbb{P}_1$ in $\mathcal{M}_\sigma(\mathcal{X})$ satisfy

- (i) $\mathbb{P}_1 = \mathbb{P}_1^a + \mathbb{P}_1^\perp$, $\mathbb{P}_1^a \ll \mathbb{P}_0$, and $\mathbb{P}_1^\perp \perp \mathbb{P}_0$;
- (ii) $\mathbb{P}_1(f) \geq \mathbb{P}_1^a(f) = L\mathbb{P}_0(f) = \mathbb{P}_0(Lf) = \mathbb{P}_1(f\mathbf{1}_{\{L<\infty\}})$ for all $f \in \overline{\mathcal{X}^+}$;
- (iii) $\mathbb{P}_0 \ll \mathbb{P}_1$ if and only if $\mathbb{P}_0(L) = 1$ if and only if $\mathbb{P}_1(d\mathbb{P}_0/d\mathbb{P}_1 = 0) = \mathbb{P}_1(L = \infty) = 0$ if and only if for all $f \in \overline{\mathcal{X}^+}$ holds $\mathbb{P}_1(f) = \mathbb{P}_0(Lf)$. \square

§04.07 **Reminder.** Consider a \mathbb{R}^k -valued statistic S defined on $(\mathcal{X}, \mathcal{X})$, i.e. $S \in \mathcal{X}^k$. If $\mathbb{P}_1 \ll \mathbb{P}_0$, then the probability measure $\mathbb{P}_1^S = \mathbb{P}_1 \circ S^{-1} \in \mathcal{W}(\mathcal{B}^k)$ induced by S under \mathbb{P}_1 can be calculated from the probability measure $\mathbb{P}_0^{(S,L)} = \mathbb{P}_0 \circ (S, L)^{-1}$ induced by the random vector (S, L) under \mathbb{P}_0 through the formula

$$\mathbb{P}_1(S \in B) = \mathbb{P}_1^S(\mathbf{1}_B) = \mathbb{P}_0(\mathbf{1}_B(S)L) = \mathbb{P}_0^{(S,L)}(\mathbf{1}_B(\Pi_S)\Pi_L) \quad \text{for all } B \in \mathcal{B}^k$$

using the coordinate maps $\Pi_L(S, L) = L$ and $\Pi_S(S, L) = S$. The formula, however, is only valid under the assumption $\mathbb{P}_1 \ll \mathbb{P}_0$, since a part of \mathbb{P}_1 orthogonal to \mathbb{P}_0 can't be recovered. \square

Here and subsequently, let $\mathbb{P}_\Theta = (\mathbb{P}_\theta)_{\theta \in \Theta}$ with $\Theta \subseteq \mathbb{R}^k$ be a family of probability measures on a measurable space $(\mathcal{X}, \mathcal{X})$, and for each $\theta_o, \theta \in \Theta$ let $L_{\theta_o}(\theta) := d\mathbb{P}_\theta/d\mathbb{P}_{\theta_o}$ denote a likelihood ratio of \mathbb{P}_θ with respect to \mathbb{P}_{θ_o} . Keep in mind, that $L_{\theta_o}(\theta_o) = 1 (= \mathbf{1}_\mathcal{X})$.

§04.08 **Definition.** Let $s \geq 1$ and $\theta_o \in \text{int}(\Theta)$. The statistical model $(\mathcal{X}, \mathcal{X}, \mathbb{P}_\Theta)$ (and the family \mathbb{P}_Θ) is called $\mathcal{L}_s(\theta_o)$ -differentiable with derivative $\dot{\ell}_{\theta_o}$, if $\dot{\ell}_{\theta_o} \in \mathcal{L}_s^k(\mathbb{P}_{\theta_o})$ and for all $\theta \rightarrow \theta_o$ hold

$$\|s(L_{\theta_o}^{1/s}(\theta) - 1) - \langle \dot{\ell}_{\theta_o}, (\theta - \theta_o) \rangle\|_{\mathcal{L}_s(\mathbb{P}_{\theta_o})} = o(\|\theta - \theta_o\|) \tag{04.03}$$

and $\mathbb{P}_\theta(L_{\theta_o}(\theta) = \infty) = o(\|\theta - \theta_o\|^s)$. \square

§04.09 **Remark.** In case $s = 1$ the defining condition $\mathbb{P}_\theta(L_{\theta_o}(\theta) = \infty) = o(\|\theta - \theta_o\|)$ follows from (04.03) (Witting [1985], Hilfssatz 1.178, p164). We note that $\mathcal{L}_1(\theta_o)$ -differentiability implies $\mathcal{L}_1(\mathbb{P}_{\theta_o})$ -continuity of $\theta \mapsto L_{\theta_o}(\theta)$ in θ_o , i.e., $\|L_{\theta_o}(\theta) - L_{\theta_o}(\theta_o)\|_{\mathcal{L}_1(\mathbb{P}_{\theta_o})} = o(1)$ as $\theta \rightarrow \theta_o$. Since $L_{\theta_o}(\theta)$ is unique up to $\mathbb{P}_\theta + \mathbb{P}_{\theta_o}$ -a.e.-equivalence $\mathcal{L}_s(\theta_o)$ -differentiability does not depend on the version $L_{\theta_o}(\theta)$ of the likelihood ratio $d\mathbb{P}_\theta/d\mathbb{P}_{\theta_o}$. \square

§04.10 **Lemma.** If \mathbb{P}_Θ is $\mathcal{L}_1(\theta_o)$ -differentiable with derivative $\dot{\ell}_{\theta_o}$, then it holds $\mathbb{P}_{\theta_o}(\dot{\ell}_{\theta_o}) = 0$. For any $s \geq r \geq 1$ if \mathbb{P}_Θ is $\mathcal{L}_s(\theta_o)$ -differentiable with derivative $\dot{\ell}_{\theta_o}$, then \mathbb{P}_Θ is also $\mathcal{L}_r(\theta_o)$ -differentiable with derivative $\dot{\ell}_{\theta_o}$.

§04.11 **Proof of Lemma §04.10.** is given in the lecture. \square

In order to avoid additional integrability conditions in Definition §04.08 the function $\theta \mapsto s(L_\mu^{1/s}(\theta) - 1)$ is considered. The next assertion formulates differentiability under additional integrability conditions.

§04.12 **Lemma.** Let $s \geq 1$ and $\theta_o \in \text{int}(\Theta)$. The family \mathbb{P}_θ is $\mathcal{L}_s(\theta_o)$ -differentiable with derivative $\dot{\ell}_{\theta_o}$, if $\dot{\ell}_{\theta_o} \in \mathcal{L}_s^k(\mathbb{P}_\theta)$, $L_{\theta_o}(\theta) \in \mathcal{L}_s(\mathbb{P}_\theta)$ for all $\theta \in U(\theta_o)$ and for all $\theta \rightarrow \theta_o$ hold

$$\|(L_{\theta_o}(\theta) - 1) - \langle \dot{\ell}_{\theta_o}, \theta - \theta_o \rangle\|_{\mathcal{L}_s(\mathbb{P}_\theta)} = o(\|\theta - \theta_o\|)$$

and $\mathbb{P}_\theta(L_{\theta_o}(\theta) = \infty) = o(\|\theta - \theta_o\|^s)$.

§04.13 **Proof of Lemma §04.12.** is given in the lecture. \square

Let us assume in addition, that the family \mathbb{P}_θ is dominated by $\mu \in \mathcal{M}_\sigma(\mathcal{X})$. For each $\theta \in \Theta$ denote by $L_\mu(\theta) := d\mathbb{P}_\theta/d\mu \in \mathcal{X}^+$ a Radon-Nikodym density of \mathbb{P}_θ with respect to μ . Keeping **Remark §04.05** in mind $L_{*,\theta_o}(\theta) = \frac{L_\mu(\theta)}{L_\mu(\theta_o)} \mathbb{1}_{\{L_\mu(\theta_o) \in \mathbb{R}_0^+\}} + \infty \mathbb{1}_{\{L_\mu(\theta_o)=0\} \cap \{L_\mu(\theta) \in \mathbb{R}_0^+\}}$ as in (04.02) is for each $\theta_o, \theta \in \Theta$ a version of the likelihood ratio $d\mathbb{P}_\theta/d\mathbb{P}_{\theta_o}$. We note that

$$\{L_{*,\theta_o}(\theta) = \infty\} = \{\{L_\mu(\theta_o) = 0\} \cap \{L_\mu(\theta) \in \mathbb{R}_0^+\}\} \subseteq \{L_\mu(\theta_o) = 0\} =: \mathcal{N}_{\theta_o},$$

where $\mathbb{P}_{\theta_o}(\mathcal{N}_{\theta_o}) = 0$, and for all $\theta \in \Theta$ holds $\frac{L_\mu(\theta)}{L_\mu(\theta_o)} \mathbb{1}_{\mathcal{N}_{\theta_o}^c} = L_{*,\theta_o}(\theta) \mathbb{1}_{\mathcal{N}_{\theta_o}^c} < \infty$ and $\mathbb{P}_\theta(\mathcal{N}_{\theta_o}) = \mathbb{P}_\theta(\mathcal{N}_{\theta_o} \cap \{L_\mu(\theta) \in \mathbb{R}_0^+\}) = \mathbb{P}_\theta(L_{*,\theta_o}(\theta) = \infty) = \mathbb{P}_\theta(d\mathbb{P}_\theta/d\mathbb{P}_{\theta_o} = \infty)$. Decomposing the integral with respect to $\mathcal{X} = \mathcal{N}_{\theta_o} \uplus \mathcal{N}_{\theta_o}^c$ it follows

$$\begin{aligned} & \|2(L_\mu^{1/2}(\theta) - L_\mu^{1/2}(\theta_o)) - \langle \dot{\ell}_{\theta_o}, (\theta - \theta_o) \rangle L_\mu^{1/2}(\theta_o)\|_{\mathcal{L}_2(\mu)}^2 \\ &= \|2(L_{*,\theta_o}^{1/2}(\theta) - 1) - \langle \dot{\ell}_{\theta_o}, (\theta - \theta_o) \rangle\|_{\mathcal{L}_2(\mathbb{P}_\theta)}^2 + \|\mathbb{1}_{\mathcal{N}_{\theta_o}} 2L_\mu^{1/2}(\theta)\|_{\mathcal{L}_2(\mu)}^2 \\ &= \|2(L_{*,\theta_o}^{1/2}(\theta) - 1) - \langle \dot{\ell}_{\theta_o}, (\theta - \theta_o) \rangle\|_{\mathcal{L}_2(\mathbb{P}_\theta)}^2 + 4\mathbb{P}_\theta(L_{*,\theta_o}(\theta) = \infty) \\ &= \|2(L_{\theta_o}^{1/2}(\theta) - 1) - \langle \dot{\ell}_{\theta_o}, (\theta - \theta_o) \rangle\|_{\mathcal{L}_2(\mathbb{P}_\theta)}^2 + 4\mathbb{P}_\theta(L_{\theta_o}(\theta) = \infty). \end{aligned} \quad (04.04)$$

Keeping **Remark §03.08** in mind for $\theta_o \in \text{int}(\Theta)$ the family \mathbb{P}_θ is *Hellinger-differentiable with derivative $\dot{\ell}_{\theta_o}$* , if $\dot{\ell}_{\theta_o} \in \mathcal{L}_2^k(\mathbb{P}_\theta)$, hence $\dot{\ell}_{\theta_o} L_\mu^{1/2}(\theta_o) \in \mathcal{L}_2^k(\mu)$, and for $\theta \rightarrow \theta_o$

$$\|L_\mu^{1/2}(\theta) - L_\mu^{1/2}(\theta_o) - \frac{1}{2} \langle \dot{\ell}_{\theta_o}, \theta - \theta_o \rangle L_\mu^{1/2}(\theta_o)\|_{\mathcal{L}_2(\mu)} = o(\|\theta - \theta_o\|).$$

Exploiting the identity (04.04) we obtain immediately the next property.

§04.14 **Property.** Let $\mathbb{P}_\theta \ll \mu \in \mathcal{M}_\sigma(\mathcal{X})$ and $\theta_o \in \text{int}(\Theta)$. The family \mathbb{P}_θ is *Hellinger-differentiable with derivative $\dot{\ell}_{\theta_o}$* if and only if \mathbb{P}_θ is $\mathcal{L}_2(\theta_o)$ -differentiable with derivative $\dot{\ell}_{\theta_o}$.

§04.15 **Proposition.** Let $\mathbb{P}_\theta \ll \mu \in \mathcal{M}_\sigma(\mathcal{X})$ with open $\Theta \subseteq \mathbb{R}^k$. If the likelihood $L_\mu(\theta) := d\mathbb{P}_\theta/d\mu$, $\theta \in \Theta$, satisfies in addition the following conditions:

(i) for each $x \in \mathcal{X}$ the map $\theta \mapsto s(\theta, x) := L_\mu^{1/2}(\theta, x)$ is continuously differentiable with derivative $\dot{s}_\theta := \frac{\partial}{\partial \theta} s$,

(ii) $\dot{s}_\theta \in \mathcal{L}_2(\mu)$ for all $\theta \in \Theta$, and hence $\mathcal{J}_\theta := 4\mu(\dot{s}_\theta \dot{s}_\theta^t) \in \mathbb{R}_{\geq}^{(k,k)}$,

(iii) the map $\theta \mapsto \mathcal{J}_\theta$ is continuous.

Then \mathbb{P}_θ is for all $\theta_o \in \Theta$ *Hellinger-differentiable with score function $\dot{\ell}_{\theta_o} = 2 \frac{\dot{s}_{\theta_o}}{s(\theta_o)} \mathbb{1}_{\{s(\theta_o) \in \mathbb{R}_0^+\}}$* .

§04.16 **Proof of Proposition §04.15.** is given in the lecture. \square

§04.17 **Example.** Consider a statistical location model $(\mathbb{R}, \mathcal{B}, \mathbb{P}_\mathbb{R})$ dominated by the Lebesgue measure $\lambda \in \mathcal{M}_\sigma(\mathcal{B})$ with likelihood for each $\theta \in \mathbb{R}$ given by $L(\theta, x) = g(x - \theta)$, $x \in \mathbb{R}$, where g is a strictly positive density. If g is continuously differentiable with derivative \dot{g} satisfying $\lambda(|\dot{g}|^2/g) < \infty$ then due to **Proposition §04.15** the family $\mathbb{P}_\mathbb{R}$ is *Hellinger-differentiable*

with score function $\dot{\ell}_\theta = -\dot{g}(x - \theta)/g(x - \theta)$. Indeed, setting $s(\theta, x) := \sqrt{g(x - \theta)}$, we have $\dot{s}_\theta(x) = \frac{\partial}{\partial \theta} \sqrt{g(x - \theta)} = -\frac{1}{2} \dot{g}(x - \theta) / \sqrt{g(x - \theta)}$ which is continuous in θ and hence condition (i) is satisfied. Moreover conditions (ii) and (iii) hold true, since $\theta \mapsto \mathcal{J}_\theta = 4\lambda(\dot{s}_\theta)^2 = \lambda(|\dot{g}|^2/g) < \infty$ is constant and thus continuous. Applying **Proposition** §04.15 the family \mathbb{P}_R is Hellinger-differentiable with score function $\dot{\ell}_{\theta_0} = 2 \frac{\dot{s}_{\theta_0}}{s(\theta_0)} \mathbb{1}_{\{s(\theta_0) \in \mathbb{R}_0^+\}} = -\dot{g}(x - \theta_0)/g(x - \theta_0)$. \square

§04|02 Contiguity

We introduce next an asymptotic version of absolute continuity. In this section we restrict our attention to probability measures $\mathbb{P}_0^n, \mathbb{P}_1^n \in \mathcal{W}(\mathcal{X}_n)$, $n \in \mathbb{N}$, in short $(\mathbb{P}_0^n)_{n \in \mathbb{N}}, (\mathbb{P}_1^n)_{n \in \mathbb{N}} \in (\mathcal{W}(\mathcal{X}_n))_{n \in \mathbb{N}}$. We aim to obtain the limiting distribution of (test) statistics $S_n \in \mathcal{X}_n^k$, $n \in \mathbb{N}$, under \mathbb{P}_1^n if its limiting distribution under \mathbb{P}_0^n is known.

§04.18 **Definition.** Let $\mathbb{P}_0^n, \mathbb{P}_1^n \in \mathcal{W}(\mathcal{X}_n)$, $n \in \mathbb{N}$. The sequence $(\mathbb{P}_1^n)_{n \in \mathbb{N}}$ is called *contiguous* with respect to $(\mathbb{P}_0^n)_{n \in \mathbb{N}}$, symbolically $\mathbb{P}_1^n \triangleleft \mathbb{P}_0^n$, if for any $(B_n)_{n \in \mathbb{N}} \in (\mathcal{X}_n)_{n \in \mathbb{N}}$ with $\lim_{n \rightarrow \infty} \mathbb{P}_0^n(B_n) = 0$ holds $\lim_{n \rightarrow \infty} \mathbb{P}_1^n(B_n) = 0$. The sequences $(\mathbb{P}_1^n)_{n \in \mathbb{N}}$ and $(\mathbb{P}_0^n)_{n \in \mathbb{N}}$ are called *mutually contiguous*, symbolically $\mathbb{P}_0^n \triangleleft \triangleright \mathbb{P}_1^n$, if both $\mathbb{P}_1^n \triangleleft \mathbb{P}_0^n$ and $\mathbb{P}_0^n \triangleleft \mathbb{P}_1^n$. \square

§04.19 **Lemma.** Let $\mathbb{P}_0^n, \mathbb{P}_1^n \in \mathcal{W}(\mathcal{X}_n)$, $n \in \mathbb{N}$.

- (i) $\mathbb{P}_1^n \triangleleft \mathbb{P}_0^n \Leftrightarrow$ for all $(S_n)_{n \in \mathbb{N}} \in (\mathcal{X}_n^k)_{n \in \mathbb{N}}$ holds: $S_n \xrightarrow{\mathbb{P}_0^n} 0 \Rightarrow S_n \xrightarrow{\mathbb{P}_1^n} 0$;
- (ii) For any statistic $S_n : (\mathcal{X}_n, \mathcal{X}_n) \rightarrow (\mathcal{S}, \mathcal{S})$, $n \in \mathbb{N}$, holds: $\mathbb{P}_1^n \triangleleft \mathbb{P}_0^n \Rightarrow \mathbb{P}_1^n \circ S_n^{-1} \triangleleft \mathbb{P}_0^n \circ S_n^{-1}$;
- (iii) For any sub-sequence $(n_k)_{k \in \mathbb{N}}$ in \mathbb{N} holds: $\mathbb{P}_1^n \triangleleft \mathbb{P}_0^n \Rightarrow \mathbb{P}_1^{n_k} \triangleleft \mathbb{P}_0^{n_k}$;
- (iv) $\mathbb{P}_1^n \triangleleft \mathbb{P}_0^n \Leftrightarrow$ for any $\varepsilon \in \mathbb{R}_0^+$ exists $\delta \in \mathbb{R}_0^+$ such that for all $(B_n)_{n \in \mathbb{N}} \in (\mathcal{X}_n)_{n \in \mathbb{N}}$ holds: $\limsup_{n \rightarrow \infty} \mathbb{P}_0^n(B_n) < \delta \Rightarrow \limsup_{n \rightarrow \infty} \mathbb{P}_1^n(B_n) < \varepsilon$;
- (v) Let $(S_n)_{n \in \mathbb{N}} \in (\mathcal{X}_n^k)_{n \in \mathbb{N}}$ and $\mathbb{P}_1^n \triangleleft \mathbb{P}_0^n$, then:
 - (v-a) $\mathbb{P}_0^n \circ S_n^{-1} \xrightarrow{d} \mathbb{P}_0$ and $\mathbb{P}_1^n \circ S_n^{-1} \xrightarrow{d} \mathbb{P}_1 \Rightarrow \mathbb{P}_1 \ll \mathbb{P}_0$;
 - (v-b) $(\mathbb{P}_0^n \circ S_n^{-1})_{n \in \mathbb{N}}$ tight $\Rightarrow (\mathbb{P}_1^n \circ S_n^{-1})_{n \in \mathbb{N}}$ tight.

§04.20 **Proof of Lemma** §04.19. is given in the lecture. \square

§04.21 **Remark.** Next we characterise contiguity in terms of the asymptotic behaviour of the likelihood ratio $L_n = d\mathbb{P}_1^n/d\mathbb{P}_0^n \in \overline{\mathcal{X}_n^+}$, $n \in \mathbb{N}$. First recall that $\mathbb{P}_1^n(L_n < \infty) = \mathbb{P}_0^n(L_n) \in [0, 1]$ and $\mathbb{P}_0^n(L_n = \infty) = \mathbb{P}_1^n(L_n = 0) = 0$ for each $n \in \mathbb{N}$. Consequently, the probability measure $\mathbb{P}_0^n \circ L_n^{-1} \in \mathcal{W}(\overline{\mathcal{B}})$ is concentrated in \mathbb{R}^+ meaning that $\mathbb{P}_0^n \circ L_n^{-1}(\mathbb{R}^+) = \mathbb{P}_0^n(L_n \in \mathbb{R}^+) = 1$ for each $n \in \mathbb{N}$. Moreover, $(\mathbb{P}_0^n \circ L_n^{-1})_{n \in \mathbb{N}}$ is tight, since for any $\varepsilon \in \mathbb{R}_0^+$ and $c > 1/\varepsilon$ holds $\mathbb{P}_0^n(L_n > c) \leq \frac{1}{c} \mathbb{P}_0^n(L_n) \leq \frac{1}{c} < \varepsilon$ by Markov's inequality. However, $\mathbb{P}_1^n \circ L_n^{-1}$ is generally not concentrated in \mathbb{R}^+ , but under $\mathbb{P}_1^n \triangleleft \mathbb{P}_0^n$ holds $\mathbb{P}_1^n(L_n = \infty) \rightarrow 0$ since $\mathbb{P}_0^n(L_n = \infty) = 0$ for all $n \in \mathbb{N}$. Thereby, the limit distribution of $\mathbb{P}_1^n \circ L_n^{-1}$ (if it exists) is concentrated in \mathbb{R}^+ . \square

Formally, we write $L_n = L_n \mathbb{1}_{\{L_n \in \mathbb{R}^+\}} + \infty \mathbb{1}_{\{L_n = \infty\}}$, where the second summand is negligible in the sense of Slutsky's lemma under contiguity $\mathbb{P}_1^n \triangleleft \mathbb{P}_0^n$.

§04.22 **Definition.** $(T_n)_{n \in \mathbb{N}} \in (\overline{\mathcal{X}_n})_{n \in \mathbb{N}}$ converges in distribution to $\mathbb{P}^T \in \mathcal{W}(\overline{\mathcal{B}})$ under \mathbb{P}^n , shortly $T_n \xrightarrow{d} \mathbb{P}^T$ under \mathbb{P}^n , if

$$\mathbb{P}^n \circ T_n^{-1} \xrightarrow{d} \mathbb{P}^T \quad :\Leftrightarrow \quad \mathbb{P}^n \circ (T_n \mathbb{1}_{\{T_n \in \mathbb{R}\}})^{-1} \xrightarrow{d} \mathbb{P}^T \quad \text{and} \quad \mathbb{P}^n(T_n \notin \mathbb{R}) \rightarrow 0. \quad (04.05)$$

We note that any family of probability measures on $(\overline{\mathbb{R}}, \overline{\mathcal{B}})$ is tight, since $\overline{\mathbb{R}}$ is compact. A non trivial formulation of tightness for probability measures provides the next definition.

§04.23 **Definition.** A sequence $(\mathbb{P}^n)_{n \in \mathbb{N}} \in \mathcal{W}(\overline{\mathcal{B}})$ is called *asymptotically tight* if for all $\varepsilon \in \mathbb{R}_0^+$ exists $M \in \mathbb{R}_0^+$ and $n_o \in \mathbb{N}$ such that for all $n > n_o$ holds $\mathbb{P}^n([-M, M]^c) < \varepsilon$. \square

§04.24 **Remark.** Asymptotic tightness of $(\mathbb{P}^n)_{n \in \mathbb{N}} \in \mathcal{W}(\overline{\mathcal{B}})$ is equivalently characterised by: for any $(M_n)_{n \in \mathbb{N}}$ in \mathbb{R} with $M_n \uparrow \infty$ holds $\mathbb{P}^n([-M_n, M_n]^c) \xrightarrow{n \rightarrow \infty} 0$. In particular, we have immediately $\mathbb{P}^n(\{-\infty, \infty\}) \xrightarrow{n \rightarrow \infty} 0$. The concept of asymptotic tightness and tightness as in **Definition** §20.21 coincide if $\mathbb{P}^n(\mathbb{R}) = 1$ for all $n \in \mathbb{N}$. Furthermore, it can be shown that the claim of Prohorov's theorem **Property** §20.24 holds also for families of asymptotically tight probability measures. \square

§04.25 **Theorem.** For each $n \in \mathbb{N}$ let $\mathbb{P}_0^n, \mathbb{P}_1^n \in \mathcal{W}(\mathcal{X}_n)$, let $L_n := d\mathbb{P}_1^n/d\mathbb{P}_0^n \in \overline{\mathcal{X}_n^+}$ be a likelihood ratio of \mathbb{P}_1^n with respect to \mathbb{P}_0^n and let $\mathbb{P}_0^L, \mathbb{P}_1^L \in \mathcal{W}(\mathcal{B})$. Then the following statements are equivalent:

- (a1) $\mathbb{P}_1^n \triangleleft \mathbb{P}_0^n$;
- (a2) $\mathbb{P}_0^n(L_n) \xrightarrow{n \rightarrow \infty} 1$ and for any $\varepsilon \in \mathbb{R}_0^+$ exists $M \in \mathbb{R}_0^+$ with $\sup_{n \in \mathbb{N}} \mathbb{P}_0^n(L_n \mathbb{1}_{\{L_n > M\}}) < \varepsilon$, i.e. $(\mathbb{P}_0^n \circ L_n^{-1})_{n \in \mathbb{N}}$ is uniformly integrable;
- (a3) $(\mathbb{P}_1^n \circ L_n^{-1})_{n \in \mathbb{N}}$ is asymptotically tight.

If in addition $L_n \xrightarrow{d} \mathbb{P}_0^L$ under \mathbb{P}_0^n , i.e. $\mathbb{P}_0^n \circ L_n^{-1} \xrightarrow{d} \mathbb{P}_0^L$, then the following statements are equivalent:

- (b1) $\mathbb{P}_1^n \triangleleft \mathbb{P}_0^n$;
- (b2) $1 = \int_{\mathbb{R}} y \mathbb{P}_0^L(dy) = \mathbb{P}_0^L(\text{id}_{\mathbb{R}}) = \mathbb{P}_0^L(\text{id}_{\mathbb{R}} \mathbb{1}_{\mathbb{R}})$;
- (b3) $L_n \xrightarrow{d} \mathbb{P}_1^L$ under \mathbb{P}_1^n with $\mathbb{P}_1^L(B) = \mathbb{P}_0^L(\text{id}_{\mathbb{R}} \mathbb{1}_B) = \int_B y \mathbb{P}_0^L(dy)$ for all $B \in \mathcal{B}$.

§04.26 **Proof of Theorem** §04.25. is given in the lecture. \square

Since $\mathbb{P}_0^n(L_n) = \mathbb{P}_1^n(L_n < \infty)$ it holds $\mathbb{P}_0^n(L_n) \rightarrow 1 \Leftrightarrow \mathbb{P}_1^n(L_n = \infty) \rightarrow 0$. Keeping (04.01) in mind the mass of the absolute continuous part of \mathbb{P}_1^n with respect to \mathbb{P}_0^n converges two 1, if and only if, the singular part vanishes.

§04.27 **Corollary.** Under the notations of **Theorem** §04.25 the following statements are equivalent:

- (i) $\mathbb{P}_1^n \triangleleft \mathbb{P}_0^n$;
- (ii) if $\mathbb{P}_0^{n_k} \circ L_{n_k}^{-1} \xrightarrow{d} \mathbb{P}_0^L \in \mathcal{W}(\mathcal{B})$ along a sub-sequence $(n_k)_{k \in \mathbb{N}}$, then $\mathbb{P}_0^L(\text{id}_{\mathbb{R}}) = 1$;
- (iii) if $\mathbb{P}_0^{n_k} \circ L_{n_k}^{-1} \xrightarrow{d} \mathbb{P}_0^L \in \mathcal{W}(\mathcal{B})$ along a sub-sequence $(n_k)_{k \in \mathbb{N}}$, then $\mathbb{P}_1^{n_k} \circ L_{n_k}^{-1} \xrightarrow{d} \mathbb{P}_1^L$, with $\mathbb{P}_1^L(B) = \mathbb{P}_0^L(\text{id}_{\mathbb{R}} \mathbb{1}_B)$ for all $B \in \mathcal{B}$.

§04.28 **Proof of Corollary** §04.27. is given in the lecture. \square

We are particularly interested in mutual contiguity $(\mathbb{P}_0^n \triangleleft \triangleright \mathbb{P}_1^n)$ of $(\mathbb{P}_0^n)_{n \in \mathbb{N}}$ and $(\mathbb{P}_1^n)_{n \in \mathbb{N}}$, which can be characterised by applying **Theorem** §04.25 and its analogous formulation switching the roles of \mathbb{P}_0^n and \mathbb{P}_1^n . However, for $n \in \mathbb{N}$ the transformation of a likelihood ratio $L_n = d\mathbb{P}_1^n/d\mathbb{P}_0^n$ into a *log-likelihood ratio* (LLR) $\ell_n := \log L_n = \log(d\mathbb{P}_1^n/d\mathbb{P}_0^n) \in \overline{\mathcal{X}}$ captures equally both orthogonal events $\{L_n = 0\}$ and $\{L_n = \infty\}$. Generally, ℓ_n takes the value $-\infty$ and $+\infty$ with positive \mathbb{P}_0^n - and \mathbb{P}_1^n -probability, respectively. In other words $\mathbb{P}_0^n \circ \ell_n^{-1}$ and $\mathbb{P}_1^n \circ \ell_n^{-1}$ is concentrated in $[-\infty, \infty)$ and $(-\infty, \infty]$, respectively, since by **Definition** §04.03 of L_n it holds

$$\mathbb{P}_0^n(\ell_n = \infty) = 0 \quad \text{and} \quad \mathbb{P}_1^n(\ell_n = -\infty) = 0 \quad \text{for all } n \in \mathbb{N}. \quad (04.06)$$

Thereby, similar to **Remark** §04.21 under mutual contiguity $\mathbb{P}_0^n \triangleleft \triangleright \mathbb{P}_1^n$ it follows

$$\mathbb{P}_1^n(\ell_n = \infty) \rightarrow 0 \quad \text{and} \quad \mathbb{P}_0^n(\ell_n = -\infty) \rightarrow 0 \quad \text{as } n \rightarrow \infty. \quad (04.07)$$

Consequently, the limit distribution of ℓ_n under both \mathbb{P}_0^n and \mathbb{P}_1^n , if it exists, is concentrated in \mathbb{R} . Keeping [Definition §04.22](#) in mind under mutual contiguity $\mathbb{P}_0^n \triangleleft \triangleright \mathbb{P}_1^n$ convergence in distribution of ℓ_n under \mathbb{P}_0^n and \mathbb{P}_1^n to $\mathbb{P}_0^\ell, \mathbb{P}_1^\ell \in \mathcal{W}(\mathcal{B})$, respectively, is equivalently characterised by

$$\begin{aligned} \mathbb{P}_0^n \circ \ell_n^{-1} \xrightarrow{d} \mathbb{P}_0^\ell &\Leftrightarrow \mathbb{P}_0^n \circ (\ell_n \mathbb{1}_{\{\ell_n > -\infty\}})^{-1} \xrightarrow{d} \mathbb{P}_0^\ell \quad \text{and} \\ \mathbb{P}_1^n \circ \ell_n^{-1} \xrightarrow{d} \mathbb{P}_1^\ell &\Leftrightarrow \mathbb{P}_1^n \circ (\ell_n \mathbb{1}_{\{\ell_n < \infty\}})^{-1} \xrightarrow{d} \mathbb{P}_1^\ell. \end{aligned} \tag{04.08}$$

If $L_n^{-1} = d\mathbb{P}_0^n/d\mathbb{P}_1^n$ is a likelihood ratio of \mathbb{P}_0^n with respect to \mathbb{P}_1^n , as for example in [Remark §04.05](#), then making use of the identity $\log L_n^{-1} = -\log L_n = -\ell_n$ the convergence in distribution of ℓ_n under \mathbb{P}_0^n respectively \mathbb{P}_1^n implies immediately the corresponding convergence of $\log L_n^{-1}$. Similar to [Theorem §04.25 \(b1\)-\(b3\)](#) the next result characterises mutual contiguity in terms of the log-likelihood ratio ℓ_n .

§04.29 **Theorem.** For each $n \in \mathbb{N}$ let $\mathbb{P}_0^n, \mathbb{P}_1^n \in \mathcal{W}(\mathcal{X}_n)$, let $\ell_n := \log L_n = \log(d\mathbb{P}_1^n/d\mathbb{P}_0^n) \in \overline{\mathcal{X}_n}$ be a log-likelihood ratio such that also $L_n^{-1} = d\mathbb{P}_0^n/d\mathbb{P}_1^n \in \overline{\mathcal{X}_n}^+$ and let $\mathbb{P}_0^\ell, \mathbb{P}_1^\ell \in \mathcal{W}(\mathcal{B})$. If in addition $\ell_n \xrightarrow{d} \mathbb{P}_0^\ell$ under \mathbb{P}_0^n , i.e. $\mathbb{P}_0^n \circ \ell_n^{-1} \xrightarrow{d} \mathbb{P}_0^\ell$, then the following statements are equivalent:

- (b'1) $\mathbb{P}_1^n \triangleleft \triangleright \mathbb{P}_0^n$;
- (b'2) $1 = \int_{\mathbb{R}} \exp(z) \mathbb{P}_0^\ell(dz) = \mathbb{P}_0^\ell(\exp) = \mathbb{P}_0^\ell(\exp \mathbb{1}_{\mathbb{R}})$
- (b'3) $\ell_n \xrightarrow{d} \mathbb{P}_1^\ell$ under \mathbb{P}_1^n with $\mathbb{P}_1^\ell(B) = \mathbb{P}_0^\ell(\exp \mathbb{1}_B) = \int_B \exp(z) \mathbb{P}_0^\ell(dz)$ for all $B \in \mathcal{B}$.

§04.30 **Proof of Theorem §04.29.** is given in the lecture. □

§04.31 **Remark.** Let \mathbb{f}_0^ℓ and \mathbb{f}_1^ℓ denote, respectively, a μ -density of \mathbb{P}_0^ℓ and \mathbb{P}_1^ℓ with respect to a measure $\mu \in \mathcal{M}_\sigma(\mathcal{B})$ dominating \mathbb{P}_0^ℓ , and hence \mathbb{P}_1^ℓ . The measure \mathbb{P}_1^ℓ in [Theorem §04.29 \(b'3\)](#) is equally defined by $\mathbb{f}_1^\ell(z) = \exp(z) \mathbb{f}_0^\ell(z)$ for μ -a.e. $z \in \mathbb{R}$. □

§04.32 **Corollary.** Under the notations of [Theorem §04.29](#) if $\mathbb{P}_0^n \circ \ell_n^{-1} \xrightarrow{d} N_{(\mu, \sigma^2)}$ for $(\mu, \sigma) \in \mathbb{R} \times \mathbb{R}^+$ then the following statements are equivalent:

- (b''1) $\mathbb{P}_1^n \triangleleft \triangleright \mathbb{P}_0^n$;
- (b''2) $\mu = -\sigma^2/2$
- (b''3) $\ell_n \xrightarrow{d} N_{(\sigma^2/2, \sigma^2)}$ under \mathbb{P}_1^n .

§04.33 **Proof of Corollary §04.32.** is given in the lecture. □

§04.34 **Example (Le Cam's first lemma).** For $n \in \mathbb{N}$ let $\mathbb{P}_0^n, \mathbb{P}_1^n \in \mathcal{W}(\mathcal{X}_n)$ and $L_n := d\mathbb{P}_1^n/d\mathbb{P}_0^n \in \overline{\mathcal{X}_n}^+$. If $\ell_n := \log L_n \xrightarrow{d} N_{(-\sigma^2/2, \sigma^2)}$ under \mathbb{P}_0^n , then $\mathbb{P}_1^n \triangleleft \triangleright \mathbb{P}_0^n$ and $\ell_n \xrightarrow{d} N_{(\sigma^2/2, \sigma^2)}$ under \mathbb{P}_1^n due to [Corollary §04.32](#). For $\sigma > 0$ from $\sigma^{-1}(\ell_n + \sigma^2/2) \xrightarrow{d} N_{(0,1)}$ under \mathbb{P}_0^n follows thus $\sigma^{-1}(\ell_n + \sigma^2/2) \xrightarrow{d} N_{(\sigma,1)}$ under \mathbb{P}_1^n . In other words in this situation there is asymptotically a location shift by σ . □

For each $n \in \mathbb{N}$ let $\mathbb{P}_0^n, \mathbb{P}_1^n \in \mathcal{W}(\mathcal{X}_n)$, let $L_n := d\mathbb{P}_1^n/d\mathbb{P}_0^n \in \overline{\mathcal{X}_n}^+$ be a likelihood ratio of \mathbb{P}_1^n with respect to \mathbb{P}_0^n , let $\ell_n := \log L_n$ and let $S_n \in \mathcal{X}_n^k$ be a \mathbb{R}^k -valued statistic defined on $(\mathcal{X}_n, \mathcal{X}_n)$. We search conditions which allow to calculate the limiting distribution of (S_n, L_n) respectively (S_n, ℓ_n) under \mathbb{P}_1^n , from the limiting distribution of (S_n, L_n) respectively (S_n, ℓ_n) under \mathbb{P}_0^n . Keeping again ([§04.22](#)) in mind under mutual contiguity $\mathbb{P}_0^n \triangleleft \triangleright \mathbb{P}_1^n$ the joint convergence in distribution of $(S_n, L_n) \xrightarrow{d} \mathbb{P}_1^{(S,L)} \in \mathcal{W}(\mathcal{B}^{k+1})$ under \mathbb{P}_1^n , $(S_n, \ell_n) \xrightarrow{d} \mathbb{P}_0^{(S,\ell)} \in \mathcal{W}(\mathcal{B}^{k+1})$

under \mathbb{P}_0^n and $(S_n, \ell_n) \xrightarrow{d} \mathbb{P}_1^{(S, \ell)} \in \mathcal{W}(\mathcal{B}^{k+1})$ under \mathbb{P}_1^n , respectively, is equally characterised by

$$\begin{aligned} \mathbb{P}_1^n \circ (S_n, L_n)^{-1} \xrightarrow{d} \mathbb{P}_1^{(S, L)} &\Leftrightarrow \mathbb{P}_1^n \circ (S_n, L_n \mathbf{1}_{\{L_n < \infty\}})^{-1} \xrightarrow{d} \mathbb{P}_1^{(S, L)}, \\ \mathbb{P}_0^n \circ (S_n, \ell_n)^{-1} \xrightarrow{d} \mathbb{P}_0^{(S, \ell)} &\Leftrightarrow \mathbb{P}_0^n \circ (S_n, \ell_n \mathbf{1}_{\{\ell_n > -\infty\}})^{-1} \xrightarrow{d} \mathbb{P}_0^{(S, \ell)} \quad \text{and} \\ \mathbb{P}_1^n \circ (S_n, \ell_n)^{-1} \xrightarrow{d} \mathbb{P}_1^{(S, \ell)} &\Leftrightarrow \mathbb{P}_1^n \circ (S_n, \ell_n \mathbf{1}_{\{\ell_n < \infty\}})^{-1} \xrightarrow{d} \mathbb{P}_1^{(S, \ell)}. \end{aligned} \quad (04.09)$$

Denote by $\Pi_L := \Pi_{k+1} \in \mathcal{B}^{k+1}$, i.e. $y = (y_i)_{i \in \llbracket k+1 \rrbracket} \mapsto \Pi_L(y) := y_{k+1}$ (respectively $\Pi_\ell := \Pi_{k+1} \in \mathcal{B}^{k+1}$) the coordinate map which allows us to write

$$\int_C y \mathbb{P}_1^{(S, L)}(ds, dy) = \int_{\mathbb{R}^{k+1}} \mathbf{1}_C(s, y) \Pi_L(s, y) \mathbb{P}_1^{(S, L)}(ds, dy) = \mathbb{P}_1^{(S, L)}(\mathbf{1}_C \Pi_L) \text{ for all } C \in \mathcal{B}^{k+1}.$$

§04.35 **Theorem.** For each $n \in \mathbb{N}$ let $\mathbb{P}_0^n, \mathbb{P}_1^n \in \mathcal{W}(\mathcal{X}_n)$, let $\ell_n = \log L_n = \log(d\mathbb{P}_1^n/d\mathbb{P}_0^n) \in \overline{\mathcal{X}_n}$ be a log-likelihood ratio, and let $S_n \in \mathcal{X}_n^k$ be a \mathbb{R}^k -valued statistic. Then, we have

- (i) If $(S_n, L_n) \xrightarrow{d} \mathbb{P}_0^{(S, L)} \in \mathcal{W}(\mathcal{B}^{k+1})$ under \mathbb{P}_0^n and $\mathbb{P}_0^{(S, L)}(\Pi_L \mathbf{1}_{\mathbb{R}^{k+1}}) = \mathbb{P}_0^{(S, L)}(\Pi_L) = 1$, then $(S_n, L_n) \xrightarrow{d} \mathbb{P}_1^{(S, L)}$ under \mathbb{P}_1^n with $\mathbb{P}_1^{(S, L)}(C) := \mathbb{P}_0^{(S, L)}(\Pi_L \mathbf{1}_C)$ for all $C \in \mathcal{B}^{k+1}$.
- (ii) If $(S_n, \ell_n) \xrightarrow{d} \mathbb{P}_0^{(S, \ell)} \in \mathcal{W}(\mathcal{B}^{k+1})$ under \mathbb{P}_0^n and $\mathbb{P}_0^{(S, \ell)}(\exp(\Pi_\ell) \mathbf{1}_{\mathbb{R}^{k+1}}) = \mathbb{P}_0^{(S, \ell)}(\exp(\Pi_\ell)) = 1$, then $(S_n, \ell_n) \xrightarrow{d} \mathbb{P}_1^{(S, \ell)}$ under \mathbb{P}_1^n with $\mathbb{P}_1^{(S, \ell)}(C) := \mathbb{P}_0^{(S, \ell)}(\exp(\Pi_\ell) \mathbf{1}_C)$ for all $C \in \mathcal{B}^{k+1}$.

§04.36 **Proof** of **Theorem** §04.35. is given in the lecture. \square

§04.37 **Example (Le Cam's third lemma).** For each $n \in \mathbb{N}$ let $\mathbb{P}_0^n, \mathbb{P}_1^n \in \mathcal{W}(\mathcal{X}_n)$, let $\ell_n = \log L_n = \log(d\mathbb{P}_1^n/d\mathbb{P}_0^n) \in \overline{\mathcal{X}_n}$ be a log-likelihood ratio, and let $S_n \in \mathcal{X}_n^k$ be a \mathbb{R}^k -valued statistic. Suppose that the limit distribution of (S_n, ℓ_n) under \mathbb{P}_0^n is multivariate normal, that is

$$\mathbb{P}_0^n \circ (S_n, \ell_n)^{-1} \xrightarrow{d} \mathbb{P}_0^{(S, \ell)} = N_{(v, M)} \quad \text{with} \quad v = \begin{pmatrix} \mu \\ -\sigma^2/2 \end{pmatrix} \text{ and } M = \begin{pmatrix} \Sigma & \tau \\ \tau^t & \sigma^2 \end{pmatrix}. \quad (04.10)$$

Then it holds $(S_n, \ell_n) \xrightarrow{d} \mathbb{P}_1^{(S, \ell)} = N_{(v', M)}$ under \mathbb{P}_1^n with $v' = (\mu + \tau, \sigma^2/2)^t$. Indeed, since $\mathbb{P}_0^{(S, \ell)}(\exp(\Pi_\ell)) = 1$ both assumptions of **Theorem** §04.35 (ii) are satisfied and hence it remains to calculate the limit distribution $\mathbb{P}_1^{(S, \ell)}(C) := \mathbb{P}_0^{(S, \ell)}(\exp(\Pi_\ell) \mathbf{1}_C)$ for all $C \in \mathcal{B}^{k+1}$. Suppose first $M > 0$, or equivalently $\Sigma > 0$ and $\sigma > 0$, then $\mathbb{P}_0^{(S, \ell)}$ has a density $f_0^{(S, \ell)}$ with respect to the Lebesgue-measure $\lambda^{k+1} \in \mathcal{M}_\sigma(\mathcal{B}^{k+1})$ and (see **Remark** §04.31) the Lebesgue-density $f_1^{(S, \ell)}$ of $\mathbb{P}_1^{(S, \ell)}$ satisfies $f_1^{(S, \ell)}(s, z) = \exp(z) f_0^{(S, \ell)}(s, z)$ for λ^{k+1} -a.e. $(s, z) \in \mathbb{R}^{k+1}$. Keeping the coordinate map Π_ℓ in mind we denote by f_0^ℓ and f_1^ℓ the marginal density of $\mathbb{P}_0^{(S, \ell)} \circ \Pi_\ell$ and $\mathbb{P}_1^{(S, \ell)} \circ \Pi_\ell$, respectively. Denoting by $f_0^{S|\ell=z}$ and $f_1^{S|\ell=z}$, respectively, a conditional density of S given $\ell = z$ under the joint distribution $\mathbb{P}_1^{(S, \ell)}$ and $\mathbb{P}_0^{(S, \ell)}$ (see **Notation** §21.11 (iv)) we have $f_1^{S|\ell=z}(s) f_1^\ell(z) = \exp(z) f_0^{S|\ell=z}(s) f_0^\ell(z)$ for λ^{k+1} -a.e. $(s, z) \in \mathbb{R}^{k+1}$. Exploiting **Theorem** §04.29 (b'3) it holds $f_1^\ell(z) = \exp(z) f_0^\ell(z)$ for λ -a.e. $z \in \mathbb{R}$ (see **Remark** §04.31). Consequently, it remains to verify that $N_{(v, M)}$ and $N_{(v', M)}$ have the same conditional distribution given $\ell = z$. Indeed, both are again multivariate normal (see **Notation** §21.11 (v)) with equal covariance matrix $\Sigma - \sigma^2 \tau \tau^t$ and conditional mean $\mathbb{P}_0^{S|\ell=z}(\text{id}_{\mathbb{R}^k}) = \mu + \sigma^{-2} \tau (z + \sigma^2/2) = \mu + \tau + \sigma^{-2} \tau (z - \sigma^2/2) = \mathbb{P}_1^{S|\ell=z}(\text{id}_{\mathbb{R}^k})$. The case of a positive semi-definite Σ and $\sigma^2 > 0$ follows by similar arguments when considering the projection onto the image of Σ . If $\sigma = 0$ the claim follows from **Lemma** §04.19 (i) together with Slutsky's lemma §20.10. In particular, note that $S_n \xrightarrow{d} N_{(\mu, \Sigma)}$ under \mathbb{P}_0^n and $S_n \xrightarrow{d} N_{(\mu + \tau, \Sigma)}$ under \mathbb{P}_1^n (see **Reminder** §04.07). \square

§05 Local asymptotic normality (LAN)

§05.01 **Aim.** For each $n \in \mathbb{N}$ let $(\mathcal{X}_n, \mathcal{X}_n, \mathbb{P}_\Theta^n = (\mathbb{P}_\theta^n)_{\theta \in \Theta})$ with $\Theta \subseteq \mathbb{R}^k$ be a statistical experiment. We aim to approximate $(\mathcal{X}_n, \mathcal{X}_n, \mathbb{P}_\Theta^n)$ in a certain sense by a Gaussian location model after suitable reparametrisation.

§05.02 **Reminder.** Consider on $(\mathbb{R}^k, \mathcal{B}^k)$ the family $N_{\mathbb{R}^k \times \{\Sigma\}} := (N_{(h, \Sigma)})_{h \in \mathbb{R}^k}$ of multivariate normal distributions with common strictly positive definite covariance matrix $\Sigma \in \mathbb{R}_{>}^{(k, k)}$ and log-likelihood ratio $\log(dN_{(h, \Sigma)}/dN_{(0, \Sigma)})(z) = \langle \Sigma^{-1}h, z \rangle - \frac{1}{2} \langle \Sigma^{-1}h, h \rangle$, $z \in \mathbb{R}^k$. Noting that for each $h \in \mathbb{R}^k$ the likelihood $L(h) = dN_{(h, \Sigma)}/d\lambda^k$ of $N_{(h, \Sigma)}$ with respect to the Lebesgue measure λ^k on \mathbb{R}^k satisfies $L(h, x) = L(0, x - h)$ for all $x \in \mathbb{R}^k$ the statistical experiment $(\mathbb{R}^k, \mathcal{B}^k, N_{\mathbb{R}^k \times \{\Sigma\}})$ is called a *Gaussian location model*. \square

Consider a localised reparametrisation centred around a parameter value $\theta_o \in \text{int}(\Theta)$ which is in the sequel regarded as fixed.

§05.03 **Definition.** Consider a sequence of statistical experiments $(\mathcal{X}_n, \mathcal{X}_n, \mathbb{P}_\Theta^n)$, $n \in \mathbb{N}$, with common parameter set $\Theta \subseteq \mathbb{R}^k$. Given a *localising rate* $(\delta_n)_{n \in \mathbb{N}}$ with $\delta_n = o(1)$ for each $n \in \mathbb{N}$ define a *local parameter set* $\Theta_o^n := \{\delta_n^{-1}(\theta - \theta_o) : \theta \in \Theta\} \subseteq \mathbb{R}^k$. For each $\theta \in \Theta$ and associated *local parameter* $h = \delta_n^{-1}(\theta - \theta_o) \in \Theta_o^n$ rewriting \mathbb{P}_θ^n as $\mathbb{P}_{\theta_o + \delta_n h}^n$ we obtain a sequence of *localised statistical experiment* $(\mathcal{X}_n, \mathcal{X}_n, \mathbb{P}_{\delta_n \Theta_o^n + \theta_o}^n := (\mathbb{P}_{\theta_o + \delta_n h}^n)_{h \in \Theta_o^n})$, $n \in \mathbb{N}$. \square

§05.04 **Remark.** In the sequel we eventually take the local parameter set Θ_o^n equal to \mathbb{R}^k which is not correct if the parameter set Θ is a strict subset of \mathbb{R}^k . However, if $\theta_o \in \text{int}(\Theta)$ is an inner point of Θ , which is assumed throughout this section, then for each $h \in \mathbb{R}^k$ the parameter $\theta = \theta_o + \delta_n h$ belongs to Θ for every sufficiently large n . In other words, the local parameter set Θ_o^n converges to the whole of \mathbb{R}^k as $n \rightarrow \infty$, i.e., $\cup_{n \in \mathbb{N}} \Theta_o^n = \mathbb{R}^k$. Thereby, we tactically may either define the probability measure $\mathbb{P}_{\theta_o + \delta_n h}^n$ arbitrarily if $\theta_o + \delta_n h$ does not belong to Θ , or assume that n is sufficiently large. \square

§05.05 **Aim.** We show, for large n , that the localised statistical experiment $(\mathcal{X}_n, \mathcal{X}_n, \mathbb{P}_{\delta_n \mathbb{R}^k + \theta_o}^n)$ and a Gaussian location model $(\mathbb{R}^k, \mathcal{B}^k, N_{\mathbb{R}^k \times \{\sigma_o^{-1}\}})$ are similar in statistical properties whenever the original experiments, i.e., $\theta \mapsto \mathbb{P}_\theta$, are “smooth”.

§05.06 **Heuristics.** Consider a statistical experiment $(\mathcal{X}, \mathcal{X}, \mathbb{P}_\Theta)$ dominated by $\mu \in \mathcal{M}_\sigma(\mathcal{X})$, i.e., $\mathbb{P}_\Theta \ll \mu$, with $\Theta \subseteq \mathbb{R}$, positive real likelihood $L(\theta) = d\mathbb{P}_\theta/d\mu \in \mathcal{X}^+$ and log-likelihood $\ell = \log L$. Assume that for all $x \in \mathcal{X}$, the map $\theta \mapsto \ell(\theta, x)$ is twice differentiable with derivatives $\dot{\ell}_\theta := \frac{\partial}{\partial \theta} \ell$ and $\ddot{\ell}_\theta := \frac{\partial^2}{\partial^2 \theta} \ell$. A *Taylor expansion* of the log-likelihood ratio leads to $\ell(\theta + h, x) - \ell(\theta, x) = h \dot{\ell}_\theta(x) + \frac{1}{2} h^2 \ddot{\ell}_\theta(x) + o_x(h^2)$ where the remainder term depends on x . Considering a product experiment $(\mathcal{X}^n, \mathcal{X}^{\otimes n}, \mathbb{P}_\Theta^{\otimes n})$ eventually it holds $\log(d\mathbb{P}_{\theta+h/\sqrt{n}}^{\otimes n}/d\mathbb{P}_\theta^{\otimes n}) = h \sqrt{n} \widehat{\mathbb{P}}_n(\dot{\ell}_\theta) + \frac{1}{2} h^2 \widehat{\mathbb{P}}_n(\ddot{\ell}_\theta) + R_n$ where the score $\dot{\ell}$ has mean zero, i.e., $\mathbb{P}_\theta(\dot{\ell}_\theta) = 0$, and the Fisher information \mathcal{J}_θ equals $-\mathbb{P}_\theta(\ddot{\ell}_\theta) = \mathbb{P}_\theta(|\dot{\ell}_\theta|^2)$. Setting $\mathcal{Z}_\theta^n := \sqrt{n} \widehat{\mathbb{P}}_n(\dot{\ell}_\theta)$ from the central limit theorem §20.13 follows $\mathcal{Z}_\theta^n \xrightarrow{d} N_{(0, \mathcal{J}_\theta)}$ under $\mathbb{P}_\theta^{\otimes n}$ while due to the law of large numbers §20.06 it holds $\widehat{\mathbb{P}}_n \ddot{\ell}_\theta = -\mathcal{J}_\theta + o_{\mathbb{P}^{\otimes n}}(1)$. If in addition the remainder term is negligible, i.e., $R_n = o_{\mathbb{P}^{\otimes n}}(1)$, then the log-likelihood ratio permits an expansion

$$\log(d\mathbb{P}_{\theta+h/\sqrt{n}}^{\otimes n}/d\mathbb{P}_\theta^{\otimes n}) = h \mathcal{Z}_\theta^n - \frac{1}{2} h^2 \mathcal{J}_\theta + o_{\mathbb{P}^{\otimes n}}(1)$$

which in the limit equals the log-likelihood ratio in a Gaussian location model. \square

§05.07 **Definition.** A sequence of statistical experiments $(\mathcal{X}_n, \mathcal{X}_n, \mathbb{P}_\theta^n)_{n \in \mathbb{N}}$ with $\Theta \subseteq \mathbb{R}^k$ is called *local asymptotic normal (LAN)* in $\theta_o \in \text{int}(\Theta)$, if there is a localising rate $(\delta_n)_{n \in \mathbb{N}}$ with $\delta_n = o(1)$, a sequence of statistics $(\mathcal{Z}_{\theta_o}^n)_{n \in \mathbb{N}} \in (\mathcal{X}_n^k)_{n \in \mathbb{N}}$ and a matrix $\mathcal{J}_{\theta_o} \in \mathbb{R}^{(k,k)}$ such that for every $h \in \mathbb{R}^k$ the following three statements hold true:

- (a) $\theta_o + \delta_n h \in \Theta$ for all sufficiently large n , i.e., $n \geq n_o(h)$;
- (b) $\mathcal{Z}_{\theta_o}^n \xrightarrow{d} N_{(0, \mathcal{J}_{\theta_o})}$ under $\mathbb{P}_{\theta_o}^n$, i.e., $\mathbb{P}_{\theta_o}^n \circ (\mathcal{Z}_{\theta_o}^n)^{-1} \xrightarrow{d} N_{(0, \mathcal{J}_{\theta_o})}$;
- (c) $\log(d\mathbb{P}_{\theta_o + \delta_n h}^n / d\mathbb{P}_{\theta_o}^n) = \langle \mathcal{Z}_{\theta_o}^n, h \rangle - \frac{1}{2} \langle \mathcal{J}_{\theta_o} h, h \rangle + R_{n,h}$ where $R_{n,h} = o_{\mathbb{P}_{\theta_o}^n}(1)$.

The matrix \mathcal{J}_{θ_o} and the sequence of statistics $(\mathcal{Z}_{\theta_o}^n)_{n \in \mathbb{N}}$ is called, respectively, *Fisher information* at θ_o and *central sequence*. \square

§05.08 **Comment.** If we assume in addition a strictly positive definite matrix $\mathcal{J}_{\theta_o} \in \mathbb{R}_>^{(k,k)}$ with inverse $\mathcal{J}_{\theta_o}^{-1}$ the sequence of statistics $(\tilde{\mathcal{Z}}_{\theta_o}^n := \mathcal{J}_{\theta_o}^{-1} \mathcal{Z}_{\theta_o}^n)_{n \in \mathbb{N}} \in (\mathcal{X}_n^k)_{n \in \mathbb{N}}$ is equally a central sequence satisfying $\tilde{\mathcal{Z}}_{\theta_o}^n \xrightarrow{d} N_{(0, \mathcal{J}_{\theta_o}^{-1})}$ under $\mathbb{P}_{\theta_o}^n$ and $\log(d\mathbb{P}_{\theta_o + \delta_n h}^n / d\mathbb{P}_{\theta_o}^n) = \langle \mathcal{J}_{\theta_o} h, \tilde{\mathcal{Z}}_{\theta_o}^n \rangle - \frac{1}{2} \langle \mathcal{J}_{\theta_o} h, h \rangle + o_{\mathbb{P}_{\theta_o}^n}(1)$. In other words the likelihood ratio $d\mathbb{P}_{\theta_o + \delta_n h}^n / d\mathbb{P}_{\theta_o}^n$ equals approximately the likelihood ratio $dN_{(h, \mathcal{J}_{\theta_o}^{-1})} / dN_{(0, \mathcal{J}_{\theta_o}^{-1})}$ as in the **Reminder** §05.02. Consequently, the localised statistical model $(\mathcal{X}_n, \mathcal{X}_n, \mathbb{P}_{\delta_n \Theta_o^n + \theta_o}^n)$ is similar to a Gaussian location model $(\mathbb{R}^k, \mathcal{B}^{\otimes k}, N_{\mathbb{R}^k \times \{\mathcal{J}_{\theta_o}^{-1}\}})$ in the sense of **Definition** §05.07. \square

§05.09 **Definition.** A LAN sequence of statistical experiments is called *uniformly local asymptotic normal (ULAN)* in $\theta_o \in \Theta$, if the condition (c) in **Definition** §05.07 is replaced by (c') for $h_n \rightarrow h$ it holds $\log(d\mathbb{P}_{\theta_o + \delta_n h_n}^n / d\mathbb{P}_{\theta_o}^n) = \langle \mathcal{Z}_{\theta_o}^n, h \rangle - \frac{1}{2} \langle \mathcal{J}_{\theta_o} h, h \rangle + o_{\mathbb{P}_{\theta_o}^n}(1)$. \square

§05.10 **Theorem.** Let $(\mathcal{X}_n, \mathcal{X}_n, \mathbb{P}_\theta^n)_{n \in \mathbb{N}}$ be LAN in $\theta_o \in \Theta \subseteq \mathbb{R}^k$ with localising rate $(\delta_n)_{n \in \mathbb{N}}$, central sequence $(\mathcal{Z}_{\theta_o}^n)_{n \in \mathbb{N}}$ and Fisher information matrix $\mathcal{J}_{\theta_o} \in \mathbb{R}^{(k,k)}$. Then for any $h, h' \in \mathbb{R}^k$ the following statements hold true:

- (i) $(\mathbb{P}_{\theta_o + \delta_n h}^n)_{n \in \mathbb{N}}$ and $(\mathbb{P}_{\theta_o + \delta_n h'}^n)_{n \in \mathbb{N}}$ are mutually contiguous, i.e., $\mathbb{P}_{\theta_o + \delta_n h}^n \triangleleft \triangleright \mathbb{P}_{\theta_o + \delta_n h'}^n$;
- (ii) $\mathcal{Z}_{\theta_o}^n \xrightarrow{d} N_{(\mathcal{J}_{\theta_o} h, \mathcal{J}_{\theta_o})}$ under $\mathbb{P}_{\theta_o + \delta_n h}^n$.

If the sequence of statistical experiments is ULAN, then for any $h_n \rightarrow h$ and $h'_n \rightarrow h'$ in \mathbb{R}^k the following statements hold true:

- (i') $(\mathbb{P}_{\theta_o + \delta_n h_n}^n)_{n \in \mathbb{N}}$ and $(\mathbb{P}_{\theta_o + \delta_n h'_n}^n)_{n \in \mathbb{N}}$ are mutually contiguous, i.e., $\mathbb{P}_{\theta_o + \delta_n h_n}^n \triangleleft \triangleright \mathbb{P}_{\theta_o + \delta_n h'_n}^n$;
- (ii') $\mathcal{Z}_{\theta_o}^n \xrightarrow{d} N_{(\mathcal{J}_{\theta_o} h, \mathcal{J}_{\theta_o})}$ under $\mathbb{P}_{\theta_o + \delta_n h_n}^n$.

§05.11 **Proof** of **Theorem** §05.10. is given in the lecture. \square

§05.12 **Theorem.** Let $\mathbb{P}_\theta \ll \mu \in \mathcal{M}_\sigma(\mathcal{X})$ with open $\Theta \subseteq \mathbb{R}^k$ be Hellinger-differentiable in $\theta_o \in \Theta$ with derivative $\dot{\ell}_{\theta_o}$ and Fisher information matrix $\mathcal{J}_{\theta_o} = \mathbb{E}_{\theta_o}(\dot{\ell}_{\theta_o} \dot{\ell}_{\theta_o}^t) \in \mathbb{R}_>^{(k,k)}$. Then the sequence of product experiments $(\mathcal{X}^n, \mathcal{X}^{\otimes n}, \mathbb{P}_\theta^{\otimes n})$ is ULAN in θ_o with localising rate $\delta_n := n^{-1/2}$ and central sequence $\mathcal{Z}_{\theta_o}^n := \sqrt{n} \widehat{\mathbb{P}}_n(\dot{\ell}_{\theta_o})$, $n \in \mathbb{N}$, that is,

- (i) $\sqrt{n} \widehat{\mathbb{P}}_n(\dot{\ell}_{\theta_o}) \xrightarrow{d} N_{(0, \mathcal{J}_{\theta_o})}$ under $\mathbb{P}_{\theta_o}^{\otimes n}$ and
- (ii) for $h_n \rightarrow h$ it holds $\log(d\mathbb{P}_{\theta_o + h_n / \sqrt{n}}^{\otimes n} / d\mathbb{P}_{\theta_o}^{\otimes n}) = \langle \mathcal{Z}_{\theta_o}^n, h \rangle - \frac{1}{2} \langle \mathcal{J}_{\theta_o} h, h \rangle + o_{\mathbb{P}_{\theta_o}^{\otimes n}}(1)$.

§05.13 **Proof** of **Theorem** §05.12. is given in the lecture. \square

§05.14 **Corollary.** Under the assumptions of **Theorem** §05.12 consider for each $n \in \mathbb{N}$ a statistical product experiment $(\mathcal{X}^n, \mathcal{X}^{\otimes n}, \mathbb{P}_\theta^{\otimes n})$ and an estimator $\widehat{\gamma}_n \in (\mathcal{X}^{\otimes n})^p$ of a parameter of interest $\gamma : \Theta \rightarrow \mathbb{R}^p$ allowing an expansion $\sqrt{n}(\widehat{\gamma}_n - \gamma(\theta_o)) = \sqrt{n} \widehat{\mathbb{P}}_n(\psi_{\theta_o}) + o_{\mathbb{P}_{\theta_o}^{\otimes n}}(1)$ for some function $\psi_{\theta_o} \in$

$\mathcal{L}_2^p(\mathbb{P}_{\theta_o})$ with $\mathbb{P}_{\theta_o}(\psi_{\theta_o}) = 0$. Then, $\sqrt{n}(\widehat{\gamma}_n - \gamma(\theta_o)) \xrightarrow{d} N_{(0, \Sigma_o)}$ under $\mathbb{P}_{\theta_o}^{\otimes n}$ with $\Sigma_o := \mathbb{P}_{\theta_o}(\psi_{\theta_o} \psi_{\theta_o}^t)$ and for each $h \in \mathbb{R}^k$ holds $\sqrt{n}(\widehat{\gamma}_n - \gamma(\theta_o)) \xrightarrow{d} N_{(\tau_h, \Sigma_o)}$ under $\mathbb{P}_{\theta_o+h/\sqrt{n}}^{\otimes n}$ with $\tau_h := \mathbb{P}_{\theta_o}(\psi_{\theta_o} \dot{\ell}_{\theta_o}^t)h$.

§05.15 **Proof** of **Corollary** §05.14. is given in the lecture. □

§05.16 **Example** (*Example §03.06 continued*). Under the assumptions of **Theorem** §05.12 let $\gamma : \theta \rightarrow \mathbb{R}^p$ be a parameter of interest. Consider $m(\gamma) \in \mathcal{L}_1(\mathbb{P}_{\theta})$ for all $\gamma \in \mathbb{R}^p$, a criterion process $\widehat{M}_n(\gamma) = \widehat{\mathbb{P}}_n(m(\gamma))$, a criterion function $M(\theta, \gamma) = \mathbb{P}_{\theta}(m(\gamma))$ and a M-estimator $\widehat{\gamma}_n \in \arg \inf_{\gamma \in \Gamma} \{\widehat{M}_n(\gamma)\}$ of $\{\gamma_o := \gamma(\theta_o)\} = \arg \inf_{\gamma \in \Gamma} \{M(\theta_o, \gamma)\}$. Under regularity conditions as in **Example** §03.06 we have $\sqrt{n}(\widehat{\gamma}_n - \gamma_o) = \sqrt{n}\widehat{\mathbb{P}}_n(\psi_{\theta_o}) + o_{\mathbb{P}_{\theta_o}^{\otimes n}}(1)$ with $\psi_{\theta_o} := -\ddot{M}_o^{-1}\dot{m}(\gamma_o)$ assuming a regular matrix $\ddot{M}_o := \mathbb{P}_{\theta_o}(\ddot{m}(\gamma_o))$. Consequently, setting $\Sigma_o = \mathbb{P}_{\theta_o}(\psi_{\theta_o} \psi_{\theta_o}^t) = \ddot{M}_o^{-1}\mathbb{P}_{\theta_o}(\dot{m}(\gamma_o)\dot{m}(\gamma_o)^t)\ddot{M}_o^{-1}$ from **Corollary** §05.14 it follows

$$\sqrt{n}(\widehat{\gamma}_n - \gamma_o) \xrightarrow{d} N_{(\tau_h, \Sigma_o)} \quad \text{under} \quad \mathbb{P}_{\theta_o+h/\sqrt{n}}^{\otimes n} \quad \text{with} \quad \tau_h = -\ddot{M}_o^{-1}\mathbb{P}_{\theta_o}(\dot{m}(\gamma_o)\dot{\ell}_{\theta_o}^t)h.$$

In the particular case of a MLE $\widehat{\theta}_n$ of θ , i.e., $(\gamma = \text{id}_{\mathbb{R}^k})$, as in **Example** §03.07 setting $m := -\log(d\mathbb{P}_{\theta}/d\mathbb{P}_o)$ we have $\dot{m}(\theta_o) = -\dot{\ell}_{\theta_o}$, $\mathcal{J}_{\theta_o} = \mathbb{P}_{\theta_o}(\dot{m}(\theta_o)\dot{m}(\theta_o)^t) = \mathbb{P}_{\theta_o}(\dot{m}(\theta_o)) = \ddot{M}_o$ and thus $\Sigma_o = \ddot{M}_o^{-1}\mathbb{P}_{\theta_o}(\dot{m}(\gamma_o)\dot{m}(\gamma_o)^t)\ddot{M}_o^{-1} = \mathcal{J}_{\theta_o}^{-1}$ and $\tau_h := -\ddot{M}_o^{-1}\mathbb{P}_{\theta_o}(\dot{m}(\theta_o)\dot{\ell}_{\theta_o}^t)h = h$. Therewith, $\sqrt{n}(\widehat{\theta}_n - \theta_o) \xrightarrow{d} N_{(h, \mathcal{J}_{\theta_o}^{-1})}$ under $\mathbb{P}_{\theta_o+h/\sqrt{n}}^{\otimes n}$. □

§05.17 **Remark**. Supposing $\sqrt{n}(\widehat{\theta}_n - \theta_o) = \sqrt{n}\widehat{\mathbb{P}}_n(\psi_{\theta_o}) + o_{\mathbb{P}_{\theta_o}^{\otimes n}}(1)$ let us further assume a transformation $A : \Theta \rightarrow \mathbb{R}^p$ that is “smooth”, and hence by employing the delta method §20.16, for instance satisfies $\sqrt{n}(A(\widehat{\theta}_n) - A(\theta_o)) = \dot{A}_{\theta_o}\sqrt{n}\widehat{\mathbb{P}}_n(\psi_{\theta_o}) + o_{\mathbb{P}_{\theta_o}^{\otimes n}}(1)$. Consequently, it follows $\sqrt{n}(A(\widehat{\theta}_n) - A(\theta_o)) \xrightarrow{d} N_{(\tau_h, \Sigma_o)}$ under $\mathbb{P}_{\theta_o+h/\sqrt{n}}^{\otimes n}$ with $\tau_h = \dot{A}_{\theta_o}\mathbb{P}_{\theta_o}(\psi_{\theta_o}\dot{\ell}_{\theta_o}^t)h$ and $\Sigma_o = \dot{A}_{\theta_o}\mathbb{P}_{\theta_o}(\psi_{\theta_o}\psi_{\theta_o}^t)\dot{A}_{\theta_o}^t$. In the special case of a MLE we have $\sqrt{n}(A(\widehat{\theta}_n) - A(\theta_o)) \xrightarrow{d} N_{(\dot{A}_{\theta_o}h, \dot{A}_{\theta_o}\mathcal{J}_{\theta_o}^{-1}\dot{A}_{\theta_o}^t)}$ under $\mathbb{P}_{\theta_o+h/\sqrt{n}}^{\otimes n}$. □

§06 Asymptotic relative efficiency

§06.01 **Heuristics** (§03.09 and §03.10 continued). Under the conditions of **Corollary** §05.14 consider the statistical testing task $H_0 : A(\theta_o) = 0$ against the alternative $H_1 : A(\theta_o) \neq 0$ for some transformation $A : \Theta \rightarrow \mathbb{R}^p$ satisfying $\sqrt{n}(A(\widehat{\theta}_n) - A(\theta_o)) = \dot{A}_{\theta_o}\sqrt{n}\widehat{\mathbb{P}}_n(\psi_{\theta_o}) + o_{\mathbb{P}_{\theta_o}^{\otimes n}}(1)$. As in §03.09 let $\widehat{W}_n := nA(\widehat{\theta}_n)^t\widehat{\Sigma}_n^{-1}A(\widehat{\theta}_n)$ where $\widehat{\Sigma}_n = \Sigma + o_{\mathbb{P}_{\theta_o}^{\otimes n}}(1)$ is a consistent estimator of $\Sigma = \dot{A}_{\theta_o}\mathbb{P}_{\theta_o}(\psi_{\theta_o}\psi_{\theta_o}^t)\dot{A}_{\theta_o}^t$, then a **Wald test** is given by $\varphi_n := \mathbf{1}_{\{\widehat{W}_n > \chi_{p,1-\alpha}^2\}}$. Thereby, under H_0 , i.e. $A(\theta_o) = 0$, we have $\sqrt{n}A(\widehat{\theta}_n) = \dot{A}_{\theta_o}\sqrt{n}\widehat{\mathbb{P}}_n(\psi_{\theta_o}) + o_{\mathbb{P}_{\theta_o}^{\otimes n}}(1)$ and $\widehat{W}_n \xrightarrow{d} \chi_p^2$ under $\mathbb{P}_{\theta_o}^{\otimes n}$ which in turn implies $\mathbb{P}_{\theta_o}^{\otimes n}(\varphi_n = 1) \xrightarrow{n \rightarrow \infty} \chi_p^2((\chi_{p,1-\alpha}^2, \infty)) = \alpha$. In other words, the Wald test is asymptotically a level α test. For each $\theta \in \Theta$ let us denote $\beta_{\varphi_n}(\theta) := \mathbb{P}_{\theta}^{\otimes n}(\varphi_n) = \mathbb{P}_{\theta}^{\otimes n}(\varphi_n = 1) = \mathbb{P}_{\theta}^{\otimes n}(\widehat{W}_n > \chi_{p,1-\alpha}^2)$ which equals the power of the Wald test φ_n under H_1 , i.e. $\theta \in \Theta$ with $A(\theta) \neq 0$. In the sequel we consider local alternatives of the form $\theta = \theta_o + h/\sqrt{n}$ and thus we are interested in $\beta_{\varphi_n}(\theta_o + h/\sqrt{n}) = \mathbb{P}_{\theta_o+h/\sqrt{n}}^{\otimes n}(\widehat{W}_n > \chi_{p,1-\alpha}^2)$. Keeping **Remark** §05.17 under $\mathbb{P}_{\theta_o+h/\sqrt{n}}^{\otimes n}$ we have $\sqrt{n}A(\widehat{\theta}_n) \xrightarrow{d} N_{(\dot{A}_{\theta_o}\mathbb{P}_{\theta_o}(\psi_{\theta_o}\dot{\ell}_{\theta_o}^t)h, \Sigma)}$, assuming additionally $\Sigma > 0$ also $\Sigma^{-1/2}\sqrt{n}A(\widehat{\theta}_n) \xrightarrow{d} N_{(a_h, \text{Id}_p)}$ with $a_h := \Sigma^{-1/2}\dot{A}_{\theta_o}\mathbb{P}_{\theta_o}(\psi_{\theta_o}\dot{\ell}_{\theta_o}^t)h$, and hence, $nA(\widehat{\theta}_n)^t\Sigma^{-1}A(\widehat{\theta}_n) \xrightarrow{d} \chi_p^2(\|a_h\|^2)$. Here $\chi_p^2(c)$ denotes a non-central χ^2 -distribution with p degrees of freedom and non-centrality parameter $c \in \mathbb{R}^+$. Moreover, $\widehat{W}_n - nA(\widehat{\theta}_n)^t\Sigma^{-1}A(\widehat{\theta}_n) = o_{\mathbb{P}_{\theta_o}^{\otimes n}}(1)$ and thus $\widehat{W}_n - nA(\widehat{\theta}_n)^t\Sigma^{-1}A(\widehat{\theta}_n) = o_{\mathbb{P}_{\theta_o+h/\sqrt{n}}^{\otimes n}}(1)$ due to **Lemma** §04.19 (ii) by employing that $\mathbb{P}_{\theta_o}^{\otimes n} \triangleleft \mathbb{P}_{\theta_o+h/\sqrt{n}}^{\otimes n}$ are mutually contiguous. Consequently, $\widehat{W}_n \xrightarrow{d} \chi_p^2(\|a_h\|^2)$ under $\mathbb{P}_{\theta_o+h/\sqrt{n}}^{\otimes n}$ and thus $\beta_{\varphi_n}(\theta_o + h/\sqrt{n}) \xrightarrow{n \rightarrow \infty} \chi_p^2(\|a_h\|^2)((\chi_{p,1-\alpha}^2, \infty))$. Note that a_h simplifies to $h^t\dot{A}_{\theta_o}^t(\dot{A}_{\theta_o}\mathcal{J}_{\theta_o}^{-1}\dot{A}_{\theta_o}^t)^{-1}\dot{A}_{\theta_o}h$ in the particular case of a MLE $\widehat{\theta}_n$. □

§06.02 **Reminder (Gauß test).** In a Gaussian location model, i.e. $Y \odot N_{\mathbb{R}^k \times \{\mathcal{J}_{\theta_o}^{-1}\}}$ with $\mathcal{J}_{\theta_o} \in \mathbb{R}_{>}^{(k,k)}$, consider the binary testing task $H_0 : \{N_{(0, \mathcal{J}_{\theta_o}^{-1})}\}$ against the alternative $H_1 : \{N_{(h, \mathcal{J}_{\theta_o}^{-1})}\}$ for some $h \in \mathbb{R}^k$. In this situation the log-likelihood ratio $\ell_h = \log(dN_{(h, \mathcal{J}_{\theta_o}^{-1})}/dN_{(0, \mathcal{J}_{\theta_o}^{-1})})$ satisfies $\ell_h(y) = \langle \mathcal{J}_{\theta_o} y, h \rangle - \frac{1}{2} \sigma_h^2$ for all $y \in \mathbb{R}^k$ with $\sigma_h^2 := \langle \mathcal{J}_{\theta_o} h, h \rangle$. Consequently, $\ell_h \sim N_{(-\sigma_h^2/2, \sigma_h^2)}$ under $N_{(0, \mathcal{J}_{\theta_o}^{-1})}$, i.e. under the hypothesis H_0 , and $\ell_h \sim N_{(\sigma_h^2/2, \sigma_h^2)}$ under $N_{(h, \mathcal{J}_{\theta_o}^{-1})}$, i.e. under the alternative H_1 . For $\alpha \in (0, 1)$ let $c_{h, 1-\alpha} \in \mathbb{R}$ satisfy $N_{(-\sigma_h^2/2, \sigma_h^2)}((c_{h, 1-\alpha}, \infty)) = \alpha$ and thus $N_{(0, \mathcal{J}_{\theta_o}^{-1})}(\ell_h > c_{h, 1-\alpha}) = N_{(-\sigma_h^2/2, \sigma_h^2)}((c_{h, 1-\alpha}, \infty)) = \alpha$. Keeping in mind that any most powerful level- α test has Neyman-Pearson form and the **Gauß test** $\varphi^* := \mathbb{1}_{\{\ell_h > c_{h, 1-\alpha}\}}$ is a Neyman-Pearson level- α test. Its power given by $\beta_{\varphi^*}(h) := N_{(h, \mathcal{J}_{\theta_o}^{-1})}(\varphi^* = 1) = N_{(h, \mathcal{J}_{\theta_o}^{-1})}(\ell_h > c_{h, 1-\alpha}) = N_{(\sigma_h^2/2, \sigma_h^2)}((c_{h, 1-\alpha}, \infty))$ is maximal in the class of all level- α tests, i.e., for any level- α test φ holds $\beta_{\varphi}(h) \leq \beta_{\varphi^*}(h)$. In other words, φ^* is a most powerful level- α test (Statistik 1, Satz §18.16, p.56). \square

§06.03 **Example (Neyman-Pearson test).** Assume local asymptotic normality as in Definition §05.07 where $\ell_{h,n} := \log(d\mathbb{P}_{\theta_o + \delta_n h}^n/d\mathbb{P}_{\theta_o}^n) \xrightarrow{d} N_{(-\sigma_h^2/2, \sigma_h^2)}$ under $\mathbb{P}_{\theta_o}^n$ with $\sigma_h^2 := \langle \mathcal{J}_{\theta_o} h, h \rangle$ for $h \in \mathbb{R}^k$. Hence by Le Cam's first lemma (Example §04.34) mutual contiguity $\mathbb{P}_{\theta_o + \delta_n h}^n \triangleleft \triangleright \mathbb{P}_{\theta_o}^n$ and $\ell_{h,n} \xrightarrow{d} N_{(\sigma_h^2/2, \sigma_h^2)}$ under $\mathbb{P}_{\theta_o + \delta_n h}^n$ hold. Consider the binary testing task of the hypothesis $H_0 : \{\mathbb{P}_{\theta_o}^n\}$ against a local alternative $H_1 : \{\mathbb{P}_{\theta_o + \delta_n h}^n\}$. In this situation $\varphi_n^* = \mathbb{1}_{\{\ell_{h,n} > c_{h,n, 1-\alpha}\}}$ is a **Neyman-Pearson test**, which is a most powerful level- α test, if $\mathbb{P}_{\theta_o}^n(\varphi_n^* = 1) = \alpha$. Keeping its power function $\beta_{\varphi_n^*}(\theta) = \mathbb{P}_{\theta}^n(\varphi_n^*) = \mathbb{P}_{\theta}^n(\varphi_n^* = 1) = \mathbb{P}_{\theta}^n(\ell_{h,n} > c_{h,n, 1-\alpha})$ evaluated at θ in mind the value $\beta_{\varphi_n^*}(\theta_o + \delta_n h)$ equals the maximal size of the power in the class of all level- α tests. Considering $c_{h, 1-\alpha} \in \mathbb{R}$ as in Reminder §06.02 under local asymptotic normality it follows $\alpha = \mathbb{P}_{\theta_o}^n(\varphi_n^*) = \mathbb{P}_{\theta_o}^n(\ell_{h,n} > c_{h,n, 1-\alpha}) \xrightarrow{n \rightarrow \infty} N_{(-\sigma_h^2/2, \sigma_h^2)}((c_{h, 1-\alpha}, \infty)) = \alpha$ which implies $c_{h,n, 1-\alpha} \xrightarrow{n \rightarrow \infty} c_{h, 1-\alpha}$, and in addition $\beta_{\varphi_n^*}(\theta_o + \delta_n h) = \mathbb{P}_{\theta_o + \delta_n h}^n(\varphi_n^*) = \mathbb{P}_{\theta_o + \delta_n h}^n(\ell_{h,n} > c_{h,n, 1-\alpha}) \xrightarrow{n \rightarrow \infty} N_{(\sigma_h^2/2, \sigma_h^2)}((c_{h, 1-\alpha}, \infty)) = \beta_{\varphi^*}(h)$ with Neyman-Pearson test φ^* in a Gaussian location model as in Reminder §06.02. \square

§06.04 **Theorem.** Let $\Theta \subseteq \mathbb{R}$. Consider a one-sided test task $H_0 : (-\infty, \theta_o]$ against $H_1 : (\theta_o, \infty)$. Suppose that $(\mathcal{X}_n, \mathcal{Z}_n, \mathbb{P}_{\theta}^n)$ is LAN in $\theta_o \in \Theta$ with localising sequence $(\delta_n)_{n \in \mathbb{N}}$, central sequence $(\mathcal{Z}_{\theta_o}^n)_{n \in \mathbb{N}} \in (\mathcal{X}_n)_{n \in \mathbb{N}}$ and strictly positive Fisher information $\mathcal{J}_{\theta_o} \in \mathbb{R}_{>_0}^+$.

(i) Given a sequence $(T_n)_{n \in \mathbb{N}} \in (\mathcal{X}_n)_{n \in \mathbb{N}}$ of test statistics satisfying $(T_n, \mathcal{Z}_{\theta_o}^n) \xrightarrow{d} N_{(0, M)}$ with $M = ((\sigma^2, \rho)^t, (\rho, \mathcal{J}_{\theta_o})^t)$ consider the randomised test $\varphi_n := \mathbb{1}_{\{T_n > c_n\}} + \gamma_n \mathbb{1}_{\{T_n = c_n\}}$ with $\gamma_n \in [0, 1]$ and $c_n \in \mathbb{R}$ such that $\beta_{\varphi_n}(\theta_o) = \mathbb{P}_{\theta_o}^n(\varphi_n) = \mathbb{P}_{\theta_o}^n(T_n > c_n) + \gamma_n \mathbb{P}_{\theta_o}^n(T_n = c_n) = \alpha_n \xrightarrow{n \rightarrow \infty} \alpha$. Choosing $z_{1-\alpha} \in \mathbb{R}$ with $1 - \mathbb{F}_{[0, 1]}(z_{1-\alpha}) := N_{(0, 1)}((z_{1-\alpha}, \infty)) = \alpha$ we have

$$\beta_{\varphi_n}(\theta_o + \delta_n h) = \mathbb{P}_{\theta_o + \delta_n h}^n(\varphi_n) \xrightarrow{n \rightarrow \infty} \mathbb{F}_{[0, 1]}(-z_{1-\alpha} + h\rho/\sigma).$$

(ii) In case $T_n = \mathcal{Z}_{\theta_o}^n$ consider $\varphi_n^* = \mathbb{1}_{\{\mathcal{Z}_{\theta_o}^n > z_{1-\alpha} \mathcal{J}_{\theta_o}^{1/2}\}}$, i.e. $\gamma_n = 0$ and $c_n = z_{1-\alpha} \mathcal{J}_{\theta_o}^{1/2}$. Then $\beta_{\varphi_n^*}(\theta_o) = \mathbb{P}_{\theta_o}^n(\varphi_n^*) = \mathbb{P}_{\theta_o}^n(\mathcal{J}_{\theta_o}^{-1/2} \mathcal{Z}_{\theta_o}^n > z_{1-\alpha}) \xrightarrow{n \rightarrow \infty} 1 - \mathbb{F}_{[0, 1]}(z_{1-\alpha}) = \alpha$ and

$$\beta_{\varphi_n^*}(\theta_o + \delta_n h) = \mathbb{P}_{\theta_o + \delta_n h}^n(\varphi_n^*) \xrightarrow{n \rightarrow \infty} \mathbb{F}_{[0, 1]}(-z_{1-\alpha} + h \mathcal{J}_{\theta_o}^{1/2}).$$

§06.05 **Proof of Theorem §06.04.** is given in the lecture. \square

§06.06 **Remark.**

- (a) By using Theorem §04.35 directly it could be possible to calculate an asymptotic power of a test if $\log(d\mathbb{P}_{\theta_o + \delta_n h}^n/d\mathbb{P}_{\theta_o}^n) \xrightarrow{d} \mathbb{P}$ under $\mathbb{P}_{\theta_o}^n$ where \mathbb{P} equals not necessarily $N_{(0, 1)}$.
- (b) Let $(Y_1, Y_2) \sim N_{(0, M)}$ with $M = ((\sigma^2, \rho)^t, (\rho, \mathcal{J}_{\theta_o})^t)$ as in Theorem §06.04 (i), then $\rho^2 = |\text{Cov}(Y_1, Y_2)|^2 \leq \text{var}_m(Y_1) \text{var}_m(Y_2) = \sigma^2 \mathcal{J}_{\theta_o}$. Consequently, the test φ_n^* given in (ii) maximises the asymptotic power when considering only a randomised test φ_n as given in part (i) of Theorem §06.04. \square

- §06.07 **Theorem.** Let the assumptions of **Theorem** §06.04 be satisfied. Any test φ_n of the one-sided testing task $H_0 : (\infty, \theta_o]$ against $H_1 : (\theta_o, \infty)$ with $\beta_{\varphi_n}(\theta_o) := \mathbb{P}_{\theta_o}^n(\varphi_n) = \alpha_n \xrightarrow{n \rightarrow \infty} \alpha$ fulfils
- (i) $\limsup_{n \rightarrow \infty} \beta_{\varphi_n}(\theta_o + \delta_n h) \leq \mathbb{F}_{[0,1]}(-z_{1-\alpha} + h\sqrt{\mathcal{J}_{\theta_o}})$ for all $h \in \mathbb{R}_0^+$;
 - (ii) $\liminf_{n \rightarrow \infty} \beta_{\varphi_n}(\theta_o - \delta_n h) \geq \mathbb{F}_{[0,1]}(-z_{1-\alpha} - h\sqrt{\mathcal{J}_{\theta_o}})$ for all $h \in \mathbb{R}_0^+$.

§06.08 **Proof of Theorem** §06.07. is given in the lecture. □

§06.09 **Remark.** Keeping **Theorem** §06.07 in mind we call the test (sequence) $(\varphi_n^*)_{n \in \mathbb{N}}$ given in **Theorem** §06.04 (ii) asymptotically uniformly most powerful level- α test (sequence) in the class of all asymptotic level- α test (sequences). Its asymptotic power function equals $\mathbb{F}_{[0,1]}(-z_{1-\alpha} + h\sqrt{\mathcal{J}_{\theta_o}})$ which is the power function of the uniformly most powerful test of $H_0 : (-\infty, 0]$ against $H_1 : (0, \infty)$ in the limit Gaussian location experiment $(\mathbb{R}, \mathcal{B}, N_{\mathbb{R} \times \{\mathcal{J}_{\theta_o}^{-1}\}})$. □

§06.10 **Asymptotic relative efficiency.** Let $(\mathcal{X}_n, \mathcal{X}_n, \mathbb{P}_{\theta_o}^n)_{n \in \mathbb{N}}$ be LAN with localising rate $\delta_n := n^{-1/2}$. Consider a test φ_n^a satisfying the conditions of **Theorem** §06.04 (i) and hence, admitting an asymptotic power function such that $\beta_{\varphi_n^a}(\theta_o + h/\sqrt{n}) \xrightarrow{n \rightarrow \infty} \mathbb{F}_{[0,1]}(-z_{1-\alpha} + h\rho_a/\sigma_a)$. Thereby, choosing $\eta = h/\sqrt{n}$ the approximation $\beta_{\varphi_n^a}(\theta_o + \eta) \approx \mathbb{F}_{[0,1]}(-z_{1-\alpha} + \eta\sqrt{n}\rho_a/\sigma_a)$ is reasonable. In analogy, if φ_n^b is another test satisfying the conditions of **Theorem** §06.04 (i) and admitting $\beta_{\varphi_n^b}(\theta_o + \eta) \approx \mathbb{F}_{[0,1]}(-z_{1-\alpha} + \eta\sqrt{n}\rho_b/\sigma_b)$. Roughly speaking, this means, that at $\theta_o + \eta$ the power of the test φ_n^a and φ_n^b with sample size n_a and n_b , respectively, is approximately equal if $n_a\rho_a^2/\sigma_a^2 = n_b\rho_b^2/\sigma_b^2$. The quantity $\text{are}(\varphi_{n_a}^a, \varphi_{n_b}^b) = (n_a/n_b) = (\rho_b^2\sigma_a^2)/(\rho_a^2\sigma_b^2)$ is called *asymptotic relative efficiency*. Meaning, that a sample of size $n_a = \text{are}(\varphi_{n_a}^a, \varphi_{n_b}^b) n_b$ is needed for the test $\varphi_{n_a}^a$ to attain at $\theta_o + \eta$ approximately the same power $\mathbb{F}_{[0,1]}(-z_{1-\alpha} + \eta\sqrt{n_a}\rho_a/\sigma_a) = \mathbb{F}_{[0,1]}(-z_{1-\alpha} + \eta\sqrt{n_b}\rho_b/\sigma_b)$ as the test $\varphi_{n_b}^b$ with sample size n_b . A comparison with the test φ_n^* as in **Theorem** §06.04 (ii) allows analogously to introduce a notion of *asymptotic absolute efficiency*. □

§07 Rank tests

§07.01 **Reminder.** Consider on the sample space $(\mathbb{R}^n, \mathcal{B}^n)$ the statistic $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$ with $x \mapsto T(x) = (T_i(x))_{i \in \llbracket n \rrbracket}$ and $T_i(x) := \min\{c \in \mathbb{R} : \sum_{j \in \llbracket n \rrbracket} \mathbb{1}_{\{x_j \leq c\}} \geq i\}$, $i \in \llbracket n \rrbracket$. Since $T_1(x) \leq T_2(x) \leq \dots \leq T_n(x)$ for all $x \in \mathbb{R}^n$ the statistic T (and any other statistic with this property) is called an *order statistic*. Denote by \mathcal{S}_n the symmetric group of order n , i.e. the set of all permutations of the set $\llbracket n \rrbracket$. We identify as usual a vector $s = (s_i)_{i \in \llbracket n \rrbracket} \in \llbracket n \rrbracket^n$ with the map $s : \llbracket n \rrbracket \rightarrow \llbracket n \rrbracket$, $i \mapsto s(i) := s_i$, and hence $\mathcal{S}_n \subseteq \llbracket n \rrbracket^n$. Let $s^- \in \mathcal{S}_n$ denote the inverse permutation of $s \in \mathcal{S}_n$, i.e. $\text{id}_{\mathcal{S}_n} = s \circ s^- = s^- \circ s$. For a permutation $s = (s_i)_{i \in \llbracket n \rrbracket} \in \mathcal{S}_n$ and a vector $x = (x_i)_{i \in \llbracket n \rrbracket} \in \mathbb{R}^n$ we write shortly $x_s := (x_{s_i})_{i \in \llbracket n \rrbracket}$. A Borel-measurable map $S := (S_i)_{i \in \llbracket n \rrbracket} : \mathbb{R}^n \rightarrow \mathcal{S}_n$, i.e. $S^{-1}(\{s\}) \in \mathcal{B}^n$ for all $s \in \mathcal{S}_n$, is called a *random permutation* on $(\mathbb{R}^n, \mathcal{B}^n)$. The associated map $S^- : \mathbb{R}^n \rightarrow \mathcal{S}_n$ satisfying $\text{id}_{\mathcal{S}_n} = S^-(x) \circ S(x) = S(x) \circ S^-(x)$ for all $x \in \mathcal{X}$ is trivially again Borel-measurable, and hence called random inverse permutation of S . Moreover the statistic $X_S : \mathbb{R}^n \rightarrow \mathbb{R}^n$ with $x \mapsto X_S(x) := (x_{S_i(x)})_{i \in \llbracket n \rrbracket} = x_{S(x)} = \sum_{s \in \mathcal{S}_n} x_s \mathbb{1}_{\{s\}}(S(x))$ (a finite sum of Borel-measurable functions $x \mapsto x_s \mathbb{1}_{S^{-1}(s)}(x)$) is called a *random arrangement*. □

§07.02 **Definition.** A random permutation $O = (O_i)_{i \in \llbracket n \rrbracket}$ on $(\mathbb{R}^n, \mathcal{B}^n)$ is called *order permutation*, if the associated random arrangement $X_O : \mathbb{R}^n \rightarrow \mathbb{R}^n$ with $x \mapsto x_{O(x)}$ is an order statistic, i.e. $x_{O_1(x)} \leq x_{O_2(x)} \leq \dots \leq x_{O_n(x)}$ for all $x \in \mathbb{R}^n$. A random permutation $R = (R_i)_{i \in \llbracket n \rrbracket}$ on $(\mathbb{R}^n, \mathcal{B}^n)$ is called *rank permutation*, if its random inverse permutation $O := R^-$ is an order permutation. For $i \in \llbracket n \rrbracket$ the i -th component $R_i(x)$ of $R(x)$ is called the *rank* of the i -th component of $x \in \mathbb{R}^n$. □

§07.03 **Comment.** An order permutation O is uniquely determined on the Borel-set $\{x_i \neq x_j\} := \{(x_i)_{i \in [n]} \in \mathbb{R}^n : x_i \neq x_j, \forall j \in [n] \setminus \{i\}, \forall i \in [n]\}$ only. However, for $x \in \mathbb{R}^n$, the permutation $o := O(x) \in \mathcal{S}_n$ and $i \in [n]$ the value at the i -th position in the ordered vector x_o equals the value at the o_i -th position in the original vector x . Conversely, for the permutation $r := R(x) \in \mathcal{S}_n$ of the rank permutation $R := O^-$ the value at the r_i -th position in the ordered vector x_o equals the value at the i -th position in the original vector x . \square

§07.04 **Remark.** The map $R^* = (R_i^*)_{i \in [n]} : \mathbb{R}^n \rightarrow \mathcal{S}_n$ with $x \mapsto R_i^*(x) := \sum_{j \in [i]} \mathbb{1}_{\{x_i = x_j\}} + \sum_{j \in [n]} \mathbb{1}_{\{x_i > x_j\}}$ for each $i \in [n]$ is a rank permutation. Indeed, for each $x \in \mathbb{R}^n$ we have $r := R^*(x) \in \mathcal{S}_n$ ($r : [n] \rightarrow [n]$ is injective and hence bijective) and its inverse permutation $o := r^-$ satisfies $x_{o_1} \leq x_{o_2} \leq \dots \leq x_{o_n}$. Furthermore, each component of R^* is \mathcal{B} - $2^{[n]}$ -measurable, and hence R^* is a rank permutation. On the Borel-set $\{x_i \neq x_j\}$ each rank permutation $R = (R_i)_{i \in [n]}$ is uniquely determined by $R_i(x) = \sum_{j \in [n]} \mathbb{1}_{\{x_j \leq x_i\}} = R_i^*(x)$, $i \in [n]$. For each $y \in \mathbb{R}$ define $\widehat{F}_n(y) := \widehat{P}_n(\mathbb{1}_{(-\infty, y]})$ with $\widehat{F}_n(y, x) := \frac{1}{n} \sum_{j \in [n]} \mathbb{1}_{\{x_j \leq y\}} \in [0, 1]$ for all $x \in \mathbb{R}^n$. \widehat{F}_n is called *empirical cumulative distribution function*. If in addition $r := R(x)$ and $o := r^-$ for $x \in \{x_i \neq x_j\}$ then $i = n\widehat{F}_n(x_{o_i}, x)$ and $r_i = n\widehat{F}_n(x_i, x)$ for each $i \in [n]$. \square

§07.05 **Comment.** We assume a product probability measure $\mathbb{P}^n = \bigotimes_{j \in [n]} \mathbb{P}_j$ on the sample space $(\mathbb{R}^n, \mathcal{B}^n)$ where for each $j \in [n]$ the marginal probability measure $\mathbb{P}_j \in \mathcal{W}(\mathcal{B})$ admits a Lebesgue density $f_j = d\mathbb{P}_j/d\lambda$ and hence $\mathbb{P}^n \ll \lambda^n \in \mathcal{M}_\sigma(\mathcal{B}^n)$ with Lebesgue density $d\mathbb{P}^n/d\lambda^n = \prod_{j \in [n]} f_j$. Noting that the complement $\{x_i = x_j\} := \{x_i \neq x_j\}^c$ of the Borel-set $\{x_i \neq x_j\}$ is a λ^n null set, and hence it is also a \mathbb{P}^n null set. Thereby, each rank permutation R on $(\mathbb{R}^n, \mathcal{B}^n)$ with corresponding order permutation $O := R^-$ satisfies $x_{O_1(x)} < x_{O_2(x)} < \dots < x_{O_n(x)}$ for \mathbb{P}^n -a.e. $x \in \mathbb{R}^n$. Moreover, for \mathbb{P}^n -a.e. $x \in \mathbb{R}^n$ the vector of ranks $R(x)$ (and the rang permutation R) is determined by $R_i(x) = \sum_{j \in [n]} \mathbb{1}_{\{x_j \leq x_i\}} = n\widehat{F}_n(x_i, x)$, $i \in [n]$. \square

§07.06 **Lemma.** Consider a product probability measure $\mathbb{P}^{\otimes n}$ on $(\mathbb{R}^n, \mathcal{B}^n)$ with identical marginal distribution $\mathbb{P} \in \mathcal{W}(\mathcal{B})$, cumulative distribution function $F(y) := \mathbb{P}(\mathbb{1}_{(-\infty, y]})$, $y \in \mathbb{R}$, and Lebesgue density $f = d\mathbb{P}/d\lambda$. Let R and X_O with $O = R^-$ be a rang permutation on $(\mathbb{R}^n, \mathcal{B}^n)$ and the corresponding order statistic, respectively.

- (i) R is under $\mathbb{P}^{\otimes n}$ uniformly distributed on the symmetric group \mathcal{S}_n , precisely, $(\mathbb{P}^{\otimes n})^R(\{s\}) = (\mathbb{P}^{\otimes n} \circ R^{-1})(\{s\}) = \mathbb{P}^{\otimes n}(R = s) = \frac{1}{n!}$, $s \in \mathcal{S}_n$, in short $R \sim (\mathbb{P}^{\otimes n})^R = \mathcal{U}_{\mathcal{S}_n}$.
- (ii) R and X_O are independent under $\mathbb{P}^{\otimes n}$.
- (iii) The distribution of X_O admits under $\mathbb{P}^{\otimes n}$ a Lebesgue density $f_{X_O}^{(x)} = n! \mathbb{1}_B(x) \prod_{i \in [n]} f(x_i)$, $x \in \mathbb{R}^n$, with $B := \{(x_i)_{i \in [n]} \in \mathbb{R}^n, x_1 < \dots < x_n\}$.
- (iv) For each $i \in [n]$ the distribution of the i -th component of X_O admits under $\mathbb{P}^{\otimes n}$ a Lebesgue density $f_i(x) = i \binom{n}{i} |F(x)|^{i-1} |1 - F(x)|^{n-i} f(x)$, $x \in \mathbb{R}$.

§07.07 **Proof of Lemma §07.06.** is given in the lecture. \square

§07.08 **Definition.** Let \mathbb{P}_o and \mathbb{P} be probability measures on $(\mathbb{R}, \mathcal{B})$. We say \mathbb{P}_o is *stochastically smaller* than \mathbb{P} , or $\mathbb{P}_o \preceq \mathbb{P}$ for short, if $\mathbb{P}_o((c, \infty)) \leq \mathbb{P}((c, \infty))$ for all $c \in \mathbb{R}$. If in addition $\mathbb{P}_o \neq \mathbb{P}$, then we write $\mathbb{P}_o \prec \mathbb{P}$. \square

§07.09 **Remark.** Roughly speaking, $\mathbb{P}_o \preceq \mathbb{P}$ says that realisations of \mathbb{P}_o are typically smaller than realisations of \mathbb{P} . \square

§07.10 **Example.** For $\sigma \in \mathbb{R}^+$ consider on $(\mathbb{R}, \mathcal{B})$ a Gaussian location family $N_{\mathbb{R} \times \{\sigma^2\}}$. Then for all $a, b \in \mathbb{R}$ holds $N_{(a, \sigma^2)} \prec N_{(b, \sigma^2)}$ if and only if $a \leq b$. More generally, given a location family \mathbb{P}_r on

$(\mathbb{R}, \mathcal{B})$ as introduced in **Example §04.17** with likelihood function $L(\theta, x) = g(x - \theta)$, $\theta, x \in \mathbb{R}$, for some strictly positive Lebesgue-density g on \mathbb{R} . Then for all $a, b \in \mathbb{R}$ holds $\mathbb{P}_a \prec \mathbb{P}_b$ if and only if $a \lesssim b$. □

§07.11 **Heuristics.** Given a sample from each distribution $\mathbb{P}_o, \mathbb{P} \in \mathcal{W}(\mathcal{B})$ we consider the testing task $H_o : \mathbb{P} = \mathbb{P}_o$ against the alternative $H_1 : \mathbb{P}_o \prec \mathbb{P}$. Loosely speaking, this means, that we aim to reject the null hypothesis if realisations of \mathbb{P}_o are *significantly* smaller than realisation of \mathbb{P} . More precisely, we assume a sample of $n = m + l$ independent real random variables $(X_i)_{i \in [n]}$ where the first m and the last l have as common marginal distribution \mathbb{P}_o and \mathbb{P} , respectively. In other words $X = (X_i)_{i \in [n]}$ takes its values in the pooled sample space $(\mathbb{R}^n, \mathcal{B}^n)$. Considering a rank permutation R on $(\mathbb{R}^n, \mathcal{B}^n)$ and an observation $x \in \mathbb{R}^n$ it seems reasonable to reject the hypothesis if the sum of ranks within the first group of m random variables, i.e. $W_o(x) := \sum_{i \in [m]} R_i(x)$, takes *sufficiently* smaller values than the sum of ranks within the second group of l random variables, i.e. $W(x) := \sum_{i \in [l]} R_{i+m}(x)$ where obviously $W_o(x) + W(x) = \sum_{i \in [n]} R_i(x) = \sum_{i \in [n]} i = \frac{n(n+1)}{2}$. □

§07.12 **Lemma.** For $m, l \in \mathbb{N}$ and $n := m + l$ let $R = (R_i)_{i \in [n]}$ be a rang permutation on $(\mathbb{R}^n, \mathcal{B}^n)$, $W_o := \sum_{i \in [m]} R_i$, $W := \sum_{i \in [l]} R_{i+m}$ and $U_{ml} : \mathbb{R}^n \rightarrow \llbracket 0, ml \rrbracket$ with $x \mapsto U_{ml}(x) := \sum_{i \in [m]} \sum_{j \in [l]} \mathbb{1}_{\{x_i > x_{j+m}\}}$. Then for each $x \in \{x_i \neq x_j\}$ it holds $W_o(x) = U_{ml}(x) + \frac{m(m+1)}{2}$ and consequently $W(x) = ml - U_{ml}(x) + \frac{l(l+1)}{2}$.

§07.13 **Proof of Lemma §07.12.** is given in the lecture. □

§07.14 **Comment.** Keeping **Lemma §07.12** in mind, we use the test statistic W_o or equivalently U_{ml} to reject the hypothesis $H_o : \mathbb{P} = \mathbb{P}_o$ against the alternative $H_1 : \mathbb{P}_o \prec \mathbb{P}$, if $U_{ml} < c$ or equivalently $W_o < c + \frac{m(m+1)}{2}$ for a certain threshold $c \in (0, ml]$. The test is called (one-sided) Mann-Whitney U-test or Wilcoxon two-sample rank sum test¹. The critical value has to be chosen according to a pre-specified level $\alpha \in (0, 1)$ which under the null hypothesis necessitates the knowledge of the distribution of U_{ml} or an asymptotic approximation. Interestingly the next proposition shows that under the null hypothesis the distribution of U_{ml} is *distribution free* in the following sense: If $\mathbb{P}_o = \mathbb{P}$ and \mathbb{P} admits a Lebesgue density, then the distribution of U_{ml} is determined and it is independent of the underlying distribution \mathbb{P} . □

§07.15 **Proposition.** For $m, l \in \mathbb{N}$ and $n := m + l$ let $\mathbb{P}^{\otimes n} \in \mathcal{W}(\mathcal{B}^n)$ with identical marginal distribution $\mathbb{P} \ll \lambda$. For all $k \in \llbracket 0, ml \rrbracket$ it holds $\mathbb{P}^{\otimes n}(U_{ml} = k) = N(k; m, l) / \binom{n}{k}$ where $N(k; m, l)$ denotes the number of all partitions $\sum_{i \in [m]} k_i = k$ of k in m increasingly ordered numbers $k_1 \leq k_2 \leq \dots \leq k_m$ taking from the set $\llbracket 0, l \rrbracket$. In particular, it holds $\mathbb{P}^{\otimes n}(U_{ml} = k) = \mathbb{P}^{\otimes n}(U_{ml} = ml - k)$.

§07.16 **Proof of Proposition §07.15.** is given in the lecture. □

§07.17 **Remark.** For small values of k the partition number $N(k; m, l)$ can be calculated by combinatorial means and there exists tables gathering certain quantiles of the U_{ml} -distribution. However, for large values of k the exact calculation of quantiles of the U_{ml} -distribution may be avoided by using an appropriate asymptotic approximation. In the sequel we let m and l and thus $n = m + l$ tend to infinity, which formally means that we consider sequences $(m_n)_{n \in \mathbb{N}}$ and $(l_n)_{n \in \mathbb{N}}$ satisfying $m_n + l_n = n$ for any $n \in \mathbb{N}$. Here and subsequently we assume that $m_n/n \xrightarrow{n \rightarrow \infty} \gamma \in (0, 1)$ and hence $l_n/n \xrightarrow{n \rightarrow \infty} 1 - \gamma$. For sake of presentation, however, we drop the additional index n and write shortly $n = m + l$ with $m/n \xrightarrow{n \rightarrow \infty} \gamma$ and hence $l/n \xrightarrow{n \rightarrow \infty} 1 - \gamma$. □

¹The version based on W_o has been proposed by Wilcoxon [1945], while the U_{ml} -version has been independently be introduced by Mann and Whitney [1947].

§07.18 **Theorem.** For $m, l \in \mathbb{N}$ and $n := m + l$ let $\mathbb{P}^{\otimes n} \in \mathcal{W}(\mathcal{B}^n)$ with identical marginal distribution $\mathbb{P} \ll \lambda$, and hence continuous cumulative distribution function \mathbb{F} . Consider $U_{ml} : \mathbb{R}^n \rightarrow \llbracket 0, ml \rrbracket$ and $T_{ml} : \mathbb{R}^n \rightarrow \mathbb{R}$ with $x \mapsto U_{ml}(x) := \sum_{i \in \llbracket m \rrbracket} \sum_{j \in \llbracket l \rrbracket} \mathbb{1}_{\{x_i > x_{j+m}\}}$ and

$$x \mapsto T_{ml}(x) := l \sum_{i \in \llbracket l \rrbracket} \mathbb{F}(x_i) - m \sum_{i \in \llbracket m \rrbracket} \mathbb{F}(x_{i+m}) = l \sum_{i \in \llbracket m \rrbracket} (\mathbb{F}(x_i) - 1/2) - m \sum_{i \in \llbracket l \rrbracket} (\mathbb{F}(x_{i+m}) - 1/2).$$

Define further $v_{ml} := ml(n+1)/12$, $T_{ml}^* := T_{ml}/\sqrt{v_{ml}}$ and $U_{ml}^* := (U_{ml} - ml/2)/\sqrt{v_{ml}}$. If in addition $m/n \rightarrow \gamma \in (0, 1)$ then $U_{ml}^* - T_{ml}^* = o_{\mathbb{P}^{\otimes n}}(1)$ and $T_{ml}^* \xrightarrow{d} N_{(0,1)}$ under $\mathbb{P}^{\otimes n}$, and thus $U_{ml}^* \xrightarrow{d} N_{(0,1)}$ under $\mathbb{P}^{\otimes n}$.

§07.19 **Proof of Theorem §07.18.** is given in the lecture. \square

§07.20 **Remark.** Considering two independent samples $(X_i)_{i \in \llbracket m \rrbracket} \sim \mathbb{P}_o^{\otimes m}$ and $(X_{i+m})_{i \in \llbracket l \rrbracket} \sim \mathbb{P}^{\otimes l}$ set $n := m + l$ and $X := (X_i)_{i \in \llbracket n \rrbracket}$. Keeping **Theorem §07.18** in mind we reject the null hypothesis $H_o : \mathbb{P}_o = \mathbb{P}$ against the alternative $H_1 : \mathbb{P}_o \prec \mathbb{P}$, if $U_{ml}(X) < ml/2 + z_\alpha \sqrt{v_{ml}}$ with $\mathbb{F}_{[0,1]}(z_\alpha) = \alpha \in (0, 1)$. This test is asymptotically a level- α test due to **Theorem §07.18** by exploiting that under the null $\mathbb{P}^{\otimes n}(U_{ml} < ml/2 + z_\alpha \sqrt{v_{ml}}) \xrightarrow{n \rightarrow \infty} \mathbb{F}_{[0,1]}(z_\alpha) = \alpha$ for $m/n \xrightarrow{n \rightarrow \infty} \gamma \in (0, 1)$. Note that we reject similarly the null hypothesis $H_o : \mathbb{P}_o = \mathbb{P}$ against the alternative $H_1 : \mathbb{P} \prec \mathbb{P}_o$ if $U_{ml} > ml/2 + z_{1-\alpha} \sqrt{v_{ml}}$. Next we study the (asymptotic) size of the power of the rank test under local alternatives where we use that under the assumptions of **Theorem §07.18** it holds

$$\begin{aligned} U_{ml}^* &= \frac{U_{ml} - ml/2}{\sqrt{v_{ml}}} = \sqrt{\frac{l}{n+1}} \frac{1}{\sqrt{m}} \sum_{i \in \llbracket m \rrbracket} \frac{\mathbb{F}(X_i) - 1/2}{\sqrt{1/12}} - \sqrt{\frac{m}{n+1}} \frac{1}{\sqrt{l}} \sum_{i \in \llbracket l \rrbracket} \frac{\mathbb{F}(X_{i+m}) - 1/2}{\sqrt{1/12}} + o_{\mathbb{P}^{\otimes n}}(1) \\ &= \sqrt{1-\gamma} \sqrt{m} \widehat{\mathbb{P}}_m(g) - \sqrt{\gamma} \sqrt{l} \widehat{\mathbb{P}}_l(g) + o_{\mathbb{P}^{\otimes n}}(1) \quad (07.01) \end{aligned}$$

setting $g := \sqrt{12}(\mathbb{F} - 1/2)$, $\widehat{\mathbb{P}}_m(g) := \frac{1}{m} \sum_{i \in \llbracket m \rrbracket} g(X_i)$ and $\widehat{\mathbb{P}}_l(g) := \frac{1}{l} \sum_{i \in \llbracket l \rrbracket} g(X_{i+m})$ where $\widehat{\mathbb{P}}_m(g)$ and $\widehat{\mathbb{P}}_l(g)$ are independent, $\mathbb{P}(g) = 0$, and $\mathbb{P}(g^2) = 1$ by construction. \square

§08 Asymptotic power of rank tests

§08.01 **Motivation.** Considering the test of the hypothesis $H_o : \mathbb{P}_o = \mathbb{P}$ against the alternative $H_1 : \mathbb{P}_o \prec \mathbb{P}$ we restrict our attention to the special case that \mathbb{P}_o and \mathbb{P} belong to a location family $\mathbb{P}_\mathbb{R}$ as introduced in **Example §04.17**. Precisely, we assume that the family $\mathbb{P}_\mathbb{R}$ of probability measures on $(\mathbb{R}, \mathcal{B})$ is dominated by the Lebesgue measure. For each $\theta \in \mathbb{R}$, \mathbb{P}_θ admits a likelihood function given by $L(\theta, x) = g(x - \theta)$, $x \in \mathbb{R}$, where g is a continuous and strictly positive density on \mathbb{R} . Recall that in this context $\mathbb{P}_a \prec \mathbb{P}_b$ holds if and only if $a \not\leq b$ (see **Example §07.10**). Observe further that we can assume that $\mathbb{P}_o = \mathbb{P}_0$ (possibly after a reparametrisation). For $m, l \in \mathbb{N}$ and $n = m + l$ supposing independent random variables $(X_i)_{i \in \llbracket n \rrbracket}$ with $(X_i)_{i \in \llbracket m \rrbracket} \odot \mathbb{P}_\mathbb{R}^{\otimes m}$ and $(X_{i+m})_{i \in \llbracket l \rrbracket} \sim \mathbb{P}_0^{\otimes l}$ their joint distribution belongs to the two sample location family $\mathbb{P}_\mathbb{R}^{m,l} := (\mathbb{P}_\theta^{m,l} := \mathbb{P}_\theta^{\otimes m} \otimes \mathbb{P}_0^{\otimes l})_{\theta \in \mathbb{R}}$. Summarising, based on the statistical two sample location experiment $(\mathbb{R}^n, \mathcal{B}^n, \mathbb{P}_\mathbb{R}^{m,l})$ the aim is to test the hypothesis $H_o : \theta = 0$ against the alternative $H_1 : 0 < \theta$. \square

§08.02 **Regular location model.** A location family $\mathbb{P}_\mathbb{R}$ of probability measures on $(\mathbb{R}, \mathcal{B})$ dominated by the Lebesgue measure $\lambda \in \mathcal{M}_\sigma(\mathcal{B})$ with likelihood for each $\theta \in \mathbb{R}$ and a strictly positive density $g \in \mathcal{B}^+$ given by $L(\theta, x) = g(x - \theta)$, $x \in \mathbb{R}$, is called *regular* if the density g is in addition continuously differentiable with derivative \dot{g} satisfying $\lambda(|\dot{g}|^2/g) < \infty$. \square

§08.03 **Reminder.** A regular location family $\mathbb{P}_{\mathbb{R}}$ is Hellinger-differentiable in each $\theta \in \mathbb{R}$ with score function $\dot{\ell}_{\theta} = -\dot{g}(x - \theta)/g(x - \theta)$ and Fisher information $\mathcal{J} := \mathcal{J}_{\theta} = \lambda(|\dot{g}|^2/g)$ (see **Example** §04.17). Due to **Theorem** §05.12 the statistical product experiment $(\mathbb{R}^m, \mathcal{B}^m, \mathbb{P}_{\mathbb{R}}^{\otimes m})$ is ULAN in $\theta_0 = 0$ with localising rate $(\delta_m := m^{-1/2})_{m \in \mathbb{N}}$ and central sequence $(\mathcal{Z}_0^m := -\sqrt{m}\widehat{\mathbb{P}}_m(\dot{g}/g))_{m \in \mathbb{N}}$. Precisely, for any sequence $h_m \rightarrow h$ as $m \rightarrow \infty$ it holds

$$\log(d\mathbb{P}_{h_m/\sqrt{m}}^{\otimes m}/d\mathbb{P}_0^{\otimes m}) = -h\sqrt{m}\widehat{\mathbb{P}}_m(\dot{g}/g) - \frac{1}{2}h^2\mathcal{J} + o_{\mathbb{P}_0^{\otimes m}}(1)$$

and $\sqrt{m}\widehat{\mathbb{P}}_m(\dot{g}/g) \xrightarrow{d} N_{(0,\mathcal{J})}$ under $\mathbb{P}_0^{\otimes m}$. Given a two sample location family $\mathbb{P}_{\mathbb{R}}^{m,l}$ for any $\theta \in \mathbb{R}$ the log of the likelihood-ratio satisfies $\log(d\mathbb{P}_{\theta}^{m,l}/d\mathbb{P}_0^{m,l}) = \log(d\mathbb{P}_{\theta}^{\otimes m}/d\mathbb{P}_0^{\otimes m})$. Thereby, if the location family is regular and $m/n \xrightarrow{n \rightarrow \infty} \gamma \in (0, 1)$, whence $h_m := h\sqrt{m/n} \xrightarrow{n \rightarrow \infty} h\sqrt{\gamma}$, it follows

$$\begin{aligned} \ell_{h,n} &:= \log(d\mathbb{P}_{h/\sqrt{n}}^{m,l}/d\mathbb{P}_0^{m,l}) = \log(d\mathbb{P}_{h_m/\sqrt{m}}^{\otimes m}/d\mathbb{P}_0^{\otimes m}) \\ &= -h\sqrt{\gamma}\sqrt{m}\widehat{\mathbb{P}}_m(\dot{g}/g) - \frac{\gamma}{2}h^2\mathcal{J} + o_{\mathbb{P}_0^{\otimes m}}(1) \quad (08.01) \end{aligned}$$

We consider in the sequel a rank test $\varphi_n = \mathbb{1}_{\{U_{ml}^* > z_{1-\alpha}\}}$ with $\mathbb{F}_{[0,1]}(-z_{1-\alpha}) = \alpha \in (0, 1)$ based on the standardised test statistic $U_{ml}^* = (U_{ml} - ml/2)/\sqrt{v_{ml}}$ and its asymptotic decomposition given in (07.01). \square

§08.04 **Theorem.** Assume a two sample regular location model $(\mathbb{R}^n, \mathcal{B}^n, \mathbb{P}_{\mathbb{R}}^{m,l})$, $n = m + l \in \mathbb{N}$. Consider a rank test $\varphi_n = \mathbb{1}_{\{U_{ml}^* > z_{1-\alpha}\}}$ with $\mathbb{F}_{[0,1]}(-z_{1-\alpha}) = \alpha \in (0, 1)$ for the testing task $H_0 : \theta = 0$ against $H_1 : \theta > 0$. If $m/n \xrightarrow{n \rightarrow \infty} \gamma \in (0, 1)$, then:

- (i) Under the null hypothesis $H_0 : \theta = 0$ we have $\mathbb{P}_0^{m,l}(\varphi_n) = \mathbb{P}_0^{\otimes m+l}(U_{ml}^* > z_{1-\alpha}) \xrightarrow{n \rightarrow \infty} \alpha$, i.e., φ_n is an asymptotic level- α test;
- (ii) The power function $\beta_{\varphi_n}(\theta) = \mathbb{P}_{\theta}^{m,l}(\varphi_n)$, $\theta \in \mathbb{R}$, of the rank test $\varphi_n = \mathbb{1}_{\{U_{ml}^* > z_{1-\alpha}\}}$ satisfies under local alternatives $\beta_{\varphi_n}(h/\sqrt{n}) = \mathbb{P}_{h/\sqrt{n}}^{m,l}(U_{ml}^* > z_{1-\alpha}) \xrightarrow{n \rightarrow \infty} \mathbb{F}_{[0,1]}(-z_{1-\alpha} + h\rho)$ with $\rho = \lambda(g^2)\sqrt{12\gamma(1-\gamma)}$ for each $h \in \mathbb{R}^+$.

§08.05 **Proof** of **Theorem** §08.04. is given in the lecture. \square

§08.06 **Comment.** Let us briefly consider a rank test $\varphi_n = \mathbb{1}_{\{U_{ml} < ml/2 + z_{\alpha}\sqrt{v_{ml}}\}} = \mathbb{1}_{\{U_{ml}^* < z_{\alpha}\}}$ with $\mathbb{F}_{[0,1]}(z_{\alpha}) = \alpha \in (0, 1)$ for the testing task of the null hypothesis $H_0 : \theta = 0$ against the alternative $H_1 : \theta > 0$. Similar to **Theorem** §08.04 φ_n is an asymptotic level- α test with power for local alternatives $\beta_{\varphi_n}(-h/\sqrt{n}) = \mathbb{P}_{-h/\sqrt{n}}^{m,l}(U_{ml}^* < z_{\alpha}) \xrightarrow{n \rightarrow \infty} \mathbb{F}_{[0,1]}(z_{\alpha} + h\rho)$, $h \in \mathbb{R}^+$. \square

§08.07 **Two sample Gaussian location model.** For $m, l \in \mathbb{N}$, $n = m + l$ and variance $\sigma^2 \in \mathbb{R}_0^+$ the joint distribution of independent random variables $(X_i)_{i \in \llbracket n \rrbracket}$ with $(X_i)_{i \in \llbracket m \rrbracket} \overset{\circ}{\sim} N_{\mathbb{R} \times \{\sigma^2\}}^{\otimes m}$ and $(X_{i+m})_{i \in \llbracket l \rrbracket} \sim N_{(0,\sigma^2)}^{\otimes l}$ belongs to a *two sample Gaussian location model* $N_{\mathbb{R} \times \{\sigma^2\}}^{m,l} := (N_{(\theta,\sigma^2)}^{m,l} := N_{(\theta,\sigma^2)}^{\otimes m} \otimes N_{(0,\sigma^2)}^{\otimes l})_{\theta \in \mathbb{R}}$. \square

§08.08 **Remark.** For $\sigma \in \mathbb{R}_0^+$ a Gaussian location family $N_{\mathbb{R} \times \{\sigma^2\}}$ on $(\mathbb{R}, \mathcal{B})$ is regular with score function $\dot{\ell}_{\theta}(x) = (x - \theta)/\sigma$, $x \in \mathbb{R}$, and Fisher information $\mathcal{J}_{\theta} = \lambda(|\dot{g}|^2/g) = N_{(\theta,\sigma^2)}(\dot{\ell}_{\theta}^2) = \int_{\mathbb{R}} (x - \theta)^2/\sigma^2 N_{(\theta,\sigma^2)}(dx) = 1$, $\theta \in \mathbb{R}$ (using the notations in **Example** §04.17). \square

§08.09 **Example.** Assume a two sample Gaussian location model $(\mathbb{R}^n, \mathcal{B}^n, N_{\mathbb{R} \times \{\sigma^2\}}^{m,l})$, $n = m + l \in \mathbb{N}$ and $\sigma \in \mathbb{R}_0^+$. Define the statistics $T_{ml} := l \sum_{i \in \llbracket m \rrbracket} X_i - m \sum_{i \in \llbracket l \rrbracket} X_{i+m} \in \mathcal{B}^n$ and $V_{ml}^2 := \frac{ml(m+l)}{(m+l)-2} (\sum_{i \in \llbracket m \rrbracket} (X_i - \frac{1}{m} \sum_{i \in \llbracket m \rrbracket} X_i)^2 + \sum_{i \in \llbracket l \rrbracket} (X_{i+m} - \frac{1}{l} \sum_{i \in \llbracket l \rrbracket} X_{i+m})^2) \in \mathcal{B}^n$, and set $V_{ml} :=$

$\sqrt{V_{ml}^2}$. Under $N_{(0,\sigma^2)}^{m,l}$ the standardised (Student-) t-statistic $T_{ml}^* := T_{ml}/V_{ml} \in \overline{\mathcal{B}}^n$ has a t_{n-2} -distribution with $n - 2$ degrees of freedom, i.e. $T_{ml}^* \sim t_{n-2}$. We denote by $t_{n-2,\alpha}$ a α -quantile of a t_{n-2} -distribution, i.e., $\mathbb{F}_{t_{n-2}}(t_{n-2,\alpha}) = \alpha \in (0, 1)$. Consider for the testing task of the null hypothesis $H_0 : \theta = 0$ against the alternative $H_1 : 0 < \theta$ (or $H_1 : 0 > \theta$) the t-test $\varphi_n^* = \mathbb{1}_{\{T_{ml}^* > t_{n-2,1-\alpha}\}}$ (or $\varphi_n^* = \mathbb{1}_{\{T_{ml}^* < t_{n-2,\alpha}\}}$), which is by construction a level- α test. Since a Gaussian location model is regular (see Remark §08.08) we can directly apply Theorem §08.04 to derive its asymptotic power function under local alternatives. However, Theorem §08.04 allows us to study a t-test in an arbitrary regular location model (Definition §08.02). More precisely, for $\theta \in \mathbb{R}$ and $\sigma \in \mathbb{R}_0^+$ define $\mathfrak{v}_{(\theta,\sigma)} \in \mathcal{B}$ with $x \mapsto \mathfrak{v}_{(\theta,\sigma)}(x) := (x - \theta)/\sigma$. As in Remark §08.08 $\dot{\ell}_\theta = \mathfrak{v}_{(\theta,\sigma)}$ and $\mathcal{J}_\theta = N_{(\theta,\sigma^2)}(\mathfrak{v}_{(\theta,\sigma)}^2) = 1$, $\theta \in \mathbb{R}$, is the score function and Fisher information, respectively, in a Gaussian location family $N_{\mathbb{R} \times \{\sigma^2\}}$ with variance $\sigma^2 \in \mathbb{R}_0^+$. Considering a regular location family $\mathbb{P}_{\mathbb{R}}$ with $\mathbb{P}_0 \in \mathcal{W}_2(\mathcal{B})$ (see Notation §19.05) and hence $\mathbb{P}_\theta \in \mathcal{W}_2(\mathcal{B})$ for all $\theta \in \mathbb{R}$ we have $\sigma^2 := \mathbb{P}_0(\mathfrak{v}_{(0,1)}^2) = \lambda(\mathfrak{v}_{(0,1)}^2 g) = \lambda(\text{id}_{\mathbb{R}}^2 g) = \int_{\mathbb{R}} x^2 g(x) \lambda(dx) < \infty$ and $\mathbb{P}_\theta(\mathfrak{v}_{(\theta,\sigma)}^2) = 1$ for all $\theta \in \mathbb{R}$ exploiting the translation invariance of the Lebesgue measure. \square

§08.10 **Regular mean location family with finite variance.** A regular location family $\mathbb{P}_{\mathbb{R}}$ of probability measures on $(\mathbb{R}, \mathcal{B})$ is said to have *finite variance* $\sigma^2 \in \mathbb{R}_0^+$, if $\mathbb{P}_0 \in \mathcal{W}_2(\mathcal{B})$ (and hence $\mathbb{P}_\theta \in \mathcal{W}_2(\mathcal{B})$ for all $\theta \in \mathbb{R}$), and $\sigma^2 = \mathbb{P}_0(\mathfrak{v}_{(0,1)}^2)$ (and hence $\mathbb{P}_\theta(\mathfrak{v}_{(\theta,\sigma)}^2) = 1$ for all $\theta \in \mathbb{R}$). We call a regular location family satisfying in addition $\mathbb{P}_0(\mathfrak{v}_{(0,1)}) = 0$ (and hence $\mathbb{P}_\theta(\mathfrak{v}_{(\theta,\sigma)}) = 0$ for all $\theta \in \mathbb{R}$) a *regular mean location family*. \square

§08.11 **Theorem.** Assume a two sample regular mean location model $(\mathbb{R}^n, \mathcal{B}^n, \mathbb{P}_{\mathbb{R}}^{m,l})$ with finite variance $\sigma^2 \in \mathbb{R}_0^+$. Consider a t-test $\varphi_n^* = \mathbb{1}_{\{T_{m,l}^* > t_{n-2,1-\alpha}\}}$ with $1 - \mathbb{F}_{t_{n-2}}(t_{n-2,1-\alpha}) = \alpha \in (0, 1)$ for the testing task $H_0 : \theta = 0$ against $H_1 : \theta > 0$. If $m/n \xrightarrow{n \rightarrow \infty} \gamma \in (0, 1)$, then:

- (i) Under the null hypothesis $H_0 : \theta = 0$ we have $\mathbb{P}_0^{m,l}(\varphi_n^*) = \mathbb{P}_0^{m,l}(T_{m,l}^* > t_{n-2,1-\alpha}) \xrightarrow{n \rightarrow \infty} \alpha$, i.e., φ_n^* is an asymptotic level- α test;
- (ii) The power function $\beta_{\varphi_n^*}(\theta) = \mathbb{P}_\theta^{m,l}(\varphi_n^*)$, $\theta \in \mathbb{R}$ of the t-test $\varphi_n^* = \mathbb{1}_{\{T_{m,l}^* > t_{n-2,1-\alpha}\}}$ satisfies under local alternatives $\beta_{\varphi_n^*}(h/\sqrt{n}) = \mathbb{P}_{n/\sqrt{n}}^{m,l}(T_{m,l}^* > t_{n-2,1-\alpha}) \xrightarrow{n \rightarrow \infty} \mathbb{F}_{[0,1]}(-z_{1-\alpha} + \rho)$ with $\rho = h\sigma^{-1}\sqrt{\gamma(1-\gamma)}$.

§08.12 **Proof** of Theorem §08.11. is given in the lecture. \square

§08.13 **Remark.** Given a two sample regular mean location model $(\mathbb{R}^n, \mathcal{B}^n, \mathbb{P}_{\mathbb{R}}^{m,l})$, $n = m + l \in \mathbb{N}$ with density $g \in \mathcal{B}^+$ and finite variance $\sigma^2 \in \mathbb{R}_0^+$ let us compare the asymptotic level- α rank-test $\varphi_n = \mathbb{1}_{\{U_{m,l}^* > z_{1-\alpha}\}}$ (see Theorem §08.04) and the t-test $\varphi_n^* = \mathbb{1}_{\{T_{m,l}^* > t_{n-2,1-\alpha}\}}$ (see Theorem §08.11). Using their asymptotic power functions the *asymptotic relative efficiency* (see Definition §06.10) between both tests equals $\text{are}(\varphi_n, \varphi_n^*) = 12\sigma^2(\lambda g^2)^2$. In the particular case of a Gaussian location model, i.e., $g(x) = \frac{1}{\sqrt{2\pi\sigma}} \exp(-x^2/(2\sigma^2))$ we have $\lambda g^2 = 1/(2\sigma\sqrt{\pi})$ and hence $\text{are}(\varphi_n, \varphi_n^*) = 3/\pi \approx 0.955$. On the other hand denoting by \mathcal{D} the class of all Lebesgue-densities $g \in \mathcal{B}^+$ satisfying $\lambda(\mathfrak{v}_{(0,\sigma)} g) = 0$ and $\lambda(\mathfrak{v}_{(0,\sigma)}^2 g) = 1$ Hodges and Lehmann [1956] have shown that $\inf_{g \in \mathcal{D}} 12\sigma^2(\lambda g^2)^2 = 0.864$ and $\sup_{g \in \mathcal{D}} 12\sigma^2(\lambda g^2)^2 = \infty$. \square

Chapter 3

Nonparametric estimation by projection

This chapter presents an introduction to nonparametric estimation of curves along the lines of the textbooks by Tsybakov [2009] and Comte [2015] where far more details, examples and further discussions can be found.

Overview

§09	Review	33
§10	Noisy version of the parameter	35
§10 01	Stochastic process	35
§10 02	Noisy parameter	38
§11	Orthogonal projection	39
§11 01	Weigthed norms and inner products	39
§11 02	Orthogonal projection	40
§11 03	Global and maximal global \mathfrak{v} -error	40
§11 04	Local and maximal local ϕ -error	41
§12	Orthogonal projection estimator	42
§12 01	Global and maximal global \mathfrak{v} -risk	43
§12 02	Local and maximal local ϕ -risk	46
§13	Minimax optimal estimation	50
§13 01	Minimax theory: a general approach	50
§13 02	Deriving a lower bound	53
§13 03	Lower bound based on two hypotheses	54
§13 04	Lower bound based on m hypotheses	56
§14	Data-driven estimation	58
§14 01	Data-driven estimation procedures	58
§14 02	Model selection	59
§14 03	GSSM: data-driven global estimation	62
§14 04	Goldenshluger and Lepskij's method	66
§14 05	GSSM: data-driven local estimation	68

§09 Review

Nonparametric density estimation. Consider for a non-empty set of parameters Θ a family \mathbb{P}_Θ of probability measures on $(\mathbb{R}, \mathcal{B})$ which contains the distribution of an observable real random variable, $X \sim \mathbb{P}_\Theta$. The family \mathbb{P}_Θ captures the prior knowledge about the distribution of the observation. For example, a family given by a set of parameters Θ containing only one singleton, i.e., $\Theta = \{\theta_o\}$, and hence $X \sim \mathbb{P}_{\theta_o}$ for some probability measure $\mathbb{P}_{\theta_o} \in \mathcal{W}(\mathcal{B})$, means, the data generating process is known to us in advance. On the contrary, a parameter set $\Theta = \mathcal{W}(\mathcal{B})$ reflects a lack of prior knowledge. A parametric model \mathbb{P}_Θ for some parameter set $\Theta \subseteq \mathbb{R}^k$ provides in a certain sense a trade-off between both extremes. In this chapter

our aim is to avoid an assumption of a finite dimensional set of parameters. For example, consider $(X_i)_{i \in [n]} \stackrel{i.i.d.}{\sim} \mathbb{P} \in \mathcal{W}(\mathcal{B})$, that is, an independent and identically distributed sample with common probability measure $\mathbb{P} \in \mathcal{W}(\mathcal{B})$. A reasonable estimator of the associated cumulative distribution function (c.d.f.) $\mathbb{F}(t) := \mathbb{P}((-\infty, t])$, $t \in \mathbb{R}$, is the empirical cumulative distribution function (e.c.d.f.) $\widehat{\mathbb{F}}_n(t) := \widehat{\mathbb{P}}_n((-\infty, t])$, $t \in \mathbb{R}$. For each $t \in \mathbb{R}$, $\widehat{\mathbb{F}}_n(t)$ is an unbiased estimator of $\mathbb{F}(t)$ with variance $\text{var}_m(\widehat{\mathbb{F}}_n(t)) = \frac{1}{n}\mathbb{F}(t)(1 - \mathbb{F}(t))$. Consequently, $\widehat{\mathbb{F}}_n(t)$ converges in probability to $\mathbb{F}(t)$, and thus it is a consistent estimator. Moreover, by the law of large numbers the convergence holds almost surely in any point and also uniformly, by Glivenko-Cantelli's theorem, i.e., $\|\widehat{\mathbb{F}}_n - \mathbb{F}\|_{\mathcal{L}_\infty} = o(1)$ \mathbb{P} -a.s.. If we assume in addition that \mathbb{P} admits a Lebesgue density then $\widehat{\mathbb{F}}_n$ is a unbiased estimator with minimal variance, by Lehman-Scheffé's theorem. However, comparing different probability measures using their associated c.d.f.'s is visually difficult and as a consequence, other measures for dissimilarities are typically used. Consider, for instance, for two probability measures \mathbb{P} and \mathbb{P}_0 on $(\mathbb{R}, \mathcal{B})$ their *total variation distance* given by $\|\mathbb{P} - \mathbb{P}_0\|_{\text{TV}} := \sup\{|\mathbb{P}(B) - \mathbb{P}_0(B)|, B \in \mathcal{B}\}$. We note that for any probability measure $\mathbb{P} \in \mathcal{W}(\mathcal{B})$ admitting a Lebesgue-density we have $\|\mathbb{P} - \widehat{\mathbb{P}}_n\|_{\text{TV}} = 1$ \mathbb{P} -a.s. for any $n \in \mathbb{N}$. As a consequence the empirical probability measure $\widehat{\mathbb{P}}_n$ is not a consistent estimator of \mathbb{P} in terms of the total variation distance. In other words, depending on the measure of accuracy (metric, topology, etc.) a different estimator of \mathbb{P} might be reasonable.

§09.01 **Lemma (Scheffé's theorem).** Let $\mathbb{P}, \mathbb{P}_0 \in \mathcal{W}(\mathcal{B})$ admit a μ -density \mathfrak{p} and \mathfrak{p}_0 , respectively, for some $\mu \in \mathcal{M}_\sigma(\mathcal{B})$. Then $\|\mathbb{P} - \mathbb{P}_0\|_{\text{TV}} = \frac{1}{2}\mu(|\mathfrak{p} - \mathfrak{p}_0|) = \frac{1}{2}\|\mathfrak{p} - \mathfrak{p}_0\|_{\mathcal{L}_1(\mu)}$.

§09.02 **Proof of Lemma §09.01.** is given in the lecture. □

In the sequel let \mathcal{D} be the set of Lebesgue densities on $(\mathbb{R}, \mathcal{B})$, and hence $\mathcal{D} \subseteq \mathcal{L}_1 = \mathcal{L}_1(\mathcal{B}, \lambda)$. $\mathbb{P}_\mathfrak{p} = \mathfrak{p}\lambda$ and $\mathbb{E}_\mathfrak{p}$ denote for each density $\mathfrak{p} \in \mathcal{D}$ the associated probability measure and expectation, respectively. We consider the statistical product experiment $(\mathbb{R}^n, \mathcal{B}^n, \mathbb{P}_\mathfrak{p}^{\otimes n} = (\mathbb{P}_\mathfrak{p}^{\otimes n})_{\mathfrak{p} \in \mathcal{D}})$ and $(X_i)_{i \in [n]} \odot \mathbb{P}_\mathfrak{p}^{\otimes n}$. Typically, for $s \geq 1$ we assess the accuracy of an estimator $\widehat{\mathfrak{p}}$ of \mathfrak{p} either by a local measure, e.g. $\mathbb{P}_\mathfrak{p}^{\otimes n}(|\widehat{\mathfrak{p}}(t) - \mathfrak{p}(t)|^s)$, for $t \in \mathbb{R}$, or by a global measure, e.g. $\mathbb{P}_\mathfrak{p}^{\otimes n}(\|\widehat{\mathfrak{p}} - \mathfrak{p}\|_{\mathcal{L}_s}^s) = \mathbb{P}_\mathfrak{p}^{\otimes n}(\lambda(|\widehat{\mathfrak{p}} - \mathfrak{p}|^s))$, with a focus on the special cases $s = 1$ and $s = 2$. For an introduction to Kernel density estimation we refer to the lecture course [Statistik 1](#) (§22 - §24).

Nonparametric regression. We describe the dependence of the variation of a real-valued random variable Y (response) on the variation of an explanatory random variable X by a functional relationship $\mathbb{E}(Y|X = x) = f(x)$ where f is an unknown functional parameter of interest. For a detailed discussion of the case of a deterministic explanatory variable we refer to Tsybakov [2009]. Here and subsequently, we restrict our attention to the special case of a real-valued explanatory variable X , and hence, a random vector (X, Y) taking values in $(\mathbb{R}^2, \mathcal{B}^2)$. The joint distribution of (X, Y) is uniquely determined by the functional parameter of interest f , the conditional distribution of the error $\xi := Y - f(X)$ given X and the marginal distribution of X which are generally all not known in advance. However, the joint distribution is typically parametrised by the regression function f only and we write shortly $(X, Y) \sim \mathbb{P}_f$. Thereby, the dependence on the marginal distribution \mathbb{P}^X of the regressor X and the conditional distribution of the error term ξ given X is usually not made explicit. For sake of simplicity, suppose in addition that the joint distribution \mathbb{P}_f of (X, Y) admits a joint Lebesgue density \mathfrak{p} . Denoting by \mathfrak{p}^X the marginal density of X we use for the conditional density $\mathfrak{p}^{Y|X}$ of Y given X the \mathbb{P}_f -a.s.

identity $\mathbb{P}^X \mathbb{P}^{Y|X} = \mathbb{P}$ which allows for \mathbb{P}_f -a.e. $x \in \mathbb{R}$ to write

$$\begin{aligned} \mathfrak{q}(x) &:= f(x) \mathbb{P}^X(x) = \mathbb{E}(Y | X = x) \mathbb{P}^X(x) \\ &= \int_{\mathbb{R}} y \mathbb{P}^{Y|X=x}(y) dy \mathbb{P}^X(x) = \int_{\mathbb{R}} y \mathbb{P}(x, y) dy. \end{aligned} \quad (09.01)$$

Consequently, the function of interest is \mathbb{P}_f -a.s. given by $f = \mathfrak{q} / \mathbb{P}^X$ which motivates the following estimation strategy. Given a sample of (X, Y) estimate separately \mathfrak{q} and \mathbb{P}^X , say by $\hat{\mathfrak{q}}$ and $\hat{\mathbb{P}}^X$, and then form a estimator $\hat{f} = \hat{\mathfrak{q}} / \hat{\mathbb{P}}^X$ (possibly in addition to be regularised). There are many different approaches including local smoothing techniques, orthogonal series estimation, penalised smoothing techniques and combinations of them, to name but a few. In the sequel let \mathcal{F} be a family of regression functions and for each $f \in \mathcal{F}$ denote by \mathbb{P}_f and \mathbb{E}_f the associated probability measure of (X, Y) and its expectation, respectively. We denote by $\mathbb{P}_{\mathcal{F}}$ the family of possible distributions of (X, Y) , but keep in mind, that the distribution \mathbb{P}_f of (X, Y) is generally not uniquely determined by $f \in \mathcal{F}$ only. If $\{(X_i, Y_i) : i \in \llbracket n \rrbracket\}$ form an independent and identically distributed (i.i.d.) sample of $(X, Y) \sim \mathbb{P}_f$ then $\mathbb{P}_f^{\otimes n} = \otimes_{j \in \llbracket n \rrbracket} \mathbb{P}_f$ denotes the joint product probability measure of the family $((X_i, Y_i))_{i \in \llbracket n \rrbracket}$. We write $((X_i, Y_i))_{i \in \llbracket n \rrbracket} \stackrel{i.i.d.}{\sim} \mathbb{P}_f$ or $((X_i, Y_i))_{i \in \llbracket n \rrbracket} \sim \mathbb{P}_f^{\otimes n}$ for short. We denote by $\mathbb{P}_{\mathcal{F}}^{\otimes n} := (\mathbb{P}_f^{\otimes n})_{f \in \mathcal{F}}$ the corresponding family of product probability measures. For $s \geq 1$ we assess also the accuracy of an estimator \hat{f} of f either by a local measure, e.g. $\mathbb{P}_f^{\otimes n}(\|\hat{f}(t) - f(t)\|^s)$, for $t \in \mathbb{R}$, or by a global measure, e.g. $\mathbb{P}_f^{\otimes n}(\|\hat{f} - f\|_{L^s}^s) = \mathbb{P}_f^{\otimes n}(\lambda(\|\hat{f} - f\|^s))$ with a focus on the special cases $s = 1$ and $s = 2$. For an introduction to smoothing techniques we refer to the lecture course [Statistik 1](#) (§22 - §24).

§10 Noisy version of the parameter

Let $(\mathbb{H}, \langle \cdot, \cdot \rangle_{\mathbb{H}})$ be a separable real Hilbert spaces. We are interested in the reconstruction of $\theta \in \mathbb{H}$ from a noisy version of it, which we formalise first in this section by introducing stochastic processes.

§10|01 Stochastic process

§10.01 Notation. Here and subsequently, a non-empty and generally non-finite subset \mathcal{J} of \mathbb{N}, \mathbb{Z} or \mathbb{R} and a subset \mathcal{U} of \mathcal{J} denote an index set. We consider the product spaces $\mathbb{R}^{\mathcal{J}} = \prod_{j \in \mathcal{J}} \mathbb{R}$ and $\mathbb{R}^{\mathcal{U}} = \prod_{u \in \mathcal{U}} \mathbb{R}$, where we identify the family $\mathbf{y} = (y_j)_{j \in \mathcal{J}} \in \mathbb{R}^{\mathcal{J}}$ and the map $y_{\cdot} : \mathcal{J} \rightarrow \mathbb{R}$ with $j \mapsto y_j$. The map $\Pi_{\mathcal{U}} : \mathbb{R}^{\mathcal{J}} \rightarrow \mathbb{R}^{\mathcal{U}}$ given by $\mathbf{y} = (y_j)_{j \in \mathcal{J}} \mapsto (y_j)_{j \in \mathcal{U}} =: \Pi_{\mathcal{U}} \mathbf{y}$, is called *canonical projection*. In particular, for each $j \in \mathcal{J}$ the *coordinate map* $\Pi_j := \Pi_{\{j\}} : \mathbb{R}^{\mathcal{J}} \rightarrow \mathbb{R}$ is given by $\mathbf{y} = (y_{j'})_{j' \in \mathcal{J}} \mapsto y_j =: \Pi_j \mathbf{y}$. Moreover, $\mathbb{R}^{\mathcal{J}}$ is equipped with the product Borel- σ -algebra $\mathcal{B}^{\otimes \mathcal{J}} := \otimes_{j \in \mathcal{J}} \mathcal{B}$. Recall that $\mathcal{B}^{\otimes \mathcal{J}}$ equals the smallest σ -algebra on $\mathbb{R}^{\mathcal{J}}$ such that all coordinate maps $\Pi_j, j \in \mathcal{J}$ are measurable. i.e., $\mathcal{B}^{\otimes \mathcal{J}} = \sigma(\Pi_j, j \in \mathcal{J})$. Moreover, let $(\mathcal{J}, \mathcal{J}, \nu)$ be a measure space with σ -finite $\nu \in \mathcal{M}_{\sigma}(\mathcal{J})$ and $\mathcal{L}_2(\nu) := \mathcal{L}_2(\mathcal{J}, \mathcal{J}, \nu)$ the usual set of square integrable real-valued functions defined on $(\mathcal{J}, \mathcal{J}, \nu)$. Define the set of equivalence classes $\mathbb{J} := \mathbb{L}_2(\nu) := \mathbb{L}_2(\mathcal{J}, \mathcal{J}, \nu)$, which forms a Hilbert space endowed with usual inner product $\langle \cdot, \cdot \rangle_{\mathbb{J}} := \langle \cdot, \cdot \rangle_{\mathbb{L}_2(\nu)}$ and induced norm $\|\cdot\|_{\mathbb{J}} := \|\cdot\|_{\mathbb{L}_2(\nu)}$. Eventually, we define arithmetic operations on elements of $\mathbb{R}^{\mathcal{J}}$ coordinate-wise, for example meaning $a \cdot b = (a_j b_j)_{j \in \mathcal{J}}$ and $r a_{\cdot} = (r a_j)_{j \in \mathcal{J}}$ for $a, b \in \mathbb{R}^{\mathcal{J}}$ and $r \in \mathbb{R}$. Let us further introduce $\mathbf{0} := (0)_{j \in \mathcal{J}}$ and $\mathbf{1} := (1)_{j \in \mathcal{J}}$. \square

§10.02 Comment. Given a measure space $(\Omega, \mathcal{A}, \mu)$, $s \in [1, \infty]$ and the usual space $\mathcal{L}_s(\Omega, \mathcal{A}, \mu)$ of $\mathcal{L}_s(\mu)$ -integrable functions introduce for each \mathcal{A} - \mathcal{B} -measurable $h : \Omega \rightarrow \mathbb{R}$, in short $h \in \mathcal{A}$, the

μ -equivalence class $\{h\}_\mu := \{h_o \in \mathcal{A} : h = h_o \mu\text{-a.e.}\}$. For $s \in \overline{\mathbb{R}}^+$ define the set of equivalence classes $\mathbb{L}_s(\mu) := \mathbb{L}_s(\mathcal{A}, \mu) := \mathbb{L}_s(\Omega, \mathcal{A}, \mu) := \{\{h\}_\mu : h \in \mathcal{L}_s(\mathcal{A}, \mu)\}$ and $\|\{h\}_\mu\|_{\mathbb{L}_s(\mu)} := \|h\|_{\mathcal{L}_s(\mu)}$ for $\{h\}_\mu \in \mathbb{L}_s(\mu)$. For $s \geq 1$, $(\mathbb{L}_s(\mu), \|\cdot\|_{\mathbb{L}_s(\mu)})$ is a normed vector space. Formally, we denote by $\{\cdot\}_\mu : \mathcal{L}_s(\mu) \rightarrow \mathbb{L}_s(\mu)$ the natural injection $h \mapsto \{h\}_\mu$. In case $s = 2$ the norm $\|\{h\}_\mu\|_{\mathbb{L}_2(\mu)} := \|h\|_{\mathcal{L}_2(\mu)} = (\mu(|h|^2))^{1/2}$ is induced by the inner product $(\{h\}_\mu, \{h_o\}_\mu) \mapsto \langle \{h\}_\mu, \{h_o\}_\mu \rangle_{\mathbb{L}_2(\mu)} := \mu(h h_o)$, and hence $(\mathbb{L}_2(\mu), \langle \cdot, \cdot \rangle_{\mathbb{L}_2(\mu)})$ is a Hilbert space. As usual we identify the equivalence class $\{h\}_\mu$ with its representative h , and write $h \in \mathbb{L}_2(\mu)$ for short. If $\lambda = \mu$ is the Lebesgue-measure then we write shortly $(\mathbb{L}_2, \langle \cdot, \cdot \rangle_{\mathbb{L}_2})$ and $\{\cdot\} : \mathcal{L}_2 \rightarrow \mathbb{L}_2$. \square

§10.03 **Stochastic process.** Let $(Y_j)_{j \in \mathcal{J}}$ be a family of real-valued random variables on a common probability space $(\Omega, \mathcal{A}, \mathbb{P})$, that is, $Y_j \in \mathcal{A}$ for each $j \in \mathcal{J}$. Consider the $\mathbb{R}^{\mathcal{J}}$ -valued random variable $Y_\cdot := (Y_j)_{j \in \mathcal{J}}$ where $Y_\cdot : \Omega \rightarrow \mathbb{R}^{\mathcal{J}}$ is a \mathcal{A} - $\mathcal{B}^{\otimes \mathcal{J}}$ -measurable map given by $\omega \mapsto (Y_j(\omega))_{j \in \mathcal{J}} =: Y_\cdot(\omega)$. Y_\cdot is called a *stochastic process*. Its *distribution* $\mathbb{P}^X := \mathbb{P} \circ Y_\cdot^{-1}$ is the image probability measure of \mathbb{P} under the map Y_\cdot , i.e. $Y_\cdot \sim \mathbb{P}^X$ for short. Further, denote by $\mathbb{P}^{Y_u} = \mathbb{P} \circ Y_u^{-1} = \mathbb{P}^X \circ \Pi_u^{-1}$ the distribution of the stochastic process $Y_u := \Pi_u Y_\cdot = (Y_u)_{u \in \mathcal{U}}$ on $\mathcal{U} \subseteq \mathcal{J}$. The family $(\mathbb{P}^{Y_u})_{\mathcal{U} \subseteq \mathcal{J} \text{ finite}}$ is called *family of finite-dimensional distributions* of Y_\cdot or \mathbb{P}^X . In particular, $\mathbb{P}^{Y_j} = \mathbb{P}^{\Pi_j} = \mathbb{P}^X \circ \Pi_j^{-1}$ denotes the distribution of $Y_j = \Pi_j Y_\cdot$. Furthermore, for $j, j_o \in \mathcal{J}$ we write $\mathbb{P}(Y_j) = \mathbb{P}^X(\Pi_j)$ and $\text{Cov}(Y_j, Y_{j_o}) := \mathbb{P}(Y_j Y_{j_o}) - \mathbb{P}(Y_j)\mathbb{P}(Y_{j_o}) = \mathbb{P}^X(\Pi_j \Pi_{j_o}) - \mathbb{P}^X(\Pi_j)\mathbb{P}^X(\Pi_{j_o})$, if it exists, for the expectation of Y_j and the covariance of Y_j and Y_{j_o} with respect to \mathbb{P}^X . \square

§10.04 **Assumption.** The stochastic process $Y_\cdot = (Y_j)_{j \in \mathcal{J}}$ on a common measurable space (Ω, \mathcal{A}) as a function $\Omega \times \mathcal{J} \rightarrow \mathbb{R}$ with $(\omega, j) \mapsto Y_j(\omega)$ is $\mathcal{A} \otimes \mathcal{J}$ - \mathcal{B} -measurable, $Y_\cdot \in \mathcal{A} \otimes \mathcal{J}$ for short. \square

§10.05 **Definition.** Let $Y_\cdot = (Y_j)_{j \in \mathcal{J}} \sim \mathbb{P}^X$ be a stochastic process satisfying Assumption §10.04. If $\mathbb{P}(|Y_j|) \in \mathbb{R}^+$, i.e. $Y_j \in \mathcal{L}_1(\mathbb{P})$ or $\Pi_j \in \mathcal{L}_1(\mathbb{P}^X)$ in equal, for each $j \in \mathcal{J}$, then $\mathbf{m}_\cdot := (\mathbf{m}_j := \mathbb{P}(Y_j))_{j \in \mathcal{J}} \in \mathbb{R}^{\mathcal{J}}$ is called *mean function* of Y_\cdot where $\mathbf{m}_\cdot \in \mathcal{J}$ due to Assumption §10.04. If in addition $\nu(\mathbf{m}_\cdot^2) \in \mathbb{R}^+$, i.e. $\mathbf{m}_\cdot \in \mathbb{J}$ then \mathbf{m}_\cdot is called (\mathbb{J} -) *mean*. If $\mathbb{P}(|Y_j|^2) \in \mathbb{R}^+$, i.e., $Y_j \in \mathcal{L}_2(\mathbb{P})$ or $\Pi_j \in \mathcal{L}_2(\mathbb{P}^X)$ in equal, for each $j \in \mathcal{J}$, then $\text{cov}_\cdot = (\text{cov}_{j,j_o} := \text{Cov}(Y_j, Y_{j_o}))_{j,j_o \in \mathcal{J}} \in \mathbb{R}^{\mathcal{J}^2}$ is called *covariance function* of Y_\cdot , where $\text{cov}_\cdot \in \mathcal{J}^2$ due to Assumption §10.04. A linear and bounded (continuous) operator from \mathbb{J} into itself, $\Gamma \in \mathbb{L}(\mathbb{J})$ for short, satisfying $\langle \Gamma x_\cdot, y_\cdot \rangle_{\mathbb{J}} = \text{Cov}(\nu(x_\cdot Y_\cdot), \nu(y_\cdot Y_\cdot)) = \int_{\mathcal{J}} \int_{\mathcal{J}} y_j \text{cov}_{j,j_o} x_{j_o} \nu(dj) \nu(dj_o)$ for all $y_\cdot, x_\cdot \in \mathbb{J} = \mathbb{L}_2(\nu)$ is called *covariance operator* of Y_\cdot or \mathbb{P}^X . If Y_\cdot admits a mean function $\mathbf{m}_\cdot \in \mathcal{J}$ (respectively mean $\mathbf{m}_\cdot \in \mathbb{J}$) and a covariance function $\text{cov}_\cdot \in \mathcal{J}^2$ (respectively covariance operator $\Gamma \in \mathbb{L}(\mathbb{J})$) then we write shortly $Y_\cdot \sim \mathbb{P}_{(\mathbf{m}_\cdot, \text{cov}_\cdot)}$ (respectively $Y_\cdot \sim \mathbb{P}_{(\mathbf{m}_\cdot, \Gamma)}$). \square

§10.06 **Remark.** A covariance operator $\Gamma \in \mathbb{L}(\mathbb{J})$ associated with a stochastic process $Y_\cdot \sim \mathbb{P}^X$ is self-adjoint and non-negative definite, $\Gamma \in \mathbb{L}^{\geq}(\mathbb{J})$ for short. If

$$\sup \{ \mathbb{P}(|\nu(y_\cdot Y_\cdot)|^2) : y_\cdot \in \mathbb{J} = \mathbb{L}_2(\nu), \|y_\cdot\|_{\mathbb{J}} \leq 1 \} \in \mathbb{R}^+,$$

which holds whenever $\mathbb{P}(\|Y_\cdot\|_{\mathbb{J}}^2) \in \mathbb{R}^+$ or in equal $\|Y_\cdot\|_{\mathbb{J}} \in \mathcal{L}_2(\mathbb{P})$ (implying $Y_\cdot \in \mathbb{J}$ \mathbb{P} -a.s.), then there exists a covariance operator $\Gamma \in \mathbb{L}^{\geq}(\mathbb{J})$ satisfying $\langle \Gamma x_\cdot, y_\cdot \rangle_{\mathbb{J}} = \text{Cov}(\nu(x_\cdot Y_\cdot), \nu(y_\cdot Y_\cdot))$. Observe that $\|Y_\cdot\|_{\mathbb{J}}^2 = \sup \{ |\nu(y_\cdot Y_\cdot)|^2 : y_\cdot \in \mathbb{J}, \|y_\cdot\|_{\mathbb{J}} \leq 1 \}$. Note that $\|Y_\cdot\|_{\mathbb{J}} \in \mathcal{L}_2(\mathbb{P})$ is a sufficient condition for the existence of a covariance operator, but it is not a necessary condition. \square

§10.07 **Empirical mean model.** Assume a probability space $(\mathcal{Z}, \mathcal{Z}, \mathbb{P})$ and a stochastic process $\psi_\cdot = (\psi_j)_{j \in \mathcal{J}} \in \mathcal{Z} \otimes \mathcal{J}$, i.e. $\mathcal{Z} \times \mathcal{J} \ni (z, j) \mapsto \psi_j(z) \in \mathbb{R}$ is $\mathcal{Z} \otimes \mathcal{J}$ - \mathcal{B} -measurable, satisfying in addition $\psi_j \in \mathcal{L}_1(\mathbb{P}) := \mathcal{L}_1(\mathcal{Z}, \mathcal{Z}, \mathbb{P})$ for each $j \in \mathcal{J}$. Consider the product probability

space $(\mathcal{Z}^n, \mathcal{Z}^{\otimes n}, \mathbb{P}^{\otimes n})$ and $Y_\cdot = (Y_j)_{j \in \mathcal{J}}$ with $Y_j := \widehat{\mathbb{P}}_n(\psi_j) \in \mathcal{Z}^{\otimes n}$ where $z^n = (z_i^n)_{i \in [n]} \mapsto Y_j(z^n) = (\widehat{\mathbb{P}}_n(\psi_j))(z^n) = n^{-1} \sum_{i \in [n]} \psi_j(z_i^n)$ for each $j \in \mathcal{J}$ and $Y_\cdot \in \mathcal{Z}^{\otimes n} \otimes \mathcal{J}$. By construction $m_\cdot = (m_j = \mathbb{P}(\psi_j))_{j \in \mathcal{J}} \in \mathcal{J}$ is the mean function of Y_\cdot . For each $j \in \mathcal{J}$ the statistic $\varepsilon_j := n^{1/2}(\widehat{\mathbb{P}}_n(\psi_j) - \mathbb{P}(\psi_j)) \in \mathcal{Z}^{\otimes n}$ is centred, i.e. $\varepsilon_j \in \mathcal{L}_1(\mathbb{P}^{\otimes n})$ with $\mathbb{P}^{\otimes n}(\varepsilon_j) = 0$, and $\varepsilon_\cdot = (\varepsilon_j)_{j \in \mathcal{J}} \in \mathcal{Z}^{\otimes n} \otimes \mathcal{J}$. Since $Y_j = m_j + n^{-1/2}\varepsilon_j$ for each $j \in \mathcal{J}$ by construction we write shortly $Y_\cdot = m_\cdot + n^{-1/2}\varepsilon_\cdot$. If for each $j \in \mathcal{J}$ in addition $\psi_j \in \mathcal{L}_2(\mathbb{P}) := \mathcal{L}_2(\mathcal{Z}, \mathcal{Z}, \mathbb{P})$ then we have $Y_j = \widehat{\mathbb{P}}_n(\psi_j) \in \mathcal{L}_2(\mathbb{P}^{\otimes n})$ and, hence $\varepsilon_j \in \mathcal{L}_2(\mathbb{P}^{\otimes n})$ by construction. The *covariance function* $\text{cov}_\cdot \in \mathcal{J}^2$ of $\varepsilon_\cdot = (\varepsilon_j)_{j \in \mathcal{J}}$ is given for each $j, j_o \in \mathcal{J}$ by

$$\text{cov}_{j,j_o} = \text{Cov}(\varepsilon_j, \varepsilon_{j_o}) = \mathbb{P}(\psi_j \psi_{j_o}) - \mathbb{P}(\psi_j)\mathbb{P}(\psi_{j_o}) = n \text{Cov}(Y_j, Y_{j_o}).$$

Consequently, we have $\varepsilon_\cdot \sim \mathbb{P}_{(0, \text{cov}_\cdot)}$ and $Y_\cdot = m_\cdot + n^{-1/2}\varepsilon_\cdot \sim \mathbb{P}_{(m_\cdot, n^{-1}\text{cov}_\cdot)}$. There exists a covariance operator $\Gamma \in \mathbb{L}(\mathbb{J})$, if in addition $\sup \{ \mathbb{P}(|\nu(y, \psi)|^2) : y \in \mathbb{J} = \mathbb{L}_2(\nu), \|y\|_{\mathbb{J}} \leq 1 \} \in \mathbb{R}^+$, which holds whenever $\|\psi_\cdot\|_{\mathbb{J}} \in \mathcal{L}_2(\mathbb{P})$ or in equal $\mathbb{P}(\|\psi_\cdot\|_{\mathbb{J}}^2) \in \mathbb{R}^+$. Observe that $\|\psi_\cdot\|_{\mathbb{J}}^2 = \sup \{ |\nu(y, \psi)|^2 : y \in \mathbb{J}, \|y\|_{\mathbb{J}} \leq 1 \}$. Note that $\|\psi_\cdot\|_{\mathbb{J}} \in \mathcal{L}_2(\mathbb{P})$ is a sufficient condition for the existence of a covariance operator, but it is not necessary. \square

§10.08 **White noise process.** A stochastic process $\dot{W}_\cdot = (\dot{W}_j)_{j \in \mathcal{J}}$ is called *white noise process*, if $(\dot{W}_j)_{j \in \mathcal{J}}$ is a family of independent and identically distributed random variables, where each \dot{W}_j has zero mean and variance one, $\dot{W}_j \sim \mathbb{P}_{(0,1)}$ and $\dot{W}_\cdot \sim \mathbb{P}_{(0,1)}^{\otimes \mathcal{J}}$ in short. \square

§10.09 **Notation.** In other words, the distribution $\mathbb{P}^{\dot{W}_\cdot}$ of a white noise process $\dot{W}_\cdot = (\dot{W}_j)_{j \in \mathcal{J}} \sim \mathbb{P}^{\dot{W}_\cdot}$ equals the product of its marginal $\mathbb{P}_{(0,1)}$ -distributions, i.e. $\mathbb{P}^{\dot{W}_\cdot} = \otimes_{j \in \mathcal{J}} \mathbb{P}^{\dot{W}_j} = \otimes_{j \in \mathcal{J}} \mathbb{P}_{(0,1)} = \mathbb{P}_{(0,1)}^{\otimes \mathcal{J}}$. \square

§10.10 **Remark.** The centred stochastic process $\varepsilon_\cdot := (\varepsilon_j)_{j \in \mathcal{J}}$ of error terms in an Empirical mean model §10.07 is in general not a white noise process. \square

§10.11 **Notation.** We denote by $\ell_2 := \mathbb{L}_2(\nu_{\mathbb{N}}) = \mathbb{L}_2(\mathbb{N}, 2^{\mathbb{N}}, \nu_{\mathbb{N}}) = \mathbb{J}$ the space of all square-summable real-valued sequences endowed with counting measure $\nu_{\mathbb{N}} := \sum_{j \in \mathbb{N}} \delta_{\{j\}}$ over the index set \mathbb{N} . \square

§10.12 **Property.** Let $\dot{W}_\cdot := (\dot{W}_j)_{j \in \mathbb{N}} \sim \mathbb{P}_{(0,1)}^{\otimes \mathbb{N}}$ be a white noise process. By assumption \dot{W}_\cdot admits $0_\cdot := (0)_{j \in \mathbb{N}}$ as ℓ_2 -mean and $\Gamma = \text{id}_{\ell_2} \in \mathbb{L}(\ell_2)$ as covariance operator, i.e. $\dot{W}_\cdot \sim \mathbb{P}_{(0, \text{id}_{\ell_2})}$, since $\langle x_\cdot, y_\cdot \rangle_{\ell_2} = \sum_{j \in \mathbb{N}} y_j x_j = \sum_{j \in \mathbb{N}} y_j \sum_{j_o \in \mathbb{N}} \text{cov}_{j,j_o} x_{j_o} = \langle \Gamma x_\cdot, y_\cdot \rangle_{\ell_2}$. \square

§10.13 **Gaussian process.** A stochastic process $Y_\cdot = (Y_j)_{j \in \mathcal{J}} \sim \mathbb{P}_{(m_\cdot, \text{cov}_\cdot)}$ satisfying Assumption §10.04 with mean function $m_\cdot \in \mathcal{J}$ and covariance function $\text{cov}_\cdot \in \mathcal{J}^2$ is called a *Gaussian process*, if the family of finite-dimensional distributions $(\mathbb{P}^{Y_\cdot})_{U \subseteq \mathcal{J} \text{ finite}}$ consists of normal distributions, that is, $Y_U = (Y_u)_{u \in U}$ is normally distributed with mean vector $(m_u)_{u \in U}$ and covariance matrix $(\text{cov}_{u,u'})_{u,u' \in U}$. We write shortly $Y_\cdot \sim \mathbb{N}_{(m_\cdot, \text{cov}_\cdot)}$ or $Y_\cdot \sim \mathbb{N}_{(m_\cdot, \Gamma_\cdot)}$, if in addition there exist a covariance operator $\Gamma_\cdot \in \mathbb{L}(\mathbb{J})$ associated with Y_\cdot . The Gaussian process $\dot{B}_\cdot \sim \mathbb{N}_{(0, \text{id}_{\mathbb{J}})}$ with \mathbb{J} -mean zero and covariance operator $\text{id}_{\mathbb{J}}$ is called *iso-Gaussian process* or *Gaussian white noise process*, which equals $\dot{B}_\cdot \sim \mathbb{N}_{(0,1)}^{\otimes \mathbb{N}}$ in the particular case $\mathbb{J} = \mathbb{L}_2(\nu_{\mathbb{N}}) = \ell_2$. \square

§10.14 **Definition (Random function).** Let $(\mathbb{H}, \langle \cdot, \cdot \rangle_{\mathbb{H}})$ be an Hilbert space equipped with its Borel- σ -algebra $\mathcal{B}_{\mathbb{H}}$, which is induced by its topology. An \mathcal{A} - $\mathcal{B}_{\mathbb{H}}$ -measurable map $Y : (\Omega, \mathcal{A}) \rightarrow (\mathbb{H}, \mathcal{B}_{\mathbb{H}})$ is called an \mathbb{H} -valued random variable or a *random function* in \mathbb{H} . \square

§10.15 **Lemma.** Consider $(\ell_2, \langle \cdot, \cdot \rangle_{\ell_2})$. There does not exist a non-zero random function $Y_\cdot = (Y_j)_{j \in \mathbb{N}}$ in ℓ_2 which is a Gaussian white noise process. \square

§10.16 **Proof of Lemma §10.15.** Exercise. \square

§10|02 Noisy parameter

§10.17 **Assumption.** The Hilbert space $\mathbb{J} = \mathbb{L}_2(\mathcal{J}, \mathcal{J}, \nu)$ and the surjective partial isometry $\mathbf{U} \in \mathbb{L}(\mathbb{H}, \mathbb{J})$, i.e. $\mathbf{U}\mathbf{U}^* = \text{id}_{\mathbb{J}}$, are fixed and presumed to be known in advance. \square

§10.18 **Notation.** Here and subsequently, we write $\theta = (\theta_j)_{j \in \mathcal{J}} := \mathbf{U}\theta \in \mathbb{J}$. Keep in mind, that we identify equivalence classes and their representatives. Our aim is the reconstruction of θ and hence $\mathbf{U}^*\theta \in \mathbb{H}$ from a noisy version of θ . \square

§10.19 **Noisy parameter.** Let $\varepsilon = (\varepsilon_j)_{j \in \mathcal{J}}$ be a stochastic process satisfying Assumption §10.04 with mean zero and let $n \in \mathbb{N}$ be a sample size. The stochastic process $\widehat{\theta} = \theta + n^{-1/2}\varepsilon$, with \mathbb{J} -mean θ is called a *noisy version* of the parameter $\theta \in \mathbb{J}$, or *noisy parameter* for short. We denote by \mathbb{P}_θ^n the distribution of $\widehat{\theta}$. If ε admits (possibly depending on θ) a covariance function, say $\text{cov}_{\cdot, \cdot} \in \mathcal{J}^2$, or a covariance operator, say $\Gamma \in \mathbb{L}(\mathbb{J})$, then we eventually write $\varepsilon \sim \mathbb{P}_{(\theta, \text{cov}_{\cdot, \cdot})}$ and $\widehat{\theta} \sim \mathbb{P}_{(\theta, n^{-1}\text{cov}_{\cdot, \cdot})}$ or $\varepsilon \sim \mathbb{P}_{(\theta, \Gamma)}$ and $\widehat{\theta} \sim \mathbb{P}_{(\theta, n^{-1}\Gamma)}$ for short. The reconstruction of $\theta \in \mathbb{J}$ (or in equal $\mathbf{U}^*\theta \in \mathbb{H}$) from a noisy version $\widehat{\theta} \sim \mathbb{P}_\theta^n$ is called a *statistical direct problem*. \square

§10.20 **Sequence space model.** Consider $\mathbb{J} = \ell_2 = \mathbb{L}_2(\nu_{\mathbb{N}})$. Let $\varepsilon = (\varepsilon_j)_{j \in \mathbb{N}}$ be a real-valued stochastic process satisfying Assumption §10.04 with mean 0 $\in \ell_2$ and let $n \in \mathbb{N}$ be a sample size. The observable noisy version $\widehat{\theta} = \theta + n^{-1/2}\varepsilon \sim \mathbb{P}_\theta^n$ with ℓ_2 -mean $\theta \in \ell_2$ as in §10.11 takes the form of a *sequence space model (SSM)*

$$\widehat{\theta}_j = \theta_j + n^{-1/2}\varepsilon_j, \quad j \in \mathbb{N}. \quad (10.01)$$

If ε admits a covariance function (possibly depending on θ), say $\text{cov}_{\cdot, \cdot} \in 2^{\mathbb{N}^2}$, then we eventually write $\widehat{\theta} \sim \mathbb{P}_{(\theta, n^{-1}\text{cov}_{\cdot, \cdot})}$ for short. If in addition ε admits a covariance operator $\Gamma \in \mathbb{L}(\ell_2)$ (an infinite matrix) then we write $\widehat{\theta} \sim \mathbb{P}_{(\theta, n^{-1}\Gamma)}$. \square

§10.21 **Gaussian sequence space model.** Let $\dot{\mathbf{B}} := (\dot{\mathbf{B}}_j)_{j \in \mathbb{N}} \sim \mathbb{N}_{(0,1)}^{\otimes \mathbb{N}}$ be a Gaussian white noise process. The observable noisy version $\widehat{\theta} = \theta + n^{-1/2}\dot{\mathbf{B}}$ with ℓ_2 -mean $\theta \in \ell_2$ takes the form of a *Gaussian sequence space model (GSSM)*

$$\widehat{\theta}_j = \theta_j + n^{-1/2}\dot{\mathbf{B}}_j, \quad j \in \mathbb{N} \quad \text{with} \quad (\dot{\mathbf{B}}_j)_{j \in \mathbb{N}} \sim \mathbb{N}_{(0,1)}^{\otimes \mathbb{N}} \quad (10.02)$$

and we denote by \mathbb{N}_θ^n the distribution of the stochastic process $\widehat{\theta}$. \square

§10.22 **Notation.** Consider the measure space $([0, 1], \mathcal{B}_{[0,1]}, \lambda_{[0,1]})$ where $\lambda_{[0,1]}$ denotes the restriction of the Lebesgue measure to the Borel- σ -algebra $\mathcal{B}_{[0,1]}$ over $[0, 1]$, and the Hilbert space $\mathbb{L}_2(\lambda_{[0,1]}) := \mathbb{L}_2([0, 1], \mathcal{B}_{[0,1]}, \lambda_{[0,1]})$. Assume that $\theta \in \mathbb{L}_2(\lambda_{[0,1]}) =: \mathbb{H}$. Consider an *orthonormal system* $(\mathbf{u}_j)_{j \in \mathbb{N}}$ in $\mathbb{L}_2(\lambda_{[0,1]})$. Then $\mathbf{U} : \mathbb{L}_2(\lambda_{[0,1]}) \rightarrow \ell_2$ with $h \mapsto \mathbf{U}h := h_\bullet = (h_j := \langle h, \mathbf{u}_j \rangle_{\mathbb{H}})_{j \in \mathbb{N}}$ is a surjective partial isometry $\mathbf{U} \in \mathbb{L}(\mathbb{L}_2(\lambda_{[0,1]}), \ell_2)$. Its adjoint operator $\mathbf{U}^* \in \mathbb{L}(\ell_2, \mathbb{L}_2(\lambda_{[0,1]}))$ satisfies $\mathbf{U}^*a_\bullet = \sum_{j \in \mathbb{N}} a_j \mathbf{u}_j$ for all $a_\bullet \in \ell_2$. We call $h_\bullet = (h_j)_{j \in \mathbb{N}}$ (*generalised*) *Fourier coefficients* and \mathbf{U} (*generalised*) *Fourier series transform*. \square

§10.23 **Nonparametric density estimation on $[0, 1]$.** Let \mathbb{D}_2 be a set of square-integrable Lebesgue densities on $([0, 1], \mathcal{B}_{[0,1]})$, and hence $\mathbb{D}_2 \subseteq \mathbb{L}_2(\lambda_{[0,1]}) =: \mathbb{H}$. We denote for each density $\mathbb{p} \in \mathbb{D}_2$ by $\mathbb{P}_\mathbb{p} := \mathbb{p} \lambda_{[0,1]}$ and $\mathbb{E}_\mathbb{p}$ the associated probability measure and expectation, respectively. Assuming an iid. sample $(X_i)_{i \in [n]}$ of size $n \in \mathbb{N}$ we consider the statistical product experiment $([0, 1]^n, \mathcal{B}_{[0,1]}^{\otimes n}, \mathbb{P}_\mathbb{p}^{\otimes n} := (\mathbb{P}_\mathbb{p}^{\otimes n})_{\mathbb{p} \in \mathbb{D}_2})$. Let $\mathbf{U} \in \mathbb{L}(\mathbb{L}_2(\lambda_{[0,1]}), \ell_2)$ be a generalised Fourier series transform (see **Notation** §10.22) which is fixed and known in advanced. Evidently, for each $\mathbb{p} \in \mathbb{D}_2 \subseteq \mathbb{L}_2(\lambda_{[0,1]})$ the generalised Fourier coefficients $\mathbb{p} = (\mathbb{p}_j)_{j \in \mathbb{N}} = \mathbf{U}\mathbb{p}$ satisfy

$\mathbb{p} = \langle \mathbb{p}, u_j \rangle_{\mathbb{H}} = \lambda_{[0,1]}(\mathbb{p} u_j) = \mathbb{p} \lambda_{[0,1]}(u_j) = \mathbb{P}_p(u_j)$, i.e. $u_j \in \mathbb{L}_1([0,1], \mathcal{B}_{[0,1]}, \mathbb{P}_p) =: \mathbb{L}_1(\mathbb{P}_p)$, for each $j \in \mathbb{N}$. Moreover, the stochastic process $(u_j)_{j \in \mathbb{N}}$ on $([0,1], \mathcal{B}_{[0,1]}, \mathbb{P}_p)$ is $\mathcal{B}_{[0,1]} \otimes 2^{\mathbb{N}}$ - \mathcal{B} -measurable. Similar to an Empirical mean model §10.07 we define $\widehat{\mathbb{p}} = (\widehat{\mathbb{p}} := \widehat{\mathbb{P}}_n(u_j))_{j \in \mathbb{N}} \in \mathcal{B}_{[0,1]}^{\otimes n} \otimes 2^{\mathbb{N}}$ where $x^n = (x_i)_{i \in [n]} \mapsto \widehat{\mathbb{p}}(x^n) = (\widehat{\mathbb{P}}_n(u_j))(x^n) = n^{-1} \sum_{i \in [n]} u_j(x_i^n)$ for each $j \in \mathbb{N}$. By construction $\mathbb{p} = (\mathbb{p} = \mathbb{P}_p(u_j))_{j \in \mathbb{N}} \in 2^{\mathbb{N}}$ is the mean function of $\widehat{\mathbb{p}}$. For each $j \in \mathbb{N}$ the statistic $\varepsilon_j := n^{1/2}(\widehat{\mathbb{P}}_n(u_j) - \mathbb{P}_p(u_j)) \in \mathcal{B}_{[0,1]}^{\otimes n}$ is centred, i.e. $\varepsilon_j \in \mathbb{L}_1([0,1]^n, \mathcal{B}_{[0,1]}^{\otimes n}, \mathbb{P}_p^{\otimes n}) =: \mathbb{L}_1(\mathbb{P}_p^{\otimes n})$ with $\mathbb{P}_p^{\otimes n}(\varepsilon_j) = 0$, and $\varepsilon = (\varepsilon_j)_{j \in \mathbb{N}} \in \mathcal{B}_{[0,1]}^{\otimes n} \otimes 2^{\mathbb{N}}$. Since $\widehat{\mathbb{p}} = \mathbb{p} + n^{-1/2} \varepsilon$ for each $j \in \mathbb{N}$ by construction $\widehat{\mathbb{p}} = \mathbb{p} + n^{-1/2} \varepsilon$ is a noisy version of \mathbb{p} . \square

§11 Orthogonal projection

§11.01 **Notation (Reminder).** Consider a measure space $(\mathcal{J}, \mathcal{J}, \nu)$ as in **Notation** §10.01. For $w \in \mathbb{R}^{\mathcal{J}}$ define the multiplication map $M_w : \mathbb{R}^{\mathcal{J}} \rightarrow \mathbb{R}^{\mathcal{J}}$ with $a \mapsto M_w a := w \cdot a := (w_j a_j)_{j \in \mathcal{J}}$. If $w \in \mathcal{J}$, i.e. w is \mathcal{J} - \mathcal{B} -measurable, then we have $M_w : \mathcal{J} \rightarrow \mathcal{J}$ too. We denote by $\mathbb{M}_{\mathcal{J}}$ the set of all multiplication maps defined on \mathcal{J} . If in addition $w \in \mathcal{L}_{\infty}(\nu)$ then we have also $M_w : \mathbb{L}_2(\nu) = \mathbb{J} \rightarrow \mathbb{J}$ identifying eventually equivalence classes and representatives. We set $\mathbb{L}(\mathbb{J}) := \{M_w \in \mathbb{M}_{\mathcal{J}} : w \in \mathcal{L}_{\infty}(\nu)\} \subseteq \mathbb{L}(\mathbb{J})$ noting that $\|M_w\|_{\mathbb{L}(\mathbb{J})} = \sup\{\|w \cdot a\|_{\mathbb{J}} : \|a\|_{\mathbb{J}} \leq 1\} \leq \|w\|_{\mathcal{L}_{\infty}(\nu)}$ for each $M_w \in \mathbb{L}(\mathbb{J})$. \square

§11.02 **Notation.** For $A \in \mathcal{J}$ we denote by $\mathbb{1}^A = (\mathbb{1}_j^A)_{j \in \mathcal{J}}$ the indicator function where for each $j \in \mathcal{J}$, $\mathbb{1}_j^A = 1$ if $j \in A$ and $\mathbb{1}_j^A = 0$ otherwise. Obviously, $\mathbb{1}^A$ is \mathcal{J} - \mathcal{B} -measurable, i.e. $\mathbb{1}^A \in \mathcal{J}$, and it belongs to $\mathbb{L}_{\infty}(\nu)$, and to $\mathbb{L}_2(\nu)$ whenever $\nu(A) \in \mathbb{R}^+$. Since $\{j\} \in \mathcal{J}$ we have $\mathbb{1}^{\{j\}} \in \mathcal{J}$ and $\mathbb{1}^{\{j\}} \in \mathbb{L}_{\infty}(\nu)$. Obviously, we have $\mathbb{1} = \mathbb{1}^{\mathcal{J}} \in \mathbb{L}_{\infty}(\nu)$ and $M_{\mathbb{1}} \in \mathbb{L}(\mathbb{J})$. For each $w \in \mathcal{L}_{\infty}(\nu)$ set $\mathbb{J}w := \{a \cdot w\}_{\nu} : a \in \mathcal{L}_2(\nu)\} = \{a \cdot w : a \in \mathbb{J} = \mathbb{L}_2(\nu)\}$ and hence in particular $\mathbb{J}\mathbb{1}^A = \{a \cdot \mathbb{1}^A : a \in \mathbb{J}\}$. Given $0 = (0)_{j \in \mathcal{J}}$ for $w \in \mathcal{J}$ we write further $\mathcal{N}_w := \{w = 0\} := \{j \in \mathcal{J} : w_j = 0\} \in \mathcal{J}$, and denote by $\text{dom}(M_w) = \{a \in \mathbb{J} : a \cdot w \in \mathbb{J}\}$, $\text{ran}(M_w) = \{a \cdot w : a \in \text{dom}(M_w) \subseteq \mathbb{J}\}$ and $\text{ker}(M_w) = \{a \in \mathbb{J} : \{a \cdot w\}_{\nu} = 0\}$, respectively, the domain, range and nullspace of $M_w : \mathbb{J} \supseteq \text{dom}(M_w) \rightarrow \mathbb{J}$. We write $w \in \mathcal{J}_0$, if $w \in \mathcal{J}$ and $\nu(\mathcal{N}_w) = 0$. Similarly, if $w \in (\mathbb{R}^+)^{\mathcal{J}}$ is \mathcal{J} - \mathcal{B}^+ -measurable, then we write $w \in \mathcal{J}^+$, and $w \in \mathcal{J}_0^+$ assuming additionally $\nu(\mathcal{N}_w) = 0$. \square

§11.03 **Property.** For each $w \in \mathcal{J}^+ \cap \mathcal{L}_{\infty}(\nu)$ the multiplication $M_w \in \mathbb{L}(\mathbb{J}) \subseteq \mathbb{L}(\mathbb{J})$ is a positive semi-definite operator. Keeping $\mathcal{N}_w = \{w = 0\} \in \mathcal{J}$ in mind its range and null space is given by $\text{ran}(M_w) = \mathbb{J}w$ and $\text{ker}(M_w) = \mathbb{J}\mathbb{1}^{\mathcal{N}_w} = \text{ran}(M_{\mathbb{1}^{\mathcal{N}_w}})$, respectively. $M_w \in \mathbb{L}(\mathbb{J})$ is consequently **injective** if and only if $w \in \mathcal{J}_0^+$, i.e. $w \in \mathcal{J}$ and $\nu(\mathcal{N}_w) = 0$. For each $A \in \mathcal{J}$ setting $A^c := \mathcal{J} \setminus A \in \mathcal{J}$ the range and null space of $M_{\mathbb{1}^A} \in \mathbb{L}(\mathbb{J}) \subseteq \mathbb{L}(\mathbb{J})$ is given by $\text{ran}(M_{\mathbb{1}^A}) = \mathbb{J}\mathbb{1}^A$ and $\text{ker}(M_{\mathbb{1}^A}) = \mathbb{J}\mathbb{1}^{A^c}$, respectively. Obviously, we have $M_{\mathbb{1}^A}^2 = M_{\mathbb{1}^A}$ and hence $M_{\mathbb{1}^A}$ is an **orthogonal projection** and $\mathbb{J} = \mathbb{J}\mathbb{1}^A \oplus \mathbb{J}\mathbb{1}^{A^c}$. Moreover, the map $M_{\mathbb{1}} = \text{id}_{\mathbb{J}}$ equals the identity on \mathbb{J} . \square

§11|01 Weighted norms and inner products

§11.04 **Notation.** Extending the real line by the points $-\infty$ and $+\infty$ we define $\overline{\mathbb{R}} := \mathbb{R} \cup \{\pm\infty\}$. We denote by $\overline{\mathcal{B}}$ the Borel- σ -field over $\overline{\mathbb{R}}$ and note that the trace of $\overline{\mathcal{B}} \cap \mathbb{R}$ over \mathbb{R} equals \mathcal{B} . Thereby, each $a \in \mathcal{J}^+$ is in a canonical way also \mathcal{J} - $\overline{\mathcal{B}}^+$ measurable, $a \in \overline{\mathcal{J}}^+$ for short. For $w \in \overline{\mathcal{J}}$ and hence $w^2 \in \overline{\mathcal{J}}^+$, consider the measure $w^2 \nu$ on $(\mathcal{J}, \mathcal{J})$, i.e., $w^2 = dw^2 \nu / d\nu$ is the Radon-Nikodym density of $w^2 \nu$ with respect to ν . We write shortly $\langle \cdot, \cdot \rangle_w := \langle \cdot, \cdot \rangle_{\mathbb{L}_2(w^2 \nu)}$ and $\|\cdot\|_w := \|\cdot\|_{\mathbb{L}_2(w^2 \nu)}$. For $w \in \mathcal{J}$ we denote its Moore-Penrose inverse by $w^\dagger := w^{-1} \mathbb{1}^{\mathcal{N}_w^c} \in \mathcal{J}$ meaning $w_j^\dagger := w_j^{-1}$ if $j \in \mathcal{N}_w^c$ and $w_j^\dagger := 0$ if $j \in \mathcal{N}_w$. Obviously, we have $w^\dagger w \cdot w^\dagger = w^\dagger$,

$w.w^\dagger w = w$ and $w.w^\dagger = w^\dagger w = \mathbb{1}_w^{\mathcal{N}_w}$. We set $\mathbb{J}^w := \mathbb{L}_2^w(\nu) := \text{dom}(M_w)$ and write $w^{2\dagger} := (w^\dagger)^2 = (w^2)^\dagger$ for short. \square

§11.05 **Property.** Let $w \in \overline{\mathcal{J}}$. Then for each $a \in \mathcal{L}_2(w^2\nu)$ we have $w^2\nu(|a|^2) = \nu(|w.a|^2)$. If $w \in \mathbb{R}$ ν -a.e., then $w^2\nu \in \mathcal{M}_\sigma(\mathcal{J})$ is a σ -finite measure and $\mathbb{L}_2(w^2\nu)$ endowed with inner product $\langle \cdot, \cdot \rangle_w = \langle \cdot, \cdot \rangle_{\mathbb{L}_2(w^2\nu)} = \langle M_w \cdot, M_w \cdot \rangle_{\mathbb{L}_2(\nu)}$ is a separable Hilbert space. If in addition $w \in \mathcal{L}_\infty(\nu)$, then

$$\mathcal{L}_2(w^{2\dagger}\nu) = \mathcal{L}_2(\nu)w + \mathcal{J}\mathbb{1}_w^{\mathcal{N}_w} = \{w.h : h \in \mathcal{L}_2(\nu)\} + \{h.\mathbb{1}_w^{\mathcal{N}_w} : h \in \mathcal{J}\}. \quad (11.01)$$

Indeed, for each $h \in \mathcal{J}$ consider the decomposition $h = w.w^\dagger h + h.\mathbb{1}_w^{\mathcal{N}_w}$. The claim follows immediately from the equivalence of $h \in \mathcal{L}_2(w^{2\dagger}\nu)$ and $w^\dagger h \in \mathcal{L}_2(\nu)$. Under $w \in \mathcal{L}_\infty(\nu)$ the map $M_w : \mathcal{L}_2(\nu) \rightarrow \mathcal{L}_2(\nu)$ is well-defined, and setting $\text{dom}(M_w) = \{h \in \mathcal{L}_2(\nu) : w^\dagger h \in \mathcal{L}_2(\nu)\} = \mathcal{L}_2(\nu)w + \mathcal{L}_2(\nu)\mathbb{1}_w^{\mathcal{N}_w} \subseteq \mathcal{L}_2(w^{2\dagger}\nu)$ (similar to (11.01)). Consequently, if in addition $\nu(\mathcal{N}_w) = 0$, then $\text{dom}(M_w) = \mathcal{L}_2(w^{2\dagger}\nu)$. If $w \in \mathbb{L}_\infty(\nu)$ then $M_w \in \mathbb{L}(\mathbb{J})$, and $M_w : \mathbb{J} \supseteq \text{dom}(M_w) \rightarrow \mathbb{J}$. Moreover, we have $\text{dom}(M_w) = \mathbb{J}$, $\text{ran}(M_w) = \mathbb{J}w$ and $\text{ker}(M_w) = \mathbb{J}\mathbb{1}_w^{\mathcal{N}_w}$ (see **Property §11.03**). Therewith, it follows $\text{dom}(M_w) = \mathbb{J}w \oplus \mathbb{J}\mathbb{1}_w^{\mathcal{N}_w}$. Consequently, if in addition $\nu(\mathcal{N}_w) = 0$, then $\mathbb{J}^w = \mathbb{L}_2^w(\nu) = \text{dom}(M_w) = \mathbb{J}w = \mathbb{L}_2(w^{2\dagger}\nu)$. The last equality follows from (11.01) since both measures $w^{2\dagger}\nu$ and ν share the same null sets (i.e. they mutually dominate each other). \square

§11|02 Orthogonal projection

§11.06 **Notation.** For a non-empty and generally non-finite subset \mathcal{J} of \mathbb{N} , \mathbb{Z} or \mathbb{R} and $m \in \mathbb{N}$ we set $\llbracket m \rrbracket := [-m, m] \cap \mathcal{J}$ and we write shortly $\mathbb{1}^m = (\mathbb{1}_j^m)_{j \in \mathcal{J}} := \mathbb{1}^{\llbracket m \rrbracket}$. Furthermore, we define $\mathbb{1}^{m\perp} := \mathbb{1} - \mathbb{1}^m$. \square

§11.07 **Property.** For each $m \in \mathbb{N}$, $M_{\mathbb{1}^m} \in \mathbb{L}(\mathbb{J})$ and $M_{\mathbb{1}^{m\perp}} \in \mathbb{L}(\mathbb{J})$ is the *orthogonal projection* onto the linear subspace $\mathbb{J}\mathbb{1}^m \subseteq \mathbb{J}$ and its orthogonal complement $\mathbb{J}\mathbb{1}^{m\perp} = (\mathbb{J}\mathbb{1}^m)^\perp \subseteq \mathbb{J}$, respectively, that is $\mathbb{J} = \mathbb{J}\mathbb{1}^m \oplus \mathbb{J}\mathbb{1}^{m\perp}$. We have point-wise $\mathbb{1}^m - \mathbb{1} = o(1)$ as $m \rightarrow \infty$ meaning that for each $j \in \mathcal{J}$ holds $\mathbb{1}_j^m - \mathbb{1}_j = o(1)$ as $m \rightarrow \infty$. Considering the orthogonal projection $M_{\mathbb{1}^m} \in \mathbb{L}(\mathbb{J})$ and the identity $M_{\mathbb{1}} = \text{id}_{\mathbb{J}} \in \mathbb{L}(\mathbb{J})$ point-wise convergence $M_{\mathbb{1}^m} - \text{id}_{\mathbb{J}} = o(1)$ as $m \rightarrow \infty$ holds too, that is, $\|(M_{\mathbb{1}^m} - \text{id}_{\mathbb{J}})a\|_{\mathbb{J}} = \|(\mathbb{1}^m - \mathbb{1})a\|_{\mathbb{J}} = \|\mathbb{1}^{m\perp}a\|_{\mathbb{J}} = o(1)$ as $m \rightarrow \infty$ for all $a \in \mathbb{J}$. \square

§11.08 **Orthogonal projection.** Given $m \in \mathbb{N}$ we define for each $\theta = U\theta \in \mathbb{J}$ its orthogonal projection $\theta^m := \theta.\mathbb{1}^m \in \mathbb{J}\mathbb{1}^m$ (and $\theta^m := U^*\theta^m \in \mathbb{H}$). \square

§11|03 Global and maximal global \mathfrak{v} -error

We shall measure first globally the accuracy of the orthogonal projection $\theta^m := \theta.\mathbb{1}^m$ of $\theta \in \mathbb{J}$.

§11.09 **Property.** If $\mathfrak{v} \in \mathcal{J}_{\nu_0}$ (i.e. $\nu(\mathcal{N}_{\mathfrak{v}}) = 0$) and $\theta \in \mathbb{L}_2(\mathfrak{v}^2\nu)$ (i.e. $\|\theta\|_{\mathfrak{v}}^2 = \mathfrak{v}^2\nu(\theta^2) \in \mathbb{R}^+$), then for each $m \in \mathbb{N}$ we have $\theta^m \in \mathbb{L}_2(\mathfrak{v}^2\nu)$ too, since $\|\theta^m\|_{\mathfrak{v}}^2 = \mathfrak{v}^2\nu(\theta^2.\mathbb{1}^m) \leq \mathfrak{v}^2\nu(\theta^2)$. Moreover, it holds $\|\theta^m - \theta\|_{\mathfrak{v}}^2 = \|\theta.\mathbb{1}^{m\perp}\|_{\mathfrak{v}}^2 = \mathfrak{v}^2\nu(\theta^2.\mathbb{1}^{m\perp}) \leq \mathfrak{v}^2\nu(\theta^2) \in \mathbb{R}^+$ and $\|\theta^m - \theta\|_{\mathfrak{v}}^2 = o(1)$ as $m \rightarrow \infty$ by dominated convergence. \square

§11.10 **Comment.** We assume throughout this chapter that the Hilbert space $\mathbb{J} = \mathbb{L}_2(\mathcal{J}, \mathcal{J}, \nu)$ and the surjective partial isometry $U \in \mathbb{L}(\mathbb{H}, \mathbb{J})$ is fixed and known in advance. Considering a \mathfrak{v} -error means the weight sequences $\mathfrak{v} \in \mathcal{J}$ is also fixed and known in advance. Consequently, the condition $\mathfrak{v} \in \mathcal{J}_{\nu_0}$ does not impose an additional restriction. \square

§11.11 **Global \mathfrak{v} -error.** Given $\mathfrak{v} \in \mathcal{J}_{\nu_0}$, $m \in \mathbb{N}$, a parameter $\theta = U\theta \in \mathbb{L}_2(\mathfrak{v}^2\nu)$ and its orthogonal projection $\theta^m = \theta.\mathbb{1}^m \in \mathbb{J}\mathbb{1}^m$ we call $\|\theta^m - \theta\|_{\mathfrak{v}} = \|\theta.\mathbb{1}^{m\perp}\|_{\mathfrak{v}} \in \mathbb{R}^+$ *global \mathfrak{v} -error*. \square

§11.12 **Assumption.** Consider weights $\alpha, \mathbf{v} \in \mathcal{J}_{\setminus 0}$, i.e. $\nu(\mathcal{N}_\alpha) = 0 = \nu(\mathcal{N}_\mathbf{v})$, such that $\alpha \in \mathbb{L}_\infty(\nu)$ and $(\alpha\mathbf{v})_\bullet := (\alpha_j \mathbf{v}_j)_{j \in \mathcal{J}} = \alpha \mathbf{v} \in \mathbb{L}_\infty(\nu)$. We write $(\alpha\mathbf{v})_{(m)} := \|(\alpha\mathbf{v})_\bullet \mathbf{1}^{m \perp}\|_{\mathbb{L}_\infty(\nu)} \in \mathbb{R}^+$ for each $m \in \mathbb{N}$. \square

§11.13 **Reminder.** Under Assumption §11.12 we have $\mathbb{J}^\alpha = \mathbb{L}_2(\nu) = \text{dom}(\mathbb{M}_\alpha) = \mathbb{J}\alpha = \mathbb{L}_2(\alpha^{2\ddagger}\nu)$ and the three measures ν , $\alpha^{2\ddagger}\nu$ and $\mathbf{v}^2\nu$ dominate mutually each other, i.e. they share the same null sets (see **Property** §11.05). Consequently, $\mathbb{J}^\alpha \subseteq \mathbb{J} = \mathbb{L}_2(\nu)$ and if $h_\bullet \in \mathbb{L}_2(\alpha^{2\ddagger}\nu)$ satisfies $\mathbf{v}^2\nu(h_\bullet^2) \in \mathbb{R}^+$, for example, then $h_\bullet \in \mathbb{L}_2(\mathbf{v}^2\nu)$ too. \square

§11.14 **Notation.** Under Assumption §11.12 and given a constant $r \in \mathbb{R}_{\setminus 0}^+$ we consider $\mathbb{J}^\alpha = \mathbb{L}_2(\nu) = \mathbb{L}_2(\alpha^{2\ddagger}\nu)$ endowed with $\|\cdot\|_{\alpha^\ddagger} := \|\cdot\|_{\mathbb{J}^\alpha} := \|\cdot\|_{\mathbb{L}_2(\alpha^{2\ddagger}\nu)}$ and the ellipsoid

$$\mathbb{J}^{\alpha,r} := \{h_\bullet \in \mathbb{J}^\alpha : \|h_\bullet\|_{\alpha^\ddagger}^2 = \alpha^{2\ddagger}\nu(h_\bullet^2) = \nu(\alpha^{2\ddagger}h_\bullet^2) \leq r^2\} \subseteq \mathbb{J}^\alpha.$$

Keep in mind that $(\alpha\mathbf{v})_\bullet \in \mathbb{L}_\infty(\nu)$ implies $(\alpha\mathbf{v})_{(m)} := \|(\alpha\mathbf{v})_\bullet \mathbf{1}^{m \perp}\|_{\mathbb{L}_\infty(\nu)} \in \mathbb{R}^+$ for each $m \in \mathbb{N}$. \square

§11.15 **Property.** Under Assumption §11.12 we have $\mathbb{J}^\alpha \subseteq \mathbb{L}_2(\mathbf{v}^2\nu)$. Indeed, for each $h_\bullet \in \mathbb{J}^\alpha$, i.e., $\|h_\bullet\|_{\alpha^\ddagger} \in \mathbb{R}^+$, follows $\|h_\bullet\|_{\mathbf{v}}^2 = \nu(h_\bullet^2 \alpha^{2\ddagger}(\alpha\mathbf{v})_\bullet^2) \leq \|h_\bullet\|_{\alpha^\ddagger}^2 \|(\alpha\mathbf{v})_\bullet\|_{\mathbb{L}_\infty(\nu)}^2 \in \mathbb{R}^+$. \square

§11.16 **Abstract smoothness condition.** Under Assumption §11.12 the parameter $\theta \in \mathbb{J}$ satisfies an *abstract smoothness condition* if there is $r \in \mathbb{R}_{\setminus 0}^+$ such that $\theta \in \mathbb{J}^{\alpha,r} \subseteq \mathbb{J}^\alpha$. \square

§11.17 **Lemma.** Under Assumption §11.12 for each $m \in \mathbb{N}$ the orthogonal projection $\theta_\bullet^m := \theta \mathbf{1}_\bullet^m \in \mathbb{J} \mathbf{1}_\bullet^m$ of $\theta \in \mathbb{J}^{\alpha,r} \subseteq \mathbb{L}_2(\mathbf{v}^2\nu)$ satisfies $\|\theta_\bullet^m - \theta\|_{\mathbf{v}} = \|\theta \mathbf{1}_\bullet^{m \perp}\|_{\mathbf{v}} \leq r (\alpha\mathbf{v})_{(m)}$. \square

§11.18 **Proof of Lemma** §11.17. is given in the lecture. \square

§11.19 **Maximal global \mathbf{v} -error.** Under Assumption §11.12 for $m \in \mathbb{N}$, a parameter $\theta = \mathbb{U}\theta \in \mathbb{J}^{\alpha,r}$ and its orthogonal projection $\theta_\bullet^m = \theta \mathbf{1}_\bullet^m \in \mathbb{J} \mathbf{1}_\bullet^m$ we call $\sup \{\|\theta_\bullet^m - \theta\|_{\mathbf{v}} : \theta \in \mathbb{J}^{\alpha,r}\}$ *maximal global \mathbf{v} -error* over the class of parameters $\mathbb{J}^{\alpha,r}$. \square

§11|04 Local and maximal local ϕ -error

Secondly, we measure locally the accuracy of the orthogonal projection $\theta_\bullet^m := \theta \mathbf{1}_\bullet^m \in \mathbb{J} \mathbf{1}_\bullet^m$ of $\theta = \mathbb{U}\theta \in \mathbb{J}$.

§11.20 **Notation.** For $\phi \in \mathcal{J}$ and $\text{dom}(\phi\nu) := \{h_\bullet \in \mathbb{J} = \mathbb{L}_2(\nu) : \phi h_\bullet \in \mathbb{L}_1(\nu)\}$ we consider the linear functional $\phi\nu : \mathbb{J} \supseteq \text{dom}(\phi\nu) \rightarrow \mathbb{R}$ given by $h_\bullet \mapsto \phi\nu(h_\bullet) := \nu(\phi h_\bullet)$ with a slight abuse of notations. \square

§11.21 **Comment.** If $\phi \in \mathbb{J} = \mathbb{L}_2(\nu)$, then it follows $\text{dom}(\phi\nu) = \mathbb{J}$ and $\|\phi\nu\|_{\mathbb{L}(\mathbb{J},\mathbb{R})} = \|\phi\|_{\mathbb{J}} \in \mathbb{R}^+$. Consequently, we have $\phi\nu \in \mathbb{L}(\mathbb{J},\mathbb{R})$ and $\phi\nu(h_\bullet) = \langle h_\bullet, \phi \rangle_{\mathbb{J}}$, in other words ϕ is a Fréchet-Riesz representative of the continuous linear functional $\phi\nu$. \square

§11.22 **Property.** If $\phi \in \mathcal{J}_{\setminus 0}$ (i.e. $\nu(\mathcal{N}_\phi) = 0$) and $\theta \in \text{dom}(\phi\nu)$ (i.e. $\theta \phi \in \mathbb{L}_1(\nu)$), then for each $m \in \mathbb{N}$ we have $\theta_\bullet^m \in \text{dom}(\phi\nu)$ too, since $\|\phi \theta_\bullet^m\|_{\mathbb{L}_1(\nu)} = \nu(|\phi \theta_\bullet^m| \mathbf{1}_\bullet^m) \leq \nu(|\phi \theta_\bullet^m|)$. Moreover, it holds $|\phi\nu(\theta) - \phi\nu(\theta_\bullet^m)| \leq |\phi| \nu(|\theta_\bullet^m - \theta|) = |\phi| \nu(|\theta| \mathbf{1}_\bullet^{m \perp}) \leq \nu(|\phi \theta|) \in \mathbb{R}^+$ and $|\phi\nu(\theta) - \phi\nu(\theta_\bullet^m)| = o(1)$ as $m \rightarrow \infty$ by dominated convergence. \square

§11.23 **Comment.** We assume throughout this chapter that the Hilbert space $\mathbb{J} = \mathbb{L}_2(\mathcal{J}, \mathcal{J}, \nu)$ and the surjective partial isometry $\mathbb{U} \in \mathbb{L}(\mathbb{H}, \mathbb{J})$ is fixed and known in advance. Considering a ϕ -error means the linear function $\phi\nu$ and hence in equal $\phi \in \mathcal{J}$ is also fixed and known in advance. Consequently, the condition $\phi \in \mathcal{J}_{\setminus 0}$ does not impose an additional restriction. \square

- §11.24 **Local ϕ -error.** Given $\phi \in \mathcal{J}_0$, $m \in \mathbb{N}$, a parameter $\theta = U\theta \in \text{dom}(\phi\nu)$ and its orthogonal projection $\theta^m = \theta \mathbb{1}_*^m \in \mathbb{J}^m$ we call $|\phi\nu(\theta) - \phi\nu(\theta^m)| = |\phi\nu(\theta \mathbb{1}_*^{m\perp})| \in \mathbb{R}^+$ *local ϕ -error*. \square
- §11.25 **Assumption.** Consider $\phi, \alpha \in \mathcal{J}_0$, i.e. $\nu(\mathcal{N}_\phi) = 0 = \nu(\mathcal{N}_\alpha)$, such that $\alpha \in \mathbb{L}_\infty(\nu)$ and $(\alpha\phi)_* := (\alpha_j \phi_j)_{j \in \mathcal{J}} = \alpha_* \phi_* \in \mathbb{L}_2(\nu)$ and hence $\|\alpha_* \mathbb{1}_*^{m\perp}\|_\phi = \|(\alpha\phi)_* \mathbb{1}_*^{m\perp}\|_{\mathbb{L}_2(\nu)} = o(1)$ as $m \rightarrow \infty$. \square
- §11.26 **Reminder.** Under Assumption §11.25 we have $\mathbb{J}^a = \mathbb{L}_2^a(\nu) = \text{dom}(M_\alpha) = \mathbb{J}\alpha_* = \mathbb{L}_2(\alpha^{2\uparrow}\nu)$ and the three measures ν , $|\phi| \nu$ and $\alpha^{2\uparrow}\nu$ dominate mutually each other (see **Property** §11.05). Consequently, $\mathbb{J}^a \subseteq \mathbb{J} = \mathbb{L}_2(\nu)$ and if $h_* \in \mathbb{L}_2(\alpha^{2\uparrow}\nu)$ satisfies $\nu(|\phi h_*|) \in \mathbb{R}^+$, for example, then $h_* \in \mathbb{L}_1(|\phi| \nu)$ too. \square
- §11.27 **Property.** Under Assumption §11.25 we have $\mathbb{J}^a \subseteq \text{dom}(\phi\nu)$. Indeed, for each $h_* \in \mathbb{J}^a$, i.e. $\|h_*\|_{\alpha^{2\uparrow}} \in \mathbb{R}^+$, we have $\|\phi h_*\|_{\mathbb{L}_1(\nu)} = \nu(|h_* \alpha^{2\uparrow}(\alpha\phi)_*|) \leq \|h_*\|_{\alpha^{2\uparrow}} \|(\alpha\phi)_*\|_{\mathbb{L}_2(\nu)} \in \mathbb{R}^+$. \square
- §11.28 **Notation (Reminder).** Under Assumption §11.25 the parameter $\theta = U\theta \in \mathbb{J}$ satisfies an abstract smoothness condition if there is $r \in \mathbb{R}_0^+$ such that $\theta_* \in \mathbb{J}^{a,r} = \{h_* \in \mathbb{J}^a : \|h_*\|_{\alpha^{2\uparrow}}^2 \leq r^2\} \subseteq \mathbb{J}^a$ where $\|\cdot\|_{\alpha^{2\uparrow}} = \|\cdot\|_{\mathbb{J}^a} := \|\cdot\|_{\mathbb{L}_2(\alpha^{2\uparrow}\nu)}$ (see **Definition** §11.16). Since $(\alpha\phi)_* \in \mathbb{L}_2(\nu)$ we have $\|\alpha_* \mathbb{1}_*^{m\perp}\|_\phi = \|(\alpha\phi)_* \mathbb{1}_*^{m\perp}\|_{\mathbb{L}_2(\nu)} = o(1)$ as $m \rightarrow \infty$ by dominated convergence. \square
- §11.29 **Lemma.** Under Assumption §11.25 for each $m \in \mathbb{N}$ the orthogonal projection $\theta^m := \theta \mathbb{1}_*^m \in \mathbb{J}^m$ of $\theta \in \mathbb{J}^{a,r} \subseteq \text{dom}(\phi\nu)$ satisfies $|\phi\nu(\theta - \theta^m)| = |\phi\nu(\theta \mathbb{1}_*^{m\perp})| \leq \nu(|\phi \theta \mathbb{1}_*^{m\perp}|) \leq r \|\alpha_* \mathbb{1}_*^{m\perp}\|_\phi$. \square
- §11.30 **Proof of Lemma** §11.29. is given in the lecture. \square
- §11.31 **Maximal local ϕ -error.** Under Assumption §11.25 for $m \in \mathbb{N}$, a parameter $\theta = U\theta \in \mathbb{J}^{a,r}$ and its orthogonal projection $\theta^m = \theta \mathbb{1}_*^m \in \mathbb{J}^m$ we call $\sup \{|\phi\nu(\theta) - \phi\nu(\theta^m)| : \theta \in \mathbb{J}^{a,r}\}$ *maximal local ϕ -error* over the class of parameters $\mathbb{J}^{a,r}$. \square

§12 Orthogonal projection estimator

- §12.01 **Notation (Reminder).** Consider a measure space $(\mathcal{J}, \mathcal{J}, \nu)$ as in **Notation** §10.01. For $w_* \in \mathbb{R}^{\mathcal{J}}$ define the multiplication map $M_{w_*} : \mathbb{R}^{\mathcal{J}} \rightarrow \mathbb{R}^{\mathcal{J}}$ with $a_* \mapsto M_{w_*} a_* := w_* a_*$. For $w_* \in \mathcal{J}$ we have $M_{w_*} : \mathcal{J} \rightarrow \mathcal{J}$ too. We denote by $\mathbb{M}_{\mathcal{J}}$ the set of all multiplication maps defined on \mathcal{J} . If in addition $w_* \in \mathbb{L}_\infty(\nu)$ then we have also $M_{w_*} : \mathbb{L}_2(\nu) = \mathbb{J} \rightarrow \mathbb{J}$ identifying eventually equivalence classes and representatives. We set $\mathbb{L}^{\mathbb{J}} := \mathbb{M}_{\mathcal{J}} := \{M_{w_*} \in \mathbb{M}_{\mathcal{J}} : w_* \in \mathbb{L}_\infty(\nu)\} \subseteq \mathbb{L}(\mathbb{J})$ noting that $\|M_{w_*}\|_{\mathbb{L}(\mathbb{J})} = \sup \{\|w_* a_*\|_{\mathbb{J}} : \|a_*\|_{\mathbb{J}} \leq 1\} \leq \|w_*\|_{\mathbb{L}_\infty(\nu)}$ for each $M_{w_*} \in \mathbb{L}^{\mathbb{J}}$. \square
- §12.02 **Reminder.** If $w_* \in \mathbb{L}_\infty(\nu)$ then $M_{w_*} \in \mathbb{L}^{\mathbb{J}}$, and $M_{w_*} : \mathbb{J} \supseteq \text{dom}(M_{w_*}) \rightarrow \mathbb{J}$. Moreover, we have $\text{dom}(M_{w_*}) = \mathbb{J}$, $\text{ran}(M_{w_*}) = \mathbb{J}w_*$ and $\ker(M_{w_*}) = \mathbb{J}\mathbb{1}_*^{\mathcal{N}_w}$ (see **Property** §11.03), and $\text{dom}(M_{w_*}) = \mathbb{J}w_* \oplus \mathbb{J}\mathbb{1}_*^{\mathcal{N}_w}$ (see **Property** §11.05). Consequently, if in addition $\nu(\mathcal{N}_w) = 0$, then $\mathbb{J}^w = \mathbb{L}_2^w(\nu) = \text{dom}(M_{w_*}) = \mathbb{J}w_* = \mathbb{L}_2(w^{2\uparrow}\nu)$. For each $m \in \mathbb{N}$ we write $\mathbb{1}_*^m = (\mathbb{1}_j^m)_{j \in \mathcal{J}} := \mathbb{1}_*^{[m]}$ and $\mathbb{1}_*^{m\perp} := \mathbb{1}_* - \mathbb{1}_*^m$ with $[m] := [-m, m] \cap \mathcal{J}$. Consequently, $M_{\mathbb{1}_*^m} \in \mathbb{L}^{\mathbb{J}}$ and $M_{\mathbb{1}_*^{m\perp}} \in \mathbb{L}^{\mathbb{J}}$ is the *orthogonal projection* onto the linear subspace $\mathbb{J}\mathbb{1}_*^m \subseteq \mathbb{J}$ and its orthogonal complement $\mathbb{J}\mathbb{1}_*^{m\perp} = (\mathbb{J}\mathbb{1}_*^m)^\perp \subseteq \mathbb{J}$, respectively, that is $\mathbb{J} = \mathbb{J}\mathbb{1}_*^m \oplus \mathbb{J}\mathbb{1}_*^{m\perp}$ (see **Property** §11.07). Finally, given $\theta = U\theta \in \mathbb{J}$ we consider the orthogonal projections $\theta^m = \theta \mathbb{1}_*^m \in \mathbb{J}\mathbb{1}_*^m$ (and $\theta^{m\perp} := U^* \theta^m \in \mathbb{H}$) (**Definition** §11.08). \square
- §12.03 **Notation (Reminder).** Consider a centred stochastic processes $\varepsilon_* = (\varepsilon_j)_{j \in \mathcal{J}}$ satisfying Assumption §10.04 and let $n \in \mathbb{N}$ be a sample size. The observable noisy version $\hat{\theta} = \theta + n^{-1/2} \varepsilon_*$ of the parameter $\theta = U\theta \in \mathbb{J}$ takes the form of a *statistical direct problem* (see **Definition** §10.19). We

denote by \mathbb{P}_q^n the distribution of $\widehat{\theta}$. We write $\varepsilon \sim \mathbb{P}_{(q,\Gamma)}$ if ε admits a covariance operator $\Gamma \in \mathbb{L}^{\geq}(\mathbb{J})$ possibly depending on θ . \square

§12.04 **Definition.** Given a noisy version $\widehat{\theta} \sim \mathbb{P}_q^n$ of the parameter $\theta = U\theta \in \mathbb{J}$ for each $m \in \mathbb{N}$ we call $\widehat{\theta}^m := \widehat{\theta} \mathbb{1}_m^m$ *orthogonal projection estimator (OPE)* of θ . \square

§12.05 **GSSM (§10.21 continued).** Considering $\ell_2 = \mathbb{L}_2(\mathbb{N}, 2^{\mathbb{N}}, \nu_n)$ we illustrate the OPE in a Gaussian sequence space model §10.21. Here the observable stochastic process $\widehat{\theta} = \theta + n^{-1/2} \dot{B}$ is a noisy version of $\theta = U\theta \in \ell_2$ and $\dot{B} \sim N_{(0,1)}^{\otimes \mathbb{N}}$. Consequently, $\widehat{\theta}$ admits a N_q^n -distribution belonging to the family $N_{\Theta}^n := (N_q^n)_{q \in \Theta}$. Summarising the observations satisfy a statistical product experiment $(\mathbb{R}^{\mathbb{N}}, \mathcal{B}^{\otimes \mathbb{N}}, N_{\Theta}^n)$ where $\Theta \subseteq \ell_2$. \square

§12|01 Global and maximal global v-risk

We measure first the accuracy of the OPE $\widehat{\theta}^m = \widehat{\theta} \mathbb{1}_m^m$ of $\theta^m = \theta \mathbb{1}_m^m \in \mathbb{J} \mathbb{1}_m^m$ with $\theta = U\theta \in \mathbb{J}$ by a global mean-v-error, i.e. v-risk.

§12.06 **Reminder.** If $v \in \mathcal{J}_0$ and $\theta \in \mathbb{L}_2(v^2\nu)$ then we have $\theta^m \in \mathbb{L}_2(v^2\nu)$ too and $\|\theta^m - \theta\|_v^2 = o(1)$ as $m \rightarrow \infty$ (**Property** §11.09). \square

§12.07 **Assumption.** Consider a noisy version $\widehat{\theta} = \theta + n^{-1/2} \varepsilon \sim \mathbb{P}_{\theta}^n$ of $\theta = U\theta \in \mathbb{J}$ satisfying Assumption §10.04, $v^{\theta} := \mathbb{P}_{\theta}^n(\varepsilon^2) := (\mathbb{P}_{\theta}^n(\varepsilon_j^2))_{j \in \mathcal{J}} \in \mathbb{L}_{\infty}(\nu)$ and $\varepsilon \mathbb{1}_m^m \in \mathbb{L}_{\infty}(\nu)$ \mathbb{P}_{θ}^n -a.s., for each $m \in \mathbb{N}$. \square

§12.08 **Comment.** Under Assumption §12.07 if $v \mathbb{1}_m^m \in \mathbb{L}_2(\nu)$ then we have $v \varepsilon \mathbb{1}_m^m \in \mathbb{L}_2(\nu)$ \mathbb{P}_{θ}^n -a.s.. If in addition $\theta \in \mathbb{L}_2(v^2\nu)$, and hence $\theta^m \in \mathbb{L}_2(v^2\nu)$ (**Property** §11.09), then it follows

$$v \widehat{\theta}^m = n^{-1/2} v \varepsilon \mathbb{1}_m^m + v \theta^m \in \mathbb{L}_2(\nu) \quad \mathbb{P}_{\theta}^n \text{-a.s.} \quad (12.01)$$

If $\mathcal{J} \subseteq \mathbb{Z}$ (at most countable) then Assumption §10.04 and $v^{\theta} = \mathbb{P}_{\theta}^n(\varepsilon^2) \in \mathbb{L}_{\infty}(\nu)$ implies the additional assumption $\varepsilon \mathbb{1}_m^m \in \mathbb{L}_{\infty}(\nu)$ \mathbb{P}_{θ}^n -a.s.. However, the last implication does generally not hold, if $\mathcal{J} \in \{\mathbb{R}, \mathbb{R}^+\}$ for example. \square

§12|01|01 Global v-risk

§12.09 **Definition.** Under Assumption §12.07, $v \in \mathcal{J}_0$, $\theta \in \mathbb{L}_2(v^2\nu)$ and $v \mathbb{1}_m^m \in \mathbb{J}$ for $m \in \mathbb{N}$ the *global v-risk* of an OPE $\widehat{\theta}^m = \widehat{\theta} \mathbb{1}_m^m \in \mathbb{L}_2(v^2\nu)$ \mathbb{P}_{θ}^n -a.s. satisfies

$$\mathbb{E}_{\theta}^n(\|\widehat{\theta}^m - \theta\|_v^2) = \mathbb{E}_{\theta}^n(\|(\widehat{\theta} - \theta) \mathbb{1}_m^m\|_v^2) + \|\theta \mathbb{1}_m^m\|_v^2 \quad (12.02)$$

with *variance* term $\mathbb{E}_{\theta}^n(\|(\widehat{\theta} - \theta) \mathbb{1}_m^m\|_v^2) = n^{-1} \mathbb{E}_{\theta}^n(\|v \varepsilon \mathbb{1}_m^m\|_j^2)$ and *bias* term $\|\theta \mathbb{1}_m^m\|_v^2$. \square

§12.10 **Property.** Under Assumption §12.07, $v \in \mathcal{J}_0$ and $\mathbb{1}_m^m \in \mathbb{L}_2(v^2\nu)$ for $m \in \mathbb{N}$ we have

$$\mathbb{E}_{\theta}^n(\|v \varepsilon \mathbb{1}_m^m\|_j^2) = \int_{\mathcal{J}} \mathbb{E}_{\theta}^n(\varepsilon_j^2) v_j^2 \mathbb{1}_j^m \nu(dj) = \nu(v^{\theta} v^2 \mathbb{1}_m^m) \quad (12.03)$$

and consequently $\mathbb{E}_{\theta}^n(\|(\widehat{\theta} - \theta) \mathbb{1}_m^m\|_v^2) \leq n^{-1} \|v^{\theta}\|_{\mathbb{L}_{\infty}(\nu)} \|\mathbb{1}_m^m\|_v^2 \in \mathbb{R}^+$. \square

§12.11 **Notation.** For $a \in (\mathbb{R})^{\mathbb{N}}$ with minimal value in $B \subseteq \mathbb{N}$ we define

$$\arg \min \{a_m : m \in B\} := \min \{m \in B : a_m \leq a_j, \forall j \in B\}. \quad \square$$

§12.12 **Proposition (Upper bound).** *Let Assumption §12.07, $\mathbf{v} \in \mathcal{J}_{\mathbf{v}_0}$, $\theta \in \mathbb{L}_2(\mathbf{v}^2\nu)$ and $\mathbf{1}^m \in \mathbb{L}_2(\mathbf{v}^2\nu)$ for all $m \in \mathbb{N}$ be satisfied. For all $n, m \in \mathbb{N}$ setting*

$$\begin{aligned} R_n^m(\theta, \mathbf{v}) &:= \|\theta \mathbf{1}^{m\perp}\|_{\mathbf{v}}^2 + n^{-1} \|\mathbf{1}^m\|_{\mathbf{v}}^2, \quad m_n^\circ := \arg \min \{R_n^m(\theta, \mathbf{v}) : m \in \mathbb{N}\} \\ \text{and } R_n^\circ(\theta, \mathbf{v}) &:= R_n^{m_n^\circ}(\theta, \mathbf{v}) = \min \{R_n^m(\theta, \mathbf{v}) : m \in \mathbb{N}\} \end{aligned} \quad (12.04)$$

we have $\mathbb{E}_\theta^n(\|\widehat{\theta}^{m_n^\circ} - \theta\|_{\mathbf{v}}^2) \leq (1 \vee \|\mathbf{v}^\theta\|_{\mathbb{L}_\infty(\nu)}) R_n^\circ(\theta, \mathbf{v})$.

§12.13 **Proof of Proposition §12.12.** is given in the lecture. \square

§12.14 **Definition.** Let $\theta \in \mathbb{L}_2(\mathbf{v}^2\nu)$ and $\widehat{\theta}^m \in \mathbb{L}_2(\mathbf{v}^2\nu)$ \mathbb{P}_θ^n -a.s. for all $m \in \mathbb{N}$. If there exist $C \in \mathbb{R}_{\mathbf{v}_0}^+$ and for each $n \in \mathbb{N}$, $R_n^\circ \in \mathbb{R}_{\mathbf{v}_0}^+$ and $m_n^\circ \in \mathbb{N}$ satisfying

$$C^{-1} R_n^\circ \leq \inf_{m \in \mathbb{N}} \mathbb{E}_\theta^n \|\widehat{\theta}^m - \theta\|_{\mathbf{v}}^2 \leq \mathbb{E}_\theta^n \|\widehat{\theta}^{m_n^\circ} - \theta\|_{\mathbf{v}}^2 \leq C R_n^\circ,$$

then we call R_n° *oracle bound*, m_n° *oracle dimension* and $\widehat{\theta}^{m_n^\circ}$ *oracle optimal* (up to the constant C). As a consequence, up to the constant C^2 the statistic $\widehat{\theta}^{m_n^\circ}$ attains the lower global \mathbf{v} -risk bound within the family of OPE's, that is, $\mathbb{E}_\theta^n \|\widehat{\theta}^{m_n^\circ} - \theta\|_{\mathbf{v}}^2 \leq C^2 \inf_{m \in \mathbb{N}} \mathbb{E}_\theta^n \|\widehat{\theta}^m - \theta\|_{\mathbf{v}}^2$. \square

§12.15 **Oracle inequality.** *Under Assumption §12.07 let $\mathbf{v} \in \mathcal{J}_{\mathbf{v}_0}$, $\theta \in \mathbb{L}_2(\mathbf{v}^2\nu)$ and $\mathbf{1}^m \in \mathbb{L}_2(\mathbf{v}^2\nu)$ for all $m \in \mathbb{N}$. If in addition $1 \leq \max(\|\mathbf{v}^\theta\|_{\mathbb{L}_\infty(\nu)}, \|(\mathbf{v}^\theta)^{-1}\|_{\mathbb{L}_\infty(\nu)}) \leq \mathbf{v}_\theta \in \mathbb{R}_{\mathbf{v}_0}^+$ then (12.04) implies*

$$\begin{aligned} \mathbf{v}_\theta^{-1} R_n^m(\theta, \mathbf{v}) &\leq \mathbb{E}_\theta^n(\|\widehat{\theta}^m - \theta\|_{\mathbf{v}}^2) = n^{-1} \nu(\mathbf{v}^\theta \mathbf{v}^2 \mathbf{1}^m) + \|\theta \mathbf{1}^{m\perp}\|_{\mathbf{v}}^2 \\ &\leq \mathbf{v}_\theta R_n^m(\theta, \mathbf{v}) \quad \text{for all } m, n \in \mathbb{N}. \end{aligned}$$

As a consequence we immediately obtain the following *oracle inequality*

$$\begin{aligned} \mathbf{v}_\theta^{-1} R_n^\circ(\theta, \mathbf{v}) &\leq \inf_{m \in \mathbb{N}} \mathbb{E}_\theta^n(\|\widehat{\theta}^m - \theta\|_{\mathbf{v}}^2) \leq \mathbb{E}_\theta^n(\|\widehat{\theta}^{m_n^\circ} - \theta\|_{\mathbf{v}}^2) \\ &\leq \mathbf{v}_\theta R_n^\circ(\theta, \mathbf{v}) \leq \mathbf{v}_\theta^2 \inf_{m \in \mathbb{N}} \mathbb{E}_\theta^n(\|\widehat{\theta}^m - \theta\|_{\mathbf{v}}^2), \end{aligned} \quad (12.05)$$

and, hence $R_n^\circ(\theta, \mathbf{v})$, m_n° and the statistic $\widehat{\theta}^{m_n^\circ}$, respectively, is an *oracle bound*, an *oracle dimension* and *oracle optimal* (up to the constant \mathbf{v}_θ^2). \square

§12.16 **Remark.** We shall emphasise that for each fixed $m \in \mathbb{N}$ with $\|\mathbf{1}^m\|_{\mathbf{v}} \in \mathbb{R}^+$ we have $n^{-1} \|\mathbf{1}^m\|_{\mathbf{v}} = o(1)$ as $n \rightarrow \infty$. As a consequence, if $\|\mathbf{1}^m\|_{\mathbf{v}} \in \mathbb{R}^+$ for all $m \in \mathbb{N}$ and $\|\theta \mathbf{1}^{m\perp}\|_{\mathbf{v}} = o(1)$ as $m \rightarrow \infty$ then we obtain $R_n^\circ(\theta, \mathbf{v}) = o(1)$ as $n \rightarrow \infty$, and thus, $R_n^\circ(\theta, \mathbf{v})$ is also called an *oracle rate*. Indeed, for all $\delta \in \mathbb{R}_{\mathbf{v}_0}^+$ there exists $m_\delta \in \mathbb{N}$ and $n_\delta \in \mathbb{N}$ such that we have both $\|\theta \mathbf{1}^{m_\delta\perp}\|_{\mathbf{v}}^2 \leq \delta/2$ and $n^{-1} \|\mathbf{1}^{m_\delta}\|_{\mathbf{v}}^2 \leq \delta/2$ for all $n \geq n_\delta$, and whence $R_n^\circ(\theta, \mathbf{v}) \leq R_n^{m_\delta}(\theta, \mathbf{v}) \leq \delta$. However, note that the oracle dimension $m_n^\circ = m_n^\circ(\theta, \mathbf{v})$ as defined in **Proposition §12.12** depends on the unknown parameter of interest θ , and thus also the oracle optimal statistic $\widehat{\theta}^{m_n^\circ}$. In other words $\widehat{\theta}^{m_n^\circ}$ is not a feasible estimator. \square

§12.17 **Corollary (GSSM §12.05 continued).** *Let $\widehat{\theta} = \theta + n^{-1/2} \dot{B} \sim N_q^n$ as in Model §12.05, where $\dot{B} \sim N_{(0,1)}^{\otimes \mathbb{N}}$ and $\theta = U\theta \in \ell_2$. For $\mathbf{v} \in (\mathbb{R}_{\mathbf{v}_0})^{\mathbb{N}}$ and $\theta \in \ell_2(\mathbf{v}^2)$ the (infeasible) OPE $\widehat{\theta}^{m_n^\circ} = \widehat{\theta} \mathbf{1}^{m_n^\circ} \in \ell_2 \mathbf{1}^{m_n^\circ} \subseteq \ell_2(\mathbf{v}^2)$ with oracle dimension m_n° as in (12.04) satisfies*

$$N_q^n(\|\widehat{\theta}^{m_n^\circ} - \theta\|_{\mathbf{v}}^2) = R_n^\circ(\theta, \mathbf{v}) = \inf_{m \in \mathbb{N}} N_q^n(\|\widehat{\theta}^m - \theta\|_{\mathbf{v}}^2),$$

and hence it is *oracle optimal* (with constant 1).

§12.18 **Proof** of **Corollary** §12.17. is given in the lecture. \square

§12.19 **Illustration**. Here and subsequently, we use for two sequences $a_n, b_n \in (\mathbb{R}_0^+)^{\mathbb{N}}$ the notation $a_n \simeq b_n$ if the sequence a_n/b_n is bounded away both from zero and infinity. We illustrate the last results considering usual behaviour for the bias and variance term. We distinguish the following two cases

(p) $\mathfrak{v} \in \mathbb{J}$ or there is $m \in \mathbb{N}$ with $\|\theta^m - \theta\|_{\mathfrak{v}}^2 = 0$,

(np) $\mathfrak{v} \notin \mathbb{J}$ and for all $m \in \mathbb{N}$ holds $\|\theta^m - \theta\|_{\mathfrak{v}}^2 \in \mathbb{R}_0^+$.

Interestingly, in case (p) the oracle bound is parametric, that is, $nR_n^\circ(\theta, \mathfrak{v}) = O(1)$, in case (np) the oracle bound is nonparametric, i.e. $\lim_{n \rightarrow \infty} nR_n^\circ(\theta, \mathfrak{v}) = \infty$. In case (np) consider the following two specifications:

Table 01 [§12]

Order of the oracle rate $R_n^\circ(\theta, \mathfrak{v})$ as $n \rightarrow \infty$

$(j \in \mathbb{N})$	$(a \in \mathbb{R}_0^+)$	(squared bias)	(variance)	m_n°	$R_n^\circ(\theta, \mathfrak{v})$
$\mathfrak{v}_j^2 = j^{2v}$	θ_j^2	$\ \theta \mathbf{1}_j^{m \perp}\ _{\mathfrak{v}}^2$	$\ \mathbf{1}_j^m\ _{\mathfrak{v}}^2$		
(o) $v \in (-1/2, a)$	j^{-2a-1}	$m^{-2(a-v)}$	m^{2v+1}	$n^{\frac{1}{2a+1}}$	$n^{-\frac{2(a-v)}{2a+1}}$
$v = -1/2$	j^{-2a-1}	m^{-2a-1}	$\log m$	$(\frac{n}{\log n})^{\frac{1}{2a+1}}$	$\frac{\log n}{n}$
(s) $v + 1/2 \in \mathbb{R}_0^+$	$e^{-j^{2a}}$	$m^{(1-2(a-v))+} e^{-m^{2a}}$	m^{2v+1}	$(\log n)^{\frac{1}{2a}}$	$\frac{(\log n)^{\frac{2v+1}{2a}}}{n}$
$v = -1/2$	$e^{-j^{2a}}$	$e^{-m^{2a}}$	$\log m$	$(\log n)^{\frac{1}{2a}}$	$\frac{\log \log n}{n}$

We note that in Table 01 [§12] the order of the oracle rate $R_n^\circ(\theta, \mathfrak{v})$ is depict for $v \geq -1/2$ only. In case $v < -1/2$ the oracle rate $R_n^\circ(\theta, \mathfrak{v})$ is parametric. \square

§12|01|02 Maximal global \mathfrak{v} -risk

§12.20 **Reminder**. Under Assumption §11.12 we have $\mathbb{J}^a = \mathbb{L}_2(\nu) = \text{dom}(M_{\mathfrak{a}}) = \mathbb{J}^a \subseteq \mathbb{J}$ and the three measures ν , $\mathfrak{a}^{2\ddagger}\nu$ and $\mathfrak{v}^2\nu$ dominate mutually each other, i.e. they share the same null sets (see **Property** §11.05). We consider \mathbb{J}^a endowed with $\|\cdot\|_{\mathfrak{a}^\ddagger} = \|M_{\mathfrak{a}} \cdot\|_{\mathbb{J}}$ and given a constant $r \in \mathbb{R}_0^+$ the ellipsoid $\mathbb{J}^{\mathfrak{a},r} := \{h \in \mathbb{J}^a : \|h\|_{\mathfrak{a}^\ddagger} \leq r\} \subseteq \mathbb{J}^a$. Since $(\mathfrak{a}\mathfrak{v}) \in \mathbb{L}_\infty(\nu)$, and hence $(\mathfrak{a}\mathfrak{v})_{(m)} := \|(\mathfrak{a}\mathfrak{v}) \mathbf{1}_m^{\perp}\|_{\mathbb{L}_\infty(\nu)} \in \mathbb{R}^+$ for each $m \in \mathbb{N}$ we have $\mathbb{J}^a \subseteq \mathbb{L}_2(\mathfrak{v}^2\nu)$ (**Property** §11.15), and $\|\theta \mathbf{1}_m^{\perp}\|_{\mathfrak{v}} \leq r (\mathfrak{a}\mathfrak{v})_{(m)}$ for all $\theta \in \mathbb{J}^{\mathfrak{a},r}$ (**Lemma** §11.17). \square

§12.21 **Proposition**. Let the Assumptions §12.07 and §11.12 and $\mathbf{1}_m^* \in \mathbb{L}_2(\mathfrak{v}^2\nu)$ for all $m \in \mathbb{N}$ be satisfied. For all $n, m \in \mathbb{N}$ setting

$$R_n^m(\mathfrak{a}, \mathfrak{v}) := [(\mathfrak{a}\mathfrak{v})_{(m)}^2 \vee n^{-1} \|\mathbf{1}_m^*\|_{\mathfrak{v}}^2], \quad m_n^* := \arg \min \{R_n^m(\mathfrak{a}, \mathfrak{v}) : m \in \mathbb{N}\}$$

$$\text{and } R_n^*(\mathfrak{a}, \mathfrak{v}) := R_n^{m_n^*}(\mathfrak{a}, \mathfrak{v}) = \min \{R_n^m(\mathfrak{a}, \mathfrak{v}) : m \in \mathbb{N}\} \quad (12.06)$$

we have $\mathbb{E}_\theta^n (\|\widehat{\theta}^{m_n^*} - \theta\|_{\mathfrak{v}}^2) \leq (\|\mathfrak{v}^\theta\|_{\mathbb{L}_\infty(\nu)} + r^2) R_n^*(\mathfrak{a}, \mathfrak{v})$ for all $\theta = U\theta \in \mathbb{J}^{\mathfrak{a},r}$ and $n \in \mathbb{N}$.

§12.22 **Proof** of **Proposition** §12.21. is given in the lecture. \square

§12.23 **Remark**. Under the assumptions of **Proposition** §12.21 if there exists in addition $\mathfrak{v} \in \mathbb{R}^+$ satisfying $\|\mathfrak{v}^\theta\|_{\mathbb{L}_\infty(\nu)} \leq \mathfrak{v}$ for all $\theta \in \mathbb{J}^{\mathfrak{a},r}$ then

$$\sup \{ \mathbb{E}_\theta^n (\|\widehat{\theta}^{m_n^*} - \theta\|_{\mathfrak{v}}^2) : \theta \in \mathbb{J}^{\mathfrak{a},r} \} \leq (\mathfrak{v} + r^2) R_n^*(\mathfrak{a}, \mathfrak{v}) \quad \text{for all } n \in \mathbb{N}.$$

Arguing similarly as in Remark §12.16 we note that $R_n^*(\alpha, \nu) = o(1)$ as $n \rightarrow \infty$, whenever $\|\mathbb{1}^m\|_{\nu} \in \mathbb{R}^+$ for all $m \in \mathbb{N}$ and $(\alpha\nu)_{(m)} = o(1)$ as $m \rightarrow \infty$. The latter is satisfied, for example, if $(\alpha\nu)_{\cdot} = \alpha \cdot \nu \in \mathbb{J}$ (in equal $\alpha \in \mathbb{L}_2(\nu^2\nu)$). Note that the dimension $m_n^* := m_n^*(\alpha, \nu)$ as defined in (12.06) does not depend on the unknown parameter of interest θ but on the class $\mathbb{J}^{\alpha, \nu}$ only, and thus also the statistic $\widehat{\theta}^{m_n^*}$. In other words, if the regularity of θ is known in advance, then the OPE $\widehat{\theta}^{m_n^*}$ is a feasible estimator. \square

§12.24 **Corollary** (GSSM §12.05 continued). Let $\widehat{\theta} = \theta + n^{-1/2}\dot{B} \sim N_q^n$ as in Model §12.05, where $\dot{B} \sim N_{(0,1)}^{\otimes \mathbb{N}}$ and $\theta = U\theta \in \ell_2$. Under Assumption §11.12 the OPE $\widehat{\theta}^{m_n^*} = \widehat{\theta} \cdot \mathbb{1}^{m_n^*} \in \ell_2 \cdot \mathbb{1}^{m_n^*} \subseteq \ell_2(\nu^2)$ with dimension m_n^* as in (12.06) satisfies

$$\sup \{N_q^n(\|\widehat{\theta}^{m_n^*} - \theta\|_{\nu}^2) : \theta \in \ell_2^{\alpha, \nu}\} \leq C R_n^*(\alpha, \nu) \quad \text{for all } n \in \mathbb{N} \tag{12.07}$$

with constant $C = 1 + \nu^2$.

§12.25 **Proof of Corollary** §12.24. is given in the lecture. \square

§12.26 **Illustration.** We illustrate the last results considering usual behaviour for $(\alpha\nu)_{\cdot}, \nu \in \mathcal{J}$. We distinguish the following two cases **(p)** $\nu \in \mathbb{J}$, and **(np)** $\nu \notin \mathbb{J}$. Interestingly, in case **(p)** the bound in Proposition §12.21 is parametric, that is, $nR_n^*(\alpha, \nu) = O(1)$, in case **(np)** the bound is nonparametric, i.e. $\lim_{n \rightarrow \infty} nR_n^*(\alpha, \nu) = \infty$. In case **(np)** consider the following two specifications:

Table 02 [§12]

Order of the rate $R_n^*(\alpha, \nu)$ as $n \rightarrow \infty$					
$(j \in \mathbb{N})$	$(\alpha \in \mathbb{R}_0^+)$	(squared bias)	(variance)		
$\nu_j = j^{\nu}$	α_j^2	$(\alpha\nu)_{(m)}^2$	$\ \mathbb{1}^m\ _{\nu}^2$	m_n^*	$R_n^*(\alpha, \nu)$
(o) $\nu \in (-1/2, a)$	j^{-2a}	$m^{-2(a-\nu)}$	$m^{2\nu+1}$	$n^{\frac{1}{2a+1}}$	$n^{-\frac{2(a-\nu)}{2a+1}}$
$\nu = -1/2$	j^{-2a}	m^{-2a-1}	$\log m$	$(\frac{n}{\log n})^{\frac{1}{2a+1}}$	$\frac{\log n}{n}$
(s) $\nu + 1/2 \in \mathbb{R}_0^+$	$e^{-j^{2a}}$	$m^{2\nu}e^{-m^{2a}}$	$m^{2\nu+1}$	$(\log n)^{\frac{1}{2a}}$	$n^{-1}(\log n)^{\frac{2\nu+1}{2a}}$
$\nu = -1/2$	$e^{-j^{2a}}$	$e^{-m^{2a}}$	$\log m$	$(\log n)^{\frac{1}{2a}}$	$\frac{\log \log n}{n}$

We note that in Table 02 [§12] the order of the rate $R_n^*(\alpha, \nu)$ is depict for $\geq -1/2$ only. In case $\nu < -1/2$ the rate $R_n^*(\alpha, \nu)$ is parametric. \square

§12|02 Local and maximal local ϕ -risk

We measure secondly the accuracy of the OPE $\widehat{\theta}^m = \widehat{\theta} \cdot \mathbb{1}^m$ of $\theta^m = \theta \cdot \mathbb{1}^m \in \mathbb{J} \cdot \mathbb{1}^m$ with $\theta = U\theta \in \mathbb{J}$ by a local mean- ϕ -error, i.e. ϕ -risk.

§12.27 **Reminder.** If $\phi \in \mathcal{J}_{\nu}$ and $\theta \in \text{dom}(\phi\nu) := \{h \in \mathbb{J} = \mathbb{L}_2(\nu) : \phi h \in \mathbb{L}_1(\nu)\}$ then we have $|\phi\nu(\theta) - \phi\nu(\theta^m)| = o(1)$ as $m \rightarrow \infty$ (Property §11.22). \square

§12.28 **Assumption.** Consider a noisy version $\widehat{\theta} = \theta + n^{-1/2}\varepsilon \sim \mathbb{P}_{\theta}^n$ of $\theta = U\theta \in \mathbb{J}$ satisfying Assumption §10.04, $\varepsilon \sim \mathbb{P}_{(0,1)}$ with $\Gamma_{\theta} \in \mathbb{L}(\mathbb{J})$ and $\varepsilon \cdot \mathbb{1}^m \in \mathbb{L}_2(\nu)$ \mathbb{P}_{θ}^n -a.s. for each $m \in \mathbb{N}$. \square

§12.29 **Comment.** Under Assumption §12.28 if $\mathbb{1}^m \in \mathbb{L}_2(\phi^2\nu)$ then \mathbb{P}_{θ}^n -a.s. we have $|\nu(|\phi\varepsilon \cdot \mathbb{1}^m|)|^2 \leq \nu(\phi^2 \mathbb{1}^m)\nu(\varepsilon^2 \mathbb{1}^m) \in \mathbb{R}^+$ and hence $\varepsilon \cdot \mathbb{1}^m \in \text{dom}(\phi\nu)$. If in addition $\theta \in \text{dom}(\phi\nu)$, and hence $\theta^m \in \text{dom}(\phi\nu)$ (Property §11.22), then it follows

$$n^{-1/2}\varepsilon \cdot \mathbb{1}^m + \theta^m = \widehat{\theta}^m \in \text{dom}(\phi\nu) \quad \mathbb{P}_{\theta}^n\text{-a.s.} \tag{12.08}$$

If $\mathcal{J} \subseteq \mathbb{Z}$ (at most countable) then Assumption §10.04 and $\Gamma_\theta \in \mathbb{L}(\mathbb{J})$ implies $\mathbb{V}_\theta^\theta = \mathbb{P}_\theta^n(\varepsilon_\cdot^2) \in \mathbb{L}_\infty(\nu)$ and hence the additional assumption $\varepsilon_\cdot \mathbb{1}_\cdot^m \in \mathbb{L}_2(\nu)$ \mathbb{P}_θ^n -a.s.. However, the last implication does generally not hold, if $\mathcal{J} \in \{\mathbb{R}, \mathbb{R}^+\}$ for example. \square

§12|02|01 Local ϕ -risk

§12.30 **Definition.** Under Assumption §12.28, $\phi \in \mathcal{J}_0$, $\theta \in \text{dom}(\phi\nu)$ and $\mathbb{1}_\cdot^m \in \mathbb{L}_2(\phi^2\nu)$ for $m \in \mathbb{N}$ the *local ϕ -risk* of an OPE $\widehat{\theta}_\cdot^m = \widehat{\theta}_\cdot \mathbb{1}_\cdot^m \in \text{dom}(\phi\nu)$ \mathbb{P}_θ^n -a.s. satisfies

$$\mathbb{P}_\theta^n(|\phi\nu(\widehat{\theta}_\cdot^m - \theta)|^2) = \mathbb{P}_\theta^n(|\phi\nu((\widehat{\theta}_\cdot - \theta)\mathbb{1}_\cdot^m)|^2) + |\phi\nu(\theta\mathbb{1}_\cdot^{m\perp})|^2. \quad (12.09)$$

with *variance* $\mathbb{P}_\theta^n(|\phi\nu((\widehat{\theta}_\cdot - \theta)\mathbb{1}_\cdot^m)|^2) = n^{-1}\mathbb{P}_{(0,\mathbb{R})}(|\phi\nu(\varepsilon_\cdot \mathbb{1}_\cdot^m)|^2)$ and *bias* $|\phi\nu(\theta\mathbb{1}_\cdot^{m\perp})|$. \square

§12.31 **Property.** Under Assumption §12.28, $\phi \in \mathcal{J}_0$ and $\mathbb{1}_\cdot^m \in \mathbb{L}_2(\phi^2\nu)$ for $m \in \mathbb{N}$ we have

$$\mathbb{P}_{(0,\mathbb{R})}(|\phi\nu(\varepsilon_\cdot \mathbb{1}_\cdot^m)|^2) = \langle \Gamma_\theta(\phi\mathbb{1}_\cdot^m), \phi\mathbb{1}_\cdot^m \rangle_{\mathbb{J}} =: \|\phi\mathbb{1}_\cdot^m\|_{\Gamma_\theta}^2 \leq \|\Gamma_\theta\|_{\mathbb{L}(\mathbb{J})} \|\mathbb{1}_\cdot^m\|_{\phi}^2 \quad (12.10)$$

and consequently $\mathbb{P}_\theta^n(|\nu(\phi(\widehat{\theta}_\cdot - \theta)\mathbb{1}_\cdot^m)|^2) \leq n^{-1}\|\Gamma_\theta\|_{\mathbb{L}(\mathbb{J})} \|\mathbb{1}_\cdot^m\|_{\phi}^2 \in \mathbb{R}^+$. \square

§12.32 **Proposition (Upper bound).** Let Assumption §12.28, $\phi \in \mathcal{J}_0$, $\theta \in \text{dom}(\phi\nu)$ and $\mathbb{1}_\cdot^m \in \mathbb{L}_2(\phi^2\nu)$ for all $m \in \mathbb{N}$ be satisfied. For all $m, n \in \mathbb{N}$ setting

$$\begin{aligned} R_n^m(\theta, \phi) &:= |\phi\nu(\theta\mathbb{1}_\cdot^{m\perp})|^2 + n^{-1}\|\mathbb{1}_\cdot^m\|_{\phi}^2, \quad m_n^\circ := \arg \min \{R_n^m(\theta, \phi) : m \in \mathbb{N}\} \\ \text{and } R_n^\circ(\theta, \phi) &:= R_n^{m_n^\circ}(\theta, \phi) := \min \{R_n^m(\theta, \phi) : m \in \mathbb{N}\} \end{aligned} \quad (12.11)$$

we have $\mathbb{P}_\theta^n(|\phi\nu(\widehat{\theta}_\cdot^{m_n^\circ} - \theta)|^2) \leq (1 \vee \|\Gamma_\theta\|_{\mathbb{L}(\mathbb{J})})R_n^\circ(\theta, \phi)$.

§12.33 **Proof of Proposition §12.32.** is given in the lecture. \square

§12.34 **Definition.** Let $\theta \in \text{dom}(\phi\nu)$ and $\widehat{\theta}_\cdot^m \in \text{dom}(\phi\nu)$ \mathbb{P}_θ^n -a.s. for all $m \in \mathbb{N}$. If there exist $C \in \mathbb{R}_0^+$ and for each $n \in \mathbb{N}$, $R_n^\circ \in \mathbb{R}_0^+$ and $m_n^\circ \in \mathbb{N}$ satisfying

$$C^{-1}R_n^\circ \leq \inf_{m \in \mathbb{N}} \mathbb{P}_\theta^n(|\phi\nu(\widehat{\theta}_\cdot^m - \theta)|^2) \leq \mathbb{P}_\theta^n(|\phi\nu(\widehat{\theta}_\cdot^{m_n^\circ} - \theta)|^2) \leq CR_n^\circ,$$

then we call R_n° *oracle bound*, m_n° *oracle dimension* and $\widehat{\theta}_\cdot^{m_n^\circ}$ *oracle optimal* (up to the constant C). As a consequence, up to the constant C^2 the statistik $\widehat{\theta}_\cdot^{m_n^\circ}$ attains the lower local ϕ -risk bound within the family of OPE's, that is, $\mathbb{P}_\theta^n(|\phi\nu(\widehat{\theta}_\cdot^{m_n^\circ} - \theta)|^2) \leq C^2 \inf_{m \in \mathbb{N}} \mathbb{P}_\theta^n(|\phi\nu(\widehat{\theta}_\cdot^m - \theta)|^2)$. \square

§12.35 **Comment.** If $\Gamma_\theta \in \mathbb{L}(\mathbb{J})$ is invertible with inverse $\Gamma_\theta^{-1} \in \mathbb{L}(\mathbb{J})$, i.e. $\Gamma_\theta \Gamma_\theta^{-1} = \text{id}_{\mathbb{J}} = \Gamma_\theta^{-1} \Gamma_\theta$, then we write shortly $\mathbb{v}_\theta := \max(\|\Gamma_\theta\|_{\mathbb{L}(\mathbb{J})}, \|\Gamma_\theta^{-1}\|_{\mathbb{L}(\mathbb{J})}) \in \mathbb{R}_0^+$. In this situation for all $a_\cdot \in \mathbb{J}$ we have $\mathbb{v}_\theta^{-1} \|a_\cdot\|_{\mathbb{J}}^2 \leq \|a_\cdot\|_{\Gamma_\theta}^2 = \langle \Gamma_\theta a_\cdot, a_\cdot \rangle_{\mathbb{J}} \leq \mathbb{v}_\theta \|a_\cdot\|_{\mathbb{J}}^2$. \square

§12.36 **Oracle inequality.** Under Assumption §12.28 let $\phi \in \mathcal{J}_0$, $\theta \in \text{dom}(\phi\nu)$ and $\mathbb{1}_\cdot^m \in \mathbb{L}_2(\phi^2\nu)$ for all $m \in \mathbb{N}$. If in addition $1 \leq \max(\|\Gamma_\theta\|_{\mathbb{L}(\mathbb{J})}, \|\Gamma_\theta^{-1}\|_{\mathbb{L}(\mathbb{J})}) \leq \mathbb{v}_\theta \in \mathbb{R}_0^+$. then (12.11) and *Comment §12.35* imply

$$\begin{aligned} \mathbb{v}_\theta^{-1} R_n^m(\theta, \phi) &\leq \mathbb{P}_\theta^n(|\phi\nu(\widehat{\theta}_\cdot^m - \theta)|^2) = n^{-1}\|\phi\mathbb{1}_\cdot^m\|_{\Gamma_\theta}^2 + |\phi\nu(\theta\mathbb{1}_\cdot^{m\perp})|^2 \\ &\leq \mathbb{v}_\theta R_n^{m_n^\circ}(\theta, \phi) \quad \text{for all } m, n \in \mathbb{N}. \end{aligned}$$

As a consequence we immediately obtain the following *oracle inequality*

$$\begin{aligned} v_\theta^{-1} R_n^\circ(\theta, \phi) &\leq \inf_{m \in \mathbb{N}} \mathbb{P}_\theta^n (|\phi\nu(\hat{\theta}^m - \theta)|^2) \leq \mathbb{P}_\theta^n (|\phi\nu(\hat{\theta}^{m_n^\circ} - \theta)|^2) \\ &\leq v_\theta R_n^\circ(\theta, \phi) \leq v_\theta^2 \inf_{m \in \mathbb{N}} \mathbb{P}_\theta^n (|\phi\nu(\hat{\theta}^m - \theta)|^2), \end{aligned} \quad (12.12)$$

and hence, $R_n^\circ(\theta, \phi)$, m_n° and the statistic $\hat{\theta}^{m_n^\circ}$, respectively, is an *oracle bound*, an *oracle dimension* and *oracle optimal* (up to the constant v_θ^2). □

§12.37 **Remark.** Arguing similarly as in Remark §12.16 we note that $R_n^\circ(\theta, \phi) = o(1)$ as $n \rightarrow \infty$, whenever $\|\mathbf{1}_\phi^m\|_\phi^2 \in \mathbb{R}^+$ for all $m \in \mathbb{N}$ and $|\phi\nu(\theta \mathbf{1}_\phi^{m|\perp})| = o(1)$ as $m \rightarrow \infty$. The latter is satisfied, for example, if $\theta \in \text{dom}(\phi\nu)$. The oracle dimension $m_n^\circ = m_n^\circ(\theta, \phi)$ as defined in (12.11) depends again on the unknown parameter of interest θ , and thus also the oracle optimal statistic $\hat{\theta}^{m_n^\circ}$. In other words $\hat{\theta}^{m_n^\circ}$ is not a feasible estimator. □

§12.38 **Corollary** (GSSM §12.05 continued). Let $\hat{\theta} = \theta + n^{-1/2} \dot{B} \sim N_q^n$ as in Model §12.05, where $\dot{B} \sim N_{(0,1)}^{\otimes \mathbb{N}}$ and $\theta = U\theta \in \ell_2$. For $\theta \in \text{dom}(\phi\nu_K)$ the (infeasible) OPE $\hat{\theta}^{m_n^\circ} = \hat{\theta} \mathbf{1}_\phi^{m_n^\circ} \in \ell_2 \mathbf{1}_\phi^{m_n^\circ} \subseteq \text{dom}(\phi\nu_K)$ with oracle dimension m_n° as in (12.11) satisfies

$$N_q^n (|\phi\nu_K(\hat{\theta}^{m_n^\circ} - \theta)|^2) = R_n^\circ(\theta, \phi) = \inf_{m \in \mathbb{N}} N_q^n (|\phi\nu_K(\hat{\theta}^m - \theta)|^2),$$

and hence it is *oracle optimal* (with constant 1).

§12.39 **Proof of Corollary §12.38.** is given in the lecture. □

§12.40 **Illustration.** We illustrate the last results considering usual behaviour for both the variance and the bias term. Similar to the two cases **(p)** and **(np)** in Illustration §12.19 we distinguish here the following two cases

(p) $\phi \in \mathbb{J}$ or there is $K \in \mathbb{N}$ with $\sup\{|\phi\nu(\theta \mathbf{1}_\phi^{m|\perp})|^2 : m \in \mathbb{N} \cap [K, \infty)\} = 0$,

(np) $\phi \notin \mathbb{J}$ and for all $m \in \mathbb{N}$ holds $\sup\{|\phi\nu(\theta \mathbf{1}_\phi^{m|\perp})|^2 : m \in \mathbb{N} \cap [K, \infty)\} \in \mathbb{R}_0^+$.

In case **(p)** the oracle bound is again parametric, i.e. $nR_n^\circ(\theta, \phi) = O(1)$, while in case **(np)** the oracle bound is nonparametric, i.e. $\lim_{n \rightarrow \infty} nR_n^\circ(\theta, \phi) = \infty$. In case **(np)** consider the following two specifications

Table 03 [§12]

Order of the oracle rate $R_n^\circ(\theta, \phi)$ as $n \rightarrow \infty$					
$(j \in \mathbb{N})$	$(a \in \mathbb{R}_0^+)$	(squared bias)	(variance)	m_n°	$R_n^\circ(\theta, \phi)$
$\phi_j = j^{v-1/2}$	θ_j	$ \phi\nu(\theta \mathbf{1}_\phi^{m \perp}) ^2$	$\ \mathbf{1}_\phi^m\ _\phi^2$		
(o) $v \in (0, a)$	$j^{-a-1/2}$	$m^{-2(a-v)}$	m^{2v}	$n^{\frac{1}{2a}}$	$n^{-\frac{(a-v)}{a}}$
$v = 0$	$j^{-a-1/2}$	m^{-2a}	$\log m$	$(\frac{n}{\log n})^{\frac{1}{2a}}$	$\frac{\log n}{n}$
(s) $v \in \mathbb{R}_0^+$	$e^{-j^{2a}}$	$m^{(1-2(a-v))_+} e^{-2m^{2a}}$	m^{2v}	$(\log n)^{\frac{1}{2a}}$	$\frac{(\log n)^{\frac{v}{a}}}{n}$
$v = 0$	$e^{-j^{2a}}$	$m^{(1-2a)_+} e^{-m^{2a}}$	$\log m$	$(\log n)^{\frac{1}{2a}}$	$\frac{\log \log n}{n}$

We note that in Table 03 [§12] the order of the oracle rate $R_n^\circ(\theta, \phi)$ is depict for $v \geq 0$ only. For $v < 0$ the oracle rate $R_n^\circ(\theta, \phi)$ is parametric. □

§12|02|02 Maximal local ϕ -risk

§12.41 **Reminder.** Under Assumption §11.25 we have $\mathbb{J}^a = \mathbb{L}_2^a(\nu) = \text{dom}(M_{\mathfrak{a}}) = \mathbb{J}\mathfrak{a} \subseteq \mathbb{J}$ and the three measures ν , $\mathfrak{a}^{\dagger}\nu$ and $|\phi|\nu$ dominate mutually each other, i.e. they share the same null sets (see **Property** §11.05). We consider \mathbb{J}^a endowed with $\|\cdot\|_{\mathfrak{a}^\dagger} = \|M_{\mathfrak{a}}\cdot\|_{\mathbb{J}}$ and given a constant $r \in \mathbb{R}_0^+$ the ellipsoid $\mathbb{J}^{a,r} := \{h_\bullet \in \mathbb{J}^a : \|h_\bullet\|_{\mathfrak{a}^\dagger} \leq r\} \subseteq \mathbb{J}^a$. Since $(\mathfrak{a}\phi)_\bullet \in \mathbb{J}$, and hence $\|\mathfrak{a}\mathbb{1}_\bullet^{m\perp}\|_\phi = \|(\mathfrak{a}\phi)_\bullet\mathbb{1}_\bullet^{m\perp}\|_{\mathbb{J}} \in \mathbb{R}^+$ for each $m \in \mathbb{N}$ ($\|\mathfrak{a}\mathbb{1}_\bullet^{m\perp}\|_\phi = o(1)$ as $m \rightarrow \infty$ by dominated convergence) we have $\mathbb{J}^a \subseteq \text{dom}(\phi\nu)$ (**Property** §11.27), and $|\phi\nu(\theta\mathbb{1}_\bullet^{m\perp})| \leq r \|\mathfrak{a}\mathbb{1}_\bullet^{m\perp}\|_\phi$ for all $\theta \in \mathbb{J}^{a,r}$ (**Lemma** §11.29). \square

§12.42 **Proposition.** Let the Assumptions §11.25 and §12.28, and $\mathbb{1}_\bullet^m \in \mathbb{L}_2(\phi^2\nu)$ for all $m \in \mathbb{N}$ be satisfied. For all $n, m \in \mathbb{N}$ setting

$$R_n^m(\mathfrak{a}, \phi) := \|\mathfrak{a}\mathbb{1}_\bullet^{m\perp}\|_\phi^2 + n^{-1}\|\mathbb{1}_\bullet^m\|_\phi^2, \quad m_n^* := \arg \min \{R_n^m(\mathfrak{a}, \phi) : m \in \mathbb{N}\}$$

$$\text{and } R_n^*(\mathfrak{a}, \phi) := R_n^{m_n^*}(\mathfrak{a}, \phi) = \min \{R_n^m(\mathfrak{a}, \phi) : m \in \mathbb{N}\} \quad (12.13)$$

we have $\mathbb{E}_\theta^n (|\phi\nu(\widehat{\theta}_\bullet^m - \theta)|^2) \leq (\|\Gamma_\theta\|_{\mathbb{L}(\mathbb{J})} \vee r^2) R_n^*(\mathfrak{a}, \phi)$ for all $\theta = U\theta \in \mathbb{J}^{a,r}$ and $n \in \mathbb{N}$.

§12.43 **Proof of Proposition** §12.42. is given in the lecture. \square

§12.44 **Remark.** Under the assumptions of **Proposition** §12.42 if there exists in addition $v \in \mathbb{R}^+$ satisfying $\|\Gamma_\theta\|_{\mathbb{L}(\mathbb{J})} \leq v$ for all $\theta \in \mathbb{J}^{a,r}$ then

$$\sup \{ \mathbb{E}_\theta^n (|\phi\nu(\widehat{\theta}_\bullet^{m_n^*} - \theta)|^2) : \theta \in \mathbb{J}^{a,r} \} \leq (v \vee r^2) R_n^*(\mathfrak{a}, \phi) \quad \text{for all } n \in \mathbb{N}.$$

Arguing similarly as in **Remark** §12.16 we note that $R_n^*(\mathfrak{a}, \phi) = o(1)$ as $n \rightarrow \infty$, whenever $\|\mathbb{1}_\bullet^m\|_\phi^2 \in \mathbb{R}^+$ for all $m \in \mathbb{N}$ and $\|\mathfrak{a}\mathbb{1}_\bullet^{m\perp}\|_\phi = o(1)$ as $m \rightarrow \infty$. The latter is satisfied, for example, if $(\mathfrak{a}\phi)_\bullet \in \mathbb{J}$ (in equal $\mathfrak{a} \in \mathbb{L}_2(\phi^2\nu)$). Note that the dimension $m_n^* := m_n^*(\mathfrak{a}, \phi)$ as defined in (12.13) does not depend on the unknown parameter of interest θ but on the class $\mathbb{J}^{a,r}$ only, and thus also the statistic $\widehat{\theta}_\bullet^{m_n^*}$. In other words, if the regularity of θ is known in advance, then the OPE $\widehat{\theta}_\bullet^{m_n^*}$ is a feasible estimator. \square

§12.45 **Corollary** (GSSM §12.05 continued). Let $\widehat{\theta}_\bullet = \theta + n^{-1/2}\dot{B}_\bullet \sim N_q^n$ as in Model §12.05, where $\dot{B}_\bullet \sim N_{(0,1)}^{\otimes N}$ and $\theta = U\theta \in \ell_2$. Under Assumption §11.25 the OPE $\widehat{\theta}_\bullet^{m_n^*} = \widehat{\theta}_\bullet\mathbb{1}_\bullet^{m_n^*} \in \ell_2\mathbb{1}_\bullet^{m_n^*} \subseteq \text{dom}(\phi\nu_N)$ with dimension m_n^* as in (12.13) satisfies

$$\sup \{ N_q^n (|\phi\nu_N(\widehat{\theta}_\bullet^{m_n^*} - \theta)|^2) : \theta \in \ell_2^{a,r} \} \leq C R_n^*(\mathfrak{a}, \phi) \quad \text{for all } n \in \mathbb{N} \quad (12.14)$$

with constant $C = 1 \vee r^2$.

§12.46 **Proof of Corollary** §12.45. is given in the lecture. \square

§12.47 **Illustration.** We illustrate the last results considering usual behaviour for $\mathfrak{a}, \phi \in \mathcal{J}$. We distinguish the following two cases **(p)** $\phi \in \mathbb{J}$, and **(np)** $\phi \notin \mathbb{J}$. Interestingly, in case **(p)** the bound in **Proposition** §12.42 is parametric, that is, $nR_n^*(\mathfrak{a}, \phi) = O(1)$, in case **(np)** the bound is nonparametric, i.e. $\lim_{n \rightarrow \infty} nR_n^*(\mathfrak{a}, \phi) = \infty$. In case **(np)** consider the following two specifications:

Table 04 [§12]

Order of the rate $R_n^*(\mathbf{a}, \phi)$ as $n \rightarrow \infty$

$(j \in \mathbb{N})$ $\phi = j^{v-1/2}$	$(a \in \mathbb{R}_0^+)$ α_j^2	(squared bias) $\ \mathbf{a} \cdot \mathbb{1}_\phi^{m \perp}\ _\phi^2$	(variance) $\ \mathbb{1}_\phi^m\ _\phi^2$	m_n^*	$R_n^*(\mathbf{a}, \phi)$
(o) $v \in (0, a)$ $v = 0$	j^{-2a}	$m^{-2(a-v)}$	m^{2v}	$n^{\frac{1}{2a}}$	$n^{-\frac{(a-v)}{a}}$
	j^{-2a}	m^{-2a}	$\log m$	$\left(\frac{n}{\log n}\right)^{\frac{1}{2a}}$	$\frac{\log n}{n}$
(s) $v \in \mathbb{R}_0^+$ $v = 0$	$e^{-j^{2a}}$	$m^{2(v-a)+} e^{-m^{2a}}$	m^{2v}	$(\log n)^{\frac{1}{2a}}$	$\frac{(\log n)^{\frac{v}{a}}}{n}$
	$e^{-j^{2a}}$	$e^{-m^{2a}}$	$\log m$	$(\log n)^{\frac{1}{2a}}$	$\frac{\log \log n}{n}$

We note that in Table 04 [§12] the order of the rate $R_n^*(\mathbf{a}, \phi)$ is depict for $v \geq 0$ only. For $v < 0$ the rate $R_n^*(\mathbf{a}, \phi)$ is parametric. □

§13 Minimax optimal estimation

§13|01 Minimax theory: a general approach

Suppose that the function of interest θ belongs to a class $\Theta \subseteq \mathbb{H}$. For each noise level $n \in \mathbb{N}$ let $\mathbb{P}_\Theta^n := (\mathbb{P}_\theta^n)_{\theta \in \Theta}$ denote a family of probability measures and let \mathbb{E}_θ^n be the expectation with respect to the measure \mathbb{P}_θ^n in \mathbb{P}_Θ^n . Moreover, we assume that the probability measure associated with an observable quantity belongs to \mathbb{P}_Θ^n .

§13.01 **GSSM (§10.21 continued)**. Considering $\ell_2 = \mathbb{L}_2(\mathbb{N}, 2^{\mathbb{N}}, \nu_N)$ and a surjective partial isometry $\mathbf{U} \in \mathbb{L}(\mathbb{H}, \ell_2)$, which is fixed and presumed to be known in advance, we illustrate the minimax approach in a Gaussian sequence space model §10.21. Here the observable stochastic process $\hat{\theta}_n = \theta_n + n^{-1/2} \dot{B}_n$ is a noisy version of $\theta = \mathbf{U}\theta \in \ell_2$ and $\dot{B}_n \sim N_{(0,1)}^{\otimes \mathbb{N}}$. Consequently, $\hat{\theta}_n$ admits a $N_{\mathbf{a}}^n$ -distribution belonging to the family $\mathbb{N}_\Theta^n := (N_{\mathbf{a}}^n)_{\mathbf{a} \in \Theta}$. Summarising the observations satisfy a statistical product experiment $(\mathbb{R}^{\mathbb{N}}, \mathcal{B}^{\otimes \mathbb{N}}, \mathbb{N}_\Theta^n)$ where $\Theta \subseteq \ell_2$. □

Assume furthermore, that an estimator $\tilde{\theta}$ of θ based on the observable quantities is available which takes its values in \mathbb{H} but does not necessarily belong to Θ . We shall measure the accuracy of any estimator $\tilde{\theta}$ of θ by its distance $\mathfrak{d}_{\text{ist}}(\tilde{\theta}, \theta)$ where $\mathfrak{d}_{\text{ist}}(\cdot, \cdot)$ is a certain semi metric to be specified below. Moreover, we call the quantity $\mathbb{P}_\theta^n(\mathfrak{d}_{\text{ist}}^2(\tilde{\theta}, \theta))$ risk of the estimator $\tilde{\theta}$ of θ .

§13.02 **Definition**. Given an estimator $\tilde{\theta}$ of a function of interest θ belonging to a class of solutions Θ based on observable quantities with probability measure $\mathbb{P}_\theta^n \in \mathbb{P}_\Theta^n$ we call

$$\mathcal{R}_n[\tilde{\theta} | \Theta] := \sup \{ \mathbb{P}_\theta^n(\mathfrak{d}_{\text{ist}}^2(\tilde{\theta}, \theta)) : \theta \in \Theta \}$$

its *maximal risk* over Θ . □

§13.03 **Remark**. An advantage of taking a maximal risk instead of a risk is that the former does not depend on the unknown function θ . Imagine we would have taken a constant estimator, say $\tilde{\theta} = h$, of θ . This would be the perfect estimator if by chance $\theta = h$, but in all other cases this estimator is likely to perform poorly. Therefore it is reasonable to consider the supremum over the whole class of possible functions in order to get consolidated findings. However, considering the maximal risk may be a very pessimistic point of view. □

§13.04 **Definition**. Consider a maximal risk $\mathcal{R}_n[\cdot | \Theta]$ over a family \mathbb{P}_Θ^n of probability measures. Let $\hat{\theta}$ be an estimator of $\theta \in \Theta$, $C \in \mathbb{R}_0^+$ and for each $n \in \mathbb{N}$ let $R_n^* \in \mathbb{R}^+$ satisfy

(lower) R_n^* is a *lower bound* up to the constant C^{-1} of the maximal risk over Θ , that is

$$\inf_{\tilde{\theta}} \mathcal{R}_n[\tilde{\theta}|\Theta] \geq C^{-1} R_n^*$$

where the infimum is taken over all possible estimators of θ ;

(upper) R_n^* is an *upper bound* up to the constant C of the maximal risk over Θ , that is

$$\mathcal{R}_n[\hat{\theta}|\Theta] \leq C R_n^*$$

Then we call R_n^* *minimax-bound* and the estimator $\hat{\theta}$ *minimax-optimal* (up to the constant C). As a consequence, up to the constant C^2 the estimator $\hat{\theta}$ attains the lower maximal risk bound that is, $\mathcal{R}_n[\hat{\theta}|\Theta] \leq C^2 \inf_{\tilde{\theta}} \mathcal{R}_n[\tilde{\theta}|\Theta]$. \square

§13.05 **Remark.** We call a minimax-bound $(R_n^*)_{n \in \mathbb{N}}$ a *minimax-optimal rate* (of convergence) if in addition $R_n^* = o(1)$ as $n \rightarrow \infty$. It is worth noting that a minimax-optimal rate is not unique since every other rate that is equivalent of order is also minimax-optimal. \square

§13.06 **Nonparametric regression with uniform design (nRu).** Let the $[0, 1] \times \mathbb{R}$ -valued random vector (X, Y) obeys \mathbb{P}^X -a.e. a nonparametric regression model $\mathbb{P}_f(Y|X) = f$ (see section §09). For convenience, in addition the regressor X is supposed to be uniformly distributed on the interval $[0, 1]$, i.e. $X \sim \mathcal{U}_{[0,1]}$. As a consequence, we have $\mathbb{P}^X = \mathbb{1}_{[0,1]}$ and $\mathbb{L}_2([0, 1], \mathcal{B}_{[0,1]}, \mathbb{P}^X) = \mathbb{L}_2([0, 1], \mathcal{B}_{[0,1]}, \lambda_{[0,1]}) = \mathbb{L}_2(\lambda_{[0,1]})$. Here and subsequently we assume that the conditional distribution $\mathbb{P}_f^{Y|X}$ of Y given X is regular, and thus $\mathbb{P}_f^{Y|X}(\text{id}_{\mathbb{R}}) = \mathbb{P}_f(Y|X) = f$ \mathbb{P}^X -a.s.. Let us denote in this situation by $\mathcal{U}_f := \mathcal{U}_{[0,1]} \odot \mathbb{P}_f^{Y|X}$ the joint distribution of (X, Y) defined on $([0, 1] \times \mathbb{R}, \mathcal{B}_{[0,1]} \otimes \mathcal{B})$, but keep in mind, that the conditional distribution $\mathbb{P}_f^{Y|X}$ of Y given X is still not specified. We assume that $f \in \mathbb{F} \subseteq \mathbb{L}_2(\lambda_{[0,1]})$. Summarising the observations satisfy a statistical product experiment $(([0, 1] \times \mathbb{R})^n, \mathcal{B}_{[0,1] \times \mathbb{R}}^n, \mathcal{U}_f^{\otimes n} = (\mathcal{U}_f^{\otimes n})_{f \in \mathbb{F}})$ where $\mathbb{F} \subseteq \mathbb{L}_2(\lambda_{[0,1]})$. Let us assume in addition that Y given X is normally distributed with conditional mean $f(X)$ and conditional variance $\sigma^2 \in \mathbb{R}_0^+$, that is $\mathbb{P}_f^{Y|X} = \mathcal{N}_{(f(X), \sigma^2)}$. In this situation we denote by $\mathcal{U}_{f, \sigma} := \mathcal{U}_{[0,1]} \odot \mathcal{N}_{(f(X), \sigma^2)}$ the joint distribution of (X, Y) . We first consider the case that the variance σ^2 is known *a priori* (i.e. $\sigma^2 = 1$), and in a second step we dismiss this information. Obviously, the distribution $\mathcal{U}_{f, \sigma}^{\otimes n}$ depends not only on the parameter of interest $f \in \mathbb{F}$ and the noise level $n \in \mathbb{N}$, but also on the variance $\sigma^2 \in \mathbb{R}_0^+$ which plays the role of a nuisance parameter. Consequently, let $\mathcal{U}_{\mathbb{F} \times \mathbb{R}_0^+} := (\mathcal{U}_{f, \sigma})_{f \in \mathbb{F}, \sigma \in \mathbb{R}_0^+}$ denote the family of possible distributions of (X, Y) . Summarising, if the variance is unknown then the observations satisfy a statistical product experiment $(([0, 1] \times \mathbb{R})^n, \mathcal{B}_{[0,1] \times \mathbb{R}}^n, \mathcal{U}_{\mathbb{F} \times \mathbb{R}_0^+}^{\otimes n})$ where $\mathbb{F} \subseteq \mathbb{L}_2(\lambda_{[0,1]})$. \square

More generally, given a class of solutions Θ , a class of nuisance parameters Ξ and a noise level $n \in \mathbb{N}$ let $\mathbb{P}_{\Theta \times \Xi}^n := (\mathbb{P}_{\theta, \xi}^n)_{\theta \in \Theta, \xi \in \Xi}$ denote a family of probability measures. Moreover, we assume that the probability measure associated with an observable quantity belongs to $\mathbb{P}_{\Theta \times \Xi}^n$. Note that dismissing in Model §13.06 the assumption of a normally distributed error the class of nuisance parameters Ξ equals the family of possible conditional distributions of the error terms.

§13.07 **Definition.** Given an estimator $\tilde{\theta}$ of a function of interest θ belonging to a class of solutions Θ based on observable quantities with probability measure $\mathbb{P}_{\theta, \xi}^n \in \mathbb{P}_{\Theta \times \Xi}^n$ we call

$$\mathcal{R}_n[\tilde{\theta}|\Theta, \Xi] := \sup \{ \mathbb{E}_{\theta, \xi}^n (\mathfrak{d}_{\text{ist}}^2(\tilde{\theta}, \theta)) : \theta \in \Theta, \xi \in \Xi \}$$

its *maximal risk* over $\Theta \times \Xi$. \square

§13.08 **Remark.** Taking the supremum over the class of nuisances parameters allows us to quantify the additional complexity due to the presence of the nuisance parameter. Moreover, if there exist an estimator $\hat{\theta}$, a constant $C \in \mathbb{R}_0^+$ and for each $n \in \mathbb{N}$ there is $R_n^* \in \mathbb{R}^+$ such that

(lower) R_n^* is a *lower bound* up to the constant C of the maximal risk over $\Theta \times \Xi$, that is

$$\inf_{\tilde{\theta}} \mathcal{R}_n[\tilde{\theta} | \Theta, \Xi] \geq C^{-1} R_n^*$$

where the infimum is taken over all possible estimators of θ ;

(upper) R_n^* is an *upper bound* up to the constant C of the maximal risk over $\Theta \times \Xi$, that is

$$\mathcal{R}_n[\hat{\theta} | \Theta, \Xi] \leq C R_n^*,$$

then we call R_n^* *minimax-bound* and the estimator $\hat{\theta}$ *minimax-optimal* (up to the constant C).

As a consequence, up to the constant C^2 the estimator $\hat{\theta}$ attains the lower maximal risk bound that is, $\mathcal{R}_n[\hat{\theta} | \Theta, \Xi] \leq C^2 \inf_{\tilde{\theta}} \mathcal{R}_n[\tilde{\theta} | \Theta, \Xi]$. Typically, we assume first that the nuisance parameter ξ is known *a priori*, and hence $\mathbb{P}_{\Theta \times \{\xi\}}^n$ is a family of probability measures associated with the observable quantities. In this situation, we consider the maximal risk $\{\mathbb{E}_{\theta, \xi}^n(\mathfrak{d}_{\text{ist}}^2(\tilde{\theta}, \theta)) : \theta \in \Theta\}$ and we seek a bound R_n^* up to a constant which depends possibly on the nuisance parameter ξ . However, if the bound R_n^* is a valid lower and upper bound up to a constant uniformly for all nuisance parameters $\xi \in \Xi$, then it is, obviously, also a bound of the maximal risk $\mathcal{R}_n[\hat{\theta} | \Theta, \Xi]$. \square

§13.09 **Reminder.** Considering a Hilbert space $\mathbb{J} = \mathbb{L}_2(\mathcal{J}, \mathcal{J}, \nu)$ and a surjective partial isometry $U \in \mathbb{L}(\mathbb{H}, \mathbb{J})$, which are fixed and presumed to be known in advance, we study *statistical direct problems* as in Definition §10.19. Given weights $\mathfrak{a} \in \mathcal{J}_0$ we introduce $\mathbb{J}^{\mathfrak{a}} = \text{dom}(M_{\mathfrak{a}}) = \mathbb{J}\mathfrak{a} = \mathbb{L}_2(\mathfrak{a}^{2\uparrow}\nu)$ endowed with $\|\cdot\|_{\mathfrak{a}^\uparrow} := \|\cdot\|_{\mathbb{L}_2(\mathfrak{a}^{2\uparrow}\nu)}$ and the ellipsoid $\mathbb{J}^{\mathfrak{a},r} := \{h \in \mathbb{J}^{\mathfrak{a}} : \|h\|_{\mathfrak{a}^\uparrow}^2 \leq r^2\} \subseteq \mathbb{J}^{\mathfrak{a}}$, where the measures ν and $\mathfrak{a}^{2\uparrow}\nu$ dominate mutually each other. We consider the following global and local measures of accuracy (compare Subsections §12|01 and §12|02).

(global) Given weights $\mathfrak{v} \in \mathcal{J}_0$ satisfying Assumption §11.12 introduce $\mathbb{L}_2(\mathfrak{v}^2\nu) = \text{dom}(M_{\mathfrak{v}}) = \mathbb{J}\mathfrak{v}^\uparrow \subseteq \mathbb{J}$ and $\|\cdot\|_{\mathfrak{v}} = \|M_{\mathfrak{v}} \cdot\|_{\mathbb{J}}$, where $\mathbb{J}^{\mathfrak{a},r} \subseteq \mathbb{L}_2(\mathfrak{v}^2\nu)$ (Property §11.15). For $\theta = U\theta \in \mathbb{J}^{\mathfrak{a},r}$ we call $\mathfrak{d}_{\text{ist}}(\tilde{\theta}, \theta) = \|\tilde{\theta} - \theta\|_{\mathfrak{v}}$ *global v-error*, $\mathbb{E}_\theta^n(\|\tilde{\theta} - \theta\|_{\mathfrak{v}}^2)$ *global v-risk* and

$$\mathcal{R}_n^{\mathfrak{v}}[\tilde{\theta} | \mathbb{J}^{\mathfrak{a},r}] := \sup \{ \mathbb{E}_\theta^n(\|\tilde{\theta} - \theta\|_{\mathfrak{v}}^2) : \theta = U\theta \in \mathbb{J}^{\mathfrak{a},r} \}$$

maximal v-risk over $\mathbb{J}^{\mathfrak{a},r}$.

(local) Given $\phi \in \mathcal{J}_0$ satisfying Assumption §11.25 introduce $\text{dom}(\phi\nu) := \{h \in \mathbb{J} : \phi h \in \mathbb{L}_1(\nu)\}$ and the linear functional $\phi\nu : \mathbb{J} \supseteq \text{dom}(\phi\nu) \rightarrow \mathbb{R}$ with $h \mapsto \phi\nu(h) := \nu(\phi h)$ where $\mathbb{J}^{\mathfrak{a},r} \subseteq \text{dom}(\phi\nu)$ (Property §11.27). For $\theta \in \mathbb{J}^{\mathfrak{a},r}$ we call $\mathfrak{d}_{\text{ist}}(\tilde{\theta}, \theta) = |\phi\nu(\tilde{\theta} - \theta)|$ *local ϕ -error*, $\mathbb{E}_\theta^n(|\phi\nu(\tilde{\theta} - \theta)|^2)$ *local ϕ -risk* and

$$\mathcal{R}_n^\phi[\tilde{\theta} | \mathbb{J}^{\mathfrak{a},r}] := \sup \{ \mathbb{E}_\theta^n(|\phi\nu(\tilde{\theta} - \theta)|^2) : \theta = U\theta \in \mathbb{J}^{\mathfrak{a},r} \}$$

maximal ϕ -risk over $\mathbb{J}^{\mathfrak{a},r}$.

We formulate the results in terms of $\theta = U\theta \in \mathbb{J}$ rather than directly for $\theta \in \mathbb{H}$. Since U is known, considering the class $\mathbb{H}^{\mathfrak{a},r} := U^* \mathbb{J}^{\mathfrak{a},r} := \{\theta \in \mathbb{H} : U\theta \in \mathbb{J}^{\mathfrak{a},r}\}$ we obtain immediately also bounds over $\mathbb{H}^{\mathfrak{a},r}$ for the maximal global risk

$$\mathcal{R}_n^{\mathfrak{v}}[\tilde{\theta} | U^* \mathbb{J}^{\mathfrak{a},r}] := \sup \{ \mathbb{E}_\theta^n(\|U(\tilde{\theta} - \theta)\|_{\mathfrak{v}}^2) : \theta \in \mathbb{H}^{\mathfrak{a},r} \}$$

and maximal local risk

$$\mathcal{R}_n^\phi[\tilde{\theta} | U^* \mathbb{J}^{\mathfrak{a},r}] := \sup \{ \mathbb{E}_\theta^n(|\phi\nu(U(\tilde{\theta} - \theta))|^2) : \theta \in \mathbb{H}^{\mathfrak{a},r} \}$$

which we do not explicitly state in the sequel. \square

§13|02 Deriving a lower bound: a general reduction scheme

For a detailed discussion of several other strategies to derive lower bounds we refer the reader, for example, to the text book by Tsybakov [2009].

§13.10 **Definition.** Let \mathbb{P}_0 and \mathbb{P}_1 be two probability measures on a measurable space $(\mathcal{X}, \mathcal{X})$.

(a) The function

$$\text{KL}(\mathbb{P}_0|\mathbb{P}_1) = \begin{cases} \mathbb{E}_0 \left(\log \frac{d\mathbb{P}_0}{d\mathbb{P}_1} \right) = \int \log \left(\frac{d\mathbb{P}_0}{d\mathbb{P}_1} \right) d\mathbb{P}_0, & \text{if } \mathbb{P}_0 \ll \mathbb{P}_1, \\ +\infty, & \text{otherwise} \end{cases}$$

is called *Kullback-Leibler-divergence* of \mathbb{P}_0 with respect to \mathbb{P}_1 .

Let $\mu \in \mathcal{M}_\sigma(\mathcal{X})$ be a \mathbb{P}_0 and \mathbb{P}_1 dominating σ -finite measure (e.g. $\mathbb{P}_0, \mathbb{P}_1 \ll \mu = \mathbb{P}_0 + \mathbb{P}_1$). We write $d\mathbb{P}_0 := d\mathbb{P}_0/d\mu$ and $d\mathbb{P}_1 := d\mathbb{P}_1/d\mu$ for short.

(b) The *Hellinger distance* between \mathbb{P}_0 and \mathbb{P}_1 is defined by

$$H(\mathbb{P}_0, \mathbb{P}_1) := \left(\int |\sqrt{d\mathbb{P}_0} - \sqrt{d\mathbb{P}_1}|^2 \right)^{1/2} := \|\sqrt{d\mathbb{P}_0} - \sqrt{d\mathbb{P}_1}\|_{\mathbb{L}_2(\mu)}$$

(c) and the *Hellinger affinity* is given by

$$\rho(\mathbb{P}_0, \mathbb{P}_1) := \int \sqrt{d\mathbb{P}_0} \sqrt{d\mathbb{P}_1} := \langle \sqrt{d\mathbb{P}_0}, \sqrt{d\mathbb{P}_1} \rangle_{\mathbb{L}_2(\mu)},$$

where both do not depend on the choice of the dominating measure μ . □

§13.11 **Remark.** The Kullback-Leibler-divergence satisfies $\text{KL}(\mathbb{P}_0|\mathbb{P}_1) \geq 0$ as well as $\text{KL}(\mathbb{P}_0|\mathbb{P}_1) = 0$ if and only if $\mathbb{P}_0 = \mathbb{P}_1$, but $\text{KL}(\cdot|\cdot)$ is not symmetric. Moreover, for product measures holds $\text{KL}(\mathbb{P}_{0,1} \otimes \mathbb{P}_{0,2}|\mathbb{P}_{1,1} \otimes \mathbb{P}_{1,2}) = \text{KL}(\mathbb{P}_{0,1}|\mathbb{P}_{1,1}) + \text{KL}(\mathbb{P}_{0,2}|\mathbb{P}_{1,2})$. □

§13.12 **Lemma.** (i) $0 \leq H^2(\mathbb{P}_0, \mathbb{P}_1) \leq 2$; (ii) $\rho(\mathbb{P}_0, \mathbb{P}_1) = 1 - \frac{1}{2}H^2(\mathbb{P}_0, \mathbb{P}_1)$; and (iii) $H^2(\mathbb{P}_0, \mathbb{P}_1) \leq \text{KL}(\mathbb{P}_0|\mathbb{P}_1)$.

§13.13 **Proof of Lemma §13.12.** Exercise. □

§13.14 **Lemma.** For $a, b \in \ell_2$ and $n \in \mathbb{N}$ we have $\text{KL}(N_a^n|N_b^n) = \frac{n}{2}\|a - b\|_{\ell_2}^2$.

§13.15 **Proof of Lemma §13.14.** Exercise. □

§13.16 **Notation.** Recall that the semi metric $\mathfrak{d}_{\text{ist}}(\cdot, \cdot)$ is symmetric and satisfies the triangular inequality. Moreover, here and subsequently we suppose that for an estimator $\tilde{\theta}$ and parameter θ^0 and θ^1 such that $\mathfrak{d}_{\text{ist}}(\theta^0, \theta^1) \in \mathbb{R}_{>0}^+$ the quantities $\mathfrak{d}_{\text{ist}}(\tilde{\theta}, \theta^0)$ and $\mathfrak{d}_{\text{ist}}(\tilde{\theta}, \theta^1)$ are measurable. □

§13.17 **Lemma.** Let \mathbb{P}_0 and \mathbb{P}_1 be two probability measures on a measurable space $(\mathcal{X}, \mathcal{X})$. Suppose that for an estimator $\tilde{\theta}$ and parameter θ^0 and θ^1 with $\mathfrak{d}_{\text{ist}}(\theta^0, \theta^1) \in \mathbb{R}_{>0}^+$ the quantities $\mathfrak{d}_{\text{ist}}(\tilde{\theta}, \theta^0)$ and $\mathfrak{d}_{\text{ist}}(\tilde{\theta}, \theta^1)$ are measurable. Then, we have

$$\mathbb{P}_0(\mathfrak{d}_{\text{ist}}^2(\tilde{\theta}, \theta^0)) + \mathbb{P}_1(\mathfrak{d}_{\text{ist}}^2(\tilde{\theta}, \theta^1)) \geq \frac{1}{2} \mathfrak{d}_{\text{ist}}^2(\theta^0, \theta^1) \rho^2(\mathbb{P}_0, \mathbb{P}_1). \quad (13.01)$$

§13.18 **Proof of Lemma §13.17.** is given in the lecture. □

§13|03 Lower bound based on two hypotheses

§13.19 **Lemma** (*Lower bound based on two hypotheses*). Given a noise level $n \in \mathbb{N}$ let $\mathbb{P}_\Theta^n := (\mathbb{P}_\theta^n)_{\theta \in \Theta}$ be a family of probability measures. If there are $\theta^0, \theta^1 \in \Theta$ with associated probability measures $\mathbb{P}_\theta := \mathbb{P}_\theta^n$ and $\mathbb{P}_\theta := \mathbb{P}_\theta^n$ such that $H(\mathbb{P}_\theta, \mathbb{P}_\theta) \leq 1$ then we have

$$\inf_{\tilde{\theta}} \mathcal{R}_n[\tilde{\theta} | \Theta] \geq \frac{1}{16} \mathfrak{d}_{\text{ist}}^2(\theta^0, \theta^1).$$

where the infimum is taken over all possible estimators.

§13.20 **Proof** of Lemma §13.19. is given in the lecture. \square

§13.21 **Remark** (*Lower bound for a local ϕ -risk*). Due to the bounded Hellinger distance in Lemma §13.19, Le Cam's general method (see Le Cam [1973]) and Pinsker's inequality allow to derive a lower bound for a local ϕ -risk as in Reminder §13.09. However, in this special setting a lower bound can be obtained elementarily from Lemma §13.19, which in this situation states

$$\inf_{\tilde{\theta}} \mathcal{R}_n^\phi[\tilde{\theta} | \Theta] \geq \frac{1}{16} |\phi\nu(\theta^0 - \theta^1)|^2.$$

If we consider furthermore candidates $\theta^0 := \theta^*$ and $\theta^1 = -\theta^*$ for some $\theta^* \in \Theta$ such that $-\theta^* \in \Theta$, then trivially $|\phi\nu(\theta^0 - \theta^1)|^2 = 4|\phi\nu(\theta^*)|^2$ which in turn implies due to the last assertion

$$\inf_{\tilde{\theta}} \mathcal{R}_n^\phi[\tilde{\theta} | \Theta] \geq \frac{1}{4} |\phi\nu(\theta^*)|^2. \quad (13.02)$$

Often a minimax-optimal lower bound can be found by constructing a candidate $\theta^* = U\theta^* \in \Theta$ that has the largest possible $|\phi\nu(\theta^*)|^2$ -value but \mathbb{P}_θ^n and $\mathbb{P}_{-\theta^*}^n$ are still statistically indistinguishable in the sense that $H(\mathbb{P}_\theta^n, \mathbb{P}_{-\theta^*}^n) \leq 1$. \square

§13.22 **Reminder** (*Maximal local ϕ -risk in GSSM §13.01*). Given Model §13.01 we consider an OPE as in Section §12. Here the observable stochastic process $\hat{\theta}_n = \theta + n^{-1/2}\dot{B}_n \sim N_q^n$ is a noisy version of $\theta = U\theta \in \Theta \subseteq \ell_2$ and $\dot{B}_n \sim N_{(0,1)}^{\otimes N}$. Consequently, $\hat{\theta}_n$ admits a N_q^n -distribution belonging to the family $N_\Theta^n := (N_q^n)_{\theta \in \Theta}$. Summarising the observations satisfy a statistical product experiment $(\mathbb{R}^N, \mathcal{B}^N, N_\Theta^n)$ where $\Theta \subseteq \ell_2$. Under Assumption §11.25 in Corollary §12.45 an upper bound for a maximal local ϕ -risk of an OPE is shown. More precisely, the performance of the OPE $\hat{\theta}_n^m = \hat{\theta}_n \mathbb{1}_n^m \in \ell_2 \mathbb{1}_n^m \subseteq \text{dom}(\phi_{\mathcal{U}_n})$ with dimension $m \in \mathbb{N}$ is measured by its maximal local ϕ -risk, that is

$$\mathcal{R}_n^\phi[\hat{\theta}_n^m | \ell_2^{\text{a.r.}}] := \sup \{ N_q^n (|\phi_{\mathcal{U}_n}(\hat{\theta}_n^m - \theta)|^2) : \theta \in \ell_2^{\text{a.r.}} \}.$$

Let us recall (12.13) where for $n, m \in \mathbb{N}$ we have defined

$$\begin{aligned} R_n^m(\mathbf{a}, \phi) &:= \|\mathbf{a} \mathbb{1}_n^{m \perp}\|_\phi^2 + n^{-1} \|\mathbb{1}_n^m\|_\phi^2, \quad m_n^* := \arg \min \{ R_n^m(\mathbf{a}, \phi) : m \in \mathbb{N} \} \\ \text{and} \quad R_n^*(\mathbf{a}, \phi) &:= R_n^{m_n^*}(\mathbf{a}, \phi) = \min \{ R_n^m(\mathbf{a}, \phi) : m \in \mathbb{N} \}. \end{aligned} \quad (13.03)$$

By Corollary §12.45 under Assumption §11.25 the maximal local ϕ -risk of an OPE $\hat{\theta}_n^{m_n^*}$ with optimally chosen dimension m_n^* as in (13.03) satisfies

$$\mathcal{R}_n^\phi[\hat{\theta}_n^{m_n^*} | \ell_2^{\text{a.r.}}] \leq C R_n^*(\mathbf{a}, \phi)$$

with $C = 1 \vee r^2$. \square

§13.23 **Notation.** For sequences $a_n, b_n \in (\mathbb{K})^{\mathbb{N}}$ taking its values in $\mathbb{K} \in \{\mathbb{R}, \mathbb{R}^+, \mathbb{R}_{\setminus 0}^+, \mathbb{Q}, \mathbb{Z}, \dots\}$ we write $a_n \in (\mathbb{K})_{\nearrow}^{\mathbb{N}}$ and $b_n \in (\mathbb{K})_{\searrow}^{\mathbb{N}}$ if a_n and b_n , respectively, is monotonically *non-decreasing* and *non-increasing*. If in addition $a_n \rightarrow \infty$ and $b_n \rightarrow 0$ as $n \rightarrow \infty$, then we write $a_n \in (\mathbb{K})_{\nearrow\infty}^{\mathbb{N}}$ and $b_n \in (\mathbb{K})_{\searrow 0}^{\mathbb{N}}$ for short. \square

§13.24 **Assumption.** Consider $\phi, \mathbf{a} \in (\mathbb{R}_{\setminus 0})^{\mathbb{N}}$ such that $\mathbf{a} \in \ell_\infty$ and $(\mathbf{a}\phi)_n \in \ell_2$ (i.e. Assumption §11.25 is satisfied and $\|\mathbf{a}\mathbf{1}_n^{m_n \perp}\|_\phi = \|(\mathbf{a}\phi)_n \mathbf{1}_n^{m_n \perp}\|_{\ell_2} = o(1)$ as $m_n \rightarrow \infty$), and in addition $\mathbf{a}_n^2 \in (\mathbb{R}_{\setminus 0}^+)_{\searrow 0}^{\mathbb{N}}$. \square

§13.25 **Comment.** Assuming $\mathbf{a}_n^2 \in (\mathbb{R}_{\setminus 0}^+)_{\searrow}^{\mathbb{N}}$ is rather weak. If we suppose in addition $\liminf_{j \rightarrow \infty} \mathbf{a}_j^2 \geq c > 0$, and hence $\mathbf{a}_n^2 \notin (\mathbb{R}_{\setminus 0}^+)_{\searrow 0}^{\mathbb{N}}$, then the assumption $(\mathbf{a}\phi)_n \in \ell_2$ implies $\phi \in \ell_2$ and hence the rate $R_n^*(\mathbf{a}, \phi)$ is parametric (**Illustration** §12.47). Since we are interested in the case of a non-parametric, the additional assumption $\mathbf{a}_n^2 \in (\mathbb{R}_{\setminus 0}^+)_{\searrow 0}^{\mathbb{N}}$ imposes a rather weak condition satisfied also in **Illustration** §12.47.

If $\mathbf{a}_n^2 > n^{-1}$ then exploiting the definition (13.03) and $\phi \in \mathbb{R}_{\setminus 0}$ we have

$$R_n^1(\mathbf{a}, \phi) = n^{-1}\phi_1^2 + (\mathbf{a}\phi)_2^2 + \|\mathbf{a}\mathbf{1}_n^{2 \perp}\|_\phi^2 > n^{-1}\phi_1^2 + n^{-1}\phi_2^2 + \|\mathbf{a}\mathbf{1}_n^{2 \perp}\|_\phi^2 = R_n^2(\mathbf{a}, \phi),$$

and consequently $m_n^* - 1 \in \mathbb{N}$. In this situation, from (definition of the arg min)

$$\begin{aligned} n^{-1}\|\mathbf{1}_n^{m_n^* - 1}\|_\phi^2 + (\mathbf{a}\phi)_{m_n^*}^2 + \|\mathbf{a}\mathbf{1}_n^{m_n^* \perp}\|_\phi^2 \\ = R_n^{m_n^* - 1}(\mathbf{a}, \phi) > R_n^{m_n^*}(\mathbf{a}, \phi) = n^{-1}\|\mathbf{1}_n^{m_n^* - 1}\|_\phi^2 + n^{-1}\phi_{m_n^*}^2 + \|\mathbf{a}\mathbf{1}_n^{m_n^* \perp}\|_\phi^2 \end{aligned}$$

follows $(\mathbf{a}\phi)_{m_n^*}^2 > n^{-1}\phi_{m_n^*}^2$, and hence $\mathbf{a}_{m_n^*}^2 > n^{-1}$ (since $\phi_{m_n^*} \in \mathbb{R}_{\setminus 0}$). On the other hand from

$$\begin{aligned} n^{-1}\|\mathbf{1}_n^{m_n^*}\|_\phi^2 + (\mathbf{a}\phi)_{m_n^* + 1}^2 + \|\mathbf{a}\mathbf{1}_n^{m_n^* + 1 \perp}\|_\phi^2 \\ = R_n^{m_n^*}(\mathbf{a}, \phi) \leq R_n^{m_n^* + 1}(\mathbf{a}, \phi) = n^{-1}\|\mathbf{1}_n^{m_n^*}\|_\phi^2 + n^{-1}\phi_{m_n^*}^2 + \|\mathbf{a}\mathbf{1}_n^{m_n^* + 1 \perp}\|_\phi^2 \end{aligned}$$

follows $(\mathbf{a}\phi)_{m_n^* + 1}^2 \leq n^{-1}\phi_{m_n^* + 1}^2$, and hence $\mathbf{a}_{m_n^* + 1}^2 \leq n^{-1}$ (since $\phi_{m_n^* + 1} \in \mathbb{R}_{\setminus 0}$). Assuming $\mathbf{a}_n^2 > n^{-1}$ we use the property $\mathbf{a}_{m_n^*}^2 > n^{-1} \geq \mathbf{a}_{m_n^* + 1}^2$ in the next proof. \square

§13.26 **Proposition** (*GSSM §13.01 continued*). Let $\hat{\theta} = \theta + n^{-1/2}\dot{B} \sim N_q^n$ as in Model §13.01 where $\dot{B} \sim N_{(0,1)}^{\otimes \mathbb{N}}$ and $\theta = U\theta \in \ell_2$. Given Assumption §13.24 and the notations in (13.03) for all $n \in \mathbb{N} \cap (\mathbf{a}_n^{-2}, \infty)$ we have

$$\inf_{\hat{q}} \mathcal{R}_n^\phi[\tilde{\theta} | \ell_2^{a_n}] \geq 8^{-1}(1 \wedge 2r^2) R_n^*(\mathbf{a}, \phi) \quad (13.04)$$

where the infimum is taken over all estimators $\tilde{\theta}$.

§13.27 **Proof** of **Proposition** §13.26. is given in the lecture. \square

§13.28 **Illustration.** Consider the two specifications (o) and (s) depict in Table 04 [§12] of the **Illustration** §12.47. In both cases Assumption §13.24 is satisfied. Consequently, due to **Proposition** §13.26 the Table 04 [§12] presents the order of the *minimax rate* $R_n^*(\mathbf{a}, \phi)$ which is attained by the *minimax-optimal* OPE $\hat{\theta}^{m_n^*} = \hat{\theta} \mathbf{1}_n^{m_n^*} \in \ell_2 \mathbf{1}_n^{m_n^*} \subseteq \text{dom}(\phi_{\nu_n})$ with optimally selected dimension m_n^* (**Corollary** §12.45). We shall stress, that the order of m_n^* given in the Table 04 [§12] depends on the parameter $a \in \mathbb{R}_{\setminus 0}^+$ characterising the (abstract) smoothness of the solution which is generally not known in advance. \square

§13|04 Lower bound based on m hypotheses

§13.29 **Notation.** For $m \in \mathbb{N}$ set $\mathcal{T}_m := \{-1, 1\}^m$ and for each $\tau := (\tau_j)_{j \in \llbracket m \rrbracket} \in \mathcal{T}_m$ and $j \in \llbracket m \rrbracket$ introduce $\tau^{(j)} \in \mathcal{T}_m$ given by $\tau_j^{(j)} := -\tau_j$ and $\tau_l^{(j)} := \tau_l$ for $l \in \llbracket m \rrbracket \setminus \{j\}$. \square

§13.30 **Lemma (Assouad's cube technique).** Given a noise level $n \in \mathbb{N}$ let $\mathbb{P}_\Theta^n := (\mathbb{P}_\theta^n)_{\theta \in \Theta}$ be a family of probability measures. Suppose there exist $m \in \mathbb{N}$ and distances $\mathfrak{d}_{\text{ist}}^{(j)}(\cdot, \cdot)$, $j \in \llbracket m \rrbracket$ such that $\mathfrak{d}_{\text{ist}}^2(\cdot, \cdot) \geq \sum_{j \in \llbracket m \rrbracket} |\mathfrak{d}_{\text{ist}}^{(j)}(\cdot, \cdot)|^2$. If for each $\tau \in \mathcal{T}_m$ there is $\theta^\tau \in \Theta$ with associated probability measure $\mathbb{P}_\tau := \mathbb{P}_{\theta^\tau}$ such that for all $\tau \in \mathcal{T}_m$ and $j \in \llbracket m \rrbracket$ we have $H(\mathbb{P}_\tau, \mathbb{P}_{\tau^{(j)}}) \leq 1$ then we obtain

$$\inf_{\tilde{\theta}} \mathcal{R}_n[\tilde{\theta} | \Theta] \geq 2^{-m} \sum_{\tau \in \mathcal{T}_m} \frac{1}{16} \sum_{j \in \llbracket m \rrbracket} |\mathfrak{d}_{\text{ist}}^{(j)}(\theta^\tau, \theta^{\tau^{(j)}})|^2$$

where the infimum is taken over all possible estimators.

§13.31 **Proof of Lemma §13.30.** is given in the lecture. \square

§13.32 **Remark (Lower bound of a global \mathfrak{v} -risk).** The last result allows to derive a lower bound for a global \mathfrak{v} -risk as in **Reminder §13.09** which in case $\mathbb{J} = \ell_2 = \mathbb{L}_2(\mathbb{N}, 2^{\mathbb{N}}, \nu_{\mathbb{N}})$ states

$$\inf_{\tilde{\theta}} \mathcal{R}_n^{\mathfrak{v}}[\tilde{\theta} | \Theta] \geq 2^{-m} \sum_{\tau \in \mathcal{T}_m} \frac{1}{16} \sum_{j \in \llbracket m \rrbracket} \mathfrak{v}_j^2 |\theta_j^\tau - \theta_j^{\tau^{(j)}}|^2.$$

If we assume furthermore candidates $\theta^\tau := (\tau_j \theta_j^* \mathbf{1}_j^m)_{j \in \mathbb{N}} \in \Theta$, $\tau \in \mathcal{T}_m$, for some $\theta^* = U\theta^* \in \Theta$, then it is easily seen that $\sum_{j \in \llbracket m \rrbracket} \mathfrak{v}_j^2 |\theta_j^\tau - \theta_j^{\tau^{(j)}}|^2 = 4 \sum_{j \in \llbracket m \rrbracket} \mathfrak{v}_j^2 |\theta_j^*|^2 = 4 \|\theta^* \mathbf{1}_m^m\|_{\mathfrak{v}}^2$ which in turn implies

$$\inf_{\tilde{\theta}} \mathcal{R}_n^{\mathfrak{v}}[\tilde{\theta} | \Theta] \geq 2^{-m} \sum_{\tau \in \mathcal{T}_m} \frac{1}{4} \|\theta^* \mathbf{1}_m^m\|_{\mathfrak{v}}^2 = \frac{1}{4} \|\theta^* \mathbf{1}_m^m\|_{\mathfrak{v}}^2. \quad (13.05)$$

Often a minimax-optimal lower bound can be found by choosing the parameter m and the function θ^* that have the largest possible $\|\theta^* \mathbf{1}_m^m\|_{\mathfrak{v}}^2$ -value although that the associated \mathbb{P}_τ , $\tau \in \mathcal{T}_m$ are still statistically indistinguishable in the sense that $H(\mathbb{P}_\tau, \mathbb{P}_{\tau^{(j)}}) \leq 1$ for all $j \in \llbracket m \rrbracket$ and $\tau \in \mathcal{T}_m$. \square

§13.33 **Reminder (Maximal global \mathfrak{v} -risk in GSSM §13.01).** Given Model §13.01 we consider an OPE as in **Section §12**. Here the observable stochastic process $\hat{\theta}_n = \theta_n + n^{-1/2} \dot{B}_n \sim N_{\mathfrak{a}}^n$ is a noisy version of $\theta = U\theta \in \Theta \subseteq \ell_2$ and $\dot{B}_n \sim N_{(0,1)}^{\otimes \mathbb{N}}$. Consequently, $\hat{\theta}_n$ admits a $N_{\mathfrak{a}}^n$ -distribution belonging to the family $N_{\Theta}^n := (N_{\mathfrak{a}}^n)_{\mathfrak{a} \in \Theta}$. Summarising the observations satisfy a statistical product experiment $(\mathbb{R}^{\mathbb{N}}, \mathcal{B}^{\mathbb{N}}, N_{\Theta}^n)$ where $\Theta \subseteq \ell_2$. Under Assumption §11.12 in **Corollary §12.24** an upper bound for a maximal global \mathfrak{v} -risk of an OPE is shown. More precisely, the performance of the OPE $\hat{\theta}_n^m = \hat{\theta}_n \mathbf{1}_m^m \in \ell_2(\mathfrak{v}^2) \subseteq \ell_2(\mathfrak{v}^2)$ with dimension $m \in \mathbb{N}$ is measured by its maximal global \mathfrak{v} -risk over the ellipsoid $\ell_2^{\mathfrak{a}, \mathfrak{r}}$, that is

$$\mathcal{R}_n^{\mathfrak{v}}[\hat{\theta}_n^m | \ell_2^{\mathfrak{a}, \mathfrak{r}}] := \sup \{ N_{\mathfrak{a}}^n(\|\hat{\theta}_n^m - \theta\|_{\mathfrak{v}}^2) : \theta \in \ell_2^{\mathfrak{a}, \mathfrak{r}} \}.$$

Let us recall (12.06) where for $n, m \in \mathbb{N}$ we have defined $(\mathfrak{a}\mathfrak{v})_{(m)}^2 := \|(\mathfrak{a}\mathfrak{v})^2 \mathbf{1}_m^m\|_{\ell_\infty}$ and

$$\begin{aligned} R_n^m(\mathfrak{a}, \mathfrak{v}) &:= (\mathfrak{a}\mathfrak{v})_{(m)}^2 \vee n^{-1} \|\mathbf{1}_m^m\|_{\mathfrak{v}}^2, \quad m_n^* := \arg \min \{ R_n^m(\mathfrak{a}, \mathfrak{v}) : m \in \mathbb{N} \} \\ &\text{and} \quad R_n^*(\mathfrak{a}, \mathfrak{v}) := R_n^{m_n^*}(\mathfrak{a}, \mathfrak{v}) = \min \{ R_n^m(\mathfrak{a}, \mathfrak{v}) : m \in \mathbb{N} \}. \end{aligned} \quad (13.06)$$

By **Corollary §12.24** under Assumption §11.12 the maximal global \mathfrak{v} -risk of an OPE $\hat{\theta}_n^{m_n^*}$ with optimally chosen dimension m_n^* as in (13.06) satisfies

$$\mathcal{R}_n^{\mathfrak{v}}[\hat{\theta}_n^{m_n^*} | \ell_2^{\mathfrak{a}, \mathfrak{r}}] \leq C R_n^*(\mathfrak{a}, \mathfrak{v})$$

with $C = 1 + \mathfrak{r}^2$. \square

§13.34 **Notation.** For $w_{\cdot} \in \ell_{\infty} \cap (\mathbb{R}_{\setminus 0})^{\mathbb{N}}$ we set $w_{(0)}^2 := \|w_{\cdot}^2\|_{\ell_{\infty}}$ and $w_{(j)}^2 = (w_{(j)}^2) := \|w_{\cdot}^2 \mathbf{1}_{\cdot}^{j+1}\|_{\ell_{\infty}}\|_{j \in \mathbb{N}}$, where by construction $w_{(j)}^2 = \sup \{w_i^2 : i \in \mathbb{N} \cap [j+1, \infty)\}$, $j \in \mathbb{N}_0$ and $w_{(\cdot)}^2 \in (\mathbb{R}_{\setminus 0}^+)^{\mathbb{N}}$. \square

§13.35 **Assumption.** Consider $\mathbf{v}_{\cdot}, \mathbf{a}_{\cdot} \in (\mathbb{R}_{\setminus 0})^{\mathbb{N}}$ such that $\mathbf{a}_{\cdot} \in \ell_{\infty}$ and $(\mathbf{a}\mathbf{v})_{\cdot} \in \ell_{\infty}$ (i.e. Assumption §11.12 is satisfied), and in addition $(\mathbf{a}\mathbf{v})_{(\cdot)}^2 \in (\mathbb{R}_{\setminus 0}^+)^{\mathbb{N}}$ and there exists $C_{(\mathbf{a}\mathbf{v})} \in (0, 1]$ such that for all $m \in \mathbb{N}$

$$(\mathbf{a}\mathbf{v})_{(m-1)}^2 \geq \min \{(\mathbf{a}\mathbf{v})_j^2 : j \in \llbracket m \rrbracket\} \geq C_{(\mathbf{a}\mathbf{v})} (\mathbf{a}\mathbf{v})_{(m-1)}^2$$

or in equal $C_{(\mathbf{a}\mathbf{v})} \|(\mathbf{a}\mathbf{v})_{\cdot}^{-2} \mathbf{1}_{\cdot}^m\|_{\ell_{\infty}} \leq (\mathbf{a}\mathbf{v})_{(m-1)}^{-2}$. \square

§13.36 **Comment.** Note that $(\mathbf{a}\mathbf{v})_{(\cdot)}^2 \in (\mathbb{R}_{\setminus 0}^+)^{\mathbb{N}}$ by definition, hence $(\mathbf{a}\mathbf{v})_{(\cdot)}^2 \in (\mathbb{R}_{\setminus 0}^+)^{\mathbb{N}}$ is satisfied if and only if $(\mathbf{a}\mathbf{v})_{(m)}^2 = o(1)$ as $m \rightarrow \infty$ (i.e. the maximal global approximation is consistent). Moreover if $(\mathbf{a}\mathbf{v})_{(\cdot)}^2 \in (\mathbb{R}_{\setminus 0}^+)^{\mathbb{N}}$ then we have trivially $(\mathbf{a}\mathbf{v})_{(\cdot)}^2 \in (\mathbb{R}_{\setminus 0}^+)^{\mathbb{N}}$ and $\|(\mathbf{a}\mathbf{v})_{\cdot}^{-2} \mathbf{1}_{\cdot}^m\|_{\ell_{\infty}} = (\mathbf{a}\mathbf{v})_m^{-2} = (\mathbf{a}\mathbf{v})_{(m-1)}^{-2}$ for all $m \in \mathbb{N}$, i.e. Assumption §13.35 is satisfied with $C_{(\mathbf{a}\mathbf{v})} = 1$.

For m_n^* and $\mathbb{R}_n^* := \mathbb{R}_n^{m_n^*}(\mathbf{a}_{\cdot}, \mathbf{v}_{\cdot})$ as in (13.06) we distinguish case i) : $\mathbb{R}_n^* = n^{-1} \|\mathbf{1}_{\cdot}^{m_n^*}\|_{\mathbf{v}}^2 > (\mathbf{a}\mathbf{v})_{(m_n^*)}^2$ and case ii) : $\mathbb{R}_n^* = (\mathbf{a}\mathbf{v})_{(m_n^*)}^2 \geq n^{-1} \|\mathbf{1}_{\cdot}^{m_n^*}\|_{\mathbf{v}}^2$. Consider case i) first. If $(\mathbf{a}\mathbf{v})_{(1)}^2 > n^{-1} \mathbf{v}_1^2$ then $\mathbb{R}_n^1(\mathbf{a}_{\cdot}, \mathbf{v}_{\cdot}) = n^{-1} \mathbf{v}_1^2 \vee (\mathbf{a}\mathbf{v})_{(1)}^2 = (\mathbf{a}\mathbf{v})_{(1)}^2$ and hence $m_n^* - 1 \in \mathbb{N}$. In this situation, from $n^{-1} \|\mathbf{1}_{\cdot}^{m_n^*-1}\|_{\mathbf{v}}^2 < n^{-1} \|\mathbf{1}_{\cdot}^{m_n^*}\|_{\mathbf{v}}^2$ (since $\mathbf{v}_{m_n^*}^2 \in \mathbb{R}_{\setminus 0}^+$), the definition (13.06) and

$$n^{-1} \|\mathbf{1}_{\cdot}^{m_n^*-1}\|_{\mathbf{v}}^2 \vee (\mathbf{a}\mathbf{v})_{(m_n^*-1)}^2 = \mathbb{R}_n^{m_n^*-1}(\mathbf{a}_{\cdot}, \mathbf{v}_{\cdot}) > \mathbb{R}_n^* = n^{-1} \|\mathbf{1}_{\cdot}^{m_n^*}\|_{\mathbf{v}}^2$$

it follows $\mathbb{R}_n^{m_n^*-1}(\mathbf{a}_{\cdot}, \mathbf{v}_{\cdot}) = (\mathbf{a}\mathbf{v})_{(m_n^*-1)}^2$ and hence $(\mathbf{a}\mathbf{v})_{(m_n^*-1)}^2 > n^{-1} \|\mathbf{1}_{\cdot}^{m_n^*}\|_{\mathbf{v}}^2$. Consider case ii). We set

$$m_n^{\diamond} := \min \{m \in \mathbb{N} \cap [m_n^* + 1, \infty) : n^{-1} \|\mathbf{1}_{\cdot}^m\|_{\mathbf{v}}^2 \geq (\mathbf{a}\mathbf{v})_{(m)}^2\} \quad (13.07)$$

where the defining set is not empty since $(\mathbf{a}\mathbf{v})_{(\cdot)}^2 \in (\mathbb{R}_{\setminus 0}^+)^{\mathbb{N}}$. We note that $(\mathbf{a}\mathbf{v})_{(m_n^*-1)}^2 = (\mathbf{a}\mathbf{v})_{(m_n^*)}^2$. Indeed, in the non trivial case $m_n^{\diamond} - 1 > m_n^*$ for each $m \in \llbracket m_n^* + 1, m_n^{\diamond} - 1 \rrbracket$ we have $\mathbb{R}_n^m(\mathbf{a}_{\cdot}, \mathbf{v}_{\cdot}) = (\mathbf{a}\mathbf{v})_{(m)}^2 > n^{-1} \|\mathbf{1}_{\cdot}^m\|_{\mathbf{v}}^2$, which together with $(\mathbf{a}\mathbf{v})_{(m_n^*)}^2 = \mathbb{R}_n^* \leq \mathbb{R}_n^m(\mathbf{a}_{\cdot}, \mathbf{v}_{\cdot}) = (\mathbf{a}\mathbf{v})_{(m)}^2 \leq (\mathbf{a}\mathbf{v})_{(m_n^*)}^2$ (since $(\mathbf{a}\mathbf{v})_{(\cdot)}^2 \in (\mathbb{R}_{\setminus 0}^+)^{\mathbb{N}}$) implies the equality $(\mathbf{a}\mathbf{v})_{(m)}^2 = (\mathbf{a}\mathbf{v})_{(m_n^*)}^2$ for all $m \in \llbracket m_n^* + 1, m_n^{\diamond} - 1 \rrbracket$. Moreover, from $(\mathbf{a}\mathbf{v})_{(m_n^*)}^2 = \mathbb{R}_n^* \leq \mathbb{R}_n^{m_n^{\diamond}}(\mathbf{a}_{\cdot}, \mathbf{v}_{\cdot}) = n^{-1} \|\mathbf{1}_{\cdot}^{m_n^{\diamond}}\|_{\mathbf{v}}^2$ it follows $(\mathbf{a}\mathbf{v})_{(m_n^*-1)}^2 = (\mathbf{a}\mathbf{v})_{(m_n^*)}^2 \leq n^{-1} \|\mathbf{1}_{\cdot}^{m_n^{\diamond}}\|_{\mathbf{v}}^2$. To summarise, assuming $(\mathbf{a}\mathbf{v})_{(1)}^2 > n^{-1} \mathbf{v}_1^2$ we use in the next proof the properties case i) $(\mathbf{a}\mathbf{v})_{(m_n^*-1)}^2 > n^{-1} \|\mathbf{1}_{\cdot}^{m_n^*}\|_{\mathbf{v}}^2$ and case ii) $(\mathbf{a}\mathbf{v})_{(m_n^*)}^2 = (\mathbf{a}\mathbf{v})_{(m_n^*-1)}^2 \leq n^{-1} \|\mathbf{1}_{\cdot}^{m_n^{\diamond}}\|_{\mathbf{v}}^2$. \square

§13.37 **Proposition** (GSSM §13.01 continued). Let $\widehat{\theta}_{\cdot} = \theta_{\cdot} + n^{-1/2} \dot{\mathbf{B}}_{\cdot} \sim N_{\mathbf{q}}^n$ as in Model §13.01 where $\dot{\mathbf{B}}_{\cdot} \sim N_{(0,1)}^{\otimes \mathbb{N}}$ and $\theta_{\cdot} = \mathbf{U}\theta \in \ell_2$. Given Assumption §13.35 and the notations in (13.06) for all $n \in \mathbb{N} \cap (\mathbf{v}_1^2(\mathbf{a}\mathbf{v})_{(1)}^{-2}, \infty)$ we have

$$\inf_{\widehat{\theta}_{\cdot}} \mathcal{R}_n^{\mathbf{v}}[\widehat{\theta}_{\cdot} | \ell_2^{\mathbf{a}, \mathbf{r}}] \geq 8^{-1} (1 \wedge 2C_{(\mathbf{a}\mathbf{v})} r^2) \mathbb{R}_n^*(\mathbf{a}_{\cdot}, \mathbf{v}_{\cdot}) \quad (13.08)$$

where the infimum is taken over all estimators $\widehat{\theta}_{\cdot}$.

§13.38 **Proof** of Proposition §13.37. is given in the lecture. \square

§13.39 **Illustration.** Consider the two specifications (o) and (s) depict in Table 02 [§12] of the Illustration §12.26. In both cases Assumption §13.35 is satisfied. Consequently, due to Proposition §13.37 the Table 02 [§12] presents the order of the *minimax rate* $\mathbb{R}_n^*(\mathbf{a}_{\cdot}, \mathbf{v}_{\cdot})$ which is attained by the *minimax-optimal* OPE $\widehat{\theta}_{\cdot}^{m_n^*} = \widehat{\theta}_{\cdot} \mathbf{1}_{\cdot}^{m_n^*} \in \ell_2(\mathbf{v}_{\cdot}^2)$ with optimally selected dimension m_n^* (Corollary §12.24). We shall stress, that the order of m_n^* given in the Table 02 [§12] depends on the parameter $a \in \mathbb{R}_{\setminus 0}^+$ characterising the (abstract) smoothness of the solution which is generally not known in advance. \square

§14 Data-driven estimation

§14|01 Data-driven estimation procedures

Considering a Hilbert space $\mathbb{J} = \mathbb{L}_2(\mathcal{J}, \mathcal{J}, \nu)$ and a surjective partial isometry $U \in \mathbb{L}(\mathbb{H}, \mathbb{J})$, which are fixed and presumed to be known in advance, we study data-driven estimation procedures in *statistical direct problems* as in Definition §10.19. Precisely, we consider the observable noisy version $\hat{\theta}_n = \theta + n^{-1/2}\varepsilon_n$ of the parameter $\theta = U\vartheta \in \mathbb{J}$ where the centred stochastic processes $\varepsilon_n = (\varepsilon_j)_{j \in \mathcal{J}}$ satisfies Assumption §10.04 and $n \in \mathbb{N}$ is a sample size. We denote by \mathbb{P}_θ^n the distribution of $\hat{\theta}_n$. Based on the noisy parameter $\hat{\theta}_n$ we consider the family $(\hat{\theta}_n^m = \hat{\theta}_n \mathbb{1}_n^m)_{m \in \mathbb{N}}$ of orthogonal projections estimators (OPE's) of θ defined in Definition §12.04. For each $m \in \mathbb{N}$ we shall measure the accuracy of the OPE $\hat{\theta}_n^m$ by its risk $\mathbb{E}_\theta^n(\mathfrak{d}_{\text{ist}}^2(\hat{\theta}_n^m, \theta))$ where $\mathfrak{d}_{\text{ist}}(\cdot, \cdot)$ is a certain semi metric such as a global \mathfrak{v} -error (Definition §12.09) or a local ϕ -error (Definition §12.30). Moreover, given $\theta = U\vartheta \in \mathbb{J}$ we consider the family of orthogonal projections (OP's) $(\theta_n^m = \theta \mathbb{1}_n^m \in \mathbb{J} \mathbb{1}_n^m)_{m \in \mathbb{N}}$ (Definition §11.08) where we tactically set $\theta_n^\infty := \theta$. Let us here assume that there exist $C \in \mathbb{R}_0^+$, and for each $n, m \in \mathbb{N}$, $\mathfrak{var}_{n,m}(\theta, \mathfrak{d}_{\text{ist}}) \in \mathbb{R}^+$ and $R_n^m(\theta, \mathfrak{d}_{\text{ist}}) = \mathfrak{d}_{\text{ist}}^2(\theta_n^m, \theta) + \mathfrak{var}_{n,m}(\theta, \mathfrak{d}_{\text{ist}})$ such that the risk of the estimator $\hat{\theta}_n^m$ satisfy

$$C^{-1}R_n^m(\theta, \mathfrak{d}_{\text{ist}}) \leq \mathbb{E}_\theta^n(\mathfrak{d}_{\text{ist}}^2(\hat{\theta}_n^m, \theta)) \leq C\{\mathfrak{d}_{\text{ist}}^2(\theta_n^m, \theta) + \mathfrak{var}_{n,m}(\theta, \mathfrak{d}_{\text{ist}})\} = CR_n^m(\theta, \mathfrak{d}_{\text{ist}}). \quad (14.01)$$

Minimising the right hand side in the last display as a function of $m \in \mathbb{N}$ leads to an optimal dimension (if it exists) and upper bound

$$m_n^\circ := m_n^\circ(\theta, \mathfrak{d}_{\text{ist}}) := \arg \min \{R_n^m(\theta, \mathfrak{d}_{\text{ist}}) = \mathfrak{d}_{\text{ist}}^2(\theta_n^m, \theta) + \mathfrak{var}_{n,m}(\theta, \mathfrak{d}_{\text{ist}}) : m \in \mathbb{N}\} \quad \text{and} \\ R_n^\circ(\theta, \mathfrak{d}_{\text{ist}}) := R_n^{m_n^\circ}(\theta, \mathfrak{d}_{\text{ist}}) = \min \{R_n^m(\theta, \mathfrak{d}_{\text{ist}}) : m \in \mathbb{N}\}. \quad (14.02)$$

Combining (14.01) and (14.02) (up to the constant C) we have that m_n° is an *oracle dimension*, $R_n^\circ(\theta, \mathfrak{d}_{\text{ist}})$ an *oracle bound* and the OPE $\hat{\theta}_n^{m_n^\circ}$ with oracle dimension m_n° is *oracle optimal*. However, the oracle dimension $m_n^\circ(\theta, \mathfrak{d}_{\text{ist}})$ given in (14.02) depends on the unknown parameter of interest θ , and thus also the oracle optimal statistic $\hat{\theta}_n^{m_n^\circ}$. In other words $\hat{\theta}_n^{m_n^\circ}$ is not a feasible estimator. We present in what follows two data-driven procedures to select a dimension \hat{m} within an admissible subset $[[M]] \subseteq \mathbb{N}$ of dimension parameters given by an integer $M \in \mathbb{N}$, which eventually leads to a feasible data-driven estimator $\hat{\theta}_n^{\hat{m}}$ depending on the observable quantities only. We call any data-driven estimator $\hat{\theta}_n$ *adaptive* for a class Θ of solutions if for all $\theta \in \Theta$ there is a constant $K_\theta \in \mathbb{R}_0^+$ possibly depending on θ such that $\mathbb{E}_\theta^n(\mathfrak{d}_{\text{ist}}^2(\hat{\theta}_n, \theta)) \leq K_\theta R_n^\circ(\theta, \mathfrak{d}_{\text{ist}})$ for all $n \in \mathbb{N}$. Each of those two different data-driven strategies involves in addition a sequence $\mathfrak{pen}_n = (\mathfrak{pen}_{n,m})_{m \in \mathbb{N}} \in (\mathbb{R}^+)^{\mathbb{N}}$ of penalties. Both, the upper bound M and the sequence of penalties, depend on the noise level n and possibly on the class Θ of solutions. However, for ease of presentation we omit the additional subscripts. We eventually show that there are constants $C_1, C_2 \in \mathbb{R}_0^+$ possibly depending on the solution $\theta \in \Theta$ and $(\mathfrak{bias}_m(\theta, \mathfrak{d}_{\text{ist}}))_{m \in \mathbb{N}}, (R_n^{\text{re}}(\theta, \mathfrak{d}_{\text{ist}}))_{n \in \mathbb{N}} \in (\mathbb{R}^+)^{\mathbb{N}}$ such that for all $n \in \mathbb{N}$ the risk of the data-driven estimator $\hat{\theta}_n^{\hat{m}}$ satisfies

$$\mathbb{E}_\theta^n(\mathfrak{d}_{\text{ist}}^2(\hat{\theta}_n^{\hat{m}}, \theta)) \leq C_1 \min \{\mathfrak{bias}_m^2(\theta, \mathfrak{d}_{\text{ist}}) + \mathfrak{pen}_m : m \in [[M]]\} + C_2 R_n^{\text{re}}(\theta, \mathfrak{d}_{\text{ist}}). \quad (14.03)$$

If in addition for all $\theta \in \Theta$ there is a constant $C_3 \in \mathbb{R}_0^+$ such that for all $n \in \mathbb{N}$ we have also

$$m_n^\circ \in [[M]], \quad \mathfrak{bias}_{m_n^\circ}(\theta, \mathfrak{d}_{\text{ist}}) \leq C_3 \mathfrak{d}_{\text{ist}}(\theta_{n,m_n^\circ}^{\text{re}}, \theta) \quad \text{and} \quad \mathfrak{pen}_{m_n^\circ} \leq C_3 \mathfrak{var}_{n,m_n^\circ}(\theta, \mathfrak{d}_{\text{ist}}). \quad (14.04)$$

Then due to (14.03) the data-driven estimator $\hat{\theta}_n^{\hat{m}}$ satisfies

$$\mathbb{E}_\theta^n(\mathfrak{d}_{\text{ist}}^2(\hat{\theta}_n^{\hat{m}}, \theta)) \leq C_1 C_3 \min \{\mathfrak{d}_{\text{ist}}^2(\theta_n^m, \theta) + \mathfrak{var}_{n,m}(\theta, \mathfrak{d}_{\text{ist}}) : m \in [[M]]\} + C_2 R_n^{\text{re}}(\theta, \mathfrak{d}_{\text{ist}}) \\ = C_1 C_3 R_n^\circ(\theta, \mathfrak{d}_{\text{ist}}) + C_2 R_n^{\text{re}}(\theta, \mathfrak{d}_{\text{ist}}).$$

and hence, if in addition $R_n^{\text{re}}(\theta, \mathfrak{d}_{\text{ist}}) \leq C_4 R_n^\circ(\theta, \mathfrak{d}_{\text{ist}})$ for a constant $C_4 \in \mathbb{R}_0^+$, then $\widehat{\theta}^m$ is *adaptive*. Indeed, we have $\mathbb{E}_q^n(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^m, \theta)) \leq KR_n^\circ(\theta, \mathfrak{d}_{\text{ist}})$ with $K := C_1 C_3 + C_2 C_4$.

§14.01 **Remark.** In order to establish a feasible method we have to select an upper bound M . Let us briefly describe heuristically the strategy we eventually apply. For each $\theta \in \Theta$ and $n, m \in \mathbb{N}$ let $\mathfrak{var}_{n,m}(\mathfrak{d}_{\text{ist}}) = \mathfrak{var}_{n,m}(\theta, \mathfrak{d}_{\text{ist}})$ do not depend on $\theta \in \Theta$ and moreover let $\mathfrak{var}_{n,\bullet}(\mathfrak{d}_{\text{ist}}) = (\mathfrak{var}_{n,m}(\mathfrak{d}_{\text{ist}}))_{m \in \mathbb{N}} \in (\mathbb{R}_0^+)^{\mathbb{N}}$ be known in advance. Furthermore, for all $n \in \mathbb{N}$ let $\mathfrak{var}_{n,1}(\mathfrak{d}_{\text{ist}}) \leq C_{\mathfrak{d}_{\text{ist}}}$ for some constant $C_{\mathfrak{d}_{\text{ist}}} \in \mathbb{R}_0^+$, which is evidently also known in advance. Consequently, the defining set of $M_n := \max \{m \in \mathbb{N} : \mathfrak{var}_{n,m}(\mathfrak{d}_{\text{ist}}) \leq C_{\mathfrak{d}_{\text{ist}}}\}$ is not empty and finite. For all $n \in \mathbb{N}$ with $R_n^\circ(\theta, \mathfrak{d}_{\text{ist}}) \leq C_{\mathfrak{d}_{\text{ist}}}$ follows then $m_n^\circ \in \llbracket M_n \rrbracket$ since $C_{\mathfrak{d}_{\text{ist}}} \geq R_n^\circ(\theta, \mathfrak{d}_{\text{ist}}) \geq \mathfrak{var}_{n,m_n^\circ}(\mathfrak{d}_{\text{ist}})$. In other words the feasible upper bound M_n satisfies the first condition in assumption (14.04). \square

§14.02 **GSSM (§10.21 continued).** Considering $\ell_2 = \mathbb{L}_2(\mathbb{N}, 2^{\mathbb{N}}, \nu_{\mathbb{N}})$ and a surjective partial isometry $\mathbf{U} \in \mathbb{L}(\mathbb{H}, \ell_2)$, which is fixed and presumed to be known in advance, we illustrate the different data-driven procedures in a Gaussian sequence space model §10.21. Here the observable stochastic process $\widehat{\theta} = \theta + n^{-1/2}\dot{B}$ is a noisy version of $\theta = \mathbf{U}\theta \in \ell_2$ and $\dot{B} \sim N_{(0,1)}^{\otimes \mathbb{N}}$. Consequently, $\widehat{\theta}$ admits a N_q^m -distribution belonging to the family $N_\Theta^m := (N_q^m)_{\theta \in \Theta}$. Summarising the observations satisfy a statistical product experiment $(\mathbb{R}^{\mathbb{N}}, \mathcal{B}^{\otimes \mathbb{N}}, N_\Theta^m)$ where $\Theta \subseteq \ell_2$. \square

§14|02 Model selection

Given a noisy version $\widehat{\theta} \sim \mathbb{E}_q^m$ in a *statistical direct problem* as in Definition §10.19. and a collection of admissible models $\llbracket M \rrbracket$ for some $M \in \mathbb{N}$ we seek to minimise the global \mathfrak{v} -risk within the family $(\widehat{\theta}^m = \widehat{\theta} \mathbb{1}_*^m)_{m \in \llbracket M \rrbracket}$ of OPE's defined in Definition §12.04. Here and subsequently, let Assumption §12.07, $\mathfrak{v} \in \mathcal{J}_0$, $\theta \in \mathbb{L}_2(\mathfrak{v}^2\nu)$ and $\mathbb{1}_*^m \in \mathbb{L}_2(\mathfrak{v}^2\nu)$ for all $m \in \mathbb{N}$ be satisfied.

§14.03 **Reminder.** For all $n, m \in \mathbb{N}$ we set

$$R_n^m(\theta, \mathfrak{v}) := \|\theta \mathbb{1}_*^{m+1}\|_{\mathfrak{v}}^2 + n^{-1} \|\mathbb{1}_*^m\|_{\mathfrak{v}}^2, \quad m_n^\circ := \arg \min \{R_n^m(\theta, \mathfrak{v}) : m \in \mathbb{N}\}$$

$$\text{and } R_n^\circ(\theta, \mathfrak{v}) := R_n^{m_n^\circ}(\theta, \mathfrak{v}) = \min \{R_n^m(\theta, \mathfrak{v}) : m \in \mathbb{N}\} \quad (14.05)$$

Since $\theta^m = \theta \mathbb{1}_*^m \in \mathbb{L}_2(\mathfrak{v}^2\nu)$ (Property §11.09) and $\widehat{\theta}^m = \widehat{\theta} \mathbb{1}_*^m \in \mathbb{L}_2(\mathfrak{v}^2\nu)$ \mathbb{E}_q^n -a.s. (Comment §12.08), for each $a \in \mathbb{L}_2(\mathfrak{v}^2\nu) \mathbb{1}_*^m$ applying the Cauchy-Schwarz inequality we have

$$\mathfrak{v}^2\nu(|a, \widehat{\theta}|) = \mathfrak{v}^2\nu(|a, \widehat{\theta} \mathbb{1}_*^m|) \leq \|a\|_{\mathfrak{v}} \|\widehat{\theta}^m\|_{\mathfrak{v}} \in \mathbb{R}^+,$$

and hence $a, \widehat{\theta} \in \mathbb{L}_1(\mathfrak{v}^2\nu)$ \mathbb{E}_q^n -a.s.. \square

The first selection method is inspired by the work of Barron et al. [1999] and for an extensive overview of model selection by penalised contrast, the reader may refer to Massart [2007]. Let us introduce a contrast function

$$\Upsilon : \mathbb{L}_2(\mathfrak{v}^2\nu) \supseteq \bigcup_{m \in \mathbb{N}} \mathbb{L}_2(\mathfrak{v}^2\nu) \mathbb{1}_*^m \rightarrow \mathbb{R} \text{ with } a \mapsto \Upsilon(a) := \|a\|_{\mathfrak{v}}^2 - 2\mathfrak{v}^2\nu(a, \widehat{\theta}) \quad (14.06)$$

where for each $m \in \mathbb{N}$ the OPE $\widehat{\theta}^m = \widehat{\theta} \mathbb{1}_*^m$ and $a \in \mathbb{L}_2(\mathfrak{v}^2\nu) \mathbb{1}_*^m \subseteq \mathbb{L}_2(\mathfrak{v}^2\nu)$ satisfy

$$\Upsilon(a) = \|a\|_{\mathfrak{v}}^2 - 2\mathfrak{v}^2\nu(a, \widehat{\theta}) = \|a\|_{\mathfrak{v}}^2 - 2\langle a, \widehat{\theta}^m \rangle_{\mathfrak{v}} = \|a - \widehat{\theta}^m\|_{\mathfrak{v}}^2 - \|\widehat{\theta}^m\|_{\mathfrak{v}}^2.$$

Consequently, for each $m \in \mathbb{N}$ the OPE $\widehat{\theta}^m = \widehat{\theta} \mathbb{1}_*^m$ minimises the contrast function, that is,

$$-\|\widehat{\theta}^m\|_{\mathfrak{v}}^2 = \Upsilon(\widehat{\theta}^m) = \inf \{ \Upsilon(a) : a \in \mathbb{L}_2(\mathfrak{v}^2\nu) \mathbb{1}_*^m \}. \quad (14.07)$$

Given an upper bound $M \in \mathbb{N}$ and penalties $\mathfrak{pen}_\bullet = (\mathfrak{pen}_m)_{m \in \mathbb{N}} \in (\mathbb{R}^+)^{\mathbb{N}}$ we select a dimension among the collection of admissible values $\llbracket M \rrbracket$ as minimiser of a penalised contrast criterion, that is

$$\widehat{m} := \arg \min \{ \Upsilon(\widehat{\theta}^m) + \mathfrak{pen}_m : m \in \llbracket M \rrbracket \}. \quad (14.08)$$

The data-driven estimator of θ is now given by $\widehat{\theta}^{\widehat{m}}$ and below we derive an upper bound for its global \mathfrak{v} -risk $\mathbb{E}_q^n (\|\widehat{\theta}^{\widehat{m}} - \theta\|_{\mathfrak{v}}^2)$. The construction of the penalty sequence \mathfrak{pen}_\bullet and the upper bound M given below is motivated by the following key arguments used in the proof of the risk bound which we present first. Moreover, both \mathfrak{pen}_\bullet and M will depend, among others, on the noise level n , however, for sake of simplicity we will omit an additional subscript. The key argument for our reasoning is the next assertion. For $a \in \mathbb{R}$ we write $(a)_+ := a \vee 0$ shortly.

§14.04 **Lemma (key argument).** *If $\mathfrak{pen}_\bullet \in (\mathbb{R}^+)^{\mathbb{N}}$ then for all $M \in \mathbb{N}$ and $m \in \llbracket M \rrbracket$ we have*

$$\|\widehat{\theta}^m - \theta\|_{\mathfrak{v}}^2 \leq 3\|\theta^m - \theta\|_{\mathfrak{v}}^2 + 4\mathfrak{pen}_m + 8 \max \{ (\|\widehat{\theta}^j - \theta^j\|_{\mathfrak{v}}^2 - \mathfrak{pen}_j/4)_+ : j \in \llbracket m, M \rrbracket \}.$$

§14.05 **Proof of Lemma §14.04.** is given in the lecture. □

Similar to m_n° as in (14.05), which realises by construction a statistical-error-squared-bias compromise, let us fix a dimension $m^\circ \in \llbracket M \rrbracket$ to be specified below. Due to the last assertion for each $\theta \in \Theta$ we have

$$\begin{aligned} \mathbb{E}_\theta^n (\|\widehat{\theta}^m - \theta\|_{\mathfrak{v}}^2) &\leq 3\|\theta^{m^\circ} - \theta\|_{\mathfrak{v}}^2 + 4\mathfrak{pen}_{m^\circ} \\ &\quad + 8\mathbb{E}_\theta^n \left(\max \{ (\|\widehat{\theta}^j - \theta^j\|_{\mathfrak{v}}^2 - \mathfrak{pen}_j/4)_+ : j \in \llbracket m^\circ, M \rrbracket \} \right) \end{aligned} \quad (14.09)$$

Keeping in mind that $m^\circ \in \llbracket M \rrbracket$ in contrast to $m_n^\circ \in \mathbb{N}$ eventually realises an optimal statistical-error-squared-bias trade-off among the collection of admissible values $\llbracket M \rrbracket$ rather than \mathbb{N} , we wish the upper bound M to be as large as possible. In contrast, in order to control the remainder term, the last term in (14.09), we are eventually forced to use a rather small upper bound M . However, we bound the remainder term by imposing the following assumption, which though holds true for a wide range of solutions $\theta \in \Theta$ under reasonable model assumptions.

§14.06 **Assumption.** There exists a constant $C := C(\theta) \in \mathbb{R}_{>0}^+$ possibly depending on the parameter $\theta \in \Theta$ and $(R_n^{\text{re}}(\theta, \mathfrak{v}))_{n \in \mathbb{N}} \in (\mathbb{R}^+)^{\mathbb{N}}$ such that for each $n \in \mathbb{N}$ the upper bound $M \in \mathbb{N}$ and $m^\circ \in \llbracket M \rrbracket$ satisfy

$$\mathbb{E}_\theta^n \left(\max \{ (\|\widehat{\theta}^j - \theta^j\|_{\mathfrak{v}}^2 - \mathfrak{pen}_j/4)_+ : j \in \llbracket m^\circ, M \rrbracket \} \right) \leq C R_n^{\text{re}}(\theta, \mathfrak{v}).$$

The next assertion provides an upper bound for the \mathfrak{v} -risk of the estimator $\widehat{\theta}^{\widehat{m}}$ with data-driven choice \widehat{m} given by (14.08).

§14.07 **Proposition.** *Let $m^\circ \in \llbracket M \rrbracket$ satisfy the Assumption §14.06 then we have*

$$\mathbb{E}_\theta^n (\|\widehat{\theta}^{\widehat{m}} - \theta\|_{\mathfrak{v}}^2) \leq 3\|\theta^{m^\circ} - \theta\|_{\mathfrak{v}}^2 + 4\mathfrak{pen}_{m^\circ} + 8 C R_n^{\text{re}}(\theta, \mathfrak{v}).$$

§14.08 **Proof of Proposition §14.07.** is given in the lecture. □

§14.09 **Corollary.** *If $m^\circ = \arg \min \{ \|\theta^m - \theta\|_{\mathfrak{v}}^2 + \mathfrak{pen}_m : m \in \llbracket M \rrbracket \}$ satisfies Assumption §14.06, then*

$$\mathbb{E}_\theta^n (\|\widehat{\theta}^{\widehat{m}} - \theta\|_{\mathfrak{v}}^2) \leq 4 \min \{ \|\theta^m - \theta\|_{\mathfrak{v}}^2 + \mathfrak{pen}_m : m \in \llbracket M \rrbracket \} + 8 C R_n^{\text{re}}(\theta, \mathfrak{v}).$$

§14.10 **Proof** of **Corollary** §14.09. is given in the lecture. \square

§14.11 **Remark**. Considering the \mathfrak{v} -risk bound of the estimator $\widehat{\theta}^{m^\circ}$ with dimension parameter m° the first rhs. term in the upper risk-bound given in **Proposition** §14.07 is strongly reminiscent of the variance-squared-bias upper bound $R_n^{m^\circ}(\theta, \mathfrak{v}) = \|\theta^{m^\circ} - \theta\|_{\mathfrak{v}}^2 + \mathfrak{v}_{\text{ar}_{n,m^\circ}}(\theta, \mathfrak{v})$ as given in (14.02). Indeed, in many cases the penalty term $\mathfrak{p}_{\text{en}_{m^\circ}}$ is in the same order as the statistical error $\mathfrak{v}_{\text{ar}_{n,m^\circ}}(\theta, \mathfrak{v})$. Consequently, provided the remainder term $R_n^{\text{re}}(\theta, \mathfrak{v})$ is negligible compared to $R_n^{m^\circ}(\theta, \mathfrak{v})$, the upper risk bound of the data-driven estimator is given by $R_n^{m^\circ}(\theta, \mathfrak{v})$ (up to a constant). Moreover, since m° realises an optimal trade-off between squared-bias and statistical error among the admissible values M , in many cases $R_n^{m^\circ}(\theta, \mathfrak{v})$ is of optimal oracle order $R_n^\circ(\theta, \mathfrak{v})$. \square

We eventually are in a situation where the sequence of penalties $\mathfrak{p}_{\text{en}} \in (\mathbb{R}^+)^{\mathbb{N}}$ satisfying the Assumption §14.06 still depends on characteristics of the unknown parameter θ and thus it is only partially known in advance. Assuming a sequence of estimators $\widehat{\mathfrak{p}}_{\text{en}} \in (\mathbb{R}^+)^{\mathbb{N}}$ we select similar to (14.08) the dimension

$$\widehat{m} := \arg \min \{ \Upsilon(\widehat{\theta}^m) + \widehat{\mathfrak{p}}_{\text{en}_m} : m \in \llbracket M \rrbracket \}. \quad (14.10)$$

The data-driven estimator of θ is now given by $\widehat{\theta}^{\widehat{m}}$ and below we derive an upper bound for its global \mathfrak{v} -risk $\mathbb{E}_\theta^n(\|\widehat{\theta}^{\widehat{m}} - \theta\|_{\mathfrak{v}}^2)$. The key argument for our reasoning is the next assertion. Its proof follows along the lines of the **Proof** §14.05.

§14.12 **Lemma (key argument)**. If $\widehat{\mathfrak{p}}_{\text{en}}, \mathfrak{p}_{\text{en}} \in (\mathbb{R}^+)^{\mathbb{N}}$ then for all $M \in \mathbb{N}$ and $m \in \llbracket M \rrbracket$ we have

$$\begin{aligned} \|\widehat{\theta}^{\widehat{m}} - \theta\|_{\mathfrak{v}}^2 &\leq 3\|\theta \mathbf{1}^{m^\perp}\|_{\mathfrak{v}}^2 + 2\mathfrak{p}_{\text{en}_m} + 2\widehat{\mathfrak{p}}_{\text{en}_m} + 2(\mathfrak{p}_{\text{en}_{\widehat{m}}} - \widehat{\mathfrak{p}}_{\text{en}_{\widehat{m}}})_+ \\ &\quad + 8 \max \{ (\|\widehat{\theta}^j - \theta^j\|_{\mathfrak{v}}^2 - \mathfrak{p}_{\text{en}_j}/4)_+ : j \in \llbracket m, M \rrbracket \}. \end{aligned}$$

§14.13 **Proof** of **Lemma** §14.12. is given in the lecture. \square

Similar to m_n° as in (14.02), which realises by construction a statistical-error-squared-bias compromise, let us fix a dimension $m^\circ \in \llbracket M \rrbracket$ to be specified below (analogously to (14.09)). Due to the last assertion for each $\theta \in \Theta$ we have

$$\begin{aligned} \mathbb{E}_\theta^n(\|\widehat{\theta}^{\widehat{m}} - \theta\|_{\mathfrak{v}}^2) &\leq 3\|\theta \mathbf{1}^{m^\perp}\|_{\mathfrak{v}}^2 + 2\mathfrak{p}_{\text{en}_m} + 8\mathbb{E}_\theta^n(\max \{ (\|\widehat{\theta}^j - \theta^j\|_{\mathfrak{v}}^2 - \mathfrak{p}_{\text{en}_j}/4)_+ : j \in \llbracket m, M \rrbracket \}) \\ &\quad + 2\mathbb{E}_\theta^n(\widehat{\mathfrak{p}}_{\text{en}_m}) + 2\mathbb{E}_\theta^n((\mathfrak{p}_{\text{en}_{\widehat{m}}} - \widehat{\mathfrak{p}}_{\text{en}_{\widehat{m}}})_+). \quad (14.11) \end{aligned}$$

We bound the first remainder term by imposing Assumption §14.06, which though hold true for a wide range of solutions $\theta = U\theta \in \Theta$ under reasonable model assumptions.

§14.14 **Proposition**. If $m^\circ \in \llbracket M \rrbracket$ satisfies the Assumption §14.06 then we have

$$\mathbb{E}_\theta^n(\|\widehat{\theta}^{\widehat{m}} - \theta\|_{\mathfrak{v}}^2) \leq 3\|\theta \mathbf{1}^{m^\perp}\|_{\mathfrak{v}}^2 + 2\mathfrak{p}_{\text{en}_m} + 8C R_n^{\text{re}}(\theta, \mathfrak{v}) + 2\mathbb{E}_\theta^n(\widehat{\mathfrak{p}}_{\text{en}_m}) + 2\mathbb{E}_\theta^n((\mathfrak{p}_{\text{en}_{\widehat{m}}} - \widehat{\mathfrak{p}}_{\text{en}_{\widehat{m}}})_+).$$

§14.15 **Proof** of **Proposition** §14.14. is given in the lecture. \square

§14.16 **Corollary**. If $m^\circ = \arg \min \{ \|\theta \mathbf{1}^{m^\perp}\|_{\mathfrak{v}}^2 + \mathfrak{p}_{\text{en}_m} : m \in \llbracket M \rrbracket \}$ satisfies Assumption §14.06 with constant $C \in [1, \infty)$, $\mathbb{E}_\theta^n(\widehat{\mathfrak{p}}_{\text{en}_{m^\circ}}) \leq K_1 \mathfrak{p}_{\text{en}_{m^\circ}}$ and $\mathbb{E}_\theta^n((\mathfrak{p}_{\text{en}_{\widehat{m}}} - \widehat{\mathfrak{p}}_{\text{en}_{\widehat{m}}})_+) \leq K_2 R_n^{\text{re}}(\theta)$ for some $K_1, K_2 \in [1, \infty)$, then

$$\mathbb{E}_\theta^n(\|\widehat{\theta}^{\widehat{m}} - \theta\|_{\mathfrak{v}}^2) \leq 4K_1 \min \{ \|\theta \mathbf{1}^{m^\perp}\|_{\mathfrak{v}}^2 + \mathfrak{p}_{\text{en}_m} : m \in \llbracket M \rrbracket \} + (2K_2 + 8C) R_n^{\text{re}}(\theta).$$

§14.17 **Proof** of **Corollary** §14.16. is given in the lecture. \square

§14|03 GSSM: data-driven global estimation

Let us first state some elementary inequalities for Gaussian random variables. There exist several results for tail bounds of sums of independent squared Gaussian random variables and we present next a version which is due to Birgé [2001] and the formulation (14.12) can be found in Lemma 1 in Laurent and Massart [2000].

§14.18 **Lemma.** Let $a_* \in (\mathbb{R})^{\mathbb{N}}$ and $\dot{B}_* = (\dot{B}_j)_{j \in \mathbb{N}} \sim N_{(0,1)}^{\otimes \mathbb{N}}$. For all $\eta \in \mathbb{R}^+$ and $m \in \mathbb{N}$ we have

$$N_{(0,1)}^{\otimes \mathbb{N}}(\|a_* \dot{B}_* \mathbf{1}_*^m\|_{\ell_2}^2 - \|a_* \mathbf{1}_*^m\|_{\ell_2}^2 \geq 2\|a_*^2 \mathbf{1}_*^m\|_{\ell_2} \sqrt{\eta} + 2\|a_*^2 \mathbf{1}_*^m\|_{\ell_\infty} \eta) \leq \exp(-\eta). \quad (14.12)$$

which for all $\zeta \in \mathbb{R}^+$ setting $\eta := \zeta(\zeta \wedge 1)\|a_* \mathbf{1}_*^m\|_{\ell_2}^2 / (4\|a_*^2 \mathbf{1}_*^m\|_{\ell_\infty}) \in \mathbb{R}^+$ implies

$$N_{(0,1)}^{\otimes \mathbb{N}}(\|a_* \dot{B}_* \mathbf{1}_*^m\|_{\ell_2}^2 \geq (1 + 3\zeta/2)\|a_* \mathbf{1}_*^m\|_{\ell_2}^2) \leq \exp(-\eta). \quad (14.13)$$

Moreover, for any $\xi \in [1, \infty)$ we have

$$\begin{aligned} N_{(0,1)}^{\otimes \mathbb{N}}(\|a_* \dot{B}_* \mathbf{1}_*^m\|_{\ell_2}^2 - (1 + 3\xi/2)\|a_* \mathbf{1}_*^m\|_{\ell_2}^2) \\ \leq 6\|a_*^2 \mathbf{1}_*^m\|_{\ell_\infty} \exp(-(\xi/4)\|a_* \mathbf{1}_*^m\|_{\ell_2}^2 \|a_*^2 \mathbf{1}_*^m\|_{\ell_\infty}^{-1}). \end{aligned} \quad (14.14)$$

§14.19 **Proof of Lemma §14.18.** Exercise. □

§14|03|01 Global \mathfrak{v} -risk

§14.20 **Reminder (Global oracle \mathfrak{v} -risk in GSSM §14.02).** Given Model §14.02 we consider an OPE as in Section §12. Here the observable noisy version $\hat{\theta}$ admits a N_q^n -distribution belonging to the family $N_\Theta^n := (N_q^n)_{\theta \in \Theta}$, $\Theta \subseteq \ell_2$. Let us recall (12.04) in Proposition §12.12 where for $\mathfrak{v} \in (\mathbb{R}_{>0})^{\mathbb{N}}$, $\theta \in \ell_2(\mathfrak{v}^2)$ and $n, m \in \mathbb{N}$ we have defined

$$\begin{aligned} R_n^m(\theta, \mathfrak{v}) &:= \|\theta \mathbf{1}_*^{m \perp}\|_{\mathfrak{v}}^2 + n^{-1}\|\mathbf{1}_*^m\|_{\mathfrak{v}}^2, \quad m_n^\circ := \arg \min \{R_n^m(\theta, \mathfrak{v}) : m \in \mathbb{N}\} \\ \text{and } R_n^\circ(\theta, \mathfrak{v}) &:= R_n^{m_n^\circ}(\theta, \mathfrak{v}) = \min \{R_n^m(\theta, \mathfrak{v}) : m \in \mathbb{N}\}. \end{aligned} \quad (14.15)$$

Due to Corollary §12.17 the (infeasible) OPE $\hat{\theta}^{m_n^\circ} = \hat{\theta} \mathbf{1}_*^{m_n^\circ} \in \ell_2(\mathfrak{v}^2)$ with oracle dimension m_n° as in (14.15) satisfies

$$N_q^n(\|\hat{\theta}^{m_n^\circ} - \theta\|_{\mathfrak{v}}^2) = R_n^\circ(\theta, \mathfrak{v}) = \inf_{m \in \mathbb{N}} N_q^n(\|\hat{\theta}^m - \theta\|_{\mathfrak{v}}^2),$$

and hence it is *oracle optimal* (with constant 1). □

§14.21 **Assumption.** Let $\mathfrak{v} \in (\mathbb{R}_{>0})^{\mathbb{N}}$ satisfy

$$C_{\mathfrak{v}} := \sum_{m \in \mathbb{N}} 4\|\mathfrak{v}^2 \mathbf{1}_*^m\|_{\ell_\infty} \exp(-\|\mathfrak{v} \mathbf{1}_*^m\|_{\ell_2}^2 / (4\|\mathfrak{v}^2 \mathbf{1}_*^m\|_{\ell_\infty})) \in \mathbb{R}^+. \quad \square$$

§14.22 **Comment.** Since $\|\mathfrak{v}^2 \mathbf{1}_*^m\|_{\ell_\infty} \in \mathbb{R}_0^+$ and $\|\mathfrak{v}^2 \mathbf{1}_*^m\|_{\ell_\infty} \geq \mathfrak{v}_1^2$ for all $m \in \mathbb{N}$, the Assumption §14.21 implies $\|\mathfrak{v}^2 \mathbf{1}_*^m\|_{\ell_\infty} \|\mathfrak{v} \mathbf{1}_*^m\|_{\ell_2}^{-2} = o(1)$ and $\|\mathfrak{v} \mathbf{1}_*^m\|_{\ell_2}^{-2} = o(1)$ as $m \rightarrow \infty$. □

§14.23 **Illustration.** Consider $\mathfrak{v}^2 = (j^a)_{j \in \mathbb{N}}$ for $a \in \mathbb{R}_0^+$. Then $\|\mathfrak{v}^2 \mathbf{1}_*^m\|_{\ell_\infty} = m^a$ and $\|\mathfrak{v} \mathbf{1}_*^m\|_{\ell_2}^2 \simeq m^{a+1}$. Consequently, we have

$$\sum_{m \in \mathbb{N}} 4\|\mathfrak{v}^2 \mathbf{1}_*^m\|_{\ell_\infty} \exp(-\|\mathfrak{v} \mathbf{1}_*^m\|_{\ell_2}^2 / (4\|\mathfrak{v}^2 \mathbf{1}_*^m\|_{\ell_\infty})) \simeq \sum_{m \in \mathbb{N}} m^a \exp(-m) \in \mathbb{R}_0^+.$$

and the assumption Assumption §14.21 is satisfied. On the contrary if $\mathbf{v}^2 = (\exp(j^a))_{j \in \mathbb{N}}$ for $a > 1$, then $\|\mathbf{v}^2 \mathbb{1}^m\|_{\ell_\infty} = \exp(m^a)$ and $\|\mathbf{v} \mathbb{1}^m\|_{\ell_2}^2 \simeq \exp(m^a)$. Consequently, we have

$$\sum_{m \in \mathbb{N}} 4 \|\mathbf{v}^2 \mathbb{1}^m\|_{\ell_\infty} \exp(-\|\mathbf{v} \mathbb{1}^m\|_{\ell_2}^2 / (4 \|\mathbf{v}^2 \mathbb{1}^m\|_{\ell_\infty})) \simeq \sum_{m \in \mathbb{N}} \exp(m^a) \simeq \infty$$

and the assumption Assumption §14.21 is not satisfied. \square

§14.24 **Corollary.** Under Assumption §14.21 we have

$$N_q^n \left(\max \left\{ \left(\|\widehat{\theta}^j - \theta^j\|_{\mathbf{v}}^2 - \frac{5}{2} \|\mathbb{1}^j\|_{\mathbf{v}}^2 n^{-1} \right)_+ : j \in \llbracket M \rrbracket \right\} \right) \leq (3/2) C_v n^{-1} \quad \text{for all } M \in \mathbb{N}. \quad (14.16)$$

§14.25 **Proof of Corollary §14.24.** is given in the lecture. \square

§14.26 **Notation.** Consider a sequence of penalties $\text{pen}_m^{\mathbf{v}} = (\text{pen}_m^{\mathbf{v}})_{m \in \mathbb{N}} \in (\mathbb{R}_0^+)^{\mathbb{N}}$ given by

$$\text{pen}_m^{\mathbf{v}} := 10 n^{-1} \|\mathbb{1}^m\|_{\mathbf{v}}^2, \quad \text{for each } m \in \mathbb{N} \quad (14.17)$$

which is obviously known in advance. Considering the data-driven OSE $\widehat{\theta}^{\widehat{m}} = \widehat{\theta} \mathbb{1}^{\widehat{m}}$ with dimension parameter \widehat{m} selected as in (14.08) with penalty sequence $\text{pen}_m^{\mathbf{v}}$ given in (14.17) and arbitrary but fixed upper bound $M \in \mathbb{N}$ we derive below an upper bound for its global \mathbf{v} -risk, $N_q^n(\|\widehat{\theta}^{\widehat{m}} - \theta\|_{\mathbf{v}}^2)$. \square

§14.27 **Proposition (GSSM (§14.02 continued)).** Let $\widehat{\theta} = \theta + n^{-1/2} \dot{B} \sim N_q^n$ as in Model §14.02 where $\theta \in \ell_2$ and $\dot{B} \sim N_{(0,1)}^{\otimes \mathbb{N}}$. Given $\mathbf{v} \in (\mathbb{R}_{>0})^{\mathbb{N}}$, $M \in \mathbb{N}$ and $\text{pen}_m^{\mathbf{v}}$ as in (14.17) consider a data-driven OPE $\widehat{\theta}^{\widehat{m}} = \widehat{\theta} \mathbb{1}^{\widehat{m}} \in \ell_2 \mathbb{1}^{\widehat{m}} \subseteq \ell_2(\mathbf{v}^2)$ of $\theta \in \ell_2(\mathbf{v}^2)$ with

$$\widehat{m} := \arg \min \left\{ -\|\widehat{\theta}^m\|_{\mathbf{v}} + \text{pen}_m^{\mathbf{v}} : m \in \llbracket M \rrbracket \right\}. \quad (14.18)$$

If Assumption §14.21 is satisfied with $C_v \in \mathbb{R}_0^+$, then for all $n, M \in \mathbb{N}$ we have

$$N_q^n \left(\|\widehat{\theta}^{\widehat{m}} - \theta\|_{\mathbf{v}}^2 \right) \leq 40 \min \left\{ R_n^m(\theta, \mathbf{v}) : m \in \llbracket M \rrbracket \right\} + 12 C_v n^{-1} \quad (14.19)$$

where $R_n^m(\theta, \mathbf{v}) := \|\theta^m - \theta\|_{\mathbf{v}}^2 + n^{-1} \|\mathbb{1}^m\|_{\mathbf{v}}^2$ is defined as in (14.15).

§14.28 **Proof of Proposition §14.27.** is given in the lecture. \square

§14.29 **Comment.** The oracle bound $R_n^{\circ}(\theta, \mathbf{v}) = R_n^{m_n^{\circ}}(\theta, \mathbf{v}) = \min \left\{ R_n^m(\theta, \mathbf{v}) : m \in \mathbb{N} \right\}$ (for details see **Reminder §14.20**) satisfies $n R_n^{\circ}(\theta, \mathbf{v}) \geq \|\mathbb{1}^{m_n^{\circ}}\|_{\mathbf{v}}^2 \geq \mathbf{v}_1^2$. Consequently, the last upper bound in (14.19) and the oracle bound $R_n^{\circ}(\theta, \mathbf{v})$ coincide up to a constant $(40 + 12 C_v \mathbf{v}_1^{-2})$ provided the oracle dimension fulfils $m_n^{\circ} \in \llbracket M \rrbracket$. Therefore, we wish the upper bound M to be as large as possible. The next assertion shows that

$$M^{\mathbf{p}} := \max \left\{ m \in \mathbb{N} : \|\mathbb{1}^m\|_{\mathbf{v}}^2 \leq n \mathbf{v}_1^2 \right\} \in \mathbb{N} \quad (14.20)$$

is a suitable choice for the upper bound, where the defining set is not empty and finite since $\|\mathbb{1}^m\|_{\mathbf{v}}^{-2} = o(1)$ as $m \rightarrow \infty$. \square

§14.30 **Corollary (GSSM (§14.02 continued)).** Given $\mathbf{v} \in (\mathbb{R}_{>0})^{\mathbb{N}}$, $M^{\mathbf{p}} \in \mathbb{N}$ as in (14.20) and $\text{pen}_m^{\mathbf{v}}$ as in (14.17) consider a data-driven OPE $\widehat{\theta}^{\widehat{m}} = \widehat{\theta} \mathbb{1}^{\widehat{m}} \in \ell_2 \mathbb{1}^{\widehat{m}} \subseteq \ell_2(\mathbf{v}^2)$ with

$$\widehat{m} := \arg \min \left\{ -\|\widehat{\theta}^m\|_{\mathbf{v}} + \text{pen}_m^{\mathbf{v}} : m \in \llbracket M^{\mathbf{p}} \rrbracket \right\}. \quad (14.21)$$

Under the assumptions of **Proposition** §14.27 for each $n \in \mathbb{N}$ such that $R_n^\circ(\theta, \mathbf{v}) \leq \mathbf{v}_1^2$ we have

$$N_q^n(\|\widehat{\theta}^m - \theta\|_{\mathbf{v}}^2) \leq 40 R_n^\circ(\theta, \mathbf{v}) + 12 C_{\mathbf{v}} n^{-1} \leq C R_n^\circ(\theta, \mathbf{v}) \quad (14.22)$$

and, hence up to the constant $C := 40 + 12 C_{\mathbf{v}} \mathbf{v}_1^{-2}$ the feasible data-driven estimator $\widehat{\theta}^m$ is **oracle optimal**.

§14.31 **Proof** of **Corollary** §14.30. is given in the lecture. \square

§14.32 **Remark**. If Assumption §14.21 is not satisfied (see **Illustration** §14.23), then we can't make use of **Corollary** §14.24. In this situation let $\mathbf{v} \in (\mathbb{R}_{>0})^{\mathbb{N}}$ and $\delta \in ([1, \infty))^{\mathbb{N}}$ satisfy

$$C_{\mathbf{v}, \delta} := \sum_{m \in \mathbb{N}} 4 \|\mathbf{v} \cdot \mathbf{1}^m\|_{\ell_\infty} \exp\left(-\delta_m \|\mathbf{v} \cdot \mathbf{1}^m\|_{\ell_2}^2 / (4 \|\mathbf{v} \cdot \mathbf{1}^m\|_{\ell_\infty})\right) \in \mathbb{R}^+. \quad (14.23)$$

Consider a sequence of penalties $\mathbf{pen}_m^{\mathbf{v}, \delta} = (\mathbf{pen}_m^{\mathbf{v}, \delta})_{m \in \mathbb{N}} \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ given by

$$\mathbf{pen}_m^{\mathbf{v}, \delta} := 4(1 + 3\delta_m/2) n^{-1} \|\mathbf{1}^m\|_{\mathbf{v}}^2, \quad \text{for each } m \in \mathbb{N} \quad (14.24)$$

which is obviously known in advance. Similar to **Corollary** §14.24 due to (14.23) for each $M \in \mathbb{N}$ we obtain

$$\begin{aligned} n N_q^n\left(\max\left\{\|\widehat{\theta}^j - \theta^j\|_{\mathbf{v}}^2 - (1 + 3\delta_j/2) \|\mathbf{1}^j\|_{\mathbf{v}}^2 n^{-1}\right\}; j \in \llbracket M \rrbracket\right) \\ \leq \sum_{m \in \llbracket M \rrbracket} 6 \|\mathbf{v} \cdot \mathbf{1}^m\|_{\ell_\infty} \exp\left(-\delta_m \|\mathbf{v} \cdot \mathbf{1}^m\|_{\ell_2}^2 / (4 \|\mathbf{v} \cdot \mathbf{1}^m\|_{\ell_\infty})\right) = (3/2) C_{\mathbf{v}, \delta}. \end{aligned}$$

Thus, the sequence of penalties $\mathbf{pen}_m^{\mathbf{v}, \delta} \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ given in (14.24) satisfies the Assumption §14.06 with $C = (3/2) C_{\mathbf{v}, \delta}$ and $R_n^{\text{re}}(\theta, \mathbf{v}) = n^{-1}$. Consequently, the data-driven OPE $\widehat{\theta}^m = \widehat{\theta} \cdot \mathbf{1}^m \in \ell_2(\mathbf{1}^m) \subseteq \ell_2(\mathbf{v}^2)$ with

$$\widehat{m} := \arg \min \left\{ -\|\widehat{\theta}^m\|_{\mathbf{v}} + \mathbf{pen}_m^{\mathbf{v}, \delta}; m \in \llbracket M \rrbracket \right\} \quad (14.25)$$

due to **Proposition** §14.07 for all $\theta \in \ell_2(\mathbf{v}^2)$ and $n, M \in \mathbb{N}$ fulfils

$$\begin{aligned} N_q^n(\|\widehat{\theta}^{\widehat{m}} - \theta\|_{\mathbf{v}}^2) &\leq \min \left\{ 3\|\theta^m - \theta\|_{\mathbf{v}}^2 + 4\mathbf{pen}_m^{\mathbf{v}, \delta}; m \in \llbracket M \rrbracket \right\} + 8(3/2) C_{\mathbf{v}, \delta} n^{-1} \\ &\leq 40 \min \left\{ \|\theta^m - \theta\|_{\mathbf{v}}^2 + n^{-1} \delta_m \|\mathbf{1}^m\|_{\mathbf{v}}^2; m \in \llbracket M \rrbracket \right\} + 12 C_{\mathbf{v}, \delta} n^{-1} \end{aligned}$$

Introduce $R_n^\circ(\theta, \mathbf{v}) := \min \left\{ \|\theta^m - \theta\|_{\mathbf{v}}^2 + n^{-1} \delta_m \|\mathbf{1}^m\|_{\mathbf{v}}^2; m \in \mathbb{N} \right\}$ where $R_n^\circ(\theta, \mathbf{v}) \geq R_n^\circ(\theta, \mathbf{v}) \geq n^{-1} \mathbf{v}_1^2$ since in general $n^{-1} \delta_m \|\mathbf{1}^m\|_{\mathbf{v}}^2 \geq n^{-1} \|\mathbf{1}^m\|_{\mathbf{v}}^2$ for all $m \in \mathbb{N}$. Consequently, if the upper bound $M \in \mathbb{N}$ satisfies $\arg \min \left\{ \|\theta^m - \theta\|_{\mathbf{v}}^2 + n^{-1} \delta_m \|\mathbf{1}^m\|_{\mathbf{v}}^2; m \in \mathbb{N} \right\} =: m^\circ \in \llbracket M \rrbracket$ then we obtain $N_q^n(\|\widehat{\theta}^{\widehat{m}} - \theta\|_{\mathbf{v}}^2) \leq C R_n^\circ(\theta, \mathbf{v})$ with $C := 40 + 12 C_{\mathbf{v}, \delta}$. However, the upper bound $R_n^\circ(\theta, \mathbf{v})$ faces a deterioration by the factor δ and thus it is generally not an oracle bound. \square

§14|03|02 Maximal global \mathbf{v} -risk

§14.33 **Reminder** (*Maximal global \mathbf{v} -risk in GSSM §14.02*). Given Model §14.02 we consider an OPE as in **Section** §12. Here the observable noisy version $\widehat{\theta}$ admits a N_q^n -distribution belonging to the family $N_\Theta^n := (N_q^n)_{\theta \in \Theta}$, $\Theta \subseteq \ell_2$. Under Assumption §11.12 in **Corollary** §12.24 an upper bound for a maximal global \mathbf{v} -risk of an OPE is shown. More precisely, the performance of the OPE

$\hat{\theta}_\cdot^m = \hat{\theta}_\cdot \mathbb{1}_\cdot^m \in \ell_2 \mathbb{1}_\cdot^m \subseteq \ell_2(\mathfrak{v}_\cdot^2)$ with dimension $m \in \mathbb{N}$ is measured by its maximal global \mathfrak{v}_\cdot -risk over the ellipsoid $\ell_2^{\mathfrak{a}, \mathfrak{r}}$, that is

$$\mathcal{R}_n^{\mathfrak{v}}[\hat{\theta}_\cdot^m | \ell_2^{\mathfrak{a}, \mathfrak{r}}] := \sup \{ \mathbb{N}_\theta^m (\|\hat{\theta}_\cdot^m - \theta\|_{\mathfrak{v}_\cdot}^2) : \theta \in \ell_2^{\mathfrak{a}, \mathfrak{r}} \}.$$

Let us recall (12.06) where for $n, m \in \mathbb{N}$ we have defined $(\mathfrak{a}\mathfrak{v})_{(m)}^2 := \|(\mathfrak{a}\mathfrak{v})_\cdot^2 \mathbb{1}_\cdot^{m \perp}\|_{\ell_\infty}$ and

$$\begin{aligned} R_n^m(\mathfrak{a}, \mathfrak{v}) &:= (\mathfrak{a}\mathfrak{v})_{(m)}^2 \vee n^{-1} \|\mathbb{1}_\cdot^m\|_{\mathfrak{v}_\cdot}^2, \quad m_n^* := \arg \min \{ R_n^m(\mathfrak{a}, \mathfrak{v}) : m \in \mathbb{N} \} \\ \text{and } R_n^*(\mathfrak{a}, \mathfrak{v}) &:= R_n^{m_n^*}(\mathfrak{a}, \mathfrak{v}) = \min \{ R_n^m(\mathfrak{a}, \mathfrak{v}) : m \in \mathbb{N} \}. \end{aligned} \quad (14.26)$$

By **Corollary** §12.24 under Assumption §11.12 the maximal global \mathfrak{v} -risk of an OPE $\hat{\theta}_\cdot^{m_n^*}$ with optimally chosen dimension m_n^* as in (14.26) satisfies

$$\mathcal{R}_n^{\mathfrak{v}}[\hat{\theta}_\cdot^{m_n^*} | \ell_2^{\mathfrak{a}, \mathfrak{r}}] \leq C R_n^*(\mathfrak{a}, \mathfrak{v})$$

with $C = 1 + \mathfrak{r}^2$. Moreover, under Assumption §13.35 due to **Proposition** §13.37 $R_n^*(\mathfrak{a}, \mathfrak{v})$ provides (up to a constant) also a lower bound of the maximal global \mathfrak{v} -risk over the ellipsoid $\ell_2^{\mathfrak{a}, \mathfrak{r}}$ for any estimator. Consequently, (up to a constant) $R_n^*(\mathfrak{a}, \mathfrak{v})$ is a minimax bound and $\hat{\theta}_\cdot^{m_n^*}$ is minimax optimal. However, the optimal dimension m_n^* depends on $\mathfrak{a} \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ characterising the ellipsoid $\ell_2^{\mathfrak{a}, \mathfrak{r}}$. \square

§14.34 **Proposition** (*GSSM (§14.02 continued)*). Let $\hat{\theta}_\cdot = \theta + n^{-1/2} \dot{B}_\cdot \sim \mathbb{N}_\theta^m$ as in Model §14.02 where $\theta \in \ell_2$ and $\dot{B}_\cdot \sim \mathbb{N}_{(0,1)}^{\otimes \mathbb{N}}$. Given $\mathfrak{v}_\cdot \in (\mathbb{R}_{>0})^{\mathbb{N}}$, $M \in \mathbb{N}$ and $\text{pen}_\cdot^{\mathfrak{v}}$ as in (14.17) consider a data-driven OPE $\hat{\theta}_\cdot^{\hat{m}} = \hat{\theta}_\cdot \mathbb{1}_\cdot^{\hat{m}} \in \ell_2 \mathbb{1}_\cdot^{\hat{m}} \subseteq \ell_2(\mathfrak{v}_\cdot^2)$ with \hat{m} as in (14.18). If Assumptions §11.12 and §14.21 (with $C_\mathfrak{v} \in \mathbb{R}_{>0}^+$) are satisfied, then for all $n, M \in \mathbb{N}$ we have

$$\mathcal{R}_n^{\mathfrak{v}}[\hat{\theta}_\cdot^{\hat{m}} | \ell_2^{\mathfrak{a}, \mathfrak{r}}] \leq (3\mathfrak{r}^2 + 40) \min \{ R_n^m(\mathfrak{a}, \mathfrak{v}) : m \in \llbracket M \rrbracket \} + 12 C_\mathfrak{v} n^{-1} \quad (14.27)$$

where $R_n^m(\theta, \mathfrak{v}) := (\mathfrak{a}\mathfrak{v})_{(m)}^2 \vee n^{-1} \|\mathbb{1}_\cdot^m\|_{\mathfrak{v}_\cdot}^2$ is defined as in (14.26).

§14.35 **Proof of Proposition** §14.34. is given in the lecture. \square

§14.36 **Comment**. The minimax bound $R_n^*(\mathfrak{a}, \mathfrak{v}) = R_n^{m_n^*}(\mathfrak{a}, \mathfrak{v}) = \min \{ R_n^m(\mathfrak{a}, \mathfrak{v}) : m \in \mathbb{N} \}$ (for details see **Reminder** §14.33) satisfies $n R_n^*(\mathfrak{a}, \mathfrak{v}) \geq \|\mathbb{1}_\cdot^{m_n^*}\|_{\mathfrak{v}_\cdot}^2 \geq \mathfrak{v}_1^2$. Consequently, the last upper bound in (14.27) and the minimax bound $R_n^*(\mathfrak{a}, \mathfrak{v})$ coincide up to a constant $(3\mathfrak{r}^2 + 40 + 12 C_\mathfrak{v} \mathfrak{v}_1^{-2})$ provided the minimax dimension fulfils $m_n^* \in \llbracket M \rrbracket$. Therefore, we wish the upper bound M to be as large as possible. The next assertion shows that $M^\mathfrak{v}$ as in (14.20) is a suitable choice for the upper bound. \square

§14.37 **Corollary** (*GSSM (§14.02 continued)*). Given $\mathfrak{v}_\cdot \in (\mathbb{R}_{>0})^{\mathbb{N}}$, $M^\mathfrak{v} \in \mathbb{N}$ as in (14.20) and $\text{pen}_\cdot^{\mathfrak{v}}$ as in (14.17) consider a data-driven OPE $\hat{\theta}_\cdot^{\hat{m}} = \hat{\theta}_\cdot \mathbb{1}_\cdot^{\hat{m}} \in \ell_2 \mathbb{1}_\cdot^{\hat{m}} \subseteq \ell_2(\mathfrak{v}_\cdot^2)$ with \hat{m} as in (14.21). Under the assumptions of **Proposition** §14.34 for each $n \in \mathbb{N}$ such that $R_n^*(\mathfrak{a}, \mathfrak{v}) \leq \mathfrak{v}_1^2$ we have

$$\mathcal{R}_n^{\mathfrak{v}}[\hat{\theta}_\cdot^{\hat{m}} | \ell_2^{\mathfrak{a}, \mathfrak{r}}] \leq (3\mathfrak{r}^2 + 40) R_n^*(\mathfrak{a}, \mathfrak{v}) + 12 C_\mathfrak{v} n^{-1} \leq C R_n^*(\mathfrak{a}, \mathfrak{v}) \quad (14.28)$$

and, hence up to the constant $C := 3\mathfrak{r}^2 + 40 + 12 C_\mathfrak{v} \mathfrak{v}_1^{-2}$ the feasible data-driven estimator $\hat{\theta}_\cdot^{\hat{m}}$ is **minimax optimal**.

§14.38 **Proof of Corollary** §14.37. is given in the lecture. \square

§14|04 Goldenshluger and Lepskij's method

The next selection method is inspired by a bandwidth selection method in kernel density estimation proposed in Goldenshluger and Lepskij [2011]. Let us consider a probability measure \mathbb{P}_θ^n for some $\theta \in \Theta$. We shall measure the accuracy of the estimator $\widehat{\theta}^m$ of θ by its risk $\mathbb{E}_\theta^n(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^m, \theta))$ where $\mathfrak{d}_{\text{ist}}(\cdot, \cdot)$ is a certain semi metric to be specified below. Inspired by Lepskij's method (which appeared in a series of papers by Lepskij [1990, 1991, 1992a,b]) given an integer $M \in \mathbb{N}$ and a sequence $\mathfrak{p}_{\text{en}, \cdot} \in (\mathbb{R}^+)^{\mathbb{N}}$ of penalties we define a contrast $\text{contr}_\cdot \in (\mathbb{R}^+)^{\llbracket M \rrbracket}$ by

$$\begin{aligned} \text{contr}_m &:= \max \left\{ \left(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^j, \widehat{\theta}^m) - \mathfrak{p}_{\text{en}_j} - \mathfrak{p}_{\text{en}_m} \right)_+ : j \in \llbracket m, M \rrbracket \right\} \\ &= \max \left\{ \left(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^{m \vee j}, \widehat{\theta}^{m \wedge j}) - \mathfrak{p}_{\text{en}_{m \vee j}} - \mathfrak{p}_{\text{en}_{m \wedge j}} \right)_+ : j \in \llbracket m, M \rrbracket \right\}, \quad m \in \llbracket M \rrbracket. \end{aligned} \quad (14.29)$$

In the spirit of Goldenshluger and Lepskij [2011] combining the contrast given in (14.29) and the penalisation approach of model selection in [Subsection §14|02](#) we select the dimension

$$\widehat{m} := \arg \min \left\{ \text{contr}_m + \mathfrak{p}_{\text{en}_m} : m \in \llbracket M \rrbracket \right\}. \quad (14.30)$$

The data-driven estimator of θ is now given by $\widehat{\theta}^{\widehat{m}}$ and below we derive an upper bound for its risk $\mathbb{E}_\theta^n(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^{\widehat{m}}, \theta))$. The construction of the penalty sequence $\mathfrak{p}_{\text{en}, \cdot}$ and the upper bound M given below is motivated by the following key argument used in the proof of the risk bound which we present first. Moreover, both $\mathfrak{p}_{\text{en}, \cdot}$ and M will depend, among others, on the noise level n , however, for sake of simplicity we will omit an additional subscript. The key argument for our reasoning is the next assertion.

§14.39 **Lemma (key argument).** Let $\text{bias}_\cdot(\theta, \mathfrak{d}_{\text{ist}}) = (\text{bias}_m(\theta, \mathfrak{d}_{\text{ist}}))_{m \in \mathbb{N}} \in (\mathbb{R}^+)^{\mathbb{N}}$ be defined by

$$\text{bias}_m(\theta, \mathfrak{d}_{\text{ist}}) := \sup \left\{ \mathfrak{d}_{\text{ist}}(\theta^j, \theta^m) : j \in \llbracket m, \infty \rrbracket := \mathbb{N} \cap [m, \infty) \cup \{\infty\} \right\}, \quad \forall m \in \mathbb{N}. \quad (14.31)$$

If $\mathfrak{p}_{\text{en}, \cdot} \in (\mathbb{R}^+)^{\mathbb{N}}$ then for all $M \in \mathbb{N}$ and $m \in \llbracket M \rrbracket$ we have

$$\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^m, \theta) \leq 16 \text{bias}_m^2(\theta, \mathfrak{d}_{\text{ist}}) + \frac{16}{3} \mathfrak{p}_{\text{en}_m} + 28 \max \left\{ \left(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^j, \theta^j) - \mathfrak{p}_{\text{en}_j}/3 \right)_+ : j \in \llbracket m, M \rrbracket \right\}.$$

§14.40 **Proof of Lemma §14.39.** is given in the lecture. □

Similar to m_n° as in (14.02), which realises by construction a statistical-error-squared-bias compromise, let us fix a dimension $m^\circ \in \llbracket M \rrbracket$ to be specified below. Due to the last assertion for each $\theta \in \Theta$ we have

$$\begin{aligned} \mathbb{E}_\theta^n \left(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^m, \theta) \right) &\leq 16 \text{bias}_{m^\circ}^2(\theta, \mathfrak{d}_{\text{ist}}) + \frac{16}{3} \mathfrak{p}_{\text{en}_{m^\circ}} \\ &\quad + 28 \mathbb{E}_\theta^n \left(\max \left\{ \left(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^j, \theta^j) - \mathfrak{p}_{\text{en}_j}/3 \right)_+ : j \in \llbracket m^\circ, M \rrbracket \right\} \right). \end{aligned} \quad (14.32)$$

Keeping in mind that m° in contrast to m_n° eventually realises an optimal statistical-error-squared-bias trade-off among the collection of admissible values $\llbracket M \rrbracket$ rather than \mathbb{N} , we wish the upper bound M to be as large as possible. In contrast, in order to control the remainder term, the last term in (14.32), we are forced to use a rather small upper bound M . However, we bound the remainder term by imposing a condition similar to [Assumption §14.06](#), which though holds true for a wide range of solutions $\theta = U\theta \in \Theta$ under reasonable model assumptions.

§14.41 **Assumption.** There exists a constant $C := C(\theta) \in \mathbb{R}_0^+$ and $(R_n^{\text{re}}(\theta, \mathfrak{d}_{\text{ist}}))_{n \in \mathbb{N}} \in (\mathbb{R}^+)^{\mathbb{N}}$ possibly depending on the parameter $\theta = U\theta \in \Theta$ such that for each $n \in \mathbb{N}$ the upper bound $M \in \mathbb{N}$ and $m^\circ \in \llbracket M \rrbracket$ satisfy

$$\mathbb{E}_\theta^n \left(\max \left\{ \left(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^m, \theta^m) - \mathfrak{p}_{\text{en}_m}/3 \right)_+ : m \in \llbracket m^\circ, M \rrbracket \right\} \right) \leq C R_n^{\text{re}}(\theta, \mathfrak{d}_{\text{ist}}). \quad (14.33)$$

The next assertion provides an upper bound for the risk of the estimator $\widehat{\theta}^m$ with data-driven choice \widehat{m} given by (14.30).

§14.42 **Proposition.** *Let $m^\circ \in \llbracket M \rrbracket$ satisfy the Assumption §14.41 then we have*

$$\mathbb{E}_\theta^n \left(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^m, \theta) \right) \leq 16 \text{bias}_m^2(\theta, \mathfrak{d}_{\text{ist}}) + \frac{16}{3} \text{pen}_{m^\circ} + 28C R_n^{\text{re}}(\theta, \mathfrak{d}_{\text{ist}}).$$

§14.43 **Proof of Proposition §14.42.** is given in the lecture. \square

§14.44 **Corollary.** *If $m^\circ = \arg \min \{ \text{bias}_m^2(\theta, \mathfrak{d}_{\text{ist}}) + \text{pen}_m : m \in \llbracket M \rrbracket \}$ satisfies Assumption §14.41, then*

$$\mathbb{E}_\theta^n \left(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^m, \theta) \right) \leq 16 \min \{ \text{bias}_m^2(\theta, \mathfrak{d}_{\text{ist}}) + \text{pen}_m : m \in \llbracket M \rrbracket \} + 28C R_n^{\text{re}}(\theta).$$

§14.45 **Proof of Corollary §14.44.** is given in the lecture. \square

§14.46 **Comment.** Considering a global \mathfrak{v} -error we note that $\text{bias}_m^2(\theta, \mathfrak{d}_{\text{ist}}) = \|\theta^m - \theta\|_{\mathfrak{v}}^2$ for all $m \in \mathbb{N}$, and hence the upper bound in **Corollary §14.44** equals up to the numerical constants the upper bound in **Corollary §14.09** using a model selection approach (**Subsection §14|02**). Consequently, when globally estimating the parameter in a *GSSM* with a Goldenshluger and Lepskij method rather than a model selection approach as in **Subsection §14|03** we eventually obtain the same upper bounds (up to the numerical constants). However, we shall stress that in opposite to model selection the method by Goldenshluger and Lepskij does not require, that the estimator minimises a contrast function. \square

We eventually are in a situation where the sequence of penalties $\text{pen}_\bullet \in (\mathbb{R}^+)^{\mathbb{N}}$ satisfying the Assumption §14.41 still depends on characteristics of the unknown parameter θ and thus it is only partially known in advance. Assuming a sequence of estimators $\widehat{\text{pen}}_\bullet \in (\mathbb{R}^+)^{\mathbb{N}}$ we define an estimated contrast $\widehat{\text{c}}_{\text{contr}} \in (\mathbb{R}^+)^{\llbracket M \rrbracket}$ by

$$\begin{aligned} \widehat{\text{c}}_{\text{contr}_m} &:= \max \left\{ \left(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^j, \widehat{\theta}^m) - \widehat{\text{pen}}_j - \widehat{\text{pen}}_m \right)_+ : j \in \llbracket m, M \rrbracket \right\} \\ &= \max \left\{ \left(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^{m \vee j}, \widehat{\theta}^{m \wedge j}) - \widehat{\text{pen}}_{m \vee j} - \widehat{\text{pen}}_{m \wedge j} \right)_+ : j \in \llbracket m, M \rrbracket \right\}, \quad m \in \llbracket M \rrbracket \end{aligned} \quad (14.34)$$

and similar to (14.30) we select the dimension

$$\widehat{m} := \arg \min \left\{ \widehat{\text{c}}_{\text{contr}_m} + \widehat{\text{pen}}_m : m \in \llbracket M \rrbracket \right\}. \quad (14.35)$$

The data-driven estimator of θ is now given by $\widehat{\theta}^{\widehat{m}}$ and below we derive an upper bound for its risk $\mathbb{E}_\theta^n \left(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^{\widehat{m}}, \theta) \right)$. The key argument for our reasoning is the next assertion. Its proof follows along the lines of the **Proof §14.40**.

§14.47 **Lemma (key argument).** *Let $\text{bias}_\bullet(\theta, \mathfrak{d}_{\text{ist}}) = (\text{bias}_m(\theta, \mathfrak{d}_{\text{ist}}))_{m \in \mathbb{N}} \in (\mathbb{R}^+)^{\mathbb{N}}$ be defined as in (14.31) (**Lemma §14.39**). If $\widehat{\text{pen}}_\bullet, \text{pen}_\bullet \in (\mathbb{R}^+)^{\mathbb{N}}$ then for all $M \in \mathbb{N}$ and $m \in \llbracket M \rrbracket$ we have*

$$\begin{aligned} \mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^m, \theta) &\leq 16 \text{bias}_m^2(\theta, \mathfrak{d}_{\text{ist}}) + \frac{4}{3} \text{pen}_m + 28 \max \left\{ \left(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^j, \theta^j) - \text{pen}_j / 3 \right)_+ : j \in \llbracket m, M \rrbracket \right\} \\ &\quad + 8 \max \left\{ \left(\text{pen}_j - \widehat{\text{pen}}_j \right)_+ : j \in \llbracket m, M \rrbracket \right\} + 4 \widehat{\text{pen}}_m. \end{aligned}$$

§14.48 **Proof of Lemma §14.47.** is given in the lecture. \square

Similar to m_n° as in (14.02), which realises by construction a statistical-error-squared-bias compromise, let us fix a dimension $m^\circ \in \llbracket M \rrbracket$ to be specified below (analogously to (14.32)). Due to the last assertion for each $\theta \in \Theta$ we have

$$\begin{aligned} \mathbb{E}_\theta^n \left(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^m, \theta) \right) &\leq 16 \text{bias}_m^2(\theta, \mathfrak{d}_{\text{ist}}) + \frac{4}{3} \text{pen}_{m^\circ} + 28 \mathbb{E}_\theta^n \left(\max \left\{ \left(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^j, \theta^j) - \text{pen}_j / 3 \right)_+ : j \in \llbracket m^\circ, M \rrbracket \right\} \right) \\ &\quad + 8 \mathbb{E}_\theta^n \left(\max \left\{ \left(\text{pen}_j - \widehat{\text{pen}}_j \right)_+ : j \in \llbracket m^\circ, M \rrbracket \right\} \right) + 4 \mathbb{E}_\theta^n \left(\widehat{\text{pen}}_{m^\circ} \right). \end{aligned} \quad (14.36)$$

We bound the remainder terms by imposing conditions including Assumption §14.41, which though hold true for a wide range of solutions $\theta = U\theta \in \Theta$ under reasonable model assumptions.

§14.49 **Assumption.** There exists a constant $C := C(\theta) \in \mathbb{R}_0^+$ and $(R_n^{\text{re}}(\theta, \mathfrak{d}_{\text{ist}}))_{n \in \mathbb{N}} \in (\mathbb{R}^+)^{\mathbb{N}}$ possibly depending on the parameter $\theta = U\theta \in \Theta$ such that for each $n \in \mathbb{N}$ the penalties $\widehat{\mathfrak{p}}_{\text{en},j}, \mathfrak{p}_{\text{en},j} \in (\mathbb{R}^+)^{\mathbb{N}}$, the upper bound $M \in \mathbb{N}$ and $m^\circ \in \llbracket M \rrbracket$ satisfy (14.33) in Assumption §14.41 and in addition

$$\mathbb{E}_\theta^n \left(\max \left\{ \left(\mathfrak{p}_{\text{en},j} - \widehat{\mathfrak{p}}_{\text{en},j} \right)_+ : j \in \llbracket m^\circ, M \rrbracket \right\} \right) \leq C R_n^{\text{re}}(\theta, \mathfrak{d}_{\text{ist}}).$$

The next assertion provides an upper bound for the risk of the estimator $\widehat{\theta}^{\widehat{m}}$ with data-driven choice \widehat{m} given by (14.35).

§14.50 **Proposition.** If $m^\circ \in \llbracket M \rrbracket$ satisfies the Assumption §14.49 then we have

$$\mathbb{E}_\theta^n \left(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^{\widehat{m}}, \theta) \right) \leq 16 \text{bias}_{m^\circ}^2(\theta, \mathfrak{d}_{\text{ist}}) + \frac{4}{3} \mathfrak{p}_{\text{en}_{m^\circ}} + 4 \mathbb{E}_\theta^n \left(\widehat{\mathfrak{p}}_{\text{en}_{m^\circ}} \right) + 36C R_n^{\text{re}}(\theta, \mathfrak{d}_{\text{ist}}).$$

§14.51 **Proof of Proposition §14.50.** is given in the lecture. □

§14.52 **Corollary.** If $m^\circ = \arg \min \left\{ \text{bias}_m^2(\theta, \mathfrak{d}_{\text{ist}}) + \mathfrak{p}_{\text{en}_m} : m \in \llbracket M \rrbracket \right\}$ satisfies Assumption §14.49 and $\mathbb{E}_\theta^n \left(\widehat{\mathfrak{p}}_{\text{en}_{m^\circ}} \right) \leq K \mathfrak{p}_{\text{en}_{m^\circ}}$ for some $K \in [1, \infty)$, then

$$\mathbb{E}_\theta^n \left(\mathfrak{d}_{\text{ist}}^2(\widehat{\theta}^{\widehat{m}}, \theta) \right) \leq (16 \vee 6K) \min \left\{ \text{bias}_m^2(\theta, \mathfrak{d}_{\text{ist}}) + \mathfrak{p}_{\text{en}_m} : m \in \llbracket M \rrbracket \right\} + 36C R_n^{\text{re}}(\theta).$$

§14.53 **Proof of Corollary §14.52.** is given in the lecture. □

§14|05 GSSM: data-driven local estimation

Lemma §14.18 in Subsection §14|03 presents tail bounds of sums of independent squared Gaussian random variables. We state next an elementary tail bound and a concentration inequality of a single Gaussian random variable.

§14.54 **Lemma.** Let $Z \sim N_{(0,1)}$. For all $\eta \in \mathbb{R}_0^+$ and $\zeta, K \in [1, \infty)$ we have

$$N_{(0,1)}(Z > \eta) \leq (2\pi\eta^2)^{-1/2} \exp(-\eta^2/2) \text{ and } N_{(0,1)}\left(\left(Z^2 - 2\zeta(1 + \log K)\right)_+\right) \leq K^{-\zeta}. \quad (14.37)$$

§14.55 **Proof of Lemma §14.54.** Exercise. □

§14|05|01 Local ϕ -risk

§14.56 **Reminder (Local oracle ϕ -risk in GSSM §14.02).** Given Model §14.02 we consider an OPE as in Section §12. Here the observable noisy version $\widehat{\theta}$ admits a N_q^n -distribution belonging to the family $N_\Theta^n := (N_q^n)_{\theta \in \Theta}$, $\Theta \subseteq \ell_2$. Let us recall (12.11) in Proposition §12.32 where $\phi \in (\mathbb{R}_0^+)^{\mathbb{N}}$, $\theta \in \text{dom}(\phi_{\nu_N})$ and $n, m \in \mathbb{N}$ we have defined

$$\begin{aligned} R_n^m(\theta, \phi) &:= |\phi_{\nu_N}(\theta, \mathbf{1}^{m \perp})|^2 + n^{-1} \|\mathbf{1}^m\|_\phi^2, \quad m_n^\circ := \arg \min \{ R_n^m(\theta, \phi) : m \in \mathbb{N} \} \\ \text{and } R_n^\circ(\theta, \phi) &:= R_n^{m_n^\circ}(\theta, \phi) = \min \{ R_n^m(\theta, \phi) : m \in \mathbb{N} \}. \end{aligned} \quad (14.38)$$

Due to Corollary §12.38 the (infeasible) OPE $\widehat{\theta}^{m_n^\circ} = \widehat{\theta} \cdot \mathbf{1}^{m_n^\circ} \in \ell_2 \cdot \mathbf{1}^{m_n^\circ} \subseteq \text{dom}(\phi_{\nu_N})$ with oracle dimension m_n° as in (14.38) satisfies

$$N_q^n \left(|\phi_{\nu_N}(\widehat{\theta}^{m_n^\circ} - \theta)|^2 \right) = R_n^\circ(\theta, \phi) = \inf_{m \in \mathbb{N}} N_q^n \left(|\phi_{\nu_N}(\widehat{\theta}^m - \theta)|^2 \right),$$

and hence it is *oracle optimal* (with constant 1). □

§14.57 **Corollary.** For $\phi \in (\mathbb{R}_{>0})^{\mathbb{N}}$ and $n, M \in \mathbb{N}$ setting $K_m := (\|\mathbf{1}^m\|_{\phi}^2 \vee 1)m^2 \geq 1$, $m \in \mathbb{N}$, we have

$$N_q^n \left(\max \left\{ \left(|\phi \nu_{\mathbb{N}}(\hat{\theta}^m - \theta^m)|^2 - 2(1 + \log K_m)n^{-1}\|\mathbf{1}^m\|_{\phi}^2 \right) : m \in \llbracket M \rrbracket \right\} \right) \leq 2n^{-1}. \quad (14.39)$$

§14.58 **Proof of Corollary §14.57.** is given in the lecture. \square

§14.59 **Notation.** Consider a sequence of penalties $\text{pen}^{\phi} = (\text{pen}_m^{\phi})_{m \in \mathbb{N}} \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ given by

$$\text{pen}_m^{\phi} := 6 \left(1 + \left(\log \|\mathbf{1}^m\|_{\phi}^2 \right) + 2 \log m \right) n^{-1} \|\mathbf{1}^m\|_{\phi}^2, \quad \text{for each } m \in \mathbb{N} \quad (14.40)$$

which is obviously known in advance. Moreover, studying a ϕ -error the bias term introduced in (14.31) becomes

$$\text{bias}_m(\theta, \phi) = \sup \left\{ |\phi \nu_{\mathbb{N}}(\theta^j - \theta^m)| = |\phi \nu_{\mathbb{N}}(\theta \mathbf{1}^{\llbracket m, j \rrbracket})| : j \in \llbracket m, \infty \rrbracket \right\} \quad \forall m \in \mathbb{N}.$$

If $\theta \in \text{dom}(\phi \nu_{\mathbb{N}})$ and hence $\nu_{\mathbb{N}}(|\phi \theta|) \in \mathbb{R}$ then $\text{bias}_m(\theta, \phi) \leq \nu_{\mathbb{N}}(|\phi \theta| \mathbf{1}^{m \perp}) = o(1)$ as $m \rightarrow \infty$ by dominated convergence. Considering the data-driven OSE $\hat{\theta}^m = \hat{\theta} \mathbf{1}^{\hat{m}}$ with dimension parameter \hat{m} selected as in (14.30) with penalty sequence pen^{ϕ} given in (14.40) and arbitrary but fixed upper bound $M \in \mathbb{N}$ we derive below an upper bound for its local ϕ -risk, $N_q^n (|\phi \nu_{\mathbb{N}}(\hat{\theta}^m - \theta^m)|^2)$. \square

§14.60 **Proposition (GSSM (§14.02 continued)).** Let $\hat{\theta} = \theta + n^{-1/2} \dot{B} \sim N_q^n$ as in Model §14.02 where $\theta \in \ell_2$ and $\dot{B} \sim N_{(0,1)}^{\otimes \mathbb{N}}$. Given $\phi \in (\mathbb{R}_{>0})^{\mathbb{N}}$, $M \in \mathbb{N}$ and pen^{ϕ} as in (14.40) consider a data-driven OPE $\hat{\theta}^m = \hat{\theta} \mathbf{1}^{\hat{m}} \in \ell_2 \mathbf{1}^{\hat{m}} \subseteq \text{dom}(\phi \nu_{\mathbb{N}})$ of $\theta \in \text{dom}(\phi \nu_{\mathbb{N}})$ with

$$\begin{aligned} \hat{m} &:= \arg \min \left\{ \text{contr}_m^{\phi} + \text{pen}_m^{\phi} : m \in \llbracket M \rrbracket \right\} \quad \text{and} \\ \text{contr}_m^{\phi} &:= \max \left\{ \left(|\phi \nu_{\mathbb{N}}(\hat{\theta}^j - \hat{\theta}^m)|^2 - \text{pen}_j^{\phi} - \text{pen}_m^{\phi} \right) : j \in \llbracket m, M \rrbracket \right\}, \quad m \in \llbracket M \rrbracket. \end{aligned} \quad (14.41)$$

Then for all $n, M \in \mathbb{N}$ we have

$$\begin{aligned} N_q^n (|\phi \nu_{\mathbb{N}}(\hat{\theta}^{\hat{m}} - \theta)|^2) &\leq 64 \min \left\{ \text{bias}_m^2(\theta, \phi) + \left(1 + \left(\log \|\mathbf{1}^m\|_{\phi}^2 \right) + \log m \right) n^{-1} \|\mathbf{1}^m\|_{\phi}^2 : m \in \llbracket M \rrbracket \right\} \\ &\quad + 56 n^{-1}. \end{aligned} \quad (14.42)$$

§14.61 **Proof of Proposition §14.60.** is given in the lecture. \square

§14.62 **Comment.** Let us compare the dominating part of the upper bound given in (14.42), that is

$$\min \left\{ \text{bias}_m^2(\theta, \phi) + \left(1 + \left(\log \|\mathbf{1}^m\|_{\phi}^2 \right) + \log m \right) n^{-1} \|\mathbf{1}^m\|_{\phi}^2 : m \in \llbracket M \rrbracket \right\} \quad (14.43)$$

with the oracle bound $R_n^{\circ}(\theta, \phi) = \min \left\{ |\phi \nu_{\mathbb{N}}(\theta^m - \theta)|^2 + n^{-1} \|\mathbf{1}^m\|_{\phi}^2 : m \in \mathbb{N} \right\}$ (for details see **Reminder §14.56**). In (14.43) we face eventually a deterioration by three sources. First, we generally have $\text{bias}_m(\theta, \phi) \geq |\phi \nu_{\mathbb{N}}(\theta^m - \theta)|$, but note that for $\theta \phi \in (\mathbb{R}^+)^{\mathbb{N}}$ equality holds, that is

$$\text{bias}_m(\theta, \phi) = \sup \left\{ \nu_{\mathbb{N}}(\phi \theta \mathbf{1}^{\llbracket m, j \rrbracket}) : j \in \llbracket m, \infty \rrbracket \right\} = \nu_{\mathbb{N}}(\phi \theta \mathbf{1}^{m \perp}) = |\phi \nu_{\mathbb{N}}(\theta^m - \theta)|$$

for all $m \in \mathbb{N}$. Secondly, the variance term features an additional factor $1 + \left(\log \|\mathbf{1}^m\|_{\phi}^2 \right) + \log m$, and finally the upper bound M might impose an additional deterioration. We note that the oracle bound $R_n^{\circ}(\theta, \phi)$ is parametric, i.e. $n R_n^{\circ}(\theta, \phi) = O(1)$ as $n \rightarrow \infty$, if $\phi \in \ell_2$ (case **(p)** in **Illustration §12.40**). In the sequel we consider only the case $\phi \notin \ell_2$, i.e. $\nu_{\mathbb{N}}(|\phi|^2) = \infty$. We set

$$M^{\phi} := \max \left\{ m \in \mathbb{N} : \|\mathbf{1}^m\|_{\phi}^2 \leq n \phi^2 \right\} \in \mathbb{N} \quad (14.44)$$

where the defining set is not empty and finite since $\|\phi\|_{\ell_2}^2 = \infty$. The next assertion shows that this is a suitable choice for the upper bound. Moreover, we estimate the bias term by $\text{bias}_m(\theta, \phi) \leq \nu(|\phi \theta| \mathbf{1}^{m \perp})$ where equality holds whenever $\theta \phi \in (\mathbb{R}^+)^{\mathbb{N}}$. \square

§14.63 **Corollary** (GSSM (§14.02 continued)). Given $\phi \in (\mathbb{R}_{\setminus 0})^{\mathbb{N}}$ with $\phi \notin \ell_2$, $M^\phi \in \mathbb{N}$ as in (14.44) and pen_m^ϕ as in (14.40) consider a data-driven OPE $\widehat{\theta}^m = \widehat{\theta} \mathbf{1}^m \in \ell_2 \mathbf{1}^m \subseteq \text{dom}(\phi_{\mathbb{N}})$ of $\theta \in \text{dom}(\phi_{\mathbb{N}})$ with

$$\widehat{m} := \arg \min \left\{ \text{contr}_m^\phi + \text{pen}_m^\phi : m \in \llbracket M^\phi \rrbracket \right\} \quad \text{and} \\ \text{contr}_m^\phi := \max \left\{ \left(|\phi_{\mathbb{N}}(\widehat{\theta}^j - \widehat{\theta}^m)|^2 - \text{pen}_j^\phi - \text{pen}_m^\phi \right)_+ : j \in \llbracket m, M^\phi \rrbracket \right\}, \quad m \in \llbracket M^\phi \rrbracket. \quad (14.45)$$

For $n, m \in \mathbb{N}$ we set

$$R_n^m(\theta, \phi) := \left(\nu_{\mathbb{N}}(|\phi \theta| \mathbf{1}^{m \perp}) \right)^2 + \left(1 + \left(\log \|\mathbf{1}^m\|_\phi^2 \right)_+ + \log m \right) n^{-1} \|\mathbf{1}^m\|_\phi^2, \\ m^\circ := \arg \min \left\{ R_n^m(\theta, \phi) : m \in \mathbb{N} \right\} \quad \text{and} \\ R_n^\circ(\theta, \phi) := R_n^{m^\circ}(\theta, \phi) = \min \left\{ R_n^m(\theta, \phi) : m \in \mathbb{N} \right\}. \quad (14.46)$$

Under the assumptions of **Proposition** §14.60 for each $n \in \mathbb{N}$ such that $R_n^\circ(\theta, \phi) \leq \phi_1^2$ we have

$$N_q^n \left(|\phi_{\mathbb{N}}(\widehat{\theta}^m - \theta)|^2 \right) \leq 64 R_n^\circ(\theta, \phi) + 56n^{-1} \leq (64 + 56\phi_1^{-2}) R_n^\circ(\theta, \phi). \quad (14.47)$$

§14.64 **Proof of Proof** §14.64. is given in the lecture. \square

§14.65 **Comment**. The data-driven bound $R_n^\circ(\theta, \phi)$ compared to the oracle bound $R_n^\circ(\theta, \phi)$ features a deterioration of the variance term at least by a logarithmic factor. The appearance of the logarithmic factor within the bound is a known fact in the context of local estimation (cf. Laurent et al. [2008] who consider model selection given direct Gaussian observations). Brown and Low [1996] show that it is unavoidable in the context of nonparametric Gaussian regression and hence it is widely considered as an acceptable price for adaptation. \square

§14.66 **Illustration**. We illustrate the last results considering the two specifications (o) and (s) given in Table 03 [§12] (**Illustration** §12.40). We restrict ourselves to the case $\phi \notin \ell_2$ only.

Table 01 [§14]

Order of the oracle rate $R_n^\circ(\theta, \phi)$ and the data-driven rate $R_n^\circ(\theta, \phi)$ as $n \rightarrow \infty$

$(j \in \mathbb{N})$	$(a \in \mathbb{R}_0^+)$	(squared bias)	(variance)	M^ϕ	m°	$R_n^\circ(\theta, \phi)$	$R_n^\circ(\theta, \phi)$	
$\phi = j^{v-1/2}$	θ_j	$(\nu_{\mathbb{N}}(\phi \theta \mathbf{1}^{m \perp}))^2$	$\ \mathbf{1}^m\ _\phi^2$					
(o)	$v \in (0, a)$	$j^{-a-1/2}$	$m^{-2(a-v)}$	m^{2v}	$n^{\frac{1}{2v}}$	$\left(\frac{n}{\log n} \right)^{\frac{1}{2a}}$	$n^{-\frac{(a-v)}{a}}$	$\left(\frac{\log n}{n} \right)^{\frac{(a-v)}{a}}$
	$v = 0$	$j^{-a-1/2}$	m^{-2a}	$\log m$	e^n	$\left(\frac{n}{(\log n)^2} \right)^{\frac{1}{2a}}$	$\frac{\log n}{n}$	$\frac{(\log n)^2}{n}$
(s)	$v \in \mathbb{R}_0^+$	$e^{-j^{2a}}$	$m^{(1-2(a-v))_+} e^{-2m^{2a}}$	m^{2v}	$n^{\frac{1}{2v}}$	$(\log n)^{\frac{1}{2a}}$	$\frac{(\log n)^v}{n}$	$\frac{(\log n)^v (\log \log n)}{n}$
	$v = 0$	$e^{-j^{2a}}$	$m^{(1-2a)_+} e^{-2m^{2a}}$	$\log m$	e^n	$(\log n)^{\frac{1}{2a}}$	$\frac{\log \log n}{n}$	$\frac{(\log \log n)^2}{n}$

We note that in Table 01 [§14] the order of the oracle rate $R_n^\circ(\theta, \phi)$ and the data-driven rate $R_n^\circ(\theta, \phi)$ is depicted for $v \geq 0$ only. In case $v < 0$ we have $\phi \in \ell_2$ and thus **Corollary** §14.63 is not applicable. \square

§14|05|02 Maximal local ϕ -risk

§14.67 **Reminder** (Maximal local ϕ -risk in GSSM §14.02). Given Model §14.02 we consider an OPE as in **Section** §12. Here the observable noisy version $\widehat{\theta}$ admits a N_q^n -distribution belonging to the family $N_\Theta^n := (N_q^n)_{\theta \in \Theta}$, $\Theta \subseteq \ell_2$. Under Assumption §11.25 in **Corollary** §12.45 an upper bound

for a maximal local ϕ -risk of an OPE is shown. More precisely, the performance of the OPE $\widehat{\theta}^m = \widehat{\theta} \mathbf{1}^m \in \ell_2 \mathbf{1}^m \subseteq \text{dom}(\phi_{\nu_n})$ with dimension $m \in \mathbb{N}$ is measured by its maximal local ϕ -risk over the ellipsoid ℓ_2^{ar} , that is

$$\mathcal{R}_n^\phi[\widehat{\theta}^m | \ell_2^{\text{ar}}] := \sup \{ N_n^m (|\phi_{\nu_n}(\widehat{\theta}^m - \theta)|^2) : \theta \in \ell_2^{\text{ar}} \}.$$

Let us recall (12.13) where for $n, m \in \mathbb{N}$ we have defined

$$\begin{aligned} R_n^m(\mathbf{a}, \phi) &:= \|\mathbf{a} \mathbf{1}^{m \perp}\|_\phi^2 + n^{-1} \|\mathbf{1}^m\|_\phi^2, \quad m_n^* := \arg \min \{ R_n^m(\mathbf{a}, \phi) : m \in \mathbb{N} \} \\ \text{and } R_n^*(\mathbf{a}, \phi) &:= R_n^{m_n^*}(\mathbf{a}, \phi) = \min \{ R_n^m(\mathbf{a}, \phi) : m \in \mathbb{N} \}. \end{aligned} \quad (14.48)$$

By **Corollary** §12.45 under Assumption §11.25 the maximal local ϕ -risk of an OPE $\widehat{\theta}^{m_n^*}$ with optimally chosen dimension m_n^* as in (14.48) satisfies

$$\mathcal{R}_n^\phi[\widehat{\theta}^{m_n^*} | \ell_2^{\text{ar}}] \leq C R_n^*(\mathbf{a}, \phi)$$

with $C = 1 \vee r^2$. Moreover, under Assumption §13.24 due to **Proposition** §13.26 $R_n^*(\mathbf{a}, \phi)$ provides (up to a constant) also a lower bound of the maximal local ϕ -risk over the ellipsoid ℓ_2^{ar} for any estimator. Consequently, (up to a constant) $R_n^*(\mathbf{a}, \phi)$ is a minimax bound and $\widehat{\theta}^{m_n^*}$ is minimax optimal. However, the optimal dimension m_n^* depends on $\mathbf{a} \in (\mathbb{R}_0^+)^{\mathbb{N}}$ characterising the ellipsoid ℓ_2^{ar} . \square

§14.68 **Proposition** (GSSM (§14.02 continued)). Let $\widehat{\theta} = \theta + n^{-1/2} \dot{B} \sim N_n^n$ as in Model §14.02 where $\theta \in \ell_2$ and $\dot{B} \sim N_{(0,1)}^{\otimes \mathbb{N}}$. Given $\phi \in (\mathbb{R}_{>0})^{\mathbb{N}}$, $M \in \mathbb{N}$ and pen_ϕ° as in (14.40) consider a data-driven OPE $\widehat{\theta}^{\widehat{m}} = \widehat{\theta} \mathbf{1}^{\widehat{m}} \in \ell_2 \mathbf{1}^{\widehat{m}} \subseteq \text{dom}(\phi_{\nu_n})$ with \widehat{m} as in (14.41). If Assumption §11.25 is satisfied, then for all $n, M \in \mathbb{N}$ we have

$$\begin{aligned} \mathcal{R}_n^\phi[\widehat{\theta}^{\widehat{m}} | \ell_2^{\text{ar}}] &\leq (16r^2 \vee 64) \min \{ \|\mathbf{a} \mathbf{1}^{m \perp}\|_\phi^2 + (1 + (\log \|\mathbf{1}^m\|_\phi^2)_+ + \log m) n^{-1} \|\mathbf{1}^m\|_\phi^2 : m \in \llbracket M \rrbracket \} \\ &\quad + 56 n^{-1}. \end{aligned} \quad (14.49)$$

§14.69 **Proof** of **Proposition** §14.68. is given in the lecture. \square

§14.70 **Corollary** (GSSM (§14.02 continued)). Given $\phi \in (\mathbb{R}_{>0})^{\mathbb{N}}$ with $\phi \notin \ell_2$, $M^\phi \in \mathbb{N}$ as in (14.44) and pen_ϕ° as in (14.40) consider a data-driven OPE $\widehat{\theta}^{\widehat{m}} = \widehat{\theta} \mathbf{1}^{\widehat{m}} \in \ell_2 \mathbf{1}^{\widehat{m}} \subseteq \text{dom}(\phi_{\nu_n})$ with \widehat{m} as in (14.45). For $n, m \in \mathbb{N}$ we set

$$\begin{aligned} R_n^m(\mathbf{a}, \phi) &:= \|\mathbf{a} \mathbf{1}^{m \perp}\|_\phi^2 + (1 + (\log \|\mathbf{1}^m\|_\phi^2)_+ + \log m) n^{-1} \|\mathbf{1}^m\|_\phi^2, \\ m^\circ &:= \arg \min \{ R_n^m(\mathbf{a}, \phi) : m \in \mathbb{N} \} \quad \text{and} \\ R_n^\circ(\mathbf{a}, \phi) &:= R_n^{m^\circ}(\mathbf{a}, \phi) = \min \{ R_n^m(\mathbf{a}, \phi) : m \in \mathbb{N} \}. \end{aligned} \quad (14.50)$$

Under Assumption §11.25 for each $n \in \mathbb{N}$ such that $R_n^\circ(\mathbf{a}, \phi) \leq \phi_1^2$ we have

$$\mathcal{R}_n^\phi[\widehat{\theta}^{\widehat{m}} | \ell_2^{\text{ar}}] \leq (16r^2 \vee 64) R_n^\circ(\mathbf{a}, \phi) + 56n^{-1} \leq (16r^2 \vee 64 + 56\phi_1^{-2}) R_n^\circ(\mathbf{a}, \phi). \quad (14.51)$$

§14.71 **Proof** of **Proposition** §14.71. is given in the lecture. \square

§14.72 **Comment**. The data-driven bound $R_n^\circ(\mathbf{a}, \phi)$ compared to the minimax bound $R_n^*(\mathbf{a}, \phi)$ features a deterioration of the variance term at least by a factor $\log n$. The appearance of the logarithmic factor within the bound is a known fact in the context of local estimation (cf. Laurent et al. [2008] who consider model selection given direct Gaussian observations). Brown and Low [1996] show that it is unavoidable in the context of nonparametric Gaussian regression and hence it is widely considered as an acceptable price for adaptation. \square

§14.73 **Illustration.** We illustrate the last results considering usual behaviour for $\mathbf{a}, \phi \in (\mathbb{R}_{\setminus 0})^{\mathbb{N}}$. As in **Illustration** §12.47 we distinguish again the following two cases **(p)** $\phi \in \ell_2$, and **(np)** $\phi \notin \ell_2$. Interestingly, in case **(p)** the minimax bound $R_n^*(\mathbf{a}, \phi)$ in **Proposition** §12.42 is parametric, that is, $nR_n^*(\mathbf{a}, \phi) = O(1)$, in case **(np)** the bound is nonparametric, i.e. $\lim_{n \rightarrow \infty} nR_n^*(\mathbf{a}, \phi) = \infty$. In case **(np)** consider the following two specifications:

Table 02 [§14]

Order of minimax rate $R_n^*(\mathbf{a}, \phi)$ and the data-driven rate $R_n^\circ(\mathbf{a}, \phi)$ as $n \rightarrow \infty$

	$(j \in \mathbb{N})$	$(a \in \mathbb{R}_0^+)$	(squared bias)	(variance)			$R_n^*(\mathbf{a}, \phi)$	$R_n^\circ(\mathbf{a}, \phi)$
	$\phi = j^{v-1/2}$	\mathbf{a}_j^2	$\ \mathbf{a} \cdot \mathbf{1}_\phi^{m \cdot}\ ^2_\phi$	$\ \mathbf{1}_\phi^m\ ^2_\phi$	M^ϕ	m°		
(o)	$v \in (0, a)$	j^{-2a}	$m^{-2(a-v)}$	m^{2v}	$n^{\frac{1}{2v}}$	$\left(\frac{n}{\log n}\right)^{\frac{1}{2a}}$	$n^{-\frac{(a-v)}{a}}$	$\left(\frac{\log n}{n}\right)^{\frac{a-v}{a}}$
	$v = 0$	j^{-2a}	m^{-2a}	$\log m$	e^n	$\left(\frac{n}{(\log n)^2}\right)^{\frac{1}{2a}}$	$\frac{\log n}{n}$	$\frac{(\log n)^2}{n}$
(s)	$v \in \mathbb{R}_0^+$	$e^{-j^{2a}}$	$m^{2(v-a)+} e^{-m^{2a}}$	m^{2v}	$n^{\frac{1}{2v}}$	$(\log n)^{\frac{1}{2a}}$	$\frac{(\log n)^{\frac{v}{a}}}{n}$	$\frac{(\log n)^{\frac{v}{a}} (\log \log n)}{n}$
	$v = 0$	$e^{-j^{2a}}$	$e^{-m^{2a}}$	$\log m$	e^n	$(\log n)^{\frac{1}{2a}}$	$\frac{\log \log n}{n}$	$\frac{(\log \log n)^2}{n}$

We note that in Table 02 [§14] the order of the minimax rate $R_n^*(\mathbf{a}, \phi)$ and the data-driven rate $R_n^\circ(\mathbf{a}, \phi)$ is depict for $v \geq 0$ only. For $v < 0$ we have $\phi \in \ell_2$ and thus **Corollary** §14.70 is not applicable. □

Chapter 4

Nonparametric density estimation

This chapter presents nonparametric density estimation along the lines of the textbooks by Tsybakov [2009] and Comte [2015] where far more details, examples and further discussions can be found.

Overview

§15	Noisy density coefficients	73
§16	Projection density estimator	75
	§16 01 Global and maximal global \mathfrak{v} -risk	76
	§16 02 Local and maximal local ϕ -risk	78
§17	Minimax optimal density estimation	80
	§17 01 Maximal local ϕ -risk	80
	§17 02 Maximal global \mathfrak{v} -risk	82
§18	Data-driven density estimation	83
	§18 01 Data-driven global estimation by model selection	83
	§18 02 Data-driven local estimation by Goldenshluger and Lepskij's method	87

§15 Noisy density coefficients

§15.01 **Notation (Reminder).** Consider the measure space $([0, 1], \mathcal{B}_{[0,1]}, \lambda_{[0,1]})$ where $\lambda_{[0,1]}$ denotes the restriction of the Lebesgue measure to the Borel- σ -algebra $\mathcal{B}_{[0,1]}$ over $[0, 1]$, and the Hilbert space $\mathbb{L}_2(\lambda_{[0,1]}) := \mathbb{L}_2([0, 1], \mathcal{B}_{[0,1]}, \lambda_{[0,1]})$ of square Lebesgue-integrable functions endowed with its usual inner product $\langle h_1, h_2 \rangle_{\mathbb{L}_2(\lambda_{[0,1]})} = \lambda_{[0,1]}(h_1 h_2)$ for all $h_1, h_2 \in \mathbb{L}_2(\lambda_{[0,1]})$. Let \mathbb{D}_2 be a set of square-integrable Lebesgue densities on $([0, 1], \mathcal{B}_{[0,1]})$, and hence $\mathbb{D}_2 \subseteq \mathbb{L}_2(\lambda_{[0,1]}) (=:\mathbb{H})$ as in Model §10.23. We denote for each density $\mathfrak{p} \in \mathbb{D}_2$ by $\mathbb{P}_{\mathfrak{p}} := \mathfrak{p} \lambda_{[0,1]}$ and $\mathbb{E}_{\mathfrak{p}}$ the associated probability measure and expectation, respectively. Keep in mind, that we identify equivalence classes and their representatives. \square

§15.02 **Assumption.** We consider the statistical product experiment $([0, 1]^n, \mathcal{B}_{[0,1]}^{\otimes n}, \mathbb{P}_{\mathfrak{p}}^{\otimes n} := (\mathbb{P}_{\mathfrak{p}}^{\otimes n})_{\mathfrak{p} \in \mathbb{D}_2})$ of size $n \in \mathbb{N}$ and for $\mathfrak{p} \in \mathbb{D}_2 \subseteq \mathbb{L}_2(\lambda_{[0,1]})$ we denote by $(X_i)_{i \in [n]} \sim \mathbb{P}_{\mathfrak{p}}^{\otimes n}$ an iid. sample of $X \sim \mathbb{P}_{\mathfrak{p}}$.

§15.03 **Notation (Reminder).** Consider an orthonormal system $(u_j)_{j \in \mathbb{N}}$ in $\mathbb{L}_2(\lambda_{[0,1]})$. Then $U : \mathbb{L}_2(\lambda_{[0,1]}) \rightarrow \ell_2$ with $h \mapsto Uh := h_{\bullet} = (h_j := \langle h, u_j \rangle_{\mathbb{L}_2(\lambda_{[0,1]})})_{j \in \mathbb{N}}$ is a surjective partial isometry $U \in \mathbb{L}(\mathbb{L}_2(\lambda_{[0,1]}), \ell_2)$. Its adjoint operator $U^* \in \mathbb{L}(\ell_2, \mathbb{L}_2(\lambda_{[0,1]}))$ satisfies $U^* a_{\bullet} = \sum_{j \in \mathbb{N}} a_j u_j =: \nu_{\mathbb{N}}(a_{\bullet}, u_{\bullet})$ for all $a_{\bullet} \in \ell_2$. We call $h_{\bullet} = (h_j)_{j \in \mathbb{N}}$ (*generalised*) *Fourier coefficients* and U (*generalised*) *Fourier series transform*. \square

§15.04 **Remark.** Let $U \in \mathbb{L}(\mathbb{L}_2(\lambda_{[0,1]}), \ell_2)$ be a generalised Fourier series transform as in Notation §15.03 where $\mathbb{L}_2(\lambda_{[0,1]}) = \ker(U) \oplus \text{ran}(U^*)$ and $\text{ran}(U^*) = \{U^* a_{\bullet} = \nu_{\mathbb{N}}(a_{\bullet}, u_{\bullet}) : a_{\bullet} \in \ell_2\}$. If U is not injective, then there exists $\mathcal{K} \subseteq \mathbb{N}$ and an orthonormal basis $(v_j)_{j \in \mathcal{K}}$ of $\ker(U)$, and each $h \in \mathbb{L}_2(\lambda_{[0,1]})$ with $h_{\bullet} := Uh$ admits an expansion $h = U^* h_{\bullet} + \sum_{j \in \mathcal{K}} \langle h, v_j \rangle_{\mathbb{L}_2(\lambda_{[0,1]})} v_j$. We denote by $\mathbb{1}_{[0,1]} \in \mathcal{B}_{[0,1]}$ the constant function with $x \mapsto \mathbb{1}_{[0,1]}(x) := 1$. If $\mathbb{1}_{[0,1]} \in \ker(U)$ then we have $\langle h, \mathbb{1}_{[0,1]} \rangle_{\mathbb{L}_2(\lambda_{[0,1]})} = 0$ for all $h \in \text{ran}(U^*)$, or in equal $\langle \nu_{\mathbb{N}}(a_{\bullet}, u_{\bullet}), \mathbb{1}_{[0,1]} \rangle_{\mathbb{L}_2(\lambda_{[0,1]})} = 0$ for all $a_{\bullet} \in \ell_2$, and in particular $\langle u_j, \mathbb{1}_{[0,1]} \rangle_{\mathbb{L}_2(\lambda_{[0,1]})} =$

0 for all $j \in \mathbb{N}$. For each density $\mathbb{p} \in \mathbb{L}_2(\lambda_{[0,1]})$ we have $\langle \mathbb{p}, \mathbb{1}_{[0,1]} \rangle_{\mathbb{L}_2(\lambda_{[0,1]})} = \lambda_{[0,1]}(\mathbb{p}) = 1$. In other words the coefficient $\langle \mathbb{p}, \mathbb{1}_{[0,1]} \rangle_{\mathbb{L}_2(\lambda_{[0,1]})}$ is always known. Therefore, we assume here and subsequently $\mathbb{1}_{[0,1]} \in \ker(U)$. Moreover we have $\mathbb{L}_2(\lambda_{[0,1]}) \subseteq \mathbb{L}_1([0,1], \mathcal{B}_{[0,1]}, \mathbb{P}_{\mathbb{p}}) =: \mathbb{L}_1(\mathbb{P}_{\mathbb{p}})$. Indeed, $h \in \mathbb{L}_2(\lambda_{[0,1]})$ satisfies $\mathbb{P}_{\mathbb{p}}(|h|) = \mathbb{p} \lambda_{[0,1]}(|h|) = \langle \mathbb{p}, |h| \rangle_{\mathbb{L}_2(\lambda_{[0,1]})} \leq \| \mathbb{p} \|_{\mathbb{L}_2(\lambda_{[0,1]})} \| h \|_{\mathbb{L}_2(\lambda_{[0,1]})} \in \mathbb{R}^+$, and hence $h \in \mathbb{L}_1(\mathbb{P}_{\mathbb{p}})$. Evidently, we have $u_j \in \mathbb{L}_1(\mathbb{P}_{\mathbb{p}})$ for all $j \in \mathbb{N}$ and the Fourier coefficients $\mathbb{p} = (\mathbb{p}_j)_{j \in \mathbb{N}} = U \mathbb{p} \in \ell_2$ of $\mathbb{p} \in \mathbb{L}_2(\lambda_{[0,1]})$ fulfil $\mathbb{p}_j = \langle \mathbb{p}, u_j \rangle_{\mathbb{L}_2(\lambda_{[0,1]})} = \lambda_{[0,1]}(\mathbb{p} u_j) = \mathbb{p} \lambda_{[0,1]}(u_j) = \mathbb{P}_{\mathbb{p}}(u_j)$ for all $j \in \mathbb{N}$. In addition we assume that for each $\mathbb{p} \in \mathbb{D}_2 \subseteq \mathbb{L}_2(\lambda_{[0,1]})$ the orthonormal system $(u_j)_{j \in \mathbb{N}}$ belongs also to $\mathbb{L}_2(\mathbb{P}_{\mathbb{p}}) := \mathbb{L}_2([0,1], \mathcal{B}_{[0,1]}, \mathbb{P}_{\mathbb{p}})$, i.e. $u_j \in \mathbb{L}_2(\mathbb{P}_{\mathbb{p}})$ for all $j \in \mathbb{N}$. \square

§15.05 **Assumption.** The orthonormal system $(u_j)_{j \in \mathbb{N}}$ in $\mathbb{L}_2(\lambda_{[0,1]})$, and its associated generalised Fourier series transform $U \in \mathbb{L}(\mathbb{L}_2(\lambda_{[0,1]}), \ell_2)$ with $h \mapsto Uh := \mathbf{h}_\bullet = (h_j := \langle h, u_j \rangle_{\mathbb{L}_2(\lambda_{[0,1]})})_{j \in \mathbb{N}}$, is fixed and known in advance. U is a partial isometry with **(os1)** $\mathbb{1}_{[0,1]} \in \ker(U)$. **(os2)** For all $\mathbb{p} \in \mathbb{D}_2 \subseteq \mathbb{L}_2(\lambda_{[0,1]})$ the orthonormal system $(u_j)_{j \in \mathbb{N}}$ belongs to $\mathbb{L}_2(\mathbb{P}_{\mathbb{p}})$. \square

§15.06 **Remark.** If in addition $\mathbb{D}_2 \subseteq \mathbb{L}_\infty(\lambda_{[0,1]})$ then for each $\mathbb{p} \in \mathbb{D}_2$ we have $\mathbb{L}_2(\lambda_{[0,1]}) \subseteq \mathbb{L}_2(\mathbb{P}_{\mathbb{p}})$. Indeed, $h \in \mathbb{L}_2(\lambda_{[0,1]})$ satisfies $\mathbb{P}_{\mathbb{p}}(|h|^2) = \lambda_{[0,1]}(\mathbb{p}|h|^2) \leq \| \mathbb{p} \|_{\mathbb{L}_\infty(\lambda_{[0,1]})} \| h \|_{\mathbb{L}_2(\lambda_{[0,1]})}^2 \in \mathbb{R}^+$. Consequently, any orthonormal system in $\mathbb{L}_2(\lambda_{[0,1]})$ belongs also to $\mathbb{L}_2(\mathbb{P}_{\mathbb{p}})$, and **(os2)** in Assumption §15.05 is satisfied. Alternatively, **(os2)** is fulfilled for arbitrary $\mathbb{D}_2 \subseteq \mathbb{L}_2(\lambda_{[0,1]})$ if $(u_j)_{j \in \mathbb{N}}$ belongs also to $\mathbb{L}_\infty(\lambda_{[0,1]})$. \square

§15.07 **Notation (Reminder).** Similar to an Empirical mean model §10.07 for each $j \in \mathbb{N}$ we define $\widehat{\mathbb{p}}_j := \widehat{\mathbb{P}}_n(u_j) \in \mathcal{B}_{[0,1]}^{\otimes n}$ with $x^n \mapsto (\widehat{\mathbb{P}}_n(u_j))(x^n) = n^{-1} \sum_{i \in [n]} u_j(x_i^n)$. Since the stochastic process $\mathbf{u}_\bullet = (u_j)_{j \in \mathbb{N}}$ on $([0,1], \mathcal{B}_{[0,1]})$ is $\mathcal{B}_{[0,1]} \otimes 2^{\mathbb{N}}$ - \mathcal{B} -measurable, the stochastic process $\widehat{\mathbb{p}}_\bullet = (\widehat{\mathbb{p}}_j := \widehat{\mathbb{P}}_n(u_j))_{j \in \mathbb{N}}$ on $([0,1]^n, \mathcal{B}_{[0,1]}^{\otimes n})$ is $\mathcal{B}_{[0,1]}^{\otimes n} \otimes 2^{\mathbb{N}}$ - \mathcal{B} -measurable, $\widehat{\mathbb{p}}_\bullet \in \mathcal{B}_{[0,1]}^{\otimes n} \otimes 2^{\mathbb{N}}$ for short. By construction $\mathbb{p}_\bullet = (\mathbb{p}_j = \mathbb{P}_{\mathbb{p}}(u_j))_{j \in \mathbb{N}} \in 2^{\mathbb{N}}$ is the ℓ_2 -mean of $\widehat{\mathbb{p}}_\bullet$. For each $j \in \mathbb{N}$ the statistic $\boldsymbol{\varepsilon}_j := n^{1/2}(\widehat{\mathbb{P}}_n(u_j) - \mathbb{P}_{\mathbb{p}}(u_j)) \in \mathcal{B}_{[0,1]}^{\otimes n}$ is centred, i.e. $\boldsymbol{\varepsilon}_j \in \mathbb{L}_1([0,1]^n, \mathcal{B}_{[0,1]}^{\otimes n}, \mathbb{P}_{\mathbb{p}}^{\otimes n}) =: \mathbb{L}_1(\mathbb{P}_{\mathbb{p}}^{\otimes n})$ with $\mathbb{P}_{\mathbb{p}}^{\otimes n}(\boldsymbol{\varepsilon}_j) = 0$, and $\boldsymbol{\varepsilon}_\bullet = (\boldsymbol{\varepsilon}_j)_{j \in \mathbb{N}} \in \mathcal{B}_{[0,1]}^{\otimes n} \otimes 2^{\mathbb{N}}$. Since $\widehat{\mathbb{p}}_j = \mathbb{p}_j + n^{-1/2} \boldsymbol{\varepsilon}_j$ for each $j \in \mathbb{N}$ by construction $\widehat{\mathbb{p}}_\bullet = \mathbb{p}_\bullet + n^{-1/2} \boldsymbol{\varepsilon}_\bullet$ is a noisy version of \mathbb{p}_\bullet (see Definition §10.19). Moreover, under Assumption §15.05 $\widehat{\mathbb{p}}_\bullet$ admits a covariance function $\text{cov}_{\bullet}^{\mathbb{p}} \in \mathbb{R}^{\mathbb{N}^2}$ given for $j, j_o \in \mathbb{N}$ by

$$n \text{Cov}(\widehat{\mathbb{p}}_j, \widehat{\mathbb{p}}_{j_o}) = \text{Cov}(\boldsymbol{\varepsilon}_j, \boldsymbol{\varepsilon}_{j_o}) = \mathbb{P}_{\mathbb{p}}^{\otimes n}(\boldsymbol{\varepsilon}_j \boldsymbol{\varepsilon}_{j_o}) = \mathbb{P}_{\mathbb{p}}(u_j u_{j_o}) - \mathbb{P}_{\mathbb{p}}(u_j) \mathbb{P}_{\mathbb{p}}(u_{j_o}) = \mathbb{P}_{\mathbb{p}}(u_j u_{j_o}) - \mathbb{p}_j \mathbb{p}_{j_o} =: \text{cov}_{j,j_o}^{\mathbb{p}}.$$

Consequently, we have $\boldsymbol{\varepsilon}_\bullet \sim P_{(0, \text{cov}_{\bullet}^{\mathbb{p}})}$ and $\widehat{\mathbb{p}}_\bullet = \mathbb{p}_\bullet + n^{-1/2} \boldsymbol{\varepsilon}_\bullet \sim P_{(\mathbb{p}_\bullet, n^{-1} \text{cov}_{\bullet}^{\mathbb{p}})}$ (see Definition §10.19). \square

§15.08 **Noisy density coefficients.** Under Assumptions §15.02 and §15.05 the stochastic process $\boldsymbol{\varepsilon}_\bullet = (\boldsymbol{\varepsilon}_j := n^{1/2}(\widehat{\mathbb{P}}_n(u_j) - \mathbb{P}_{\mathbb{p}}(u_j)))_{j \in \mathbb{N}}$ satisfies Assumption §10.04, i.e. $\boldsymbol{\varepsilon}_\bullet \in \mathcal{B}_{[0,1]}^{\otimes n} \otimes 2^{\mathbb{N}}$, and $\boldsymbol{\varepsilon}_\bullet$ has mean zero under $\mathbb{P}_{\mathbb{p}}^{\otimes n}$. The stochastic process $\widehat{\mathbb{p}}_\bullet = \mathbb{p}_\bullet + n^{-1/2} \boldsymbol{\varepsilon}_\bullet$ with ℓ_2 -mean \mathbb{p}_\bullet is called a *noisy version* of the density coefficients $\mathbb{p}_\bullet = U \mathbb{p} \in \ell_2$, or *noisy density coefficients* for short. Moreover $\boldsymbol{\varepsilon}_\bullet$ admits under $\mathbb{P}_{\mathbb{p}}^{\otimes n}$ a covariance function $\text{cov}_{\bullet}^{\mathbb{p}} \in \mathbb{R}^{\mathbb{N}^2}$ given for $j, j_o \in \mathbb{N}$ by $\text{cov}_{j,j_o}^{\mathbb{p}} = \mathbb{P}_{\mathbb{p}}(u_j u_{j_o}) - \mathbb{P}_{\mathbb{p}}(u_j) \mathbb{P}_{\mathbb{p}}(u_{j_o})$. We eventually write $\boldsymbol{\varepsilon}_\bullet \sim P_{(0, \text{cov}_{\bullet}^{\mathbb{p}})}$ and $\widehat{\mathbb{p}}_\bullet \sim P_{(\mathbb{p}_\bullet, n^{-1} \text{cov}_{\bullet}^{\mathbb{p}})}$. If in addition $\boldsymbol{\varepsilon}_\bullet$ admits a covariance operator $\Gamma_{\mathbb{p}} \in \mathbb{L}^{\mathbb{Z}}(\ell_2)$ then we write $\boldsymbol{\varepsilon}_\bullet \sim P_{(0, \Gamma_{\mathbb{p}})}$ and $\widehat{\mathbb{p}}_\bullet \sim P_{(\mathbb{p}_\bullet, n^{-1} \Gamma_{\mathbb{p}})}$ for short. \square

§15.09 **Remark.** The centred stochastic process $\boldsymbol{\varepsilon}_\bullet := (\boldsymbol{\varepsilon}_j)_{j \in \mathbb{N}}$ of error terms in Definition §15.08 is in general not a white noise process. \square

§15.10 **Lemma.** Under Assumptions §15.02 and §15.05 consider the stochastic process $\boldsymbol{\varepsilon}_\bullet \in \mathcal{B}_{[0,1]}^{\otimes n} \otimes 2^{\mathbb{N}}$ as in Definition §15.08.

(i) If $\mathbb{p} \in \mathbb{L}_\infty(\lambda_{[0,1]})$ then under $\mathbb{P}_{\mathbb{p}}^{\otimes n}$, $\boldsymbol{\varepsilon}_\bullet \sim P_{(0, \text{cov}_{\bullet}^{\mathbb{p}})}$ admits a covariance operator $\Gamma_{\mathbb{p}} \in \mathbb{L}^{\mathbb{Z}}(\ell_2)$ given by

$$a_\bullet \mapsto \Gamma_{\mathbb{p}} a_\bullet = (\nu_{\mathbb{N}}(\text{cov}_{j,j_o}^{\mathbb{p}} a_\bullet))_{j_o \in \mathbb{N}} = \sum_{j_o \in \mathbb{N}} \text{cov}_{j,j_o}^{\mathbb{p}} a_{j_o}$$

where $\|\Gamma_{\mathbb{p}}\|_{\mathbb{L}(\ell_2)} \leq \|\mathbb{p}\|_{\mathbb{L}_{\infty}(\lambda_{[0,1]})}$.

(ii) If $\mathbb{p} \in \mathbb{L}_{\infty}(\lambda_{[0,1]})$ and $\mathbb{p}^{-1} := \frac{1}{\mathbb{p}} \in \mathbb{L}_{\infty}(\lambda_{[0,1]})$ then $\Gamma_{\mathbb{p}} \in \mathbb{L}(\ell_2)$ is invertible with inverse $\Gamma_{\mathbb{p}}^{-1} \in \mathbb{L}(\ell_2)$ where $\|\Gamma_{\mathbb{p}}^{-1}\|_{\mathbb{L}(\ell_2)} \leq \|\mathbb{p}^{-1}\|_{\mathbb{L}_{\infty}(\lambda_{[0,1]})}$.

Consequently, if $v_{\mathbb{p}} := \max(\|\mathbb{p}\|_{\mathbb{L}_{\infty}(\lambda_{[0,1]})}, \|\mathbb{p}^{-1}\|_{\mathbb{L}_{\infty}(\lambda_{[0,1]})}) \in \mathbb{R}_{>0}^+$ then for all $a_{\bullet} \in \ell_2$ we have

$$v_{\mathbb{p}}^{-1} \|a_{\bullet}\|_{\ell_2}^2 \leq \|a_{\bullet}\|_{\Gamma_{\mathbb{p}}}^2 = \langle \Gamma_{\mathbb{p}} a_{\bullet}, a_{\bullet} \rangle_{\ell_2} \leq v_{\mathbb{p}} \|a_{\bullet}\|_{\ell_2}^2.$$

§15.11 **Proof of Lemma §15.10.** is given in the lecture. □

§15.12 **Remark.** For each $j \in \mathbb{N}$ consider $\mathbf{1}_{\cdot}^{(j)} = (\mathbf{1}_{\{j\}}(l))_{l \in \mathbb{N}} \in (\mathbb{R})^{\mathbb{N}}$ where $(\mathbf{1}_{\cdot}^{(j)})_{j \in \mathbb{N}}$ forms an orthonormal basis in ℓ_2 . If $\mathbb{p} \in \mathbb{L}_{\infty}(\lambda_{[0,1]})$ from **Lemma §15.10 (i)** for each $j \in \mathbb{N}$ we obtain

$$\mathbb{P}_{\mathbb{p}}^{\otimes n}(\varepsilon_j^2) = \mathbb{P}_{\mathbb{p}}^{\otimes n}(|\nu_{\mathbb{N}}(\mathbf{1}_{\cdot}^{(j)} \varepsilon_{\cdot})|^2) = \langle \Gamma_{\mathbb{p}} \mathbf{1}_{\cdot}^{(j)}, \mathbf{1}_{\cdot}^{(j)} \rangle_{\ell_2} \leq \|\mathbb{p}\|_{\mathbb{L}_{\infty}(\lambda_{[0,1]})} \|\mathbf{1}_{\cdot}^{(j)}\|_{\ell_2}^2 = \|\mathbb{p}\|_{\mathbb{L}_{\infty}(\lambda_{[0,1]})}$$

Keeping the last identities in mind if $v_{\mathbb{p}} := \max(\|\mathbb{p}\|_{\mathbb{L}_{\infty}(\lambda_{[0,1]})}, \|\mathbb{p}^{-1}\|_{\mathbb{L}_{\infty}(\lambda_{[0,1]})}) \in \mathbb{R}_{>0}^+$ then due to **Lemma §15.10** for all $j \in \mathbb{N}$ we have $v_{\mathbb{p}}^{-1} \leq \mathbb{P}_{\mathbb{p}}^{\otimes n}(\varepsilon_j^2) \leq v_{\mathbb{p}}$. □

§16 Projection density estimator

§16.01 **Notation (Reminder).** Consider the measure space $(\mathbb{N}, 2^{\mathbb{N}}, \nu_{\mathbb{N}})$ as in **Notation §10.11**. For $w_{\bullet} \in \mathbb{R}^{\mathbb{N}}$ define the multiplication map $M_{w_{\bullet}} : \mathbb{R}^{\mathbb{N}} \rightarrow \mathbb{R}^{\mathbb{N}}$ with $a_{\bullet} \mapsto M_{w_{\bullet}} a_{\bullet} := w_{\bullet} a_{\bullet}$. Note that each $w_{\bullet} \in \mathbb{R}^{\mathbb{N}}$ is $2^{\mathbb{N}}$ - \mathcal{B} -measurable. We denote by $\mathbb{M}_{\mathbb{R}^{\mathbb{N}}}$ the set of all multiplication maps defined on $\mathbb{R}^{\mathbb{N}}$. If in addition $w_{\bullet} \in \ell_{\infty} = \mathbb{L}_{\infty}(\mathbb{N}, 2^{\mathbb{N}}, \nu_{\mathbb{N}})$ then we have also $M_{w_{\bullet}} : \ell_2 \rightarrow \ell_2$. We set $\mathbb{L}^{\mathbb{M}(\ell_2)} = \{M_{w_{\bullet}} \in \mathbb{M}_{\mathbb{R}^{\mathbb{N}}} : w_{\bullet} \in \ell_{\infty}\} \subseteq \mathbb{L}(\ell_2)$ noting that $\|M_{w_{\bullet}}\|_{\mathbb{L}(\ell_2)} = \sup\{\|w_{\bullet} a_{\bullet}\|_{\ell_2} : \|a_{\bullet}\|_{\ell_2} \leq 1\} \leq \|w_{\bullet}\|_{\ell_{\infty}}$ for each $M_{w_{\bullet}} \in \mathbb{L}^{\mathbb{M}(\ell_2)}$. □

§16.02 **Reminder.** If $w_{\bullet} \in \ell_{\infty}$ then $M_{w_{\bullet}} \in \mathbb{L}^{\mathbb{M}(\ell_2)}$, and $M_{w_{\bullet}^{\dagger}} : \ell_2 \supseteq \text{dom}(M_{w_{\bullet}}) \rightarrow \ell_2$. Moreover, we have $\text{dom}(M_{w_{\bullet}}) = \ell_2$, $\text{ran}(M_{w_{\bullet}}) = \ell_2 w_{\bullet}$ and $\ker(M_{w_{\bullet}}) = \ell_2 \mathbf{1}_{\cdot}^{\mathcal{N}_w}$ with $\mathcal{N}_w = \{j \in \mathbb{N} : w_j = 0\} \in 2^{\mathbb{N}}$ (see **Property §11.03**), and $\text{dom}(M_{w_{\bullet}^{\dagger}}) = \ell_2 w_{\bullet} \oplus \ell_2 \mathbf{1}_{\cdot}^{\mathcal{N}_w}$ (see **Property §11.05**). Consequently, if in addition $\nu_{\mathbb{N}}(\mathcal{N}_w) = 0$ or in equal $w_{\bullet} \in (\mathbb{R}_{>0})^{\mathbb{N}}$, then $w_{\bullet}^{\dagger} = w_{\bullet}^{-1} \in (\mathbb{R}_{>0})^{\mathbb{N}}$, hence $w_{\bullet}^{2\dagger} = w_{\bullet}^{-2} \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$, and $\ell_2^w = \text{dom}(M_{w_{\bullet}^{-1}}) = \ell_2 w_{\bullet} = \mathbb{L}_2(w_{\bullet}^{-2} \nu_{\mathbb{N}}) =: \ell_2(w_{\bullet}^{-2})$. For each $m \in \mathbb{N}$ we write $\mathbf{1}_{\cdot}^m = (\mathbf{1}_j^m)_{j \in \mathbb{N}} := \mathbf{1}_{[m]}$ and $\mathbf{1}_{\cdot}^{m\perp} := \mathbf{1}_{\cdot} - \mathbf{1}_{\cdot}^m$ with $[m] := [-m, m] \cap \mathbb{N}$. Consequently, $M_{\mathbf{1}_{\cdot}^m} \in \mathbb{L}(\ell_2)$ and $M_{\mathbf{1}_{\cdot}^{m\perp}} \in \mathbb{L}(\ell_2)$ is the *orthogonal projection* onto the linear subspace $\ell_2 \mathbf{1}_{\cdot}^m \subseteq \ell_2$ and its orthogonal complement $\ell_2 \mathbf{1}_{\cdot}^{m\perp} = (\ell_2 \mathbf{1}_{\cdot}^m)^{\perp} \subseteq \ell_2$, respectively, that is $\ell_2 = \ell_2 \mathbf{1}_{\cdot}^m \oplus \ell_2 \mathbf{1}_{\cdot}^{m\perp}$ (see **Property §11.07**). Finally, given $h_{\bullet} = U h \in \ell_2$ for $h \in \mathbb{L}_2(\lambda_{[0,1]})$ we consider the orthogonal projections $h_{\bullet}^m = h_{\bullet} \mathbf{1}_{\cdot}^m \in \ell_2 \mathbf{1}_{\cdot}^m$ and $h_{\bullet}^m := U^* h_{\bullet}^m \in \mathbb{L}_2(\lambda_{[0,1]})$ (**Definition §11.08**). □

§16.03 **Notation (Reminder).** Consider the stochastic processes $\varepsilon_{\cdot} = (\varepsilon_j := n^{1/2}(\widehat{\mathbb{P}}_n(u_j) - \mathbb{P}_{\mathbb{p}}(u_j)))_{j \in \mathbb{N}}$ given in **Definition §15.08**. The observable noisy density coefficients $\widehat{\mathbb{p}} = \mathbb{p} + n^{-1/2} \varepsilon_{\cdot}$ of the density coefficients $\mathbb{p} = U \mathbb{p} \in \ell_2$ take the form of a *statistical direct problem* (see **Definition §10.19**). Under Assumptions §15.02 and §15.05 ε_{\cdot} is centred and admits a covariance function $\text{cov}_{\cdot}^{\mathbb{p}} \in \mathbb{R}^{\mathbb{N}^2}$ given in **Definition §15.08**, i.e. $\varepsilon_{\cdot} \sim P_{(0, \text{cov}_{\cdot}^{\mathbb{p}})}$ and $\widehat{\mathbb{p}} \sim P_{(\mathbb{p}, n^{-1/2} \text{cov}_{\cdot}^{\mathbb{p}})}$. If in addition $\mathbb{p} \in \mathbb{L}_{\infty}(\lambda_{[0,1]})$ then ε_{\cdot} admits a covariance operator $\Gamma_{\mathbb{p}} \in \mathbb{L}(\ell_2)$ given in **Lemma §15.10**, i.e. $\varepsilon_{\cdot} \sim P_{(0, \Gamma_{\mathbb{p}})}$ and $\widehat{\mathbb{p}} \sim P_{(\mathbb{p}, n^{-1/2} \Gamma_{\mathbb{p}})}$. □

§16.04 **Definition.** Given a noisy version $\widehat{\mathbb{p}} = \mathbb{p} + n^{-1/2} \varepsilon_{\cdot}$ of the density coefficients $\mathbb{p} = U \mathbb{p} \in \ell_2$ for each $m \in \mathbb{N}$ we call $\widehat{\mathbb{p}}^m := \widehat{\mathbb{p}} \mathbf{1}_{\cdot}^m$ *orthogonal projection estimator (OPE)* of \mathbb{p} . □

§16.05 **Remark.** Under Assumptions §15.02 and §15.05 we consider the function (with random coefficients) $\widehat{\mathbb{p}}^m := \mathbb{1}_{[0,1]} + U^* \widehat{\mathbb{p}}^m$ which belongs to $\mathbb{L}_2(\lambda_{0,1})$, integrates to one, but may take on negative values. Fortunately, there is a simple remedy — its $\mathbb{L}_2(\lambda_{0,1})$ -projection onto a class of nonnegative densities,

$$\widehat{\mathbb{p}} := (\widehat{\mathbb{p}}^m - c)_+ \quad \text{with } c \in \mathbb{R}^+ \text{ such that } \lambda_{0,1}(\widehat{\mathbb{p}}) = 1.$$

We call $\widehat{\mathbb{p}}^m$ an *orthogonal projection density estimator* of \mathbb{p} . If $\mathbb{p} = \mathbb{1}_{[0,1]} + U^* \mathbb{p}$ (for example $\ker(U)$ is spanned by $u_0 := \mathbb{1}_{[0,1]}$ or in equal $(u_j)_{j \in \mathbb{N}_0}$ is an orthonormal basis of $\mathbb{L}_2(\lambda_{0,1})$), then we have

$$\|\widehat{\mathbb{p}}^m - \mathbb{p}\|_{\mathbb{L}_2(\lambda_{0,1})}^2 = \|\widehat{\mathbb{p}}^m - \mathbb{p}\|_{\ell_2}^2.$$

In this situation all results for the OPE $\widehat{\mathbb{p}}^m$ of the density coefficients immediately transfer onto the orthogonal projection density estimator $\widehat{\mathbb{p}}^m$ of the density \mathbb{p} . \square

§16|01 Global and maximal global \mathfrak{v} -risk

We measure first the accuracy of the OPE $\widehat{\mathbb{p}}^m = \widehat{\mathbb{p}} \mathbb{1}^m$ of $\mathbb{p}^m = \mathbb{p} \mathbb{1}^m \in \ell_2 \mathbb{1}^m$ with $\mathbb{p} = U \mathbb{p} \in \ell_2$ by a global mean- \mathfrak{v} -error, i.e. \mathfrak{v} -risk.

§16.06 **Reminder.** If $\mathfrak{v} \in (\mathbb{R}_{>0})^{\mathbb{N}}$ and $\mathbb{p} \in \ell_2(\mathfrak{v}^2)$ then we have $\mathbb{p}^m = \mathbb{p} \mathbb{1}^m \in \ell_2(\mathfrak{v}^2)$ too and $\|\mathbb{p}^m - \mathbb{p}\|_{\mathfrak{v}}^2 = o(1)$ as $m \rightarrow \infty$ (**Property** §11.09). Moreover, $\boldsymbol{\varepsilon} \in \mathcal{B}_{[0,1]}^{\otimes n} \otimes 2^{\mathbb{N}}$ given in **Definition** §15.08 satisfies $\mathfrak{v} \boldsymbol{\varepsilon} \mathbb{1}^m \in \ell_2$ (note that $\mathbb{1}^m \in \ell_2$ and $\mathfrak{v} \mathbb{1}^m, \boldsymbol{\varepsilon} \mathbb{1}^m \in \ell_\infty$) and thus also

$$n^{-1/2} \mathfrak{v} \boldsymbol{\varepsilon} \mathbb{1}^m + \mathfrak{v} \widehat{\mathbb{p}}^m = \mathfrak{v} \widehat{\mathbb{p}}^m \in \ell_2. \quad (16.01)$$

Finally, under Assumptions §15.02 and §15.05 and $\mathbb{p} \in \mathbb{L}_\infty(\lambda_{0,1})$ due to **Lemma** §15.10 we have $\mathbb{P}_p^{\otimes n}(\boldsymbol{\varepsilon}^2) \in \ell_\infty$, more precisely, $\|\mathbb{P}_p^{\otimes n}(\boldsymbol{\varepsilon}^2)\|_{\ell_\infty} \leq \|\mathbb{p}\|_{\mathbb{L}_\infty(\lambda_{0,1})}$ (see **Remark** §15.12). \square

§16|01|01 Global \mathfrak{v} -risk

§16.07 **Proposition (Upper bound).** Let Assumptions §15.02 and §15.05, $\mathfrak{v} \in (\mathbb{R}_{>0})^{\mathbb{N}}$ and $\mathbb{p} \in \ell_2(\mathfrak{v}^2)$ be satisfied and for all $n, m \in \mathbb{N}$ set

$$\begin{aligned} R_n^m(\mathbb{p}, \mathfrak{v}) &:= \|\mathbb{p} \mathbb{1}^{m \perp}\|_{\mathfrak{v}}^2 + n^{-1} \|\mathbb{1}^m\|_{\mathfrak{v}}^2, \quad m_n^\circ := \arg \min \{R_n^m(\mathbb{p}, \mathfrak{v}) : m \in \mathbb{N}\} \\ \text{and } R_n^\circ(\mathbb{p}, \mathfrak{v}) &:= R_{m_n^\circ}^m(\mathbb{p}, \mathfrak{v}) = \min \{R_n^m(\mathbb{p}, \mathfrak{v}) : m \in \mathbb{N}\}. \end{aligned} \quad (16.02)$$

If $\mathbb{p} \in \mathbb{L}_\infty(\lambda_{0,1})$ then we have $\mathbb{P}_p^{\otimes n}(\|\widehat{\mathbb{p}}^{m_n^\circ} - \mathbb{p}\|_{\mathfrak{v}}^2) \leq \|\mathbb{p}\|_{\mathbb{L}_\infty(\lambda_{0,1})} R_n^\circ(\mathbb{p}, \mathfrak{v})$.

§16.08 **Proof** of **Proposition** §16.07. is given in the lecture. \square

§16.09 **Oracle inequality.** Under Assumptions §15.02 and §15.05 let $\mathfrak{v} \in (\mathbb{R}_{>0})^{\mathbb{N}}$ and $\mathbb{p} \in \ell_2(\mathfrak{v}^2)$. If in addition $\mathfrak{v}_p := \max(\|\mathbb{p}\|_{\mathbb{L}_\infty(\lambda_{0,1})}, \|\mathbb{p}^{-1}\|_{\mathbb{L}_\infty(\lambda_{0,1})}) \in \mathbb{R}_{>0}^+$ then $\mathfrak{v}_p^{-1} \leq \mathfrak{v}_p^p := \mathbb{P}_p^{\otimes n}(\boldsymbol{\varepsilon}_j^2) \leq \mathfrak{v}_p$ for all $j \in \mathbb{N}$ (see **Remark** §15.12), and hence **Property** §12.15 implies

$$\begin{aligned} \mathfrak{v}_p^{-1} R_n^m(\mathbb{p}, \mathfrak{v}) &\leq \mathbb{P}_p^{\otimes n}(\|\widehat{\mathbb{p}}^m - \mathbb{p}\|_{\mathfrak{v}}^2) = n^{-1} \nu_n(\mathfrak{v}^p \mathfrak{v}^2 \mathbb{1}^m) + \|\mathbb{p} \mathbb{1}^{m \perp}\|_{\mathfrak{v}}^2 \\ &\leq \mathfrak{v}_p R_n^m(\mathbb{p}, \mathfrak{v}) \quad \text{for all } m, n \in \mathbb{N}. \end{aligned}$$

As a consequence we immediately obtain the following *oracle inequality* (see **Definition** §12.14)

$$\begin{aligned} \mathfrak{v}_p^{-1} R_n^\circ(\mathbb{p}, \mathfrak{v}) &\leq \inf_{m \in \mathbb{N}} \mathbb{P}_p^{\otimes n}(\|\widehat{\mathbb{p}}^m - \mathbb{p}\|_{\mathfrak{v}}^2) \leq \mathbb{P}_p^{\otimes n}(\|\widehat{\mathbb{p}}^{m_n^\circ} - \mathbb{p}\|_{\mathfrak{v}}^2) \\ &\leq \mathfrak{v}_p R_n^\circ(\mathbb{p}, \mathfrak{v}) \leq \mathfrak{v}_p^2 \inf_{m \in \mathbb{N}} \mathbb{P}_p^{\otimes n}(\|\widehat{\mathbb{p}}^m - \mathbb{p}\|_{\mathfrak{v}}^2), \end{aligned} \quad (16.03)$$

and, hence $R_n^\circ(\mathbb{p}, \mathbf{v})$, m_n° and the statistic $\widehat{\mathbb{p}}^{m_n^\circ}$, respectively, is an *oracle bound*, an *oracle dimension* and *oracle optimal* (up to the constant v_p^2). We observe that $R_n^\circ(\mathbb{p}, \mathbf{v}) = o(1)$ as $n \rightarrow \infty$ (*Remark §12.16*), and thus, $R_n^\circ(\mathbb{p}, \mathbf{v})$ is an *oracle rate*. However, note that the oracle dimension $m_n^\circ = m_n^\circ(\mathbb{p}, \mathbf{v})$ depends on the unknown density coefficients \mathbb{p} , and thus also the oracle optimal statistic $\widehat{\mathbb{p}}^{m_n^\circ}$. In other words $\widehat{\mathbb{p}}^{m_n^\circ}$ is not a feasible estimator. \square

§16.10 **Illustration.** We illustrate the last results considering usual behaviour for the bias and variance term. We distinguish the following two cases

(p) $\mathbf{v} \in \ell_2$ or there is $m \in \mathbb{N}$ with $\|\mathbb{p}^m - \mathbb{p}\|_{\mathbf{v}}^2 = 0$,

(np) $\mathbf{v} \notin \ell_2$ and for all $m \in \mathbb{N}$ holds $\|\mathbb{p}^m - \mathbb{p}\|_{\mathbf{v}}^2 \in \mathbb{R}_0^+$.

Interestingly, in case **(p)** the oracle bound is parametric, that is, $nR_n^\circ(\mathbb{p}, \mathbf{v}) = O(1)$, in case **(np)** the oracle bound is nonparametric, i.e. $\lim_{n \rightarrow \infty} nR_n^\circ(\mathbb{p}, \mathbf{v}) = \infty$. In case **(np)** consider the following two specifications:

Table 01 [§16]

Order of the oracle rate $R_n^\circ(\mathbb{p}, \mathbf{v})$ as $n \rightarrow \infty$

$(j \in \mathbb{N})$	$(a \in \mathbb{R}_0^+)$	(squared bias)	(variance)	m_n°	$R_n^\circ(\mathbb{p}, \mathbf{v})$
$\mathbf{v}_j^2 = j^{2v}$	\mathbb{p}^2	$\ \mathbb{p} \mathbf{1}^{m \perp}\ _{\mathbf{v}}^2$	$\ \mathbf{1}^m\ _{\mathbf{v}}^2$		
(o) $v \in (-1/2, a)$	j^{-2a-1}	$m^{-2(a-v)}$	m^{2v+1}	$n^{\frac{1}{2a+1}}$	$n^{-\frac{2(a-v)}{2a+1}}$
$v = -1/2$	j^{-2a-1}	m^{-2a-1}	$\log m$	$\left(\frac{n}{\log n}\right)^{\frac{1}{2a+1}}$	$\frac{\log n}{n}$
(s) $v + 1/2 \in \mathbb{R}_0^+$	$e^{-j^{2a}}$	$m^{(1-2(a-v))+} e^{-m^{2a}}$	m^{2v+1}	$(\log n)^{\frac{1}{2a}}$	$\frac{(\log n)^{\frac{2v+1}{2a}}}{n}$
$v = -1/2$	$e^{-j^{2a}}$	$e^{-m^{2a}}$	$\log m$	$(\log n)^{\frac{1}{2a}}$	$\frac{\log \log n}{n}$

We note that in Table 01 [§16] the order of the oracle rate $R_n^\circ(\mathbb{p}, \mathbf{v})$ is depict for $v \geq -1/2$ only. In case $v < -1/2$ the oracle rate $R_n^\circ(\mathbb{p}, \mathbf{v})$ is parametric. \square

§16|01|02 Maximal global v-risk

§16.11 **Assumption.** Consider weights $\mathbf{a}, \mathbf{v} \in (\mathbb{R}_0^+)^{\mathbb{N}}$ with $\mathbf{a} \in \ell_\infty$ and $(\mathbf{a}\mathbf{v}) := (\mathbf{a}_j \mathbf{v}_j)_{j \in \mathbb{N}} = \mathbf{a} \cdot \mathbf{v} \in \ell_\infty$.

We write $(\mathbf{a}\mathbf{v})_{(m)} := \|(\mathbf{a}\mathbf{v}) \cdot \mathbf{1}^{m|\perp}\|_{\ell_\infty} \in \mathbb{R}^+$ for each $m \in \mathbb{N}$. The orthonormal system $(\mathbf{u}_j)_{j \in \mathbb{N}}$ in $\mathbb{L}_2(\lambda_{[0,1]})$ and $\mathbf{u}_0 := \mathbf{1}_{[0,1]}$ form an **(os1')** orthonormal basis $(\mathbf{u}_j)_{j \in \mathbb{N}_0}$ in $\mathbb{L}_2(\lambda_{[0,1]})$ and as process $\mathbf{u}^2 = (\mathbf{u}_j^2)_{j \in \mathbb{N}}$ on $([0, 1], \mathcal{B}_{[0,1]})$ satisfies **(os2')** $\|\nu_{\mathbb{N}}(\mathbf{a}^2 \mathbf{u}^2)\|_{\mathbb{L}_\infty(\lambda_{[0,1]})} =: \tau_{\mathbf{a}, \mathbf{u}}^2 \in \mathbb{R}^+$. \square

§16.12 **Reminder.** Under Assumption §16.11 we have $\ell_2^{\mathbf{a}} = \text{dom}(M_{\mathbf{a}^{-1}}) = \ell_2 \mathbf{a} \subseteq \ell_2$ and the three measures $\nu_{\mathbb{N}}$, $\mathbf{a}^{-2} \nu_{\mathbb{N}}$ and $\mathbf{v}^2 \nu_{\mathbb{N}}$ dominate mutually each other, i.e. they share the same null sets (see *Property §11.05*). We consider $\ell_2^{\mathbf{a}}$ endowed with $\|\cdot\|_{\mathbf{a}^{-1}} = \|M_{\mathbf{a}^{-1}} \cdot\|_{\ell_2}$ and given a constant $r \in \mathbb{R}_0^+$ the ellipsoid $\ell_2^{\mathbf{a}, r} := \{\mathbf{a} \in \ell_2^{\mathbf{a}} : \|\mathbf{a}\|_{\mathbf{a}^{-1}} \leq r\} \subseteq \ell_2^{\mathbf{a}}$. Since $(\mathbf{a}\mathbf{v}) \in \ell_\infty$, and hence $(\mathbf{a}\mathbf{v})_{(m)} := \|(\mathbf{a}\mathbf{v}) \cdot \mathbf{1}^{m|\perp}\|_{\ell_\infty} \in \mathbb{R}^+$ for each $m \in \mathbb{N}$ we have $\ell_2^{\mathbf{a}} \subseteq \ell_2(\mathbf{v}^2)$ (*Property §11.15*), and $\|\mathbf{a} \cdot \mathbf{1}^{m|\perp}\|_{\mathbf{v}} \leq r (\mathbf{a}\mathbf{v})_{(m)}$ for all $\mathbf{a} \in \ell_2^{\mathbf{a}, r}$ (*Lemma §11.17*). \square

§16.13 **Remark.** We replace Assumption §15.05 **(os1)** and **(os2)**, respectively, by the stronger Assumption §16.11 **(os1')** and **(os2')**. Indeed, under **(os1')** we have **(os1)** $\mathbf{1}_{[0,1]} \in \ker(U)$. Furthermore, $(\mathbf{u}_j)_{j \in \mathbb{N}}$ belongs to $\mathbb{L}_\infty(\lambda_{[0,1]})$ due to **(os2')** (and $\mathbf{a} \in (\mathbb{R}_0^+)^{\mathbb{N}}$), and hence **(os2)** is fulfilled (see also *Remark §15.06*). \square

§16.14 **Lemma.** Under Assumption §16.11 set

$$\mathbb{D}_2^{\mathbf{a}, r} := \{\mathbb{p} \in \mathbb{L}_2(\lambda_{[0,1]}): \mathbb{p} \text{ is a density and } \mathbb{p} = U \mathbb{p} \in \ell_2^{\mathbf{a}, r}\}. \quad (16.04)$$

Then we have $\sup \{ \|\mathbb{p}\|_{\mathbb{L}_\infty(\lambda_{0,1})} : \mathbb{p} \in \mathbb{D}_2^{a,r} \} \leq 1 + r\tau_{a,u}$.

§16.15 **Proof** of **Lemma** §16.14. is given in the lecture. \square

§16.16 **Proposition (Upper bound)**. Let Assumptions §15.02 and §16.11 be satisfied. For $n \in \mathbb{N}$ considering $m_n^* \in \mathbb{N}$ and $\mathbb{R}_n^*(\mathbf{a}, \mathbf{v}) \in \mathbb{R}^+$ as in (12.06) (**Proposition** §12.21) we have

$$\sup \{ \mathbb{P}_p^{\otimes n} (\|\widehat{\mathbb{p}}^{m_n^*} - \mathbb{p}\|_{\mathbb{V}}^2) : \mathbb{p} \in \mathbb{D}_2^{a,r} \} \leq C \mathbb{R}_n^*(\mathbf{a}, \mathbf{v}).$$

with constant $C = 1 + r\tau_{a,u} + r^2$.

§16.17 **Proof** of **Proposition** §16.16. is given in the lecture. \square

§16.18 **Illustration**. The *trigonometric basis* given for $x \in [0, 1]$ by

$$u_0 := \mathbb{1}_{[0,1]}, \quad u_{2k}(x) := \sqrt{2} \cos(2\pi kt), \quad u_{2k-1}(x) := \sqrt{2} \sin(2\pi kt), \quad k \in \mathbb{N},$$

is an orthonormal basis of $\mathbb{L}_2(\lambda_{0,1})$. It satisfies Assumption §16.11 (**os1'**), since $\|u_j\|_{\mathbb{L}_2(\lambda_{0,1})}^2 = 2$ for all $j \in \mathbb{N}$. Consequently, for all $\mathbf{a} \in \ell_2$ also the Assumption §16.11 (**os2'**) is satisfied with $\tau_{a,u}^2 \leq 2\|\mathbf{a}\|_{\ell_2}^2$.

(o) If $\mathbf{a}_{2j} = \mathbf{a}_{2j-1} = j^{-a}$, $a \in \mathbb{N}$, $j \in \mathbb{N}$, then $\{h \in \mathbb{L}_2(\lambda_{0,1}) : Uh \in \ell_2^a\}$ is a subset of the *Sobolev space* of a -times differentiable periodic functions. Moreover, up to a constant, for any function $h \in \mathbb{L}_2(\lambda_{0,1})$ the weighted norm $\|h\|_{\ell_2^a}^2$ equals the \mathbb{L}_2 -norm of its a -th weak derivative $h^{(a)}$ (Tsybakov [2009]).

(s) If $\mathbf{a}_j = \exp(-j^{2a})$, $a > 1/2$, $j \in \mathbb{N}$, then $\{h \in \mathbb{L}_2(\lambda_{0,1}) : Uh \in \ell_2^a\}$ is a *class of analytic functions* (Kawata [1972]).

In Table 02 [§12] (**Illustration** §12.26) the order of the rate $\mathbb{R}_n^*(\mathbf{a}, \mathbf{v})$ is depicted for the two cases (o) and (s). We note that we have $\mathbf{a} \in \ell_2$ in case (o) for $a > 1/2$ while in case (s) for $a \in \mathbb{R}_0^+$. \square

§16|02 Local and maximal local ϕ -risk

We measure secondly the accuracy of the OPE $\widehat{\mathbb{p}}^m = \widehat{\mathbb{p}} \mathbb{1}^m$ of $\mathbb{p}^m = \mathbb{p} \mathbb{1}^m \in \ell_2 \mathbb{1}^m$ with $\mathbb{p} = U \mathbb{p} \in \ell_2$ by a local mean- ϕ -error, i.e. ϕ -risk.

§16.19 **Reminder**. If $\phi \in (\mathbb{R}_{\setminus 0})^{\mathbb{N}}$ and $\mathbb{p} \in \text{dom}(\phi_{\mathbb{N}}) := \{a \in \ell_2 : \phi a \in \ell_1\}$, then we have $\mathbb{p}^m = \mathbb{p} \mathbb{1}^m \in \text{dom}(\phi_{\mathbb{N}})$ too and $|\phi_{\mathbb{N}}(\mathbb{p} - \mathbb{p}^m)| = o(1)$ as $m \rightarrow \infty$ (**Property** §11.22). Moreover, $\boldsymbol{\varepsilon} \in \mathcal{B}_{[0,1]}^{\otimes n} \otimes 2^{\mathbb{N}}$ given in **Definition** §15.08 satisfies $\boldsymbol{\varepsilon} \mathbb{1}^m \in \text{dom}(\phi_{\mathbb{N}})$ (note that $\phi \mathbb{1}^m, \boldsymbol{\varepsilon} \mathbb{1}^m \in \ell_2$) and thus also

$$n^{-1/2} \boldsymbol{\varepsilon} \mathbb{1}^m + \mathbb{p}^m = \widehat{\mathbb{p}}^m \in \text{dom}(\phi_{\mathbb{N}}). \quad (16.05)$$

Finally, under Assumptions §15.02 and §15.05 and $\mathbb{p} \in \mathbb{L}_\infty(\lambda_{0,1})$ due to **Lemma** §15.10 (i) the process $\boldsymbol{\varepsilon} \in \mathcal{B}_{[0,1]}^{\otimes n} \otimes 2^{\mathbb{N}}$ admits a covariance operator $\Gamma_{\mathbb{p}} \in \mathbb{L}(\ell_2)$, i.e. $\boldsymbol{\varepsilon} \sim P_{(0,\Gamma)}$, satisfying $\|\Gamma_{\mathbb{p}}\|_{\mathbb{L}(\ell_2)} \leq \|\mathbb{p}\|_{\mathbb{L}_\infty(\lambda_{0,1})}$. \square

§16|02|01 Local ϕ -risk

§16.20 **Proposition (Upper bound)**. Let Assumptions §15.02 and §15.05, $\phi \in (\mathbb{R}_{\setminus 0})^{\mathbb{N}}$ and $\mathbb{p} \in \text{dom}(\phi_{\mathbb{N}})$ be satisfied and for all $n, m \in \mathbb{N}$ set

$$\begin{aligned} \mathbb{R}_n^m(\mathbb{p}, \phi) &:= |\phi_{\mathbb{N}}(\mathbb{p} \mathbb{1}^{m\perp})|^2 + n^{-1} \|\mathbb{1}^m\|_{\phi}^2, \quad m_n^\circ := \arg \min \{ \mathbb{R}_n^m(\mathbb{p}, \phi) : m \in \mathbb{N} \} \\ \text{and} \quad \mathbb{R}_n^\circ(\mathbb{p}, \phi) &:= \mathbb{R}_n^{m_n^\circ}(\mathbb{p}, \phi) := \min \{ \mathbb{R}_n^m(\mathbb{p}, \phi) : m \in \mathbb{N} \}. \end{aligned} \quad (16.06)$$

If $\mathbb{p} \in \mathbb{L}_\infty(\lambda_{0,1})$ then we have $\mathbb{P}_p^{\otimes n} (|\phi_{\mathbb{N}}(\widehat{\mathbb{p}}^{m_n^\circ} - \mathbb{p})|^2) \leq \|\mathbb{p}\|_{\mathbb{L}_\infty(\lambda_{0,1})} \mathbb{R}_n^\circ(\mathbb{p}, \phi)$.

§16.21 **Proof of Proposition §16.20.** is given in the lecture. □

§16.22 **Oracle inequality.** Under Assumptions §15.02 and §15.05 let $\phi \in (\mathbb{R}_{\lambda_0})^{\mathbb{N}}$ and $\mathbb{p} \in \text{dom}(\phi_{\mathcal{N}})$. If in addition $v_{\mathbb{p}} := \max(\|\mathbb{p}\|_{\mathbb{L}_{\infty}(\lambda_{0,1})}, \|\mathbb{p}^{-1}\|_{\mathbb{L}_{\infty}(\lambda_{0,1})}) \in \mathbb{R}_{\lambda_0}^+$ then $\max(\|\Gamma_{\mathbb{p}}\|_{\mathbb{L}(\ell_2)}, \|\Gamma_{\mathbb{p}}^{-1}\|_{\mathbb{L}(\ell_2)}) \leq v_{\mathbb{p}}$ (see Lemma §15.10), and hence **Property §12.36 implies**

$$\begin{aligned} v_{\mathbb{p}}^{-1} R_n^m(\mathbb{p}, \phi) &\leq \mathbb{P}_{\mathbb{p}}^{\otimes n} (|\phi_{\mathcal{N}}(\widehat{\mathbb{p}}^m - \mathbb{p})|^2) = n^{-1} \|\phi \mathbf{1}^m\|_{\Gamma}^2 + |\phi_{\mathcal{N}}(\mathbb{p} \mathbf{1}^{m\perp})|^2 \\ &\leq v_{\mathbb{p}} R_n^m(\mathbb{p}, \phi) \quad \text{for all } m, n \in \mathbb{N}. \end{aligned}$$

As a consequence we immediately obtain the following **oracle inequality** (see Definition §12.34)

$$\begin{aligned} v_{\mathbb{p}}^{-1} R_n^{\circ}(\mathbb{p}, \phi) &\leq \inf_{m \in \mathbb{N}} \mathbb{P}_{\mathbb{p}}^{\otimes n} (|\phi_{\mathcal{N}}(\widehat{\mathbb{p}}^m - \mathbb{p})|^2) \leq \mathbb{P}_{\mathbb{p}}^{\otimes n} (|\phi_{\mathcal{N}}(\widehat{\mathbb{p}}^{m_n^{\circ}} - \mathbb{p})|^2) \\ &\leq v_{\mathbb{p}} R_n^{\circ}(\mathbb{p}, \phi) \leq v_{\mathbb{p}}^2 \inf_{m \in \mathbb{N}} \mathbb{P}_{\mathbb{p}}^{\otimes n} (|\phi_{\mathcal{N}}(\widehat{\mathbb{p}}^m - \mathbb{p})|^2), \quad (16.07) \end{aligned}$$

and hence, $R_n^{\circ}(\mathbb{p}, \phi)$, m_n° and the statistic $\widehat{\mathbb{p}}^{m_n^{\circ}}$, respectively, is an **oracle bound**, an **oracle dimension** and **oracle optimal** (up to the constant $v_{\mathbb{p}}^2$). We observe that $R_n^{\circ}(\mathbb{p}, \phi) = o(1)$ as $n \rightarrow \infty$ (Remark §12.37), and thus, $R_n^{\circ}(\mathbb{p}, \phi)$ is an **oracle rate**. However, note that the oracle dimension $m_n^{\circ} = m_n^{\circ}(\mathbb{p}, \phi)$ depends on the unknown density coefficients \mathbb{p} , and thus also the oracle optimal statistic $\widehat{\mathbb{p}}^{m_n^{\circ}}$. In other words $\widehat{\mathbb{p}}^{m_n^{\circ}}$ is not a feasible estimator. □

§16.23 **Illustration.** We illustrate the last results considering usual behaviour for both the variance and the bias term. Similar to the two cases **(p)** and **(np)** in Illustration §16.10 we distinguish here the following two cases

(p) $\phi \in \ell_2$ or there is $K \in \mathbb{N}$ with $\sup\{|\phi_{\mathcal{N}}(\mathbb{p} \mathbf{1}^{m\perp})|^2 : m \in \mathbb{N} \cap [K, \infty)\} = 0$,

(np) $\phi \notin \ell_2$ and for all $m \in \mathbb{N}$ holds $\sup\{|\phi_{\mathcal{N}}(\mathbb{p} \mathbf{1}^{m\perp})|^2 : m \in \mathbb{N} \cap [K, \infty)\} \in \mathbb{R}_{\lambda_0}^+$.

In case **(p)** the oracle bound is again parametric, i.e. $n R_n^{\circ}(\mathbb{p}, \phi) = O(1)$, while in case **(np)** the oracle bound is nonparametric, i.e. $\lim_{n \rightarrow \infty} n R_n^{\circ}(\mathbb{p}, \phi) = \infty$. In case **(np)** consider the following two specifications

Table 02 [§16]

Order of the oracle rate $R_n^{\circ}(\mathbb{p}, \phi)$ as $n \rightarrow \infty$

$(j \in \mathbb{N})$	$(a \in \mathbb{R}_{\lambda_0}^+)$	(squared bias)	(variance)	m_n°	$R_n^{\circ}(\mathbb{p}, \phi)$
$\phi_j = j^{v-1/2}$	\mathbb{p}	$ \phi_{\mathcal{N}}(\mathbb{p} \mathbf{1}^{m\perp}) ^2$	$\ \mathbf{1}^m\ _{\phi}^2$		
(o) $v \in (0, a)$	$j^{-a-1/2}$	$m^{-2(a-v)}$	m^{2v}	$n^{\frac{1}{2a}}$	$n^{-\frac{(a-v)}{a}}$
$v = 0$	$j^{-a-1/2}$	m^{-2a}	$\log m$	$(\frac{n}{\log n})^{\frac{1}{2a}}$	$\frac{\log n}{n}$
(s) $v \in \mathbb{R}_{\lambda_0}^+$	$e^{-j^{2a}}$	$m^{(1-2(a-v))_+} e^{-2m^{2a}}$	m^{2v}	$(\log n)^{\frac{1}{2a}}$	$\frac{(\log n)^{\frac{v}{a}}}{n}$
$v = 0$	$e^{-j^{2a}}$	$m^{(1-2a)_+} e^{-m^{2a}}$	$\log m$	$(\log n)^{\frac{1}{2a}}$	$\frac{\log \log n}{n}$

We note that in Table 02 [§16] the order of the oracle rate $R_n^{\circ}(\mathbb{p}, \phi)$ is depict for $v \geq 0$ only. For $v < 0$ the oracle rate $R_n^{\circ}(\mathbb{p}, \phi)$ is parametric. □

§16|02|02 Maximal local ϕ -risk

§16.24 **Assumption.** Consider $\phi, \mathbf{a} \in (\mathbb{R}_{\lambda_0})^{\mathbb{N}}$ with $\mathbf{a} \in \ell_{\infty}$ and $(\mathbf{a}\phi)_j := (\mathbf{a}_j \phi_j)_{j \in \mathbb{N}} = \mathbf{a} \cdot \phi \in \ell_2$, and hence $\|\mathbf{a} \mathbf{1}^{m\perp}\|_{\phi} = \|(\mathbf{a}\phi) \cdot \mathbf{1}^{m\perp}\|_{\ell_2} = o(1)$ as $m \rightarrow \infty$. The orthonormal system $(u_j)_{j \in \mathbb{N}}$ in $\mathbb{L}_2(\lambda_{0,1})$ and $u_0 := \mathbf{1}_{[0,1]}$ form an **(os1')** orthonormal basis $(u_j)_{j \in \mathbb{N}_0}$ in $\mathbb{L}_2(\lambda_{0,1})$ and as process $u_j^2 = (u_j^2)_{j \in \mathbb{N}}$ on $([0, 1], \mathcal{B}_{[0,1]})$ satisfies **(os2')** $\|\nu_{\mathcal{N}}(\mathbf{a}_j^2 u_j^2)\|_{\mathbb{L}_{\infty}(\lambda_{0,1})} \leq \tau_{\mathbf{a}, u}^2$ for $\tau_{\mathbf{a}, u} \in [1, \infty)$. □

§16.25 **Reminder.** Under Assumption §16.24 we have $\ell_2^{\mathfrak{a}} = \text{dom}(M_{\mathfrak{a},\cdot}) = \ell_2 \mathfrak{a} \subseteq \ell_2$ and the three measures ν_N , $\mathfrak{a}^{-2}\nu_N$ and $|\phi|_{\nu_N}$ dominate mutually each other, i.e. they share the same null sets (see **Property** §11.05). We consider $\ell_2^{\mathfrak{a}}$ endowed with $\|\cdot\|_{\mathfrak{a}^{-1}} = \|M_{\mathfrak{a},\cdot}\|_{\ell_2}$ and given a constant $r \in \mathbb{R}_0^+$ the ellipsoid $\ell_2^{\mathfrak{a},r} := \{\mathfrak{a} \in \ell_2^{\mathfrak{a}} : \|\mathfrak{a}\|_{\mathfrak{a}^{-1}} \leq r\} \subseteq \ell_2^{\mathfrak{a}}$. Since $(\mathfrak{a}\phi)_\cdot \in \ell_2$, and hence $\|\mathfrak{a}\mathbb{1}^{m\perp}\|_\phi = \|(\mathfrak{a}\phi)_\cdot \mathbb{1}^{m\perp}\|_{\ell_2} \in \mathbb{R}^+$ for each $m \in \mathbb{N}$ ($\|\mathfrak{a}\mathbb{1}^{m\perp}\|_\phi = o(1)$ as $m \rightarrow \infty$ by dominated convergence) we have $\ell_2^{\mathfrak{a}} \subseteq \text{dom}(\phi\nu_N)$ (**Property** §11.27), and $|\phi\nu_N(\mathfrak{p}\mathbb{1}^{m\perp})| \leq r \|\mathfrak{a}\mathbb{1}^{m\perp}\|_\phi$ for all $\mathfrak{p} \in \ell_2^{\mathfrak{a},r}$ (**Lemma** §11.29). \square

§16.26 **Remark.** We replace Assumption §15.05 (os1) and (os2), respectively, by the stronger Assumption §16.24 (os1') and (os2') (see **Remark** §16.13). Moreover, considering the set $\mathbb{D}_2^{\mathfrak{a},r}$ of densities in $\mathbb{L}_2(\lambda_{[0,1]})$ defined in (16.04) we have $\|\mathfrak{p}\|_{\mathbb{L}_\infty(\lambda_{[0,1]})} \leq 1 + r\tau_{\mathfrak{a},u}$ for all $\mathfrak{p} \in \mathbb{D}_2^{\mathfrak{a},r}$ due to **Lemma** §16.14. \square

§16.27 **Proposition (Upper bound).** Let Assumptions §15.02 and §16.24 be satisfied. For $n \in \mathbb{N}$ considering $m_n^* \in \mathbb{N}$ and $R_n^*(\mathfrak{a}, \phi) \in \mathbb{R}^+$ as in (12.13) (**Proposition** §12.42) we have

$$\sup \{ \mathbb{P}_p^{\otimes n} (|\phi\nu_N(\widehat{\mathfrak{p}}^m - \mathfrak{p})|^2) : \mathfrak{p} \in \mathbb{D}_2^{\mathfrak{a},r} \} \leq C R_n^*(\mathfrak{a}, \phi).$$

with constant $C = (1 + r\tau_{\mathfrak{a},u}) \vee r^2$.

§16.28 **Proof** of **Proposition** §16.27. is given in the lecture. \square

§16.29 **Illustration.** Consider the *trigonometric basis* as in **Illustration** §16.18 which satisfies Assumption §16.24 for all $\mathfrak{a} \in \ell_2$. In Table 04 [§12] the order of the rate $R_n^*(\mathfrak{a}, \phi)$ is depicted for the two cases (o) and (s) introduced in **Illustration** §12.47. We note that we have $\mathfrak{a} \in \ell_2$ in case (o) for $a > 1/2$ while in case (s) for $a \in \mathbb{R}_0^+$. \square

§17 Minimax optimal density estimation

§17|01 Maximal local ϕ -risk

§17.01 **Reminder (Maximal local ϕ -risk).** Under Assumptions §15.02 and §16.24 the observable noisy density coefficients $\widehat{\mathfrak{p}} = \mathfrak{p} + n^{-1/2}\varepsilon$, of the density coefficients $\mathfrak{p} = \mathbb{U}\mathfrak{p} \in \ell_2$ take the form of a *statistical direct problem* (see **Definition** §10.19) where the stochastic processes $\varepsilon \in \mathcal{B}_{[0,1]}^{\otimes n} \otimes 2^{\mathbb{N}}$ is given in **Definition** §15.08. Under Assumptions §15.02 and §16.24 in **Proposition** §16.27 an upper bound for a maximal local ϕ -risk of an OPE over the class $\mathbb{D}_2^{\mathfrak{a},r}$ of densities in $\mathbb{L}_2(\lambda_{[0,1]})$ defined in (16.04) is shown. More precisely, the performance of the OPE $\widehat{\mathfrak{p}}^m = \widehat{\mathfrak{p}}\mathbb{1}^m \in \ell_2\mathbb{1}^m \subseteq \text{dom}(\phi\nu_N)$ with dimension $m \in \mathbb{N}$ is measured by its maximal local ϕ -risk, that is

$$\mathcal{R}_n^\phi[\widehat{\mathfrak{p}}^m | \mathbb{D}_2^{\mathfrak{a},r}] := \sup \{ \mathbb{P}_p^{\otimes n} (|\phi\nu_N(\widehat{\mathfrak{p}}^m - \mathfrak{p})|^2) : \mathfrak{p} \in \mathbb{D}_2^{\mathfrak{a},r} \}.$$

Let us recall (12.13) (**Proposition** §12.42) where for $n, m \in \mathbb{N}$ we have defined

$$\begin{aligned} R_n^m(\mathfrak{a}, \phi) &:= \|\mathfrak{a}\mathbb{1}^{m\perp}\|_\phi^2 + n^{-1}\|\mathbb{1}^m\|_\phi^2, & m_n^* &:= \arg \min \{ R_n^m(\mathfrak{a}, \phi) : m \in \mathbb{N} \} \\ \text{and } R_n^*(\mathfrak{a}, \phi) &:= R_n^{m_n^*}(\mathfrak{a}, \phi) = \min \{ R_n^m(\mathfrak{a}, \phi) : m \in \mathbb{N} \}. \end{aligned} \quad (17.01)$$

By **Proposition** §16.27 under Assumptions §15.02 and §16.24 the maximal local ϕ -risk of an OPE $\widehat{\mathfrak{p}}^{m_n^*}$ with optimally chosen dimension m_n^* as in (17.01) satisfies

$$\mathcal{R}_n^\phi[\widehat{\mathfrak{p}}^{m_n^*} | \mathbb{D}_2^{\mathfrak{a},r}] \leq C R_n^*(\mathfrak{a}, \phi)$$

with $C = (1 + r\tau_{\mathfrak{a},u}) \vee r^2$. \square

§17.02 **Lemma (Lower bound based on two hypotheses).** *If there are $\mathbb{p}^0, \mathbb{p}^1 \in \mathbb{D}_2^{\text{a,r}}$ with associated probability measures $\mathbb{P}_0 := \mathbb{P}_{\mathbb{p}^0}$ and $\mathbb{P}_1 := \mathbb{P}_{\mathbb{p}^1}$ such that $H^2(\mathbb{P}_0, \mathbb{P}_1) \leq 2n^{-1}$ then for all $n \geq 2$ we have*

$$\inf_{\tilde{\mathbb{p}}} \mathcal{R}_n^{\phi}[\tilde{\mathbb{p}} | \mathbb{D}_2^{\text{a,r}}] \geq \frac{1}{64} |\phi_{\mathcal{U}_N}(\mathbb{p}^0 - \mathbb{p}^1)|^2.$$

where the infimum is taken over all possible estimators.

§17.03 **Proof of Lemma §17.02.** is given in the lecture. □

§17.04 **Remark.** If we consider furthermore candidate densities $\mathbb{p}^0 := \mathbb{1}_{[0,1]} + U^* \mathbb{p}^*$ and $\mathbb{p}^1 = \mathbb{1}_{[0,1]} - U^* \mathbb{p}^*$ for some $\mathbb{p}^* \in \ell_2^{\text{a,r}}$, and hence by definition $\mathbb{p}^0, \mathbb{p}^1 \in \mathbb{D}_2^{\text{a,r}}$, then trivially $|\phi_{\mathcal{V}}(\mathbb{p}^0 - \mathbb{p}^1)|^2 = 4|\phi_{\mathcal{U}_N}(\mathbb{p}^*)|^2$. If the associated probability measures $\mathbb{P}_0 := \mathbb{P}_{\mathbb{p}^0}$ and $\mathbb{P}_1 := \mathbb{P}_{\mathbb{p}^1}$ satisfy $H^2(\mathbb{P}_0, \mathbb{P}_1) \leq 2n^{-1}$ then due to **Lemma §17.02** for all $n \geq 2$ we have

$$\inf_{\tilde{\mathbb{p}}} \mathcal{R}_n^{\phi}[\tilde{\mathbb{p}} | \mathbb{D}_2^{\text{a,r}}] \geq \frac{1}{16} |\phi_{\mathcal{U}_N}(\mathbb{p}^*)|^2. \quad (17.02)$$

We find a minimax-optimal lower bound by constructing a candidate $\mathbb{p}^* \in \ell_2^{\text{a,r}}$ that has the largest possible $|\phi_{\mathcal{U}_N}(\mathbb{p}^*)|^2$ -value but $\mathbb{P}_{\mathbb{p}^0}^{\otimes n}$ and $\mathbb{P}_{\mathbb{p}^1}^{\otimes n}$ are still statistically indistinguishable in the sense that $H^2(\mathbb{P}_{\mathbb{p}^0}, \mathbb{P}_{\mathbb{p}^1}) \leq 2n^{-1}$. □

§17.05 **Lemma.** *Under Assumption §16.24 let $\mathbb{p}^* \in \ell_2^{\text{a,r}}$ with $\|\mathbb{p}^*\|_{\alpha^{-1}} \leq 1/(2\tau_{\text{a,u}})$. Then $\mathbb{p}^0 := \mathbb{1}_{[0,1]} + U^* \mathbb{p}^*$ and $\mathbb{p}^1 := \mathbb{1}_{[0,1]} - U^* \mathbb{p}^*$ belong to $\mathbb{D}_2^{\text{a,r}}$, and the associated probability measures $\mathbb{P}_0 := \mathbb{P}_{\mathbb{p}^0}$ and $\mathbb{P}_1 := \mathbb{P}_{\mathbb{p}^1}$ satisfy $H^2(\mathbb{P}_0, \mathbb{P}_1) \leq 2\|\mathbb{p}^*\|_{\ell_2}^2$.*

§17.06 **Proof of Lemma §17.05.** is given in the lecture. □

§17.07 **Reminder.** Under Assumption §16.24 let in addition $\alpha^2 \in (\mathbb{R}_{\setminus 0}^+)^{\mathbb{N}}$ (see **Notation §13.23**), then Assumption §13.24 is satisfied. If $\alpha_m^2 > n^{-1}$ then exploiting the definition (17.01) of m_n^* we have $\alpha_{m_n^*}^2 > n^{-1} \geq \alpha_{m_n^*+1}^2$ (see **Comment §13.25**) which we use in the next proof. □

§17.08 **Proposition (Lower bound).** *Let Assumptions §15.02 and §16.24 be satisfied. If $\alpha^2 \in (\mathbb{R}_{\setminus 0}^+)^{\mathbb{N}}$ then for all $n \in \mathbb{N} \cap (1 \vee \alpha^{-2}, \infty)$ we have*

$$\inf_{\tilde{\mathbb{p}}} \mathcal{R}_n^{\phi}[\tilde{\mathbb{p}} | \mathbb{D}_2^{\text{a,r}}] \geq C R_n^*(\alpha, \phi) \quad (17.03)$$

with constant $C := 16^{-1}(\mathfrak{r}^2 \wedge 1/(4\tau_{\text{a,u}}) \wedge 1)$ and infimum taken over all estimators.

§17.09 **Proof of Proposition §17.08.** is given in the lecture. □

§17.10 **Illustration.** Consider the *trigonometric basis* as in **Illustration §16.18** which satisfies Assumption §16.24 for all $\alpha \in \ell_2$ (see **Illustration §16.29**). In Table 04 [§12] the order of the rate $R_n^*(\alpha, \phi)$ is depicted for the two cases (o) and (s) introduced in **Illustration §16.29**. We note that we have $\alpha \in \ell_2$ in case (o) for $a > 1/2$ while in case (s) for $a \in \mathbb{R}_{\setminus 0}^+$. In both cases the additional assumption $\alpha^2 \in (\mathbb{R}_{\setminus 0}^+)^{\mathbb{N}}$ is satisfied. Consequently, due to **Proposition §17.08** the Table 04 [§12] presents the order of the *minimax rate* $R_n^*(\alpha, \phi)$ which is attained by the *minimax-optimal OPE* $\hat{\mathbb{p}}^{m_n^*} = \hat{\mathbb{p}} \mathbb{1}^{m_n^*} \in \ell_2 \mathbb{1}^{m_n^*} \subseteq \text{dom}(\phi_{\mathcal{U}_N})$ with optimally selected dimension m_n^* (**Proposition §16.27**). We shall stress, that the order of m_n^* given in the Table 04 [§12] depends on the parameter $a \in \mathbb{R}_{\setminus 0}^+$ characterising the (abstract) smoothness of the density of interest which is generally not known in advance. □

§17|02 Maximal global v-risk

§17.11 **Reminder** (*Maximal global v-risk*). Under Assumptions §15.02 and §16.11 the observable noisy density coefficients $\widehat{\mathfrak{p}} = \mathfrak{p} + n^{-1/2}\boldsymbol{\varepsilon}$ of the density coefficients $\mathfrak{p} = \mathbb{U}\mathfrak{p} \in \ell_2$ take the form of a *statistical direct problem* (see Definition §10.19) where the stochastic processes $\boldsymbol{\varepsilon} \in \mathcal{B}_{[0,1]}^{\otimes n} \otimes 2^{\mathbb{N}}$ is given in Definition §15.08. Under Assumptions §15.02 and §16.11 in Proposition §16.16 an upper bound for a maximal global v-risk of an OPE over the class $\mathbb{D}_2^{\text{a,r}}$ of densities in $\mathbb{L}_2(\lambda_{[0,1]})$ defined in (16.04) is shown. More precisely, the performance of the OPE $\widehat{\mathfrak{p}}^m = \widehat{\mathfrak{p}} \mathbb{1}^m \in \ell_2(\mathfrak{v}^2)\mathbb{1}^m \subseteq \ell_2(\mathfrak{v}^2)$ with dimension $m \in \mathbb{N}$ is measured by its maximal global v-risk, that is

$$\mathcal{R}_n^{\mathfrak{v}}[\widehat{\mathfrak{p}}^m | \mathbb{D}_2^{\text{a,r}}] := \sup \{ \mathbb{P}_{\mathfrak{p}}^{\otimes n} (\|\widehat{\mathfrak{p}}^m - \mathfrak{p}\|_{\mathfrak{v}}^2) : \mathfrak{p} \in \mathbb{D}_2^{\text{a,r}} \}.$$

Let us recall (12.06) (Proposition §12.21) where for $n, m \in \mathbb{N}$ we have defined $(\mathfrak{a}\mathfrak{v})_{(m)}^2 = \|(\mathfrak{a}\mathfrak{v})\mathbb{1}^{m\perp}\|_{\ell_\infty}^2$ and

$$\begin{aligned} R_n^m(\mathfrak{a}, \mathfrak{v}) &:= [(\mathfrak{a}\mathfrak{v})_{(m)}^2 \vee n^{-1}\|\mathbb{1}^m\|_{\mathfrak{v}}^2], \quad m_n^* := \arg \min \{ R_n^m(\mathfrak{a}, \mathfrak{v}) : m \in \mathbb{N} \} \\ &\text{and} \quad R_n^*(\mathfrak{a}, \mathfrak{v}) := R_n^{m_n^*}(\mathfrak{a}, \mathfrak{v}) = \min \{ R_n^m(\mathfrak{a}, \mathfrak{v}) : m \in \mathbb{N} \}. \end{aligned} \quad (17.04)$$

By Proposition §16.16 under Assumptions §15.02 and §16.11 the maximal global v-risk of an OPE $\widehat{\mathfrak{p}}^{m_n^*}$ with optimally chosen dimension m_n^* as in (17.04) satisfies

$$\mathcal{R}_n^{\mathfrak{v}}[\widehat{\mathfrak{p}}^{m_n^*} | \mathbb{D}_2^{\text{a,r}}] \leq C R_n^*(\mathfrak{a}, \mathfrak{v})$$

with $C = 1 + r\tau_{\mathfrak{a},\mathfrak{u}} + r^2$. Furthermore, as in Notation §13.29 for $m \in \mathbb{N}$ we set $\mathcal{J}_m := \{-1, 1\}^m$ and for each $\tau := (\tau_j)_{j \in \llbracket m \rrbracket} \in \mathcal{J}_m$ and $j \in \llbracket m \rrbracket$ we introduce $\tau^{(j)} \in \mathcal{J}_m$ given by $\tau_j^{(j)} := -\tau_j$ and $\tau_l^{(j)} := \tau_l$ for $l \in \llbracket m \rrbracket \setminus \{j\}$. \square

§17.12 **Lemma** (*Assouad's cube technique*). If for each $\tau \in \mathcal{J}_m$ there is $\mathfrak{p}^\tau \in \mathbb{D}_2^{\text{a,r}}$ with associated probability measure $\mathbb{P}_\tau := \mathbb{P}_{\mathfrak{p}^\tau}$ such that for all $\tau \in \mathcal{J}_m$ and $j \in \llbracket m \rrbracket$ we have $H(\mathbb{P}_\tau, \mathbb{P}_{\tau^{(j)}}) \leq 2n^{-1}$ then for all $n \geq 2$

$$\inf_{\widehat{\mathfrak{p}}} \mathcal{R}_n^{\mathfrak{v}}[\widehat{\mathfrak{p}} | \mathbb{D}_2^{\text{a,r}}] \geq 2^{-m} \sum_{\tau \in \mathcal{J}_m} \frac{1}{64} \sum_{j \in \llbracket m \rrbracket} (\mathfrak{v}_j^2 |\mathfrak{p}_j^\tau - \mathfrak{p}_j^{\tau^{(j)}}|^2)$$

where the infimum is taken over all possible estimators.

§17.13 **Proof** of Lemma §17.12. is given in the lecture. \square

§17.14 **Remark**. If we assume furthermore candidate densities $\mathfrak{p}^\tau := \mathbb{1}_{[0,1]} + \mathbb{U}^* \mathfrak{p}^\tau$ with $\mathfrak{p}^\tau := (\tau_j \mathfrak{p}_j^* \mathbb{1}_j^m)_{j \in \mathbb{N}}$, $\tau \in \mathcal{J}_m$, for some $\mathfrak{p}^* \in \ell_2^{\text{a,r}}$, where evidently $\mathfrak{p}^\tau \in \ell_2^{\text{a,r}}$ too and hence $\mathfrak{p}^\tau \in \mathbb{D}_2^{\text{a,r}}$, then it is easily seen that $\sum_{j \in \llbracket m \rrbracket} (\mathfrak{v}_j^2 |\mathfrak{p}_j^\tau - \mathfrak{p}_j^{\tau^{(j)}}|^2) = 4\|\mathfrak{p}^* \mathbb{1}^m\|_{\mathfrak{v}}^2$. If for all $\tau \in \mathcal{J}_m$ and $j \in \llbracket m \rrbracket$ the associated probability measures $\mathbb{P}_\tau := \mathbb{P}_{\mathfrak{p}^\tau}$ and $\mathbb{P}_{\tau^{(j)}} := \mathbb{P}_{\mathfrak{p}^{\tau^{(j)}}}$ satisfy $H(\mathbb{P}_\tau, \mathbb{P}_{\tau^{(j)}}) \leq 2n^{-1}$ then due to Lemma §17.12 for all $n \geq 2$ we have

$$\inf_{\widehat{\mathfrak{p}}} \mathcal{R}_n^{\mathfrak{v}}[\widehat{\mathfrak{p}} | \mathbb{D}_2^{\text{a,r}}] \geq 2^{-m} \sum_{\tau \in \mathcal{J}_m} \frac{1}{16} \|\mathfrak{p}^* \mathbb{1}^m\|_{\mathfrak{v}}^2 = \frac{1}{16} \|\mathfrak{p}^* \mathbb{1}^m\|_{\mathfrak{v}}^2. \quad (17.05)$$

We find a minimax-optimal lower bound by choosing the parameter m and the function \mathfrak{p}^* that have the largest possible $\|\mathfrak{p}^* \mathbb{1}^m\|_{\mathfrak{v}}^2$ -value although that the associated \mathbb{P}_τ , $\tau \in \mathcal{J}_m$ are still statistically indistinguishable in the sense that $H^2(\mathbb{P}_\tau, \mathbb{P}_{\tau^{(j)}}) \leq 2n^{-1}$ for all $j \in \llbracket m \rrbracket$ and $\tau \in \mathcal{J}_m$. \square

§17.15 **Lemma**. Under Assumption §16.11 let $\mathfrak{p}^* \in \ell_2^{\text{a,r}}$ with $\|\mathfrak{p}^*\|_{\mathfrak{a}^{-1}} \leq 1/(2\tau_{\mathfrak{a},\mathfrak{u}})$. Then for each $\tau \in \mathcal{J}_m$, $\mathfrak{p}^\tau := \mathbb{1}_{[0,1]} + \mathbb{U}^* \mathfrak{p}^\tau$ with $\mathfrak{p}^\tau := (\tau_j \mathfrak{p}_j^* \mathbb{1}_j^m)_{j \in \mathbb{N}}$ belongs to $\mathbb{D}_2^{\text{a,r}}$, and for each $j \in \llbracket m \rrbracket$ the associated probability measures $\mathbb{P}_\tau := \mathbb{P}_{\mathfrak{p}^\tau}$ and $\mathbb{P}_{\tau^{(j)}} := \mathbb{P}_{\mathfrak{p}^{\tau^{(j)}}}$ satisfy $H^2(\mathbb{P}_\tau, \mathbb{P}_{\tau^{(j)}}) \leq 2\|\mathfrak{p}^* \mathbb{1}^m\|_{\ell_\infty}^2$.

§17.16 **Proof of Lemma §17.15.** is given in the lecture. \square

§17.17 **Reminder.** For $w_{\cdot} \in \ell_{\infty}$ we set $w_{(0)}^2 := \|w_{\cdot}^2\|_{\ell_{\infty}}$ and $w_{(j)}^2 = (w_{(j)}^2 := \|w_{\cdot}^2 \mathbb{1}_{\cdot}^{j+1}\|_{\ell_{\infty}})_{j \in \mathbb{N}}$ (**Notation §13.34**) where by construction $w_{(j)}^2 = \sup \{w_i^2: i \in \mathbb{N} \cap [j+1, \infty)\}$, $j \in \mathbb{N}_0$ and $w_{(0)}^2 \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$. Under Assumption §16.11 let in addition $(\mathbf{av})_{(\cdot)}^2 \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ and there exists $C_{(\mathbf{av})} \in (0, 1]$ such that $C_{(\mathbf{av})} \|(\mathbf{av})_{\cdot}^{-2} \mathbb{1}_{\cdot}^m\|_{\ell_{\infty}} \leq (\mathbf{av})_{(m-1)}^{-2}$ or in equal

$$(\mathbf{av})_{(m-1)}^2 \geq \min \{(\mathbf{av})_j^2: j \in \llbracket m \rrbracket\} \geq C_{(\mathbf{av})} (\mathbf{av})_{(m-1)}^2$$

for all $m \in \mathbb{N}$, then Assumption §13.35 is satisfied. For m_n^* and $R_n^* := R_n^{m_n^*}(\mathbf{a}, \mathbf{v})$ as in (17.04) we distinguish case i) : $R_n^* = n^{-1} \|\mathbb{1}_{\cdot}^{m_n^*}\|_{\mathbf{v}}^2 > (\mathbf{av})_{(m_n^*)}^2$ and case ii) : $R_n^* = (\mathbf{av})_{(m_n^*)}^2 \geq n^{-1} \|\mathbb{1}_{\cdot}^{m_n^*}\|_{\mathbf{v}}^2$. Due to **Comment §13.36** if $(\mathbf{av})_{(1)}^2 > n^{-1} \mathbf{v}_1^2$ then in case i) $(\mathbf{av})_{(m_n^*-1)}^2 \geq n^{-1} \|\mathbb{1}_{\cdot}^{m_n^*}\|_{\mathbf{v}}^2$, while in case ii) setting (the defining set is not empty since $(\mathbf{av})_{(\cdot)}^2 \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$)

$$m_n^{\circ} := \min \{m \in \mathbb{N} \cap [m_n^* + 1, \infty): n^{-1} \|\mathbb{1}_{\cdot}^m\|_{\mathbf{v}}^2 \geq (\mathbf{av})_{(m)}^2\} \quad (17.06)$$

we have $(\mathbf{av})_{(m_n^{\circ})}^2 = (\mathbf{av})_{(m_n^{\circ}-1)}^2 \leq n^{-1} \|\mathbb{1}_{\cdot}^{m_n^{\circ}}\|_{\mathbf{v}}^2$. We use those estimates in the next proof. \square

§17.18 **Proposition (Lower bound).** Let Assumptions §15.02 and §16.11 be satisfied. If $(\mathbf{av})_{(\cdot)}^2 \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ and there exists $C_{(\mathbf{av})} \in (0, 1]$ such that $C_{(\mathbf{av})} \|(\mathbf{av})_{\cdot}^{-2} \mathbb{1}_{\cdot}^m\|_{\ell_{\infty}} \leq (\mathbf{av})_{(m-1)}^{-2}$ for all $m \in \mathbb{N}$, then for all $n \in \mathbb{N} \cap (1 \vee \mathbf{v}_1^2 (\mathbf{av})_{(1)}^{-2}, \infty)$ we have

$$\inf_{\hat{\mathbb{P}}} \mathcal{R}_n^{\mathfrak{D}}[\hat{\mathbb{P}} | \mathbb{D}_2^{\mathbf{a}, \mathbf{v}}] \geq C R_n^*(\mathbf{a}, \mathbf{v}) \quad (17.07)$$

with constant $C := (C_{(\mathbf{av})}/16)(r^2 \wedge 1/(4\tau_{\mathbf{a}, \mathbf{u}}^2) \wedge 1)$ and infimum taken over all estimators.

§17.19 **Proof of Proposition §17.18.** is given in the lecture. \square

§17.20 **Illustration.** Consider the *trigonometric basis* as in **Illustration §16.18** which satisfies Assumption §16.11 for all $\mathbf{a}_{\cdot} \in \ell_2$ (see **Illustration §16.18**). In Table 02 [§12] the order of the rate $R_n^*(\mathbf{a}, \mathbf{v})$ is depict for the two cases **(o)** and **(s)** introduced in **Illustration §16.18**. We note that we have $\mathbf{a}_{\cdot} \in \ell_2$ in case **(o)** for $a > 1/2$ while in case **(s)** for $a \in \mathbb{R}_0^+$. In both cases the additional assumptions, $(\mathbf{av})_{(\cdot)}^2 \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ and there exists $C_{(\mathbf{av})} \in (0, 1]$ such that $C_{(\mathbf{av})} \|(\mathbf{av})_{\cdot}^{-2} \mathbb{1}_{\cdot}^m\|_{\ell_{\infty}} \leq (\mathbf{av})_{(m-1)}^{-2}$ for all $m \in \mathbb{N}$, are satisfied. Consequently, due to **Proposition §17.18** the Table 02 [§12] presents the order of the *minimax rate* $R_n^*(\mathbf{a}, \mathbf{v})$ which is attained by the *minimax-optimal* OPE $\hat{\mathbb{P}}^{m_n^*} = \hat{\mathbb{P}} \mathbb{1}_{\cdot}^{m_n^*} \in \ell_2 \mathbb{1}_{\cdot}^{m_n^*} \subseteq \ell_2(\mathbf{v}^2)$ with optimally selected dimension m_n^* (**Proposition §16.16**). We shall stress, that the order of m_n^* given in the Table 02 [§12] depends on the parameter $a \in \mathbb{R}_0^+$ characterising the (abstract) smoothness of the density of interest which is generally not known in advance. \square

§18 Data-driven density estimation

§18|01 Data-driven global estimation by model selection

The next assertion provides our key argument in order to control the deviations of the reminder term. The inequality is due to Talagrand [1996] and in this form for example given in Klein and Rio [2005].

§18.01 **Lemma (Talagrand's inequality).** Let $(Z_i)_{i \in \llbracket n \rrbracket}$ be independent $(\mathcal{Z}, \mathcal{Z})$ -valued random variables and let $\{\mathfrak{r}_t: t \in \mathcal{T}\} \subseteq \mathcal{Z}$ be a countable class of Borel-measurable functions. For $t \in \mathcal{T}$ setting $\bar{\mathfrak{r}}_t = n^{-1} \sum_{i \in \llbracket n \rrbracket} \{\mathfrak{r}_t(Z_i) - \mathbb{E}(\mathfrak{r}_t(Z_i))\}$ we have

$$\mathbb{E} \left(\left(\sup \{|\bar{\mathfrak{r}}_t|^2: t \in \mathcal{T}\} - 6H^2 \right)_+ \right) \leq C_{\text{tal}} \left\{ \frac{\mathbf{v}}{n} \exp \left(\frac{-nH^2}{6\mathbf{v}} \right) + \frac{h^2}{n^2} \exp \left(\frac{-nH}{100h} \right) \right\} \quad (18.01)$$

for some universal numerical constant $C_{\text{tal}} \in [1, \infty)$ and where

$$\sup \{ |r_t(z)| : t \in \mathcal{T}, z \in \mathcal{Z} \} \leq h, \quad \mathbb{E} \left(\sup \{ |\bar{r}_t| : t \in \mathcal{T} \} \right) \leq H, \quad \sup \{ n \mathbb{E} (|\bar{r}_t|^2) : t \in \mathcal{T} \} \leq v. \quad (18.02)$$

§18.02 **Remark.** Let us briefly reconsider the OPE $\hat{\mathbb{p}}^m = \hat{\mathbb{p}} \mathbb{1}^m \in \ell_2 \mathbb{1}^m$ with dimension $m \in \mathbb{N}$ (**Definition** §16.04) where $\hat{\mathbb{p}} = \hat{\mathbb{P}}_n \mathbf{u} = (\hat{\mathbb{P}}_n \mathbf{u}_j)_{j \in \mathbb{N}}$ are noisy versions (**Definition** §15.08) of the density coefficients $\mathbb{p} = \mathbb{U} \mathbb{p} = \mathbb{P}_p(\mathbf{u}) = (\mathbb{P}_p \mathbf{u}_j = \lambda_{[0,1]}(\mathbb{p} \mathbf{u}_j))_{j \in \mathbb{N}}$. For $m \in \mathbb{N}$ introduce the unit ball $\mathbb{B}_m := \{ a_\bullet \in \ell_2(v^2) \mathbb{1}^m : \|a_\bullet\|_v \leq 1 \}$ contained in the linear subspace $\ell_2(v^2) \mathbb{1}^m$ spanned by $(\mathbb{1}^{\{j\}})_{j \in [m]}$. Clearly, for each $a_\bullet \in \ell_2(v^2) \mathbb{1}^m$ we have $r_{a_\bullet} := \sum_{j \in [m]} v_j^2 a_j \mathbf{u}_j = \nu_N(v^2 a_\bullet \mathbf{u}) \in \mathcal{B}_{[0,1]}$, i.e. it is a $\mathcal{B}_{[0,1]}$ - \mathcal{B} -measurable function, where $\hat{\mathbb{P}}_n(r_{a_\bullet}) = \nu_N(v^2 a_\bullet \hat{\mathbb{P}}_n \mathbf{u}) = \nu_N(v^2 a_\bullet \hat{\mathbb{p}})$, $\mathbb{P}_p(r_{a_\bullet}) = \nu_N(v^2 a_\bullet \mathbb{P}_p \mathbf{u}) = \nu_N(v^2 a_\bullet \mathbb{p})$ and hence $\bar{r}_{a_\bullet} = \hat{\mathbb{P}}_n(r_{a_\bullet}) - \mathbb{P}_p(r_{a_\bullet}) = \nu_N(v^2 a_\bullet (\hat{\mathbb{p}} - \mathbb{p})) = \langle \hat{\mathbb{p}} - \mathbb{p}, a_\bullet \rangle_v$. Consequently, we obtain

$$\|\hat{\mathbb{p}}^m - \mathbb{p}^m\|_v^2 = \sup \{ |\langle \hat{\mathbb{p}} - \mathbb{p}, a_\bullet \rangle_v|^2 : a_\bullet \in \mathbb{B}_m \} = \sup \{ |\bar{r}_{a_\bullet}|^2 : a_\bullet \in \mathbb{B}_m \}$$

The last identity provides the necessary argument to apply below Talagrand's inequality (§18.01). Note that, the unit ball \mathbb{B}_m is not a countable set, however, it contains a countable dense subset, say \mathcal{B}_m , since $\ell_2(v^2)$ is separable. Exploiting the continuity of the inner product it is straightforward to see that $\sup \{ |\langle b, a_\bullet \rangle_v|^2 : a_\bullet \in \mathbb{B}_m \} = \sup \{ |\langle b, a_\bullet \rangle_v|^2 : a_\bullet \in \mathcal{B}_m \}$ for all $b \in \ell_2(v^2)$. Consequently, provided that

$$\begin{aligned} \sup \{ \|\mathbf{u}(x) \mathbb{1}^m\|_v : x \in [0, 1] \} &= \sup \{ |r_{a_\bullet}(x)| : a_\bullet \in \mathcal{B}_m, x \in [0, 1] \} \leq h, \\ \mathbb{P}_p^{\otimes n} (\|\hat{\mathbb{p}}^m - \mathbb{p}^m\|_v^2) &= \mathbb{P}_p^{\otimes n} (\sup \{ |\bar{r}_{a_\bullet}|^2 : a_\bullet \in \mathcal{B}_m \}) \leq H^2, \\ \sup \{ \mathbb{P}_p (\|\nu_N(v^2 a_\bullet (\mathbf{u} - \mathbb{P}_p \mathbf{u}))\|^2) : a_\bullet \in \mathcal{B}_m \} &= \sup \{ n \mathbb{P}_p^{\otimes n} (|\bar{r}_{a_\bullet}|^2) : a_\bullet \in \mathcal{B}_m \} \leq v. \end{aligned} \quad (18.03)$$

due to Talagrand's inequality (§18.01) we have

$$\mathbb{P}_p^{\otimes n} (\|\hat{\mathbb{p}}^m - \mathbb{p}^m\|_v^2 - 6H^2)_+ \leq C_{\text{tal}} \left\{ \frac{v}{n} \exp\left(\frac{-nH^2}{6v}\right) + \frac{H^2}{n^2} \exp\left(\frac{-nH}{100h}\right) \right\} \quad (18.04)$$

for some universal numerical constant $C_{\text{tal}} \in [1, \infty)$. \square

§18|01|01 Global v -risk

§18.03 **Assumption.** The weights $v_\bullet \in (\mathbb{R}_{>0})^{\mathbb{N}}$ satisfy

$$\forall x \in \mathbb{R}_{>0}^+ : \sum_{m \in \mathbb{N}} \{ x \|v^2 \mathbb{1}^m\|_{\ell_\infty} \exp(-\|v \mathbb{1}^m\|_{\ell_2}^2 / (x \|v^2 \mathbb{1}^m\|_{\ell_\infty})) \} =: C_v(x) \in \mathbb{R}^+. \quad (18.05)$$

The orthonormal system $(\mathbf{u}_j)_{j \in \mathbb{N}}$ in $\mathbb{L}_2(\lambda_{[0,1]})$ and $\mathbf{u}_0 := \mathbb{1}_{[0,1]}$ form an **(os1')** orthonormal basis $(\mathbf{u}_j)_{j \in \mathbb{N}_0}$ in $\mathbb{L}_2(\lambda_{[0,1]})$ and as process $\mathbf{u}_\bullet = (\mathbf{u}_j)_{j \in \mathbb{N}}$ on $([0, 1], \mathcal{B}_{[0,1]})$ for all $m \in \mathbb{N}$ satisfies **(os2'')** $\sup \{ \|\mathbf{u}(x) \mathbb{1}^m\|_v^2 : x \in [0, 1] \} \leq \tau_{v,\mathbf{u}}^2 \|\mathbb{1}^m\|_v^2 \in \mathbb{R}^+$ for $\tau_{v,\mathbf{u}} \in [1, \infty)$. \square

§18.04 **Remark.** We replace Assumption §15.05 **(os1)** and **(os2)**, respectively, by the stronger Assumption §18.03 **(os1')** and **(os2'')**. Indeed, under **(os1')** we have **(os1)** $\mathbb{1}_{[0,1]} \in \ker(\mathbb{U})$. Furthermore, $(\mathbf{u}_j)_{j \in \mathbb{N}}$ belongs to $\mathbb{L}_\infty(\lambda_{[0,1]})$ due to **(os2'')** (and $v_\bullet \in (\mathbb{R}_{>0})^{\mathbb{N}}$), and hence **(os2)** is fulfilled (see also **Remark** §15.06). Under Assumption §18.03 (18.05) we have $\|v \mathbb{1}^m\|_{\ell_2}^{-2} = o(1)$ as $m \rightarrow \infty$ (**Comment** §14.22), see also **Illustration** §14.23 for an example when (18.05) is not satisfied. \square

§18.05 **Reminder (Global oracle v -risk).** Given Assumptions §15.02 and §18.03 we consider an OPE as in **Definition** §16.04. Here the observable noisy density coefficients $\hat{\mathbb{p}} = \mathbb{p} + n^{-1/2} \boldsymbol{\varepsilon}$ of

the density coefficients $\mathbb{p} = \cup \mathbb{p} \in \ell_2$ take the form of a *statistical direct problem* (see [Definition §10.19](#)) where the stochastic processes $\varepsilon \in \mathcal{B}_{[0,1]}^{\otimes n} \otimes 2^{\mathbb{N}}$ is given in [Definition §15.08](#). Under Assumptions §15.02 and §18.03, (and hence Assumption §15.05 and $\mathbf{v} \in (\mathbb{R}_{\mathbf{v}_0})^{\mathbb{N}}$ see [Remark §18.04](#)) and $\mathbb{p} \in \ell_2(\mathbf{v}^2)$ in §16.09 an *oracle inequality* for the global \mathbf{v} -risk of the OPE's is shown. More precisely, as in (16.02) ([Proposition §16.07](#)) for all $n, m \in \mathbb{N}$ setting

$$\begin{aligned} R_n^m(\mathbb{p}, \mathbf{v}) &:= \|\mathbb{p} \mathbf{1}^{m \perp}\|_{\mathbf{v}}^2 + n^{-1} \|\mathbf{1}^m\|_{\mathbf{v}}^2, \quad m_n^{\circ} := \arg \min \{R_n^m(\mathbb{p}, \mathbf{v}) : m \in \mathbb{N}\} \\ \text{and } R_n^{\circ}(\mathbb{p}, \mathbf{v}) &:= R_n^{m_n^{\circ}}(\mathbb{p}, \mathbf{v}) = \min \{R_n^m(\mathbb{p}, \mathbf{v}) : m \in \mathbb{N}\}. \end{aligned} \quad (18.06)$$

and assuming $\mathbf{v}_p := \max(\|\mathbb{p}\|_{\mathbb{L}_{\infty}(\lambda_{\mathbf{v}_0,1})}, \|\mathbb{p}^{-1}\|_{\mathbb{L}_{\infty}(\lambda_{\mathbf{v}_0,1})}) \in \mathbb{R}_{\mathbf{v}_0}^+$ due to [Property §16.09](#) the (infeasible) OPE $\widehat{\mathbb{p}}^{m_n^{\circ}} = \widehat{\mathbb{p}} \mathbf{1}^{m_n^{\circ}} \in \ell_2(\mathbf{v}^2) \mathbf{1}^{m_n^{\circ}} \subseteq \ell_2(\mathbf{v}^2)$ with oracle dimension m_n° as in (18.06) satisfies

$$\begin{aligned} \mathbf{v}_p^{-1} R_n^{\circ}(\mathbb{p}, \mathbf{v}) &\leq \inf_{m \in \mathbb{N}} \mathbb{P}_p^{\otimes n} (\|\widehat{\mathbb{p}}^m - \mathbb{p}\|_{\mathbf{v}}^2) \leq \mathbb{P}_p^{\otimes n} (\|\widehat{\mathbb{p}}^{m_n^{\circ}} - \mathbb{p}\|_{\mathbf{v}}^2) \\ &\leq \mathbf{v}_p R_n^{\circ}(\mathbb{p}, \mathbf{v}) \leq \mathbf{v}_p^2 \inf_{m \in \mathbb{N}} \mathbb{P}_p^{\otimes n} (\|\widehat{\mathbb{p}}^m - \mathbb{p}\|_{\mathbf{v}}^2), \end{aligned}$$

and hence it is *oracle optimal* (with constant \mathbf{v}_p^2). □

§18.06 **Notation.** Consider a sequence of penalties $\mathbf{pen}_{\bullet}^{\mathbf{v}} = (\mathbf{pen}_m^{\mathbf{v}})_{m \in \mathbb{N}} \in (\mathbb{R}_{\mathbf{v}_0}^+)^{\mathbb{N}}$ given by

$$\mathbf{pen}_m^{\mathbf{v}} := 24 \tau_{\mathbf{v},u}^2 n^{-1} \|\mathbf{1}^m\|_{\mathbf{v}}^2, \quad \text{for each } m \in \mathbb{N} \quad (18.07)$$

and the upper bound (where the defining set is not empty)

$$\mathbb{M}^{\mathbf{v}} := \max \{m \in \mathbb{N} : \|\mathbf{1}^m\|_{\mathbf{v}}^2 \leq n \mathbf{v}_t^2, m \leq \exp(\frac{n^{1/2}}{100})\} \quad (18.08)$$

which are obviously known in advance. Considering the data-driven OSE $\widehat{\mathbb{p}}^{\widehat{m}} = \widehat{\mathbb{p}} \mathbf{1}^{\widehat{m}}$ with dimension parameter

$$\widehat{m} := \arg \min \{-\|\widehat{\mathbb{p}}^m\|_{\mathbf{v}}^2 + \mathbf{pen}_m^{\mathbf{v}} : m \in \llbracket \mathbb{M}^{\mathbf{v}} \rrbracket\} \quad (18.09)$$

we derive below an upper bound for its global \mathbf{v} -risk, $\mathbb{P}_p^{\otimes n} (\|\widehat{\mathbb{p}}^{\widehat{m}} - \mathbb{p}\|_{\mathbf{v}}^2)$. □

§18.07 **Lemma.** Under Assumptions §15.02 and §18.03 and $\mathbb{p} \in \mathbb{L}_{\infty}(\lambda_{\mathbf{v}_0,1})$ for $\mathbf{pen}_{\bullet}^{\mathbf{v}} \in (\mathbb{R}_{\mathbf{v}_0}^+)^{\mathbb{N}}$ as in (18.07) and $\mathbb{M}^{\mathbf{v}} \in \mathbb{N}$ as in (18.08) we have

$$\mathbb{P}_p^{\otimes n} (\max \{(\|\widehat{\mathbb{p}}^m - \mathbb{p}^m\|_{\mathbf{v}}^2 - \mathbf{pen}_m^{\mathbf{v}}/4) : m \in \llbracket \mathbb{M}^{\mathbf{v}} \rrbracket\}) \leq C_{\text{tal}} \tau_{\mathbf{v},u}^2 (C_{\mathbf{v}}(x_p) + \mathbf{v}_t^2) n^{-1} \quad (18.10)$$

for some universal numerical constant $C_{\text{tal}} \in [1, \infty)$ and $x_p := 6 \|\mathbb{p}\|_{\mathbb{L}_{\infty}(\lambda_{\mathbf{v}_0,1})} \tau_{\mathbf{v},u}^{-2} \in \mathbb{R}^+$.

§18.08 **Proof of Lemma §18.07.** is given in the lecture. □

§18.09 **Proposition (Upper bound).** Under Assumptions §15.02 and §18.03 and $\mathbb{p} \in \mathbb{L}_{\infty}(\lambda_{\mathbf{v}_0,1})$ for $\mathbb{M}^{\mathbf{v}} \in \mathbb{N}$ as in (18.08) and $\mathbf{pen}_{\bullet}^{\mathbf{v}} \in (\mathbb{R}_{\mathbf{v}_0}^+)^{\mathbb{N}}$ as in (18.07) the data-driven OPE $\widehat{\mathbb{p}}^{\widehat{m}} = \widehat{\mathbb{p}} \mathbf{1}^{\widehat{m}} \in \ell_2(\mathbf{v}^2) \mathbf{1}^{\widehat{m}} \subseteq \ell_2(\mathbf{v}^2)$ of $\mathbb{p} \in \ell_2(\mathbf{v}^2)$ with data-driven dimension $\widehat{m} \in \llbracket \mathbb{M}^{\mathbf{v}} \rrbracket$ as in (18.09) satisfies

$$\mathbb{P}_p^{\otimes n} (\|\widehat{\mathbb{p}}^{\widehat{m}} - \mathbb{p}\|_{\mathbf{v}}^2) \leq 96 \tau_{\mathbf{v},u}^2 \min \{R_n^m(\mathbb{p}, \mathbf{v}) : m \in \llbracket \mathbb{M}^{\mathbf{v}} \rrbracket\} + C \tau_{\mathbf{v},u}^2 (C_{\mathbf{v}}(x_p) + \mathbf{v}_t^2) n^{-1} \quad (18.11)$$

for some universal numerical constant $C = 8C_{\text{tal}} \in [1, \infty)$ and $x_p := 6 \|\mathbb{p}\|_{\mathbb{L}_{\infty}(\lambda_{\mathbf{v}_0,1})} \tau_{\mathbf{v},u}^{-2} \in \mathbb{R}^+$.

§18.10 **Proof of Proposition §18.09.** is given in the lecture. □

§18|01|02 Maximal global \mathfrak{v} -risk

§18.11 **Assumption.** Consider weights $\mathfrak{a}_\cdot, \mathfrak{v}_\cdot \in (\mathbb{R}_{>0})^{\mathbb{N}}$ with $\mathfrak{a}_\cdot \in \ell_\infty$ and $(\mathfrak{a}\mathfrak{v})_\cdot := (\mathfrak{a}_j \mathfrak{v}_j)_{j \in \mathbb{N}} = \mathfrak{a}_\cdot \mathfrak{v}_\cdot \in \ell_\infty$. We write $(\mathfrak{a}\mathfrak{v})_{(m)} := \|(\mathfrak{a}\mathfrak{v})_\cdot \mathbf{1}_\cdot^{m \perp}\|_{\ell_\infty} \in \mathbb{R}^+$ for each $m \in \mathbb{N}$. The weights $\mathfrak{v}_\cdot \in (\mathbb{R}_{>0})^{\mathbb{N}}$ satisfy (18.05). The orthonormal system $(\mathfrak{u}_j)_{j \in \mathbb{N}}$ in $\mathbb{L}_2(\lambda_{\mathfrak{v},1})$ and $\mathfrak{u}_0 := \mathbf{1}_{[0,1]}$ form an **(os1')** orthonormal basis $(\mathfrak{u}_j)_{j \in \mathbb{N}_0}$ in $\mathbb{L}_2(\lambda_{\mathfrak{v},1})$ and as process $\mathfrak{u}_\cdot^2 = (\mathfrak{u}_j^2)_{j \in \mathbb{N}}$ on $([0,1], \mathcal{B}_{[0,1]})$ satisfies **(os2')** $\|\nu_{\mathbb{N}}(\mathfrak{a}_\cdot^2 \mathfrak{u}_\cdot^2)\|_{\mathbb{L}_\infty(\lambda_{\mathfrak{v},1})} \leq \tau_{\mathfrak{a},\mathfrak{u}}^2$ and **(os2'')** $\sup \{ \|\mathfrak{u}_\cdot(x) \mathbf{1}_\cdot^m\|_{\mathfrak{v}}^2 : x \in [0,1] \} \leq \tau_{\mathfrak{v},\mathfrak{u}} \|\mathbf{1}_\cdot^m\|_{\mathfrak{v}}^2 \in \mathbb{R}^+$ for $\tau_{\mathfrak{a},\mathfrak{u}}, \tau_{\mathfrak{v},\mathfrak{u}} \in [1, \infty)$. \square

§18.12 **Reminder (Maximal global \mathfrak{v} -risk).** Given Assumptions §15.02 and §18.11 we consider an OPE as in Definition §16.04. Here the observable noisy density coefficients $\widehat{\mathfrak{p}} = \mathfrak{p} + n^{-1/2} \boldsymbol{\varepsilon}_\cdot$ of the density coefficients $\mathfrak{p} = \cup \mathfrak{p} \in \ell_2$ take the form of a *statistical direct problem* (see Definition §10.19) where the stochastic processes $\boldsymbol{\varepsilon}_\cdot \in \mathcal{B}_{[0,1]}^{\otimes \mathbb{N}} \otimes 2^{\mathbb{N}}$ is given in Definition §15.08. Under Assumptions §15.02 and §18.11 in Proposition §16.16 an upper bound for a maximal global \mathfrak{v} -risk of an OPE is shown over the set $\mathbb{D}_2^{\mathfrak{a},\mathfrak{r}}$ given in (16.04) (Lemma §16.14). More precisely, the performance of the OPE $\widehat{\mathfrak{p}}^m = \widehat{\mathfrak{p}} \mathbf{1}_\cdot^m \in \ell_2(\mathfrak{v}^2) \mathbf{1}_\cdot^m \subseteq \ell_2(\mathfrak{v}^2)$ with dimension $m \in \mathbb{N}$ is measured by its maximal global \mathfrak{v} -risk over the ellipsoid $\mathbb{D}_2^{\mathfrak{a},\mathfrak{r}}$, that is

$$\mathcal{R}_n^{\mathfrak{v}}[\widehat{\mathfrak{p}}^m | \mathbb{D}_2^{\mathfrak{a},\mathfrak{r}}] := \sup \{ \mathbb{P}_{\mathfrak{p}}^{\otimes n} (\|\widehat{\mathfrak{p}}^m - \mathfrak{p}\|_{\mathfrak{v}}^2) : \mathfrak{p} \in \mathbb{D}_2^{\mathfrak{a},\mathfrak{r}} \}.$$

As in (12.06) for $n, m \in \mathbb{N}$ setting $(\mathfrak{a}\mathfrak{v})_{(m)}^2 := \|(\mathfrak{a}\mathfrak{v})_\cdot^2 \mathbf{1}_\cdot^{m \perp}\|_{\ell_\infty}$ and

$$\begin{aligned} R_n^m(\mathfrak{a}_\cdot, \mathfrak{v}_\cdot) &:= (\mathfrak{a}\mathfrak{v})_{(m)}^2 \vee n^{-1} \|\mathbf{1}_\cdot^m\|_{\mathfrak{v}}^2, & m_n^* &:= \arg \min \{ R_n^m(\mathfrak{a}_\cdot, \mathfrak{v}_\cdot) : m \in \mathbb{N} \} \\ \text{and } R_n^*(\mathfrak{a}_\cdot, \mathfrak{v}_\cdot) &:= R_n^{m_n^*}(\mathfrak{a}_\cdot, \mathfrak{v}_\cdot) = \min \{ R_n^m(\mathfrak{a}_\cdot, \mathfrak{v}_\cdot) : m \in \mathbb{N} \} \end{aligned} \quad (18.12)$$

by Proposition §16.16 under Assumptions §15.02 and §18.03 the maximal global \mathfrak{v} -risk of an OPE $\widehat{\mathfrak{p}}^{m_n^*}$ with optimally chosen dimension m_n^* as in (18.12) satisfies

$$\mathcal{R}_n^{\mathfrak{v}}[\widehat{\mathfrak{p}}^{m_n^*} | \mathbb{D}_2^{\mathfrak{a},\mathfrak{r}}] \leq C R_n^*(\mathfrak{a}_\cdot, \mathfrak{v}_\cdot)$$

with $C = 1 + r\tau_{\mathfrak{a},\mathfrak{u}} + r^2$. Moreover, due to Proposition §17.18 $R_n^*(\mathfrak{a}_\cdot, \mathfrak{v}_\cdot)$ provides (up to a constant) also a lower bound of the maximal global \mathfrak{v} -risk over the ellipsoid $\mathbb{D}_2^{\mathfrak{a},\mathfrak{r}}$ for any estimator. Consequently, (up to a constant) $R_n^*(\mathfrak{a}_\cdot, \mathfrak{v}_\cdot)$ is a minimax bound and $\widehat{\mathfrak{p}}^{m_n^*}$ is minimax optimal. However, the optimal dimension m_n^* depends on $\mathfrak{a}_\cdot \in (\mathbb{R}_{>0})^{\mathbb{N}}$ characterising the ellipsoid $\mathbb{D}_2^{\mathfrak{a},\mathfrak{r}}$. \square

§18.13 **Proposition (Upper bound).** Under Assumptions §15.02 and §18.03 for $M^{\mathfrak{v}} \in \mathbb{N}$ as in (18.08) and $\text{pen}_\cdot^{\mathfrak{v}} \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ as in (18.07) the data-driven OPE $\widehat{\mathfrak{p}}^{\widehat{m}} = \widehat{\mathfrak{p}} \mathbf{1}_\cdot^{\widehat{m}} \in \ell_2(\mathfrak{v}^2) \mathbf{1}_\cdot^{\widehat{m}} \subseteq \ell_2(\mathfrak{v}^2)$ with data-driven dimension $\widehat{m} \in \llbracket M^{\mathfrak{v}} \rrbracket$ as in (18.09) satisfies

$$\mathcal{R}_n^{\mathfrak{v}}[\widehat{\mathfrak{p}}^{\widehat{m}} | \mathbb{D}_2^{\mathfrak{a},\mathfrak{r}}] \leq (3r^2 + 96\tau_{\mathfrak{v},\mathfrak{u}}^2) \min \{ R_n^m(\mathfrak{a}_\cdot, \mathfrak{v}_\cdot) : m \in \llbracket M^{\mathfrak{v}} \rrbracket \} + C\tau_{\mathfrak{v},\mathfrak{u}}^2 (C_{\mathfrak{v}}(\xi) + \mathfrak{v}_1^2) n^{-1} \quad (18.13)$$

for some universal numerical constant $C = 8C_{\text{tal}} \in [1, \infty)$ and $\xi := 6(1 + r\tau_{\mathfrak{a},\mathfrak{u}})\tau_{\mathfrak{v},\mathfrak{u}}^{-2} \in \mathbb{R}^+$.

§18.14 **Proof of Proposition §18.13.** is given in the lecture. \square

§18.15 **Comment.** The minimax bound $R_n^*(\mathfrak{a}_\cdot, \mathfrak{v}_\cdot) = R_n^{m_n^*}(\mathfrak{a}_\cdot, \mathfrak{v}_\cdot) = \min \{ R_n^m(\mathfrak{a}_\cdot, \mathfrak{v}_\cdot) : m \in \mathbb{N} \}$ (for details see Reminder §18.12) satisfies $nR_n^*(\mathfrak{a}_\cdot, \mathfrak{v}_\cdot) \geq \|\mathbf{1}_\cdot^{m_n^*}\|_{\mathfrak{v}}^2 \geq \mathfrak{v}_1^2$. Consequently, the last upper bound in (18.13) and the minimax bound $R_n^*(\mathfrak{a}_\cdot, \mathfrak{v}_\cdot)$ coincide up to a constant $(3r^2 + 96\tau_{\mathfrak{v},\mathfrak{u}}^2 + C\tau_{\mathfrak{v},\mathfrak{u}}^2 (C_{\mathfrak{v}}(x_{\mathfrak{a}})\mathfrak{v}_1^{-2} + 1))$ provided the minimax dimension fulfils $m_n^* \in \llbracket M^{\mathfrak{v}} \rrbracket$. Therefore, we wish the upper bound $M^{\mathfrak{v}}$ to be as large as possible. The next assertion shows that $M^{\mathfrak{v}}$ as in (18.08) is a suitable choice for the upper bound. \square

§18.16 **Corollary.** Under the assumptions of *Proposition §18.13* for each $n \in \mathbb{N}$ such that $R_n^*(\mathbf{a}, \mathbf{v}) \leq \mathbf{v}_1^2$ and $m_n^* \leq \exp(\frac{n^{1/2}}{100})$ we have

$$\begin{aligned} \mathcal{R}_n^{\mathbf{v}}[\widehat{\mathbb{P}}^{\widehat{m}} | \mathbb{D}_2^{\mathbf{a}, \mathbf{r}}] &\leq (3r^2 + 96\tau_{\mathbf{v}, \mathbf{u}}^2) \min \{R_n^m(\mathbf{a}, \mathbf{v}): m \in \llbracket M^{\mathbf{v}} \rrbracket\} + C\tau_{\mathbf{v}, \mathbf{u}}^2 (C_{\mathbf{v}}(x_a) + \mathbf{v}_1^2) n^{-1} \\ &\leq KR_n^*(\mathbf{a}, \mathbf{v}) \end{aligned} \quad (18.14)$$

and, hence up to the constant $K := 3r^2 + 96\tau_{\mathbf{v}, \mathbf{u}}^2 + C\tau_{\mathbf{v}, \mathbf{u}}^2 (C_{\mathbf{v}}(\xi)\mathbf{v}_1^{-2} + 1)$ the feasible data-driven estimator $\widehat{\mathbb{P}}^{\widehat{m}}$ is *minimax optimal*.

§18.17 **Proof of Corollary §18.16.** is given in the lecture. □

§18.18 **Illustration.** Consider the *trigonometric basis* as in *Illustration §16.18* which satisfies Assumption §18.11 (*os1'*), (*os2'*) for all $\mathbf{a}_i \in \ell_2$ and (*os2''*). In Table 02 [§12] (*Illustration §12.26*) the order of the rate $R_n^*(\mathbf{a}, \mathbf{v})$ is depict for the two specifications (*o*) and (*s*). We note that we have $\mathbf{a}_i \in \ell_2$ in case (*o*) for $a > 1/2$ while in case (*s*) for $a \in \mathbb{R}_0^+$. The sequence \mathbf{v} satisfies Assumption §18.11, i.e. (18.05), for $v \geq -1/2$. Moreover, the optimal dimension m_n^* given in Table 02 [§12] satisfies $m_n^* \leq \exp(\frac{n^{1/2}}{100})$, and thus (under the above restrictions) the adaptive density estimator attains the minimax optimal rate $R_n^*(\mathbf{a}, \mathbf{v})$ up to the constant given in *Corollary §18.16*. □

§18|02 Data-driven local estimation by Goldenshluger and Lepskij's method

The next assertion provides our key argument in order to control the deviations of the reminder term. The Bernstein inequality in the formulation (18.15) *Exercise* is for example given in Comte [2015], Appendix B, Lemma B.2.

§18.19 **Lemma (Bernstein inequality).** Let $(Z_i)_{i \in \llbracket n \rrbracket}$ be independent random variables with $\mathbb{P}(Z_i) = 0$, $\mathbb{P}(Z_i^2) \leq v^2 \in \mathbb{R}^+$ and $|Z_i| \leq 2b \in \mathbb{R}^+$ for all $i \in \llbracket n \rrbracket$. Then for any $x \in \mathbb{R}^+$ we have

$$\begin{aligned} \mathbb{P}\left(\frac{1}{n} \sum_{i \in \llbracket n \rrbracket} Z_i \geq x\right) &\leq \max \left\{ \exp\left(-\frac{nx^2}{4v^2}\right), \exp\left(-\frac{nx}{4b}\right) \right\} \text{ and} \\ \mathbb{P}\left(\left|\frac{1}{\sqrt{n}} \sum_{i \in \llbracket n \rrbracket} Z_i\right| \geq x\right) &\leq 2 \max \left\{ \exp\left(-\frac{x^2}{4v^2}\right), \exp\left(-\frac{n^{1/2}x}{4b}\right) \right\}. \end{aligned} \quad (18.15)$$

Moreover, for any $K \in [1, \infty)$ we have

$$\mathbb{P}\left(\left(\left|\frac{1}{\sqrt{n}} \sum_{i \in \llbracket n \rrbracket} Z_i\right|^2 - (4v^2 + 32b^2(\log K)n^{-1}) \log K\right)_+\right) \leq 8K^{-1}\{v^2 + 16b^2n^{-1}\}. \quad (18.16)$$

§18.20 **Proof of Lemma §18.19.** Exercise □

§18.21 **Remark.** Let us briefly reconsider the OPE $\widehat{\mathbb{P}}^m = \widehat{\mathbb{P}} \mathbf{1}^m \in \ell_2 \mathbf{1}^m$ with dimension $m \in \mathbb{N}$ (*Definition §16.04*) where $\widehat{\mathbb{P}} = \widehat{\mathbb{P}}_n \mathbf{u} = (\widehat{\mathbb{P}}_n \mathbf{u}_j)_{j \in \mathbb{N}}$ are noisy versions (*Definition §15.08*) of the density coefficients $\mathbb{P} = \mathbb{U} \mathbb{p} = \mathbb{P}_p \mathbf{u} = (\mathbb{P}_p \mathbf{u}_j = \lambda_{[0,1]}(\mathbb{p} \mathbf{u}_j))_{j \in \mathbb{N}}$. Clearly, $\phi_{\mathbb{U}_N}((\mathbf{u} - \mathbb{P}) \mathbf{1}^m)$ is a $\mathcal{B}_{[0,1]}$ - \mathcal{B} -measurable function. Therefore, given $(X_i)_{i \in \llbracket n \rrbracket} \sim \mathbb{P}_p^{\otimes n}$ for $i \in \llbracket n \rrbracket$ setting $Z_i := \phi_{\mathbb{U}_N}((\mathbf{u}(X_i) - \mathbb{P}) \mathbf{1}^m)$ we have $\mathbb{P}_p(Z_i) = 0$ exploiting $\mathbb{P} = \mathbb{P}_p \mathbf{u}$, and

$$\phi_{\mathbb{U}_N}((\widehat{\mathbb{P}} - \mathbb{P}) \mathbf{1}^m) = \widehat{\mathbb{P}}_n(\phi_{\mathbb{U}_N}((\mathbf{u} - \mathbb{P}) \mathbf{1}^m)) = n^{-1} \sum_{i \in \llbracket n \rrbracket} Z_i.$$

Consequently, provided that

$$\begin{aligned} \mathbb{P}_p(Z_i^2) &= \mathbb{P}_p(|\phi_{\mathcal{U}_N}((\mathbf{u} - \mathbb{p})\mathbf{1}^m)|^2) \leq v_{p,m}^2 \in \mathbb{R}^+, \\ \sup \{|\phi_{\mathcal{U}_N}(\mathbf{u}(x)\mathbf{1}^m)| : x \in [0, 1]\} &\leq b_m \in \mathbb{R}^+, \text{ and hence } |Z_i| \leq 2b_m, \forall i \in \llbracket n \rrbracket, \end{aligned} \quad (18.17)$$

due to the Bernstein inequality ([Lemma §18.19 \(18.16\)](#)) we have

$$\begin{aligned} \mathbb{P}((|n^{1/2}\phi_{\mathcal{U}_N}((\widehat{\mathbb{p}} - \mathbb{p})\mathbf{1}^m)|^2 - (4v_{p,m}^2 + 32b_m^2(\log K)n^{-1})\log K)_+) \\ \leq 8K^{-1}\{v_{p,m}^2 + 16b_m^2n^{-1}\}. \end{aligned} \quad (18.18)$$

for any $K \in [1, \infty)$. \square

§18|02|01 Local ϕ -risk

§18.22 **Assumption.** Let $\phi \in (\mathbb{R}_{\setminus 0})^{\mathbb{N}}$ and the orthonormal system $(\mathbf{u}_j)_{j \in \mathbb{N}}$ in $\mathbb{L}_2(\lambda_{0,1})$ and $\mathbf{u}_0 := \mathbf{1}_{[0,1]}$ form an **(os1')** orthonormal basis $(\mathbf{u}_j)_{j \in \mathbb{N}_0}$ in $\mathbb{L}_2(\lambda_{0,1})$ and as process $\mathbf{u} = (\mathbf{u}_j)_{j \in \mathbb{N}}$ on $([0, 1], \mathcal{B}_{[0,1]})$ for all $m \in \mathbb{N}$ satisfies **(os2'')** $\sup \{\|\mathbf{u}(x)\mathbf{1}^m\|_{\ell_2}^2 : x \in [0, 1]\} \leq \tau_u^2 m \in \mathbb{R}^+$ for $\tau_u \in [1, \infty)$. \square

§18.23 **Remark.** We replace Assumption §15.05 **(os1)** and **(os2)**, respectively, by the stronger Assumption §18.22 **(os1')** and **(os2'')**. Indeed, under **(os1')** we have **(os1)** $\mathbf{1}_{[0,1]} \in \ker(U)$. Furthermore, $(\mathbf{u}_j)_{j \in \mathbb{N}}$ belongs to $\mathbb{L}_\infty(\lambda_{0,1})$ due to **(os2'')**, and hence **(os2)** is fulfilled (see also [Remark §15.06](#)). We use in the sequel that under Assumption §18.22 **(os2'')** for each $m \in \mathbb{N}$

$$\begin{aligned} \sup \{|\phi_{\mathcal{U}_N}(\mathbf{u}(x))\mathbf{1}^m|^2 : x \in [0, 1]\} \\ \leq \|\mathbf{1}^m\|_\phi^2 \sup \{\|\mathbf{u}(x)\mathbf{1}^m\|_{\ell_2}^2 : x \in [0, 1]\} \leq \tau_u^2 m \|\mathbf{1}^m\|_\phi^2 =: b_m^2 \end{aligned} \quad (18.19)$$

by applying the Cauchy Schwarz inequality and moreover (see [Proof §15.11](#))

$$\mathbb{P}_p(|\phi_{\mathcal{U}_N}((\mathbf{u} - \mathbb{p})\mathbf{1}^m)|^2) \leq \mathbb{P}_p(|\phi_{\mathcal{U}_N}(\mathbf{u}\mathbf{1}^m)|^2) =: v_{p,m}^2 \leq \|\mathbb{p}\|_{\mathbb{L}_\infty(\lambda_{0,1})} \|\mathbf{1}^m\|_\phi^2 \quad (18.20)$$

exploiting [Lemma §15.10 \(i\)](#). Combining (18.19), (18.20) and (18.18) ([Remark §18.21](#)) we obtain

$$\begin{aligned} \mathbb{P}((|n^{1/2}\phi_{\mathcal{U}_N}((\widehat{\mathbb{p}} - \mathbb{p})\mathbf{1}^m)|^2 - (4v_{p,m}^2 + 32b_m^2(\log K)n^{-1})\log K)_+) \\ \leq 8K^{-1}\{\|\mathbb{p}\|_{\mathbb{L}_\infty(\lambda_{0,1})} + 16\tau_u^2 mn^{-1}\} \|\mathbf{1}^m\|_\phi^2 \end{aligned} \quad (18.21)$$

for any $m \in \mathbb{N}$ and $K \in [1, \infty)$. \square

§18.24 **Reminder (Local oracle ϕ -risk).** Given Assumptions §15.02 and §18.22 we consider an OPE as in [Definition §16.04](#). Here the observable noisy density coefficients $\widehat{\mathbb{p}} = \mathbb{p} + n^{-1/2}\boldsymbol{\varepsilon}$ of the density coefficients $\mathbb{p} = U\mathbb{p} \in \ell_2$ take the form of a *statistical direct problem* (see [Definition §10.19](#)) where the stochastic processes $\boldsymbol{\varepsilon} \in \mathcal{B}_{[0,1]}^{\otimes n} \otimes 2^{\mathbb{N}}$ is given in [Definition §15.08](#). Under Assumptions §15.02 and §18.22, (and hence Assumption §15.05 and $\phi \in (\mathbb{R}_{\setminus 0})^{\mathbb{N}}$ see [Remark §18.23](#)) and $\mathbb{p} \in \text{dom}(\phi_{\mathcal{U}_N})$ in §16.22 an *oracle inequality* for the local ϕ -risk of the OPE's is shown. More precisely, as in (16.06) ([Proposition §16.20](#)) for all $n, m \in \mathbb{N}$ setting

$$\begin{aligned} R_n^m(\mathbb{p}, \phi) &:= |\phi_{\mathcal{U}_N}(\mathbb{p}\mathbf{1}^{m\perp})|^2 + n^{-1}\|\mathbf{1}^m\|_\phi^2, \quad m_n^\circ := \arg \min \{R_n^m(\mathbb{p}, \phi) : m \in \mathbb{N}\} \\ \text{and } R_n^\circ(\mathbb{p}, \phi) &:= R_n^{m_n^\circ}(\mathbb{p}, \phi) = \min \{R_n^m(\mathbb{p}, \phi) : m \in \mathbb{N}\}. \end{aligned} \quad (18.22)$$

and assuming $v_p := \max(\|p\|_{\mathbb{L}_\infty(\lambda_{0,1})}, \|p^{-1}\|_{\mathbb{L}_\infty(\lambda_{0,1})}) \in \mathbb{R}_{>0}^+$, and hence $\max(\|\Gamma_p\|_{\mathbb{L}(\ell_2)}, \|\Gamma_p^{-1}\|_{\mathbb{L}(\ell_2)}) \leq v_p$ (see Lemma §15.10), due to Property §16.22 the (infeasible) OPE $\widehat{p}^{m_n^c} = \widehat{p} \mathbb{1}^{m_n^c} \in \ell_2 \mathbb{1}^{m_n^c} \subseteq \text{dom}(\phi\nu)$ with oracle dimension m_n^c as in (18.22) satisfies

$$\begin{aligned} v_p^{-1} R_n^c(p, \phi) &\leq \inf_{m \in \mathbb{N}} \mathbb{P}_p^{\otimes n} (|\phi\nu_N(\widehat{p}^m - p)|^2) \leq \mathbb{P}_p^{\otimes n} (|\phi\nu_N(\widehat{p}^{m_n^c} - p)|^2) \\ &\leq v_p R_n^c(p, \phi) \leq v_p^2 \inf_{m \in \mathbb{N}} \mathbb{P}_p^{\otimes n} (|\phi\nu_N(\widehat{p}^m - p)|^2), \end{aligned}$$

and hence it is *oracle optimal* (with constant v_p^2). \square

Partially known penalty sequence

§18.25 **Notation.** Consider first a sequence of penalties $\text{pen}_m^{p,\phi} = (\text{pen}_m^{p,\phi})_{m \in \mathbb{N}} \in (\mathbb{R}_0^+)^{\mathbb{N}}$ given by

$$\begin{aligned} \text{pen}_m^{p,\phi} &:= 12n^{-1} (v_{p,m}^2 + 8b_m^2 (\log K_m) n^{-1}) (\log K_m) \quad \text{with } v_{p,m}^2 := \mathbb{P}_p (|\phi\nu_N(u, \mathbb{1}^m)|^2), \\ b_m^2 &:= \tau_u^2 m \|\mathbb{1}^m\|_\phi^2, \quad \text{and } K_m := (1 \vee \|\mathbb{1}^m\|_\phi^2) m^3 \in [1, \infty) \quad \text{for each } m \in \mathbb{N}, \end{aligned} \quad (18.23)$$

which is obviously only *partially known* in advance, and arbitrary but fixed upper bound $M \in \mathbb{N}$. Considering the data-driven OSE $\widehat{p}^{\widehat{m}} = \widehat{p} \mathbb{1}^{\widehat{m}}$ with dimension parameter selected by Goldenshluger and Lepskij's method

$$\begin{aligned} \widehat{m} &:= \arg \min \{ \text{contr}_m^{p,\phi} + \text{pen}_m^{p,\phi} : m \in \llbracket M \rrbracket \} \quad \text{and} \\ \text{contr}_m^{p,\phi} &:= \max \{ (|\phi\nu_N(\widehat{p}^j - \widehat{p}^m)|^2 - \text{pen}_j^{p,\phi} - \text{pen}_m^{p,\phi})_+ : j \in \llbracket m, M \rrbracket \}, \quad m \in \llbracket M \rrbracket. \end{aligned} \quad (18.24)$$

Moreover, studying a ϕ -error the bias term introduced in (14.31) becomes

$$\text{bias}_m(p, \phi) = \sup \{ |\phi\nu_N(p^j - p^m)| = |\phi\nu_N(p \mathbb{1}^{m,j})| : j \in \llbracket m, \infty \rrbracket \} \quad \forall m \in \mathbb{N}.$$

If $p \in \text{dom}(\phi\nu_N)$ and hence $\nu_N(|\phi p|) \in \mathbb{R}$ then $\text{bias}_m(p, \phi) \leq \nu_N(|\phi p| \mathbb{1}^{m,\perp}) = o(1)$ as $m \rightarrow \infty$ by dominated convergence. Considering the data-driven OSE $\widehat{p}^{\widehat{m}} = \widehat{p} \mathbb{1}^{\widehat{m}}$ with dimension parameter \widehat{m} selected as in (18.24) with penalty sequence $\text{pen}_m^{p,\phi}$ given in (18.23) and arbitrary upper bound $M \in \mathbb{N}$ we derive below an upper bound for its local ϕ -risk, $\mathbb{P}_p^{\otimes n} (|\phi\nu_N(\widehat{p}^{\widehat{m}} - p)|^2)$. \square

§18.26 **Lemma.** Under Assumptions §15.02 and §18.22 and $p \in \mathbb{L}_\infty(\lambda_{0,1})$ for $\text{pen}_m^{p,\phi} \in (\mathbb{R}_0^+)^{\mathbb{N}}$ as in (18.23) and for any $M \in \mathbb{N}$ we have

$$\mathbb{P}_p^{\otimes n} (\max \{ (|\phi\nu_N(\widehat{p}^m - p^m)|^2 - \text{pen}_m^{p,\phi}/3)_+ : m \in \llbracket M \rrbracket \}) \leq 14 \{ \|p\|_{\mathbb{L}_\infty(\lambda_{0,1})} + 16\tau_u^2 n^{-1} \} n^{-1}. \quad (18.25)$$

§18.27 **Proof of Lemma §18.26.** is given in the lecture. \square

§18.28 **Proposition (Upper bound).** Under Assumptions §15.02 and §18.22 and $p \in \mathbb{L}_\infty(\lambda_{0,1})$ for $\text{pen}_m^{p,\phi} \in (\mathbb{R}_0^+)^{\mathbb{N}}$ as in (18.23) the data-driven OPE $\widehat{p}^{\widehat{m}} = \widehat{p} \mathbb{1}^{\widehat{m}} \in \ell_2 \mathbb{1}^{\widehat{m}} \subseteq \text{dom}(\phi\nu_N)$ of $p \in \text{dom}(\phi\nu_N)$ with data-driven dimension $\widehat{m} \in \llbracket M \rrbracket$ as in (18.24) satisfies for all $n, M \in \mathbb{N}$

$$\begin{aligned} \mathbb{P}_p^{\otimes n} (|\phi\nu_N(\widehat{p}^{\widehat{m}} - p)|^2) &\leq 64 (\|p\|_{\mathbb{L}_\infty(\lambda_{0,1})} + 8\tau_u^2) \\ &\quad \times \min \{ \text{bias}_m^2(p, \phi) + n^{-1} \|\mathbb{1}^m\|_\phi^2 (\log K_m) (1 \vee (\log K_m) m n^{-1}) : m \in \llbracket M \rrbracket \} \\ &\quad + 392 (\|p\|_{\mathbb{L}_\infty(\lambda_{0,1})} + 16\tau_u^2 n^{-1}) n^{-1}. \end{aligned} \quad (18.26)$$

§18.29 **Proof of Proposition §18.28.** is given in the lecture. \square

§18.30 **Comment.** Let us compare the dominating part of the upper bound given in (18.26), that is

$$\min \left\{ \text{bias}_m^2(\mathbb{p}, \phi) + n^{-1} \|\mathbb{1}^m\|_\phi^2 (\log K_m) (1 \vee (\log K_m) m n^{-1}) : m \in \llbracket M \rrbracket \right\} \quad (18.27)$$

with the oracle bound $R_n^\circ(\theta, \phi) = \min \left\{ |\phi \nu_{\mathbb{N}}(\mathbb{p}^m - \mathbb{p})|^2 + n^{-1} \|\mathbb{1}^m\|_\phi^2 : m \in \mathbb{N} \right\}$ (for details see **Reminder** §18.24). In (18.27) we face eventually a deterioration by three sources. First, we generally have $\text{bias}_m(\mathbb{p}, \phi) \geq |\phi \nu_{\mathbb{N}}(\mathbb{p}^m - \mathbb{p})|$, but note that for $\mathbb{p}, \phi \in (\mathbb{R}^+)^{\mathbb{N}}$ equality holds, that is

$$\text{bias}_m(\mathbb{p}, \phi) = \sup \left\{ \nu_{\mathbb{N}}(\phi \mathbb{p} \mathbb{1}_*^{m,j}) : j \in \llbracket m, \infty \rrbracket \right\} = \nu_{\mathbb{N}}(\phi \mathbb{p} \mathbb{1}_*^{m,\perp}) = |\phi \nu_{\mathbb{N}}(\mathbb{p}^m - \mathbb{p})|$$

for all $m \in \mathbb{N}$. Secondly, the variance term features an additional factor $(\log K_m) (1 \vee (\log K_m) m n^{-1})$, and finally the upper bound M might impose an additional deterioration. We note that the oracle bound $R_n^\circ(\mathbb{p}, \phi)$ is parametric, i.e. $n R_n^\circ(\mathbb{p}, \phi) = O(1)$ as $n \rightarrow \infty$, if $\phi \in \ell_2$ (case **(p)** in **Illustration** §12.40). In the sequel we consider only the case $\phi \notin \ell_2$, i.e. $\nu_{\mathbb{N}}(|\phi|^2) = \infty$. We set

$$M^\phi := \max \left\{ m \in \mathbb{N} : \|\mathbb{1}^m\|_\phi^2 \leq n \phi_1^2 \right\} \in \mathbb{N} \quad (18.28)$$

where the defining set is not empty and finite since $\|\phi\|_{\ell_2}^2 = \infty$. The next assertion shows that this is a suitable choice for the upper bound. Moreover, we estimate the bias term by $\text{bias}_m(\mathbb{p}, \phi) \leq \nu(|\phi \mathbb{p} \mathbb{1}_*^{m,\perp}|)$ where equality holds whenever $\mathbb{p}, \phi \in (\mathbb{R}^+)^{\mathbb{N}}$. \square

§18.31 **Corollary.** Given $\phi \in (\mathbb{R}_{>0})^{\mathbb{N}}$ with $\phi \notin \ell_2$, $M^\phi \in \mathbb{N}$ as in (18.28) and $\text{pen}^{\mathbb{p},\phi}$ as in (18.23) consider a data-driven OPE $\widehat{\mathbb{p}}^{\widehat{m}} = \widehat{\mathbb{p}} \mathbb{1}_*^{\widehat{m}} \in \ell_2 \mathbb{1}_*^{\widehat{m}} \subseteq \text{dom}(\phi \nu_{\mathbb{N}})$ of $\mathbb{p} \in \text{dom}(\phi \nu_{\mathbb{N}})$ with

$$\begin{aligned} \widehat{m} &:= \arg \min \left\{ \text{contr}_m^\phi + \text{pen}_m^{\mathbb{p},\phi} : m \in \llbracket M^\phi \rrbracket \right\} \quad \text{and} \\ \text{contr}_m^{\mathbb{p},\phi} &:= \max \left\{ (|\phi \nu_{\mathbb{N}}(\widehat{\theta}^j - \widehat{\theta}^m)|^2 - \text{pen}_j^{\mathbb{p},\phi} - \text{pen}_m^{\mathbb{p},\phi})_+ : j \in \llbracket m, M^\phi \rrbracket \right\}, \quad m \in \llbracket M^\phi \rrbracket. \end{aligned} \quad (18.29)$$

For $n, m \in \mathbb{N}$ we set

$$\begin{aligned} R_n^m(\mathbb{p}, \phi) &:= \left(\nu_{\mathbb{N}}(|\phi \mathbb{p} \mathbb{1}_*^{m,\perp}|) \right)^2 \\ &\quad + (1 + (\log \|\mathbb{1}^m\|_\phi^2)_+ + \log m) (1 + ((\log \|\mathbb{1}^m\|_\phi^2)_+ + \log m) m n^{-1}) n^{-1} \|\mathbb{1}^m\|_\phi^2, \\ m^\circ &:= \arg \min \left\{ R_n^m(\mathbb{p}, \phi) : m \in \mathbb{N} \right\} \quad \text{and} \\ R_n^\circ(\mathbb{p}, \phi) &:= R_n^{m^\circ}(\mathbb{p}, \phi) = \min \left\{ R_n^m(\mathbb{p}, \phi) : m \in \mathbb{N} \right\}. \end{aligned} \quad (18.30)$$

Under the assumptions of **Proposition** §18.28 for each $n \in \mathbb{N}$ such that $R_n^\circ(\mathbb{p}, \phi) \leq \phi_1^2$ we have

$$\begin{aligned} \mathbb{P}_p^{\otimes n} (|\phi \nu_{\mathbb{N}}(\widehat{\mathbb{p}}^{\widehat{m}} - \mathbb{p})|^2) &\leq 576 (\|\mathbb{p}\|_{\mathbb{L}_\infty(\lambda_{0,i})} + 8\tau_u^2) R_n^\circ(\mathbb{p}, \phi) + 392 (\|\mathbb{p}\|_{\mathbb{L}_\infty(\lambda_{0,i})} + 16\tau_u^2 n^{-1}) n^{-1} \\ &\leq (576 + 784\phi_1^{-2}) (\|\mathbb{p}\|_{\mathbb{L}_\infty(\lambda_{0,i})} + 8\tau_u^2) R_n^\circ(\mathbb{p}, \phi). \end{aligned} \quad (18.31)$$

§18.32 **Proof** of **Corollary** §18.31. is given in the lecture. \square

§18.33 **Comment.** The data-driven bound $R_n^\circ(\mathbb{p}, \phi)$ compared to the oracle bound $R_n^\circ(\mathbb{p}, \phi)$ features a deterioration of the variance term at least by a logarithmic factor. The appearance of the logarithmic factor within the bound is a known fact in the context of local estimation (cf. Laurent et al. [2008] who consider model selection given direct Gaussian observations). Brown and Low [1996] show that it is unavoidable in the context of nonparametric Gaussian regression and hence it is widely considered as an acceptable price for adaptation. \square

§18.34 **Illustration.** We illustrate the last results considering the two specifications **(o)** and **(s)** given in Table 03 [§12] (**Illustration** §12.40). We restrict ourselves to the case $\phi \notin \ell_2$ only.

Table 01 [§18]

Order of the oracle rate $R_n^\circ(\mathbb{p}, \phi)$ and the data-driven rate $R_n^\diamond(\mathbb{p}, \phi)$ as $n \rightarrow \infty$

$(j \in \mathbb{N})$	$(a \in \mathbb{R}_{\setminus 0}^+)$	(squared bias)	(variance)	M^ϕ	m°	$R_n^\circ(\mathbb{p}, \phi)$	$R_n^\diamond(\mathbb{p}, \phi)$	
$\phi = j^{v-1/2}$	\mathbb{p}	$(\nu_{\mathbb{N}}(\phi _{\mathbb{p}} \mathbf{1}^{m ^\perp}))^2$	$\ \mathbf{1}^m\ _\phi^2$					
(o)	$v \in (0, a)$	$j^{-a-1/2}$	$m^{-2(a-v)}$	m^{2v}	$n^{\frac{1}{2v}}$	$n^{-\frac{(a-v)}{a}}$	$\left(\frac{\log n}{n}\right)^{\frac{(a-v)}{a}}$	
	$a \in (1/2, \infty)$				$\left(\frac{n}{\log n}\right)^{\frac{1}{2a}}$		$\left(\frac{\log n}{n}\right)^{\frac{2(a-v)}{a+1/2}}$	
	$a \in (0, 1/2]$				$\left(\frac{n}{\log n}\right)^{\frac{1}{a+1/2}}$			
	$v = 0$	$j^{-a-1/2}$	m^{-2a}	$\log m$	e^n	$\frac{\log n}{n}$	$\frac{(\log n)^2}{n}$	
	$a \in (1/2, \infty)$				$\left(\frac{n}{(\log n)^2}\right)^{\frac{1}{2a}}$		$\frac{(\log n)^2}{n}$	
	$a \in (0, 1/2]$				$\left(\frac{n^2}{(\log n)^3}\right)^{\frac{1}{2a+1}}$		$\left(\frac{(\log n)^3}{n^2}\right)^{\frac{a}{a+1/2}}$	
(s)	$v \in \mathbb{R}_{\setminus 0}^+$	$e^{-j^{2a}}$	$m^{(1-2(a-v))+} e^{-2m^{2a}}$	m^{2v}	$n^{\frac{1}{2v}}$	$(\log n)^{\frac{1}{2a}}$	$\frac{(\log n)^{\frac{v}{a}}}{n}$	$\frac{(\log n)^{\frac{v}{a}} (\log \log n)}{n}$
	$v = 0$	$e^{-j^{2a}}$	$m^{(1-2a)+} e^{-2m^{2a}}$	$\log m$	e^n	$(\log n)^{\frac{1}{2a}}$	$\frac{\log \log n}{n}$	$\frac{(\log \log n)^2}{n}$

We note that in Table 01 [§18] the order of the oracle rate $R_n^\circ(\mathbb{p}, \phi)$ and the data-driven rate $R_n^\diamond(\mathbb{p}, \phi)$ is depicted for $v \geq 0$ only. In case $v < 0$ we have $\phi \in \ell_2$ and thus **Corollary** §18.31 is not applicable. Moreover, in case **(s)** for $a \in \mathbb{R}_{\setminus 0}^+$ and **(o)** for $a \in (1/2, \infty)$ the rate $R_n^\circ(\mathbb{p}, \phi)$ features only an additional logarithmic factor compared with the oracle rate $R_n^\circ(\mathbb{p}, \phi)$. \square

Estimated penalty sequence

§18.35 **Notation.** The penalty sequence $\text{pen}_{\cdot, v}^{\mathbb{p}, \phi} \in (\mathbb{R}_0^+)^{\mathbb{N}}$ given in (18.23) still depends on characteristics of the unknown density \mathbb{p} . More precisely, for $m \in \mathbb{N}$ the term $\text{pen}_m^{\mathbb{p}, \phi}$ involves the quantity $\mathbf{v}_{\mathbb{p}, m}^2 = \mathbb{P}_{\mathbb{p}}(|\phi \nu_{\mathbb{N}}(\mathbf{u}, \mathbf{1}^m)|^2)$ which we eventually estimate without bias by $\widehat{\mathbf{v}}_m^2 := \widehat{\mathbb{P}}_{\mathbb{p}}(|\phi \nu_{\mathbb{N}}(\mathbf{u}, \mathbf{1}^m)|^2)$. Based on this estimator let us introduce a fully data-driven sequence of penalties $\widehat{\text{pen}}_{\cdot}^{\phi} = (\widehat{\text{pen}}_m^{\phi})_{m \in \mathbb{N}} \in (\mathbb{R}_{\setminus 0}^+)^{\mathbb{N}}$ given by

$$\begin{aligned} \widehat{\text{pen}}_m^{\phi} &:= 12n^{-1} (2\widehat{\mathbf{v}}_m^2 + 3 \times 8b_m^2 (\log K_m) n^{-1}) (\log K_m) \quad \text{with} \quad \widehat{\mathbf{v}}_m^2 := \widehat{\mathbb{P}}_{\mathbb{p}}(|\phi \nu_{\mathbb{N}}(\mathbf{u}, \mathbf{1}^m)|^2), \\ b_m^2 &:= \tau_u^2 m \|\mathbf{1}^m\|_\phi^2, \quad \text{and} \quad K_m := (1 \vee \|\mathbf{1}^m\|_\phi^2) m^3 \in [1, \infty) \quad \text{for each } m \in \mathbb{N}, \end{aligned} \quad (18.32)$$

which is now *fully known* in advance, and arbitrary but fixed upper bound $M \in \mathbb{N}$. Considering the data-driven OSE $\widehat{\mathbb{p}}^{\widehat{m}} = \widehat{\mathbb{p}} \mathbf{1}^{\widehat{m}}$ with dimension parameter selected by Goldenshluger and Lepskij's method

$$\begin{aligned} \widehat{m} &:= \arg \min \{ \widehat{\text{contr}}_m^{\phi} + \widehat{\text{pen}}_m^{\phi} : m \in \llbracket M \rrbracket \} \quad \text{and} \\ \widehat{\text{contr}}_m^{\phi} &:= \max \{ (|\phi \nu_{\mathbb{N}}(\widehat{\mathbb{p}}^j - \widehat{\mathbb{p}}^m)|^2 - \widehat{\text{pen}}_j^{\phi} - \widehat{\text{pen}}_m^{\phi})_+ : j \in \llbracket m, M \rrbracket \}, \quad m \in \llbracket M \rrbracket \end{aligned} \quad (18.33)$$

we derive below an upper bound for its local ϕ -risk, $\mathbb{P}_{\mathbb{p}}^{\otimes n}(|\phi \nu_{\mathbb{N}}(\widehat{\mathbb{p}}^{\widehat{m}} - \mathbb{p})|^2)$. \square

§18.36 **Lemma.** Under Assumptions §15.02 and §18.22 and $\mathbb{p} \in \mathbb{L}_{\infty}(\lambda_{[0,1]})$ for $\text{pen}_{\cdot, v}^{\mathbb{p}, \phi}, \widehat{\text{pen}}_{\cdot}^{\phi} \in (\mathbb{R}_{\setminus 0}^+)^{\mathbb{N}}$ as in (18.23) and (18.32), respectively, and for any $M \in \mathbb{N}$ we have

$$\mathbb{P}_{\mathbb{p}}^{\otimes n} \left(\max \{ (\text{pen}_j^{\mathbb{p}, \phi} - \widehat{\text{pen}}_j^{\phi})_+ : j \in \llbracket M \rrbracket \} \right) \leq 40 \{ \|\mathbb{p}\|_{\mathbb{L}_{\infty}(\lambda_{[0,1]})} + 6\tau_u^2 n^{-1} \} n^{-1}. \quad (18.34)$$

§18.37 **Proof of Lemma** §18.36. is given in the lecture. \square

§18.38 **Proposition (Upper bound).** Under Assumptions §15.02 and §18.22 and $\mathbb{p} \in \mathbb{L}_\infty(\lambda_{[0,1]})$ for $\widehat{\mathbb{p}}_{\text{en}}^\phi \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ as in (18.32) the data-driven OPE $\widehat{\mathbb{p}}^{\widehat{m}} = \widehat{\mathbb{p}} \mathbb{1}^{\widehat{m}} \in \ell_2 \mathbb{1}^{\widehat{m}} \subseteq \text{dom}(\phi_{\nu_N})$ of $\mathbb{p} \in \text{dom}(\phi_{\nu_N})$ with data-driven dimension $\widehat{m} \in \llbracket M \rrbracket$ as in (18.33) satisfies for all $n, M \in \mathbb{N}$

$$\begin{aligned} \mathbb{E}_p^{\otimes n} (|\phi_{\nu_N}(\widehat{\mathbb{p}}^{\widehat{m}} - \mathbb{p})|^2) &\leq 112(\|\mathbb{p}\|_{\mathbb{L}_\infty(\lambda_{[0,1]})} + 12\tau_u^2) \\ &\quad \times \min \left\{ \text{bias}_m^2(\mathbb{p}, \phi) + n^{-1} \|\mathbb{1}^m\|_\phi^2 (\log K_m)(1 \vee (\log K_m)mn^{-1}): m \in \llbracket M \rrbracket \right\} \\ &\quad + 1440(\|\mathbb{p}\|_{\mathbb{L}_\infty(\lambda_{[0,1]})} + 16\tau_u^2 n^{-1})n^{-1}. \end{aligned} \quad (18.35)$$

§18.39 **Proof of Proposition §18.38.** is given in the lecture. \square

§18.40 **Comment.** We shall stress that the last upper bound (18.35) in **Proposition §18.38** (for the fully data-driven procedure) and the upper bound (18.26) in **Proposition §18.28** (for the partially data-driven procedure) differ only in the numerical constants. Thus, thus the proof of the next results follows line by line their counterparts above. \square

§18.41 **Corollary.** Given $\phi \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ with $\phi \notin \ell_2$, $M^\phi \in \mathbb{N}$ as in (18.28) and $\widehat{\mathbb{p}}_{\text{en}}^\phi$ as in (18.32) consider a data-driven OPE $\widehat{\mathbb{p}}^{\widehat{m}} = \widehat{\mathbb{p}} \mathbb{1}^{\widehat{m}} \in \ell_2 \mathbb{1}^{\widehat{m}} \subseteq \text{dom}(\phi_{\nu_N})$ of $\mathbb{p} \in \text{dom}(\phi_{\nu_N})$ with

$$\begin{aligned} \widehat{m} &:= \arg \min \left\{ \widehat{\text{contr}}_m^\phi + \widehat{\text{pen}}_m^\phi : m \in \llbracket M^\phi \rrbracket \right\} \quad \text{and} \\ \widehat{\text{contr}}_m^\phi &:= \max \left\{ (|\phi_{\nu_N}(\widehat{\theta}^j - \widehat{\theta}^m)|^2 - \widehat{\text{pen}}_j^\phi - \widehat{\text{pen}}_m^\phi)_+ : j \in \llbracket m, M^\phi \rrbracket \right\}, \quad m \in \llbracket M^\phi \rrbracket. \end{aligned} \quad (18.36)$$

For $n, m \in \mathbb{N}$ let m^\diamond and $R_n^\diamond(\mathbb{p}, \phi)$ defined as in (18.30). Under the assumptions of **Proposition §18.38** for each $n \in \mathbb{N}$ such that $R_n^\diamond(\mathbb{p}, \phi) \leq \phi_1^2$ we have

$$\begin{aligned} \mathbb{E}_p^{\otimes n} (|\phi_{\nu_N}(\widehat{\mathbb{p}}^{\widehat{m}} - \mathbb{p})|^2) &\leq 1008(\|\mathbb{p}\|_{\mathbb{L}_\infty(\lambda_{[0,1]})} + 8\tau_u^2)R_n^\diamond(\mathbb{p}, \phi) + 1440(\|\mathbb{p}\|_{\mathbb{L}_\infty(\lambda_{[0,1]})} + 16\tau_u^2 n^{-1})n^{-1} \\ &\leq (1008 + 1920\phi_1^{-2})(\|\mathbb{p}\|_{\mathbb{L}_\infty(\lambda_{[0,1]})} + 12\tau_u^2)R_n^\diamond(\mathbb{p}, \phi). \end{aligned} \quad (18.37)$$

§18.42 **Proof of Proof §18.42.** is given in the lecture. \square

§18.43 **Comment.** The fullay data-driven bound $R_n^\diamond(\mathbb{p}, \phi)$ equals exactly the bound in the partially known case. Therefore, the **Comment §18.33** and the **Illustration §18.34** apply here equally. \square

§18|02|02 Maximal local ϕ -risk

§18.44 **Assumption.** Consider $\phi, \alpha \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ with $\alpha \in \ell_\infty$ and $(\alpha\phi)_\cdot := (\alpha_j \phi_j)_{j \in \mathbb{N}} = \alpha \cdot \phi \in \ell_2$, and hence $\|\alpha \cdot \mathbb{1}^{m \perp}\|_\phi = \|(\alpha\phi) \cdot \mathbb{1}^{m \perp}\|_{\ell_2} = o(1)$ as $m \rightarrow \infty$. The orthonormal system $(u_j)_{j \in \mathbb{N}}$ in $\mathbb{L}_2(\lambda_{[0,1]})$ and $u_0 := \mathbb{1}_{[0,1]}$ form an **(os1')** orthonormal basis $(u_j)_{j \in \mathbb{N}_0}$ in $\mathbb{L}_2(\lambda_{[0,1]})$ and as process $u_\cdot^2 = (u_j^2)_{j \in \mathbb{N}}$ on $([0, 1], \mathcal{B}_{[0,1]})$ satisfies **(os2')** $\|u_\cdot(\alpha_\cdot^2 u_\cdot^2)\|_{\mathbb{L}_\infty(\lambda_{[0,1]})} \leq \tau_{\alpha, u}^2$ and **(os2'')** $\sup \{ \|u_\cdot(x) \mathbb{1}^m\|_{\ell_2}^2 : x \in [0, 1] \} \leq \tau_u m \in \mathbb{R}^+$ for $\tau_{\alpha, u}, \tau_u \in [1, \infty)$. \square

§18.45 **Remark.** Assumption §18.44 contains Assumption §18.22 and thus Assumption §15.05 **(os1)** and **(os2)** are satisfied (see **Remark §18.23**). Moreover, considering the set \mathbb{D}_2^{ar} of densities in $\mathbb{L}_2(\lambda_{[0,1]})$ defined in (16.04) we have $\|\mathbb{p}\|_{\mathbb{L}_\infty(\lambda_{[0,1]})} \leq 1 + \tau_{\alpha, u}$ for all $\mathbb{p} \in \mathbb{D}_2^{\text{ar}}$ due to **(os2')** which allows us to apply **Lemma §16.14**. Consequently, given in addition Assumption §15.02 all assumptions of **Proposition §18.38** are satisfied. \square

§18.46 **Reminder (Maximal local ϕ -risk).** Given Assumptions §15.02 and §18.44 we consider an OPE as in **Definition §16.04**. Here the observable noisy density coefficients $\widehat{\mathbb{p}} = \mathbb{p} + n^{-1/2} \varepsilon_\cdot$ of the density coefficients $\mathbb{p} = \mathbb{U} \mathbb{p} \in \ell_2$ take the form of a *statistical direct problem* (see **Definition §10.19**) where the stochastic processes $\varepsilon_\cdot \in \mathcal{B}_{[0,1]}^{\otimes \mathbb{N}} \otimes 2^{\mathbb{N}}$ is given in **Definition §15.08**.

Under Assumptions §15.02 and §18.44 (and hence Assumption §16.24) in **Proposition** §16.27 an upper bound for a maximal local ϕ -risk of an OPE is shown over the set $\mathbb{D}_2^{\text{a,r}}$ given in (16.04) (**Lemma** §16.14). More precisely, as in (12.13) (**Proposition** §12.42) for all $n, m \in \mathbb{N}$ setting

$$\begin{aligned} R_n^m(\mathbf{a}, \phi) &:= \|\mathbf{a} \cdot \mathbf{1}^{m \perp}\|_\phi^2 + n^{-1} \|\mathbf{1}^m\|_\phi^2, \quad m_n^* := \arg \min \{R_n^m(\mathbf{a}, \phi) : m \in \mathbb{N}\} \\ \text{and } R_n^*(\mathbf{a}, \phi) &:= R_n^{m_n^*}(\mathbf{a}, \phi) = \min \{R_n^m(\mathbf{a}, \phi) : m \in \mathbb{N}\}. \end{aligned} \quad (18.38)$$

by **Proposition** §16.27 under Assumptions §15.02 and §18.44 the maximal local ϕ -risk of an OPE $\widehat{\mathbb{P}}^{m_n^*}$ with optimally chosen dimension m_n^* as in (18.38) satisfies

$$\mathcal{R}_n^\phi[\widehat{\mathbb{P}}^{m_n^*} | \mathbb{D}_2^{\text{a,r}}] \leq C R_n^*(\mathbf{a}, \phi)$$

with $C = (1 + r\tau_{\text{a,u}}) \vee r^2$. Moreover, due to **Proposition** §17.08 $R_n^*(\mathbf{a}, \phi)$ provides (up to a constant) also a lower bound of the maximal global ϕ -risk over the ellipsoid $\mathbb{D}_2^{\text{a,r}}$ for any estimator. Consequently, (up to a constant) $R_n^*(\mathbf{a}, \phi)$ is a minimax bound and $\widehat{\mathbb{P}}^{m_n^*}$ is minimax optimal. However, the optimal dimension m_n^* depends on $\mathbf{a} \in (\mathbb{R}_0^+)^{\mathbb{N}}$ characterising the ellipsoid $\mathbb{D}_2^{\text{a,r}}$. \square

§18.47 **Proposition (Upper bound)**. Under Assumptions §15.02 and §18.44 for $\widehat{\mathbb{P}}_{\text{pen}}^\phi \in (\mathbb{R}_0^+)^{\mathbb{N}}$ as in (18.32) the OPE $\widehat{\mathbb{P}}^{\widehat{m}} = \widehat{\mathbb{P}} \cdot \mathbf{1}^{\widehat{m}} \in \ell_2 \mathbf{1}^{\widehat{m}} \subseteq \text{dom}(\phi_{\nu_{\mathbb{N}}})$ with fully data-driven dimension $\widehat{m} \in \llbracket \mathbb{M} \rrbracket$ as in (18.33) satisfies for all $n, \mathbb{M} \in \mathbb{N}$

$$\begin{aligned} \mathcal{R}_n^\phi[\widehat{\mathbb{P}}^{\widehat{m}} | \mathbb{D}_2^{\text{a,r}}] &\leq 168(r^2 + r\tau_{\text{a,u}} + 9\tau_u^2) \\ &\quad \times \min \left\{ \|\mathbf{a} \cdot \mathbf{1}^{m \perp}\|_\phi^2 + n^{-1} \|\mathbf{1}^m\|_\phi^2 (\log K_m) (1 \vee (\log K_m) m n^{-1}) : m \in \llbracket \mathbb{M} \rrbracket \right\} \\ &\quad + 1440(1 + r\tau_{\text{a,u}} + 16\tau_u^2 n^{-1}) n^{-1}. \end{aligned} \quad (18.39)$$

§18.48 **Proof of Proposition** §18.47. is given in the lecture. \square

§18.49 **Corollary**. Under Assumptions §15.02 and §18.44 and $\phi \notin \ell_2$ given $\mathbb{M}^\phi \in \mathbb{N}$ as in (18.28) and $\widehat{\mathbb{P}}_{\text{pen}}^\phi \in (\mathbb{R}_0^+)^{\mathbb{N}}$ as in (18.32) consider a data-driven OPE $\widehat{\mathbb{P}}^{\widehat{m}} = \widehat{\mathbb{P}} \cdot \mathbf{1}^{\widehat{m}} \in \ell_2 \mathbf{1}^{\widehat{m}} \subseteq \text{dom}(\phi_{\nu_{\mathbb{N}}})$ with data-driven dimension $\widehat{m} \in \llbracket \mathbb{M}_n \rrbracket$ as in (18.36). For $n, m \in \mathbb{N}$ we set

$$\begin{aligned} R_n^m(\mathbf{a}, \phi) &:= \|\mathbf{a} \cdot \mathbf{1}^{m \perp}\|_\phi^2 \\ &\quad + (1 + (\log \|\mathbf{1}^m\|_\phi^2)_+ + \log m) (1 + ((\log \|\mathbf{1}^m\|_\phi^2)_+ + \log m) m n^{-1}) n^{-1} \|\mathbf{1}^m\|_\phi^2, \\ m^\diamond &:= \arg \min \{R_n^m(\mathbf{a}, \phi) : m \in \mathbb{N}\} \quad \text{and} \\ R_n^\diamond(\mathbf{a}, \phi) &:= R_n^{m^\diamond}(\mathbf{a}, \phi) = \min \{R_n^m(\mathbf{a}, \phi) : m \in \mathbb{N}\}. \end{aligned} \quad (18.40)$$

For each $n \in \mathbb{N}$ such that $R_n^\diamond(\mathbf{a}, \phi) \leq \phi_1^2$ we have

$$\begin{aligned} \mathcal{R}_n^\phi[\widehat{\mathbb{P}}^{\widehat{m}} | \mathbb{D}_2^{\text{a,r}}] &\leq 1512(r^2 + r\tau_{\text{a,u}} + 9\tau_u^2) R_n^\diamond(\mathbf{a}, \phi) + 1440(1 + r\tau_{\text{a,u}} + 16\tau_u^2 n^{-1}) n^{-1} \\ &\leq (1512 + 1440\phi^{-2})(r^2 + r\tau_{\text{a,u}} + 17\tau_u^2) R_n^\diamond(\mathbf{a}, \phi). \end{aligned} \quad (18.41)$$

§18.50 **Proof of Corollary** §18.49. is given in the lecture. \square

§18.51 **Comment**. The data-driven bound $R_n^\diamond(\mathbf{a}, \phi)$ compared to the minimax bound $R_n^*(\mathbf{a}, \phi)$ features a deterioration of the variance term at least by a logarithmic factor. The appearance of the logarithmic factor within the bound is a known fact in the context of local estimation (cf. Laurent et al. [2008] who consider model selection given direct Gaussian observations). Brown and Low [1996] show that it is unavoidable in the context of nonparametric Gaussian regression and hence it is widely considered as an acceptable price for adaptation. \square

§18.52 **Illustration.** We illustrate the last results considering the two specifications **(o)** and **(s)** given in Table 04 [§12] (**Illustration** §12.47). We restrict ourselves again to the case $\phi \notin \ell_2$ only.

Table 02 [§18]

Order of the minimax rate $R_n^*(\mathbf{a}, \phi)$ and the data-driven rate $R_n^\circ(\mathbf{a}, \phi)$ as $n \rightarrow \infty$

$(j \in \mathbb{N})$ $\phi = j^{v-1/2}$	$(a \in \mathbb{R}_{\setminus 0}^+)$ \mathbf{a}_j^2	(squared bias) $\ \mathbf{a} \cdot \mathbf{1}^{m \cdot} \ ^2_\phi$	(variance) $\ \mathbf{1}^m \ ^2_\phi$	M^ϕ	m°	$R_n^*(\mathbf{a}, \phi)$	$R_n^\circ(\mathbf{a}, \phi)$	
(o)	$v \in (0, a)$	j^{-a}	$m^{-2(a-v)}$	m^{2v}	$n^{\frac{1}{2v}}$	$n^{-\frac{(a-v)}{a}}$	$n^{-\frac{(a-v)}{a}}$	
	$a \in (1/2, \infty)$				$\left(\frac{n}{\log n}\right)^{\frac{1}{2a}}$		$\left(\frac{\log n}{n}\right)^{\frac{(a-v)}{a}}$	
	$a \in (0, 1/2]$				$\left(\frac{n}{\log n}\right)^{\frac{1}{a+1/2}}$		$\left(\frac{\log n}{n}\right)^{\frac{2(a-v)}{a+1/2}}$	
	$v = 0$	j^{-a}	m^{-2a}	$\log m$	e^n		$\frac{\log n}{n}$	$\frac{\log n}{n}$
	$a \in (1/2, \infty)$				$\left(\frac{n}{(\log n)^2}\right)^{\frac{1}{2a}}$		$\frac{(\log n)^2}{n}$	$\frac{(\log n)^2}{n}$
	$a \in (0, 1/2]$				$\left(\frac{n^2}{(\log n)^3}\right)^{\frac{1}{2a+1}}$		$\left(\frac{(\log n)^3}{n^2}\right)^{\frac{a}{a+1/2}}$	$\left(\frac{(\log n)^3}{n^2}\right)^{\frac{a}{a+1/2}}$
(s)	$v \in \mathbb{R}_{\setminus 0}^+$	$e^{-j^{2a}}$	$m^{2(v-a)+} e^{-m^{2a}}$	m^{2v}	$n^{\frac{1}{2v}}$	$(\log n)^{\frac{1}{2a}}$	$\frac{(\log n)^{\frac{v}{a}}}{n}$	$\frac{(\log n)^{\frac{v}{a}} (\log \log n)}{n}$
	$v = 0$	$e^{-j^{2a}}$	$e^{-m^{2a}}$	$\log m$	e^n	$(\log n)^{\frac{1}{2a}}$	$\frac{\log \log n}{n}$	$\frac{(\log \log n)^2}{n}$

We note that in Table 02 [§18] the order of the minimax rate $R_n^*(\mathbf{a}, \phi)$ and the data-driven rate $R_n^\circ(\mathbf{a}, \phi)$ is depict for $v \geq 0$ only. In case $v < 0$ we have $\phi \in \ell_2$ and thus **Corollary** §18.49 is not applicable. Moreover, in case **(s)** for $a \in \mathbb{R}_{\setminus 0}^+$ and **(o)** for $a \in (1/2, \infty)$ the rate $R_n^\circ(\mathbf{a}, \phi)$ features only an additional logarithmic factor compared with the minimax rate $R_n^*(\mathbf{a}, \phi)$. □

Chapter 5

Nonparametric regression

This chapter presents nonparametric regression with uniform design along the lines of the textbooks by Tsybakov [2009] and Comte [2015] where far more details, examples and further discussions can be found.

Overview

§19	Noisy regression coefficients	95
§20	Projection regression estimator	98
	§20 01 Global and maximal global \mathfrak{v} -risk	98
	§20 02 Local and maximal local ϕ -risk	101
§21	Minimax optimal regression	103
	§21 01 Maximal local ϕ -risk	103
	§21 02 Maximal global \mathfrak{v} -risk	105
§22	Data-driven regression	107
	§22 01 Data-driven global estimation by model selection	107
	§22 02 Data-driven local estimation by Goldenshluger and Lepskij's method	112

§19 Noisy regression coefficients

§19.01 **Notation (Reminder).** Consider the measure space $([0, 1], \mathcal{B}_{[0,1]}, \lambda_{[0,1]})$ where $\lambda_{[0,1]}$ denotes the restriction of the Lebesgue measure to the Borel- σ -algebra $\mathcal{B}_{[0,1]}$ over $[0, 1]$, and the Hilbert space $\mathbb{L}_2(\lambda_{[0,1]}) := \mathbb{L}_2([0, 1], \mathcal{B}_{[0,1]}, \lambda_{[0,1]})$ of square Lebesgue-integrable functions. Let (X, Y) be a $[0, 1] \times \mathbb{R}$ -valued random vector. We denote by $\mathbb{P}^X \in \mathcal{W}(\mathcal{B}_{[0,1]})$ the marginal distribution of X , by $\mathbb{P}^{Y|X}$ a regular conditional distribution of Y given X , and by $\mathbb{P}^{X,Y} = \mathbb{P}^X \odot \mathbb{P}^{Y|X} \in \mathcal{W}(\mathcal{B}_{[0,1]} \otimes \mathcal{B})$ the joint distribution of (X, Y) . We tactically identify X and Y with the coordinate map $\Pi_{[0,1]}$ and $\Pi_{\mathbb{R}}$, respectively, and thus (X, Y) with the identity $\text{id}_{[0,1] \times \mathbb{R}}$ such that $\mathbb{P} = \mathbb{P}^{X,Y} \in \mathcal{W}(\mathcal{B}_{[0,1]} \otimes \mathcal{B})$. If in addition $Y \in \mathcal{L}_1(\mathbb{P}) = \mathcal{L}_1([0, 1] \times \mathbb{R}, \mathcal{B}_{[0,1]} \otimes \mathcal{B}, \mathbb{P})$ then $\mathbb{P}^{Y|X}(\text{id}_{\mathbb{R}}) = \mathbb{P}(Y|X) =: f \in \mathcal{B}_{[0,1]}$ is unique up to \mathbb{P}^X -a.s. equality. Moreover, we have $f \in \mathcal{L}_1(\mathbb{P}^X) = \mathcal{L}_1([0, 1], \mathcal{B}_{[0,1]}, \mathbb{P}^X)$ and the error term $\xi := Y - f(X)$ satisfies $\xi \in \mathcal{L}_1(\mathbb{P})$ with $\mathbb{P}(\xi) = 0$. Let us denote in this situation by $\mathbb{P}_f^{Y|X}$ and $\mathbb{P}_f := \mathbb{P}^X \odot \mathbb{P}_f^{Y|X} \in \mathcal{W}(\mathcal{B}_{[0,1]} \otimes \mathcal{B})$, respectively, a regular conditional distribution of Y given X and the joint distribution of (X, Y) . Keep however in mind, that even if $f \in \mathcal{L}_1(\mathbb{P}^X)$ is fixed the conditional distribution $\mathbb{P}_f^{Y|X}$ is still not fully specified. In what follows we assume that the error term ξ has in addition a finite second moment and its distribution does not depend on the regression function, that is $\xi \sim \mathbb{P}^\xi \in \mathbb{P}_{\{0\} \times \mathbb{R}_0^+}$ where $\mathbb{P}_{\{0\} \times \mathbb{R}_0^+} \subseteq \mathcal{W}(\mathcal{B})$ is the subset of all probability distributions over $(\mathbb{R}, \mathcal{B})$ with finite second moment and mean zero. For $a \in \mathbb{R}$ denote by $\mathbb{P}_a^\xi \in \mathcal{W}(\mathcal{B})$ the distribution of $\xi + a$. If ξ and X are *independent*, which is assumed throughout this chapter, then there exists a \mathbb{P}^X -null set $\mathcal{N} \in \mathcal{B}_{[0,1]}$ such that $\mathbb{P}_f^{Y|X=x}(B) = \mathbb{P}_{f(x)}^\xi(B)$ for all $B \in \mathcal{B}$ and $x \in \mathcal{N}^c$ (Witting [1985], Satz 129, p.130). In other words, $(x, B) \mapsto \mathbb{P}_{f(x)}^\xi(B)$ is a version of the conditional distributions of Y given X . Evidently, if for each $B \in \mathcal{B}$ the map $\mathbb{P}_a^\xi(B) : \mathbb{R} \rightarrow [0, 1]$ with $a \mapsto \mathbb{P}_a^\xi(B)$ is Borel-measurable, $\mathbb{P}_a^\xi(B) \in \mathcal{B}$ for short, then $\mathbb{P}_a^\xi : \mathbb{R} \times \mathcal{B} \rightarrow [0, 1]$ with $(a, B) \mapsto \mathbb{P}_a^\xi(B)$ is a Markov kernel from $(\mathbb{R}, \mathcal{B})$ to $(\mathbb{R}, \mathcal{B})$. In this situation, for any $f \in \mathcal{B}_{[0,1]}$

the map $\mathbb{P}_{f(X)}^\xi : [0, 1] \times \mathcal{B} \rightarrow [0, 1]$ with $(x, B) \mapsto \mathbb{P}_{f(x)}^\xi(B)$ is a Markov kernel from $([0, 1], \mathcal{B}_{[0,1]})$ to $(\mathbb{R}, \mathcal{B})$, and hence it is a regular version of the conditional distribution of Y given X , in symbols $\mathbb{P}_f^{Y|X} = \mathbb{P}_{f(X)}^\xi$. Consequently, we have $\mathbb{P}_f = \mathbb{P}^X \odot \mathbb{P}_{f(X)}^\xi \in \mathcal{W}(\mathcal{B}_{[0,1]} \otimes \mathcal{B})$. We assume in what follows that $f \in \mathbb{F}_2 \subseteq \mathbb{L}_2(\mathbb{P}^X)$ identifying again equivalence classes and their representatives. \square

§19.02 **Assumption.** The $[0, 1] \times \mathbb{R}$ -valued random vector $(X, Y) \sim \mathbb{P}_f = \mathbb{P}^X \odot \mathbb{P}_f^{Y|X} \in \mathcal{W}(\mathcal{B}_{[0,1]} \otimes \mathcal{B})$ satisfies $Y \in \mathcal{L}_1(\mathbb{P}_f)$ and $\mathbb{P}_f(Y|X) = f \mathbb{P}^X$ -a.s. with regression function $f \in \mathbb{F}_2 \subseteq \mathbb{L}_2(\lambda_{[0,1]})$.

(NR1) The error term $\xi = Y - f(X) \sim \mathbb{P}^\xi \in \mathbb{P}_{\{0\} \times \mathbb{R}_+^*}$ has a finite second moment, mean zero and its distribution does not depend on the regression function f . We set $\sigma_\xi^2 := \mathbb{P}^\xi(\text{id}_{\mathbb{R}}^2) = \mathbb{P}_f(\xi^2)$.

(NR2) The error term ξ and the explanatory variable X are *independent*.

(NR3) The map $\mathbb{P}^\xi : \mathbb{R} \times \mathcal{B} \rightarrow [0, 1]$ with $(a, B) \mapsto \mathbb{P}_a^\xi(B)$ is a *Markov kernel* from $(\mathbb{R}, \mathcal{B})$ to $(\mathbb{R}, \mathcal{B})$. Consequently, under (NR2) the Markov kernel $\mathbb{P}_{f(X)}^\xi$ is a regular version of the conditional distribution of Y given X , i.e. $\mathbb{P}_f^{Y|X} = \mathbb{P}_{f(X)}^\xi$.

(NR4) The regressor X is uniformly distributed on the interval $[0, 1]$, i.e. $X \sim \mathcal{U}_{[0,1]}$, and thus $\mathbb{P}^X = \mathcal{U}_{[0,1]} = \lambda_{[0,1]}$. Denote by $\mathcal{U}_f := \mathcal{U}_{[0,1]} \odot \mathbb{P}_f^{Y|X}$ the joint distribution of (X, Y) .

Under (NR1)-(NR4) given $f \in \mathbb{F}_2$ and $\mathbb{P}^\xi \in \mathbb{P}_{\{0\} \times \mathbb{R}_+^*}$ the joint distribution $\mathcal{U}_f = \mathcal{U}_{[0,1]} \odot \mathbb{P}_{f(X)}^\xi$ of (X, Y) is fully specified and we set $\mathcal{U}_{\mathbb{F}_2 \times \mathbb{P}_{\{0\} \times \mathbb{R}_+^*}} := (\mathcal{U}_f)_{f \in \mathbb{F}_2, \mathbb{P}^\xi \in \mathbb{P}_{\{0\} \times \mathbb{R}_+^*}}$. We consider the statistical product experiment $(([0, 1] \times \mathbb{R})^n, (\mathcal{B}_{[0,1]} \otimes \mathcal{B})^{\otimes n}, \mathcal{U}_{\mathbb{F}_2 \times \mathbb{P}_{\{0\} \times \mathbb{R}_+^*}}^{\otimes n} := (\mathcal{U}_f^{\otimes n})_{f \in \mathbb{F}_2, \mathbb{P}^\xi \in \mathbb{P}_{\{0\} \times \mathbb{R}_+^*}})$ of size $n \in \mathbb{N}$ and for $f \in \mathbb{F}_2$ and $\mathbb{P}^\xi \in \mathbb{P}_{\{0\} \times \mathbb{R}_+^*}$ we denote by $((X_i, Y_i))_{i \in [n]} \sim \mathcal{U}_f^{\otimes n}$ an iid. sample of $(X, Y) \sim \mathcal{U}_f = \mathcal{U}_{[0,1]} \odot \mathbb{P}_{f(X)}^\xi$. \square

§19.03 **Notation (Reminder).** Consider an *orthonormal system* $(\mathbf{u}_j)_{j \in \mathbb{N}}$ in $\mathbb{L}_2(\lambda_{[0,1]})$. Then $\mathbf{U} : \mathbb{L}_2(\lambda_{[0,1]}) \rightarrow \ell_2$ with $h \mapsto \mathbf{U}h := \mathbf{h}_\bullet = (h_j := \langle h, \mathbf{u}_j \rangle_{\mathbb{L}_2(\lambda_{[0,1]})})_{j \in \mathbb{N}}$ is a surjective partial isometry $\mathbf{U} \in \mathbb{L}(\mathbb{L}_2(\lambda_{[0,1]}), \ell_2)$. Its adjoint operator $\mathbf{U}^* \in \mathbb{L}(\ell_2, \mathbb{L}_2(\lambda_{[0,1]}))$ satisfies $\mathbf{U}^* \mathbf{a}_\bullet = \sum_{j \in \mathbb{N}} a_j \mathbf{u}_j =: \nu_{\mathbb{N}}(\mathbf{a}_\bullet, \mathbf{u}_\bullet)$ for all $\mathbf{a}_\bullet \in \ell_2$. We call $\mathbf{h}_\bullet = (h_j)_{j \in \mathbb{N}}$ (*generalised*) *Fourier coefficients* and \mathbf{U} (*generalised*) *Fourier series transform*. \square

§19.04 **Remark.** Let $\mathbf{U} \in \mathbb{L}(\mathbb{L}_2(\lambda_{[0,1]}), \ell_2)$ be a generalised Fourier series transform as in Notation §19.03. For $f, h \in \mathcal{L}_2(\mathbb{P}^X) \subseteq \mathcal{L}_1(\mathbb{P}^X)$ we have $fh \in \mathcal{L}_1(\mathbb{P}^X)$ and thus $\mathbb{P}^X(fh) \in \mathbb{R}$. Keeping in mind that X and Y equals the coordinate map $\Pi_{[0,1]}$ and $\Pi_{\mathbb{R}}$, respectively, due to Assumption §19.02 (NR1), i.e., $\xi \in \mathcal{L}_2(\mathbb{P}_f)$, hence $\xi h(X) \in \mathcal{L}_1(\mathbb{P}_f)$, and (NR2) we have $\mathbb{P}_f(\xi h(X)) = \mathbb{P}_f(\xi) \mathbb{P}^X(h) = 0$. Consequently, we obtain $Yh(X) = (f(X) + \xi)h(X) \in \mathcal{L}_1(\mathbb{P}_f)$ and $\mathbb{P}_f(Yh(X)) = \mathbb{P}^X(fh) \in \mathbb{R}$. Moreover, if in addition $f \in \mathcal{L}_\infty(\mathbb{P}^X)$ then we have also $fh \in \mathcal{L}_2(\mathbb{P}^X)$ which together with $\mathbb{P}_f(\xi^2 h^2(X)) = \mathbb{P}_f(\xi^2) \mathbb{P}^X(h^2) = \sigma_\xi^2 \mathbb{P}^X(h^2) \in \mathbb{R}^+$ implies $Yh(X) \in \mathcal{L}_2(\mathbb{P}_f)$. Since $\lambda_{[0,1]} = \mathbb{P}^X$ and $\mathcal{U}_f = \mathbb{P}_f$ under Assumption §19.02 (NR4) for all $f, h \in \mathbb{L}_2(\lambda_{[0,1]})$ it follows immediately $\mathcal{U}_f(Yh(X)) = \lambda_{[0,1]}(fh) = \langle f, h \rangle_{\mathbb{L}_2(\lambda_{[0,1]})}$ identifying again equivalence classes and their representatives. Evidently, we have $\mathbf{u}_j \in \mathbb{L}_2(\lambda_{[0,1]})$ for all $j \in \mathbb{N}$ and the (generalised) Fourier coefficients $\mathbf{f}_\bullet = (f_j)_{j \in \mathbb{N}} = \mathbf{U}f \in \ell_2$ of $f \in \mathbb{L}_2(\lambda_{[0,1]})$ fulfil $f_j = \langle f, \mathbf{u}_j \rangle_{\mathbb{L}_2(\lambda_{[0,1]})} = \mathcal{U}_f(Y \mathbf{u}_j(X))$ for all $j \in \mathbb{N}$. \square

§19.05 **Assumption.** The *orthonormal system* $(\mathbf{u}_j)_{j \in \mathbb{N}}$ in $\mathbb{L}_2(\lambda_{[0,1]})$, and its associated *generalised Fourier series transform* $\mathbf{U} \in \mathbb{L}(\mathbb{L}_2(\lambda_{[0,1]}), \ell_2)$ with $h \mapsto \mathbf{U}h := \mathbf{h}_\bullet = (h_j := \langle h, \mathbf{u}_j \rangle_{\mathbb{L}_2(\lambda_{[0,1]})})_{j \in \mathbb{N}}$, is fixed and known in advance. \square

§19.06 **Remark.** Under Assumptions §19.02 and §19.05 we impose in the sequel that $f \in \mathbb{F}_2 \subseteq \mathbb{L}_\infty(\lambda_{[0,1]})$ which in turn for all $j \in \mathbb{N}$ implies $Y \mathbf{u}_j(X) \in \mathcal{L}_2(\mathcal{U}_f)$ with

$$\begin{aligned} \mathcal{U}_f(Y^2 \mathbf{u}_j^2(X)) &= \mathcal{U}_f(\xi^2 \mathbf{u}_j^2(X)) + \mathcal{U}_f(f^2(X) \mathbf{u}_j^2(X)) = \mathcal{U}_f(\xi^2) \mathbb{P}^X(\mathbf{u}_j^2) + \mathbb{P}^X(f^2 \mathbf{u}_j^2) \\ &= \sigma_\xi^2 + \lambda_{[0,1]}(f^2 \mathbf{u}_j^2) \leq \sigma_\xi^2 + \|f\|_{\mathbb{L}_\infty(\lambda_{[0,1]})}^2 \in \mathbb{R}^+. \end{aligned}$$

Alternatively, if $u_j \in \mathbb{L}_\infty(\lambda_{[0,1]})$ then for all $f \in \mathbb{L}_2(\lambda_{[0,1]})$ it follows that $Y u_j(X) \in \mathcal{L}_2(u_j)$ with $\mathcal{U}_f(Y^2 u_j^2(X)) \leq \sigma_\xi^2 + \|u_j^2\|_{\mathbb{L}_\infty(\lambda_{[0,1]})} \|f\|_{\mathbb{L}_2(\lambda_{[0,1]})}^2 \in \mathbb{R}^+$. \square

§19.07 **Notation.** Setting $\mathcal{Z} := [0, 1] \times \mathbb{R}$, $\mathcal{Z} := \mathcal{B}_{[0,1]} \otimes \mathcal{B}$ and $\psi_j(X, Y) := Y u_j(X) \in \mathcal{Z}$ for each $j \in \mathbb{N}$ under the Assumptions §19.02 and §19.05 the stochastic process $\psi_\bullet = (\psi_j)_{j \in \mathbb{N}} \in \mathcal{Z} \otimes 2^{\mathbb{N}}$ satisfies $\psi_j \in \mathcal{L}_1(u_j)$ for each $j \in \mathbb{N}$. Similar to an Empirical mean model §10.07 we define $\hat{f}_\bullet := \hat{\mathbb{P}}_n(\psi_\bullet) \in \mathcal{Z}^{\otimes \mathbb{N}} \otimes 2^{\mathbb{N}}$ with $z^n = ((x_i, y_i))_{i \in [n]} \mapsto \hat{f}_\bullet(z^n) = (\hat{\mathbb{P}}_n(\psi_j))(z^n) = n^{-1} \sum_{i \in [n]} \psi_j(x_i, y_i) = n^{-1} \sum_{i \in [n]} y_i u_j(x_i)$ for each $j \in \mathbb{N}$. For $f \in \mathbb{L}_2(\lambda_{[0,1]})$ by construction $f_\bullet = (f_j = \mathcal{U}_f(\psi_j))_{j \in \mathbb{N}} \in 2^{\mathbb{N}}$ is the ℓ_2 -mean of \hat{f}_\bullet . Consequently, $\varepsilon_\bullet := n^{1/2}(\hat{\mathbb{P}}_n - \mathcal{U}_f)(\psi_\bullet) = (\varepsilon_j = n^{1/2}(\hat{\mathbb{P}}_n(\psi_j) - f_j))_{j \in \mathbb{N}} \in \mathcal{Z}^{\otimes \mathbb{N}} \otimes 2^{\mathbb{N}}$ is centred, i.e. $\varepsilon_j \in \mathbb{L}_1(u_j^{\otimes n})$ with $\mathcal{U}_f^{\otimes n}(\varepsilon_j) = 0$. Evidently, $\hat{f}_\bullet = f_\bullet + n^{-1/2} \varepsilon_\bullet$ is a noisy version of f_\bullet (see Definition §10.19). Moreover, if $f \in \mathbb{L}_\infty(\lambda_{[0,1]})$ then \hat{f}_\bullet admits a covariance function $\text{cov}_\bullet^f \in \mathbb{R}^{\mathbb{N}^2}$ given for $j, j_o \in \mathbb{N}$ by

$$\begin{aligned} n \text{Cov}(\hat{f}_j, \hat{f}_{j_o}) &= \text{Cov}(\varepsilon_j, \varepsilon_{j_o}) = \mathcal{U}_f^{\otimes n}(\varepsilon_j \varepsilon_{j_o}) = \mathcal{U}_f(\psi_j \psi_{j_o}) - \mathcal{U}_f(\psi_j) \mathcal{U}_f(\psi_{j_o}) \\ &= \mathcal{U}_f(Y^2 u_j(X) u_{j_o}(X)) - f_j f_{j_o} =: \text{cov}_{j, j_o}^f. \end{aligned}$$

Consequently, we have $\varepsilon_\bullet \sim P_{(0, \text{cov}_\bullet^f)}$ and $\hat{f}_\bullet = f_\bullet + n^{-1/2} \varepsilon_\bullet \sim P_{(f_\bullet, n^{-1} \text{cov}_\bullet^f)}$ (see Definition §10.19). \square

§19.08 **Noisy regression coefficients.** Under Assumptions §19.02 and §19.05 the stochastic process $\varepsilon_\bullet = n^{1/2}(\hat{\mathbb{P}}_n - \mathcal{U}_f)(\psi_\bullet)$ satisfies Assumption §10.04, i.e. $\varepsilon_\bullet \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B})^{\otimes \mathbb{N}} \otimes 2^{\mathbb{N}}$, and ε_\bullet has mean zero under $\mathcal{U}_f^{\otimes n}$. The stochastic process $\hat{f}_\bullet = f_\bullet + n^{-1/2} \varepsilon_\bullet$ with ℓ_2 -mean f_\bullet is called a *noisy version* of the regression coefficients $f_\bullet = \mathcal{U} f \in \ell_2$, or *noisy regression coefficients* for short. Moreover, if $f \in \mathbb{L}_\infty(\lambda_{[0,1]})$ then ε_\bullet admits under $\mathcal{U}_f^{\otimes n}$ a covariance function $\text{cov}_\bullet^f \in \mathbb{R}^{\mathbb{N}^2}$ given for $j, j_o \in \mathbb{N}$ by $\text{cov}_{j, j_o}^f = \mathcal{U}_f(Y^2 u_j(X) u_{j_o}(X)) - f_j f_{j_o}$. We eventually write $\varepsilon_\bullet \sim P_{(0, \text{cov}_\bullet^f)}$ and $\hat{f}_\bullet \sim P_{(f_\bullet, n^{-1} \text{cov}_\bullet^f)}$. If in addition ε_\bullet admits a covariance operator $\Gamma_f \in \mathbb{L}(\ell_2)$ then we write $\varepsilon_\bullet \sim P_{(0, \Gamma)}$ and $\hat{f}_\bullet \sim P_{(f_\bullet, n^{-1} \Gamma)}$ for short. \square

§19.09 **Remark.** The centred stochastic process $\varepsilon_\bullet := (\varepsilon_j)_{j \in \mathbb{N}}$ of error terms in Definition §19.08 is in general not a white noise process. \square

§19.10 **Lemma.** Under Assumptions §19.02 and §19.05 consider $\varepsilon_\bullet \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B})^{\otimes \mathbb{N}} \otimes 2^{\mathbb{N}}$ as in Definition §19.08.

(i) If $f \in \mathbb{L}_\infty(\lambda_{[0,1]})$ then under $\mathcal{U}_f^{\otimes n}$, $\varepsilon_\bullet \sim P_{(0, \text{cov}_\bullet^f)}$ admits a covariance operator $\Gamma_f \in \mathbb{L}(\ell_2)$ given by

$$a_\bullet \mapsto \Gamma_f a_\bullet = (\nu_{\mathbb{N}}(\text{cov}_{j, j_o}^f a_\bullet))_{j, j_o \in \mathbb{N}}$$

where $\|\Gamma_f\|_{\mathbb{L}(\ell_2)} \leq \sigma_\xi^2 + \|f\|_{\mathbb{L}_\infty(\lambda_{[0,1]})}^2$.

(ii) If $f \in \mathbb{L}_\infty(\lambda_{[0,1]})$ and $\sigma_\xi^2 \in \mathbb{R}_0^+$ then $\Gamma_f \in \mathbb{L}(\ell_2)$ is invertible with inverse $\Gamma_f^{-1} \in \mathbb{L}(\ell_2)$ where $\|\Gamma_f^{-1}\|_{\mathbb{L}(\ell_2)} \leq \sigma_\xi^{-2}$.

Consequently, if $\nu_f := \max(\sigma_\xi^{-2}, \sigma_\xi^2 + \|f\|_{\mathbb{L}_\infty(\lambda_{[0,1]})}^2) \in \mathbb{R}_0^+$ then for all $a_\bullet \in \ell_2$ we have

$$\nu_f^{-1} \|a_\bullet\|_{\ell_2}^2 \leq \|a_\bullet\|_{\Gamma_f}^2 = \langle \Gamma_f a_\bullet, a_\bullet \rangle_{\ell_2} \leq \nu_f \|a_\bullet\|_{\ell_2}^2.$$

§19.11 **Proof of Lemma §19.10.** is given in the lecture. \square

§19.12 **Reminder.** Consider the orthonormal basis $(\mathbf{1}_\bullet^{[j]})_{j \in \mathbb{N}}$ in ℓ_2 (compare Remark §15.12). If $f \in \mathbb{L}_\infty(\lambda_{[0,1]})$ from Lemma §19.10 (i) for each $j \in \mathbb{N}$ we obtain

$$\mathcal{U}_f^{\otimes n}(\varepsilon_j^2) = \mathcal{U}_f^{\otimes n}(|\nu_{\mathbb{N}}(\mathbf{1}_\bullet^{[j]} \varepsilon_\bullet)|^2) = \langle \Gamma_f \mathbf{1}_\bullet^{[j]}, \mathbf{1}_\bullet^{[j]} \rangle_{\ell_2} \leq (\sigma_\xi^2 + \|f\|_{\mathbb{L}_\infty(\lambda_{[0,1]})}^2) \|\mathbf{1}_\bullet^{[j]}\|_{\ell_2}^2 = \sigma_\xi^2 + \|f\|_{\mathbb{L}_\infty(\lambda_{[0,1]})}^2$$

Keeping the last identities in mind if $v_f := \max(\sigma_\xi^{-2}, \sigma_\xi^2 + \|f\|_{\mathbb{L}_\infty(\lambda_{0,1})}^2) \in \mathbb{R}_{>0}^+$ then for all $j \in \mathbb{N}$ we have $v_f^{-1} \leq \mathcal{U}_f^{\otimes n}(\varepsilon_j^2) \leq v_f$ due to **Lemma** §19.10. \square

§20 Projection regression estimator

§20.01 **Notation (Reminder)**. Consider the measure space $(\mathbb{N}, 2^{\mathbb{N}}, \nu_{\mathbb{N}})$ as in **Notation** §10.11. For $w \in \mathbb{R}^{\mathbb{N}}$ define the multiplication map $M_w : \mathbb{R}^{\mathbb{N}} \rightarrow \mathbb{R}^{\mathbb{N}}$ with $a \mapsto M_w a := w \cdot a$. Note that each $w \in \mathbb{R}^{\mathbb{N}}$ is $2^{\mathbb{N}}$ - \mathcal{B} -measurable. We denote by $\mathbb{M}_{\mathbb{R}^{\mathbb{N}}}$ the set of all multiplication maps defined on $\mathbb{R}^{\mathbb{N}}$. If in addition $w \in \ell_\infty = \mathbb{L}_\infty(\mathbb{N}, 2^{\mathbb{N}}, \nu_{\mathbb{N}})$ then we have also $M_w : \ell_2 \rightarrow \ell_2$. We set $\mathbb{L}(\ell_2) = \{M_w \in \mathbb{M}_{\mathbb{R}^{\mathbb{N}}} : w \in \ell_\infty\} \subseteq \mathbb{L}(\ell_2)$ noting that $\|M_w\|_{\mathbb{L}(\ell_2)} = \sup\{\|w \cdot a\|_{\ell_2} : \|a\|_{\ell_2} \leq 1\} \leq \|w\|_{\ell_\infty}$ for each $M_w \in \mathbb{L}(\ell_2)$. \square

§20.02 **Reminder**. If $w \in \ell_\infty$ then $M_w \in \mathbb{L}(\ell_2)$, and $M_{w^\dagger} : \ell_2 \supseteq \text{dom}(M_{w^\dagger}) \rightarrow \ell_2$. Moreover, we have $\text{dom}(M_w) = \ell_2$, $\text{ran}(M_w) = \ell_2 w$, and $\ker(M_w) = \ell_2 \mathbb{1}^{\mathcal{N}_w}$ with $\mathcal{N}_w = \{j \in \mathbb{N} : w_j = 0\} \in 2^{\mathbb{N}}$ (see **Property** §11.03), and $\text{dom}(M_{w^\dagger}) = \ell_2 w \oplus \ell_2 \mathbb{1}^{\mathcal{N}_w}$ (see **Property** §11.05). Consequently, if in addition $\nu_{\mathbb{N}}(\mathcal{N}_w) = 0$ or in equal $w \in (\mathbb{R}_{>0})^{\mathbb{N}}$, then $w^\dagger = w^{-1} \in (\mathbb{R}_{>0})^{\mathbb{N}}$, hence $w^{2\dagger} = w^{-2} \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$, and $\ell_2^w = \text{dom}(M_{w^\dagger}) = \ell_2 w = \mathbb{L}_2(w^{-2} \nu_{\mathbb{N}}) =: \ell_2(w^{-2})$. For each $m \in \mathbb{N}$ we write $\mathbb{1}^m = (\mathbb{1}_j)_{j \in \mathbb{N}} := \mathbb{1}^{\llbracket m \rrbracket}$ and $\mathbb{1}^{m\perp} := \mathbb{1} - \mathbb{1}^m$ with $\llbracket m \rrbracket := [-m, m] \cap \mathbb{N}$. Consequently, $M_{\mathbb{1}^m} \in \mathbb{L}(\ell_2)$ and $M_{\mathbb{1}^{m\perp}} \in \mathbb{L}(\ell_2)$ is the *orthogonal projection* onto the linear subspace $\ell_2 \mathbb{1}^m \subseteq \ell_2$ and its orthogonal complement $\ell_2 \mathbb{1}^{m\perp} = (\ell_2 \mathbb{1}^m)^\perp \subseteq \ell_2$, respectively, that is $\ell_2 = \ell_2 \mathbb{1}^m \oplus \ell_2 \mathbb{1}^{m\perp}$ (see **Property** §11.07). Finally, given $h \cdot = U h \in \ell_2$ for $h \in \mathbb{L}_2(\lambda_{0,1})$ we consider the orthogonal projections $h \cdot^m = h \cdot \mathbb{1}^m \in \ell_2 \mathbb{1}^m$ and $h \cdot^m := U^* h \cdot^m \in \mathbb{L}_2(\lambda_{0,1})$ (**Definition** §11.08). \square

§20.03 **Notation (Reminder)**. Consider the stochastic processes $\varepsilon \cdot = n^{1/2}(\widehat{\mathbb{P}}_n - \mathcal{U}_f)(\psi) \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B})^{\otimes n} \otimes 2^{\mathbb{N}}$ given in **Definition** §19.08. The observable noisy version $\widehat{f} \cdot = f \cdot + n^{-1/2} \varepsilon \cdot$ of the regression coefficients $f \cdot = U f \in \ell_2$ take the form of a *statistical direct problem* (see **Definition** §10.19). Under Assumptions §19.02 and §19.05 $\varepsilon \cdot$ is centred and admits a covariance function $\text{cov}_{\varepsilon \cdot}^f \in \mathbb{R}^{\mathbb{N}^2}$ given in **Definition** §19.08, i.e. $\varepsilon \cdot \sim P_{(0, \text{cov}_{\varepsilon \cdot}^f)}$ and $\widehat{f} \cdot \sim P_{(f \cdot, n^{-1} \text{cov}_{\varepsilon \cdot}^f)}$. If in addition $f \in \mathbb{L}_\infty(\lambda_{0,1})$ then $\varepsilon \cdot$ admits a covariance operator $\Gamma_f \in \mathbb{L}(\ell_2)$ given in **Lemma** §19.10, i.e. $\varepsilon \cdot \sim P_{(0, \Gamma_f)}$ and $\widehat{f} \cdot \sim P_{(f \cdot, n^{-1} \Gamma_f)}$. \square

§20.04 **Definition**. Given a noisy version $\widehat{f} \cdot = f \cdot + n^{-1/2} \varepsilon \cdot$ of the regression coefficients $f \cdot = U f \in \ell_2$ for each $m \in \mathbb{N}$ we call $\widehat{f} \cdot^m := \widehat{f} \cdot \mathbb{1}^m$ *orthogonal projection estimator (OPE)* of $f \cdot$. \square

§20.05 **Remark**. If $f \cdot = U^* f \cdot$ (for example $(u_j)_{j \in \mathbb{N}}$ is an orthonormal basis of $\mathbb{L}_2(\lambda_{0,1})$), then we have

$$\|U^* \widehat{f} \cdot^m - f \cdot\|_{\mathbb{L}_2(\lambda_{0,1})}^2 = \|\widehat{f} \cdot^m - f \cdot\|_{\ell_2}^2.$$

In this situation all results for the OPE $\widehat{f} \cdot^m$ of the regression coefficients immediately transfer onto the *orthogonal projection regression estimator* $\widehat{f} \cdot^m := U^* \widehat{f} \cdot^m$ of the regression function $f \cdot$. \square

§20|01 Global and maximal global v-risk

We measure first the accuracy of the OPE $\widehat{f} \cdot^m = \widehat{f} \cdot \mathbb{1}^m$ of $f \cdot^m = f \cdot \mathbb{1}^m \in \ell_2 \mathbb{1}^m$ with $f \cdot = U f \in \ell_2$ by a global mean-v-error, i.e. v-risk.

§20.06 **Reminder**. If $v \cdot \in (\mathbb{R}_{>0})^{\mathbb{N}}$ and $f \cdot \in \ell_2(v \cdot)$ then we have $f \cdot^m = f \cdot \mathbb{1}^m \in \ell_2(v \cdot)$ too and $\|f \cdot^m - f \cdot\|_v^2 = o(1)$ as $m \rightarrow \infty$ (**Property** §11.09). Moreover, $\varepsilon \cdot \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B})^{\otimes n} \otimes 2^{\mathbb{N}}$ given in **Definition** §19.08 satisfies $v \cdot \varepsilon \cdot \mathbb{1}^m \in \ell_2$ (note that $\mathbb{1}^m \in \ell_2$ and $v \cdot \mathbb{1}^m, \varepsilon \cdot \mathbb{1}^m \in \ell_\infty$) and thus also

$$n^{-1/2} v \cdot \varepsilon \cdot \mathbb{1}^m + v \cdot f \cdot^m = v \cdot \widehat{f} \cdot^m \in \ell_2. \quad (20.01)$$

Finally, under Assumptions §19.02 and §19.05 and $f \in \mathbb{L}_\infty(\lambda_{[0,1]})$ due to **Lemma** §19.10 we have $\mathcal{U}_f^{\otimes n}(\boldsymbol{\varepsilon}^2) \in \ell_\infty$, more precisely, $\|\mathcal{U}_f^{\otimes n}(\boldsymbol{\varepsilon}^2)\|_{\ell_\infty} \leq \sigma_\xi^2 + \|f\|_{\mathbb{L}_\infty(\lambda_{[0,1]})}^2$ (see **Reminder** §19.12). \square

§20|01|01 Global \mathfrak{v} -risk

§20.07 **Proposition (Upper bound)**. *Let Assumptions §19.02 and §19.05, $\mathfrak{v}_\bullet \in (\mathbb{R}_{\setminus 0})^{\mathbb{N}}$ and $f \in \ell_2(\mathfrak{v}_\bullet^2)$ be satisfied and for all $n, m \in \mathbb{N}$ set*

$$\begin{aligned} R_n^m(f, \mathfrak{v}_\bullet) &:= \|f \cdot \mathbf{1}^{m \perp}\|_{\mathfrak{v}_\bullet}^2 + n^{-1} \|\mathbf{1}^m\|_{\mathfrak{v}_\bullet}^2, \quad m_n^\circ := \arg \min \{R_n^m(f, \mathfrak{v}_\bullet) : m \in \mathbb{N}\} \\ \text{and } R_n^\circ(f, \mathfrak{v}_\bullet) &:= R_n^{m_n^\circ}(f, \mathfrak{v}_\bullet) = \min \{R_n^m(f, \mathfrak{v}_\bullet) : m \in \mathbb{N}\}. \end{aligned} \quad (20.02)$$

If $f \in \mathbb{L}_\infty(\lambda_{[0,1]})$ then we have $\mathcal{U}_f^{\otimes n}(\|\widehat{f}_\bullet^{m_n^\circ} - f\|_{\mathfrak{v}_\bullet}^2) \leq 1 \vee (\sigma_\xi^2 + \|f\|_{\mathbb{L}_\infty(\lambda_{[0,1]})}^2) R_n^\circ(f, \mathfrak{v}_\bullet)$.

§20.08 **Proof of Proposition** §20.07. is given in the lecture. \square

§20.09 **Oracle inequality**. *Under Assumptions §19.02 and §19.05 let $\mathfrak{v}_\bullet \in (\mathbb{R}_{\setminus 0})^{\mathbb{N}}$ and $f \in \ell_2(\mathfrak{v}_\bullet^2)$. If in addition $\mathfrak{v}_f := \max(\sigma_\xi^{-2}, \sigma_\xi^2 + \|f\|_{\mathbb{L}_\infty(\lambda_{[0,1]})}^2) \in \mathbb{R}_{\setminus 0}^+$ then $\mathfrak{v}_f^{-1} \leq \mathfrak{v}_f^j := \mathcal{U}_f^{\otimes n}(\boldsymbol{\varepsilon}_j^2) \leq \mathfrak{v}_f$ for all $j \in \mathbb{N}$ (see **Reminder** §19.12), and hence **Property** §12.15 implies*

$$\begin{aligned} \mathfrak{v}_f^{-1} R_n^m(f, \mathfrak{v}_\bullet) &\leq \mathcal{U}_f^{\otimes n}(\|\widehat{f}_\bullet^m - f\|_{\mathfrak{v}_\bullet}^2) = n^{-1} \mathcal{U}_\mathbb{N}(\mathfrak{v}_f^j \mathfrak{v}_\bullet^2 \mathbf{1}^m) + \|f \cdot \mathbf{1}^{m \perp}\|_{\mathfrak{v}_\bullet}^2 \\ &\leq \mathfrak{v}_f R_n^m(f, \mathfrak{v}_\bullet) \quad \text{for all } m, n \in \mathbb{N}. \end{aligned}$$

As a consequence we immediately obtain the following **oracle inequality** (see **Definition** §12.14)

$$\begin{aligned} \mathfrak{v}_f^{-1} R_n^\circ(f, \mathfrak{v}_\bullet) &\leq \inf_{m \in \mathbb{N}} \mathcal{U}_f^{\otimes n}(\|\widehat{f}_\bullet^m - f\|_{\mathfrak{v}_\bullet}^2) \leq \mathcal{U}_f^{\otimes n}(\|\widehat{f}_\bullet^{m_n^\circ} - f\|_{\mathfrak{v}_\bullet}^2) \\ &\leq \mathfrak{v}_f R_n^\circ(f, \mathfrak{v}_\bullet) \leq \mathfrak{v}_f^2 \inf_{m \in \mathbb{N}} \mathcal{U}_f^{\otimes n}(\|\widehat{f}_\bullet^m - f\|_{\mathfrak{v}_\bullet}^2), \end{aligned} \quad (20.03)$$

and, hence $R_n^\circ(f, \mathfrak{v}_\bullet)$, m_n° and the statistic $\widehat{f}_\bullet^{m_n^\circ}$, respectively, is an **oracle bound**, an **oracle dimension** and **oracle optimal** (up to the constant \mathfrak{v}_f^2). We observe that $R_n^\circ(f, \mathfrak{v}_\bullet) = o(1)$ as $n \rightarrow \infty$ (**Remark** §12.16), and thus, $R_n^\circ(f, \mathfrak{v}_\bullet)$ is an **oracle rate**. However, note that the oracle dimension $m_n^\circ = m_n^\circ(f, \mathfrak{v}_\bullet)$ depends on the unknown regression coefficients f , and thus also the oracle optimal statistic $\widehat{f}_\bullet^{m_n^\circ}$. In other words $\widehat{f}_\bullet^{m_n^\circ}$ is not a feasible estimator. \square

§20.10 **Illustration**. We illustrate the last results considering usual behaviour for the bias and variance term. We distinguish the following two cases

(p) $\mathfrak{v}_\bullet \in \ell_2$ or there is $m \in \mathbb{N}$ with $\|f^m - f\|_{\mathfrak{v}_\bullet}^2 = 0$,

(np) $\mathfrak{v}_\bullet \notin \ell_2$ and for all $m \in \mathbb{N}$ holds $\|f^m - f\|_{\mathfrak{v}_\bullet}^2 \in \mathbb{R}_{\setminus 0}^+$.

Interestingly, in case **(p)** the oracle bound is parametric, that is, $nR_n^\circ(f, \mathfrak{v}_\bullet) = O(1)$, in case **(np)** the oracle bound is nonparametric, i.e. $\lim_{n \rightarrow \infty} nR_n^\circ(f, \mathfrak{v}_\bullet) = \infty$. In case **(np)** consider the following two specifications:

Table 01 [§20]

Order of the oracle rate $R_n^\circ(f, \mathbf{v})$ as $n \rightarrow \infty$

$(j \in \mathbb{N})$	$(a \in \mathbb{R}_0^+)$	(squared bias)	(variance)	m_n°	$R_n^\circ(f, \mathbf{v})$
$\mathbf{v}_j^2 = j^{2v}$	f_j^2	$\ f \cdot \mathbf{1}_\bullet^{m \perp}\ _{\mathbf{v}}^2$	$\ \mathbf{1}_\bullet^m\ _{\mathbf{v}}^2$		
(o) $v \in (-1/2, a)$	j^{-2a-1}	$m^{-2(a-v)}$	m^{2v+1}	$n^{\frac{1}{2a+1}}$	$n^{-\frac{2(a-v)}{2a+1}}$
$v = -1/2$	j^{-2a-1}	m^{-2a-1}	$\log m$	$(\frac{n}{\log n})^{\frac{1}{2a+1}}$	$\frac{\log n}{n}$
(s) $v + 1/2 \in \mathbb{R}_0^+$	$e^{-j^{2a}}$	$m^{(1-2(a-v))_+} e^{-m^{2a}}$	m^{2v+1}	$(\log n)^{\frac{1}{2a}}$	$\frac{(\log n)^{\frac{2v+1}{2a}}}{n}$
$v = -1/2$	$e^{-j^{2a}}$	$e^{-m^{2a}}$	$\log m$	$(\log n)^{\frac{1}{2a}}$	$\frac{\log \log n}{n}$

We note that in Table 01 [§20] the order of the oracle rate $R_n^\circ(f, \mathbf{v})$ is depict for $v \geq -1/2$ only. In case $v < -1/2$ the oracle rate $R_n^\circ(f, \mathbf{v})$ is parametric. □

§20|01|02 Maximal global \mathbf{v} -risk

§20.11 **Assumption.** Consider weights $\mathbf{a}, \mathbf{v} \in (\mathbb{R}_0^+)^{\mathbb{N}}$ with $\mathbf{a} \in \ell_\infty$ and $(\mathbf{a}\mathbf{v})_\bullet := (\mathbf{a}_j \mathbf{v}_j)_{j \in \mathbb{N}} = \mathbf{a} \cdot \mathbf{v} \in \ell_\infty$. We write $(\mathbf{a}\mathbf{v})_{(m)} := \|(\mathbf{a}\mathbf{v}) \cdot \mathbf{1}_\bullet^{m \perp}\|_{\ell_\infty} \in \mathbb{R}^+$ for each $m \in \mathbb{N}$. The orthonormal system $(\mathbf{u}_j)_{j \in \mathbb{N}}$ in $\mathbb{L}_2(\lambda_{[0,1]})$ is **(os1)** complete, i.e an orthonormal basis in $\mathbb{L}_2(\lambda_{[0,1]})$ and as process $\mathbf{u}^2 = (\mathbf{u}_j^2)_{j \in \mathbb{N}}$ on $([0, 1], \mathcal{B}_{[0,1]})$ satisfies **(os2)** $\|\nu_{\mathbb{N}}(\mathbf{a}^2 \mathbf{u}^2)\|_{\mathbb{L}_\infty(\lambda_{[0,1]})} \leq \tau_{\mathbf{a}, \mathbf{u}}^2$ for $\tau_{\mathbf{a}, \mathbf{u}} \in [1, \infty)$. □

§20.12 **Reminder.** Under Assumption §20.11 we have $\ell_2^{\mathbf{a}} = \text{dom}(M_{\mathbf{a}, \cdot}) = \ell_2 \cdot \mathbf{a} \subseteq \ell_2$ and the three measures $\nu_{\mathbb{N}}, \mathbf{a}^{-2} \nu_{\mathbb{N}}$ and $\mathbf{v}^2 \nu_{\mathbb{N}}$ dominate mutually each other, i.e. they share the same null sets (see **Property** §11.05). We consider $\ell_2^{\mathbf{a}}$ endowed with $\|\cdot\|_{\mathbf{a}^{-1}} = \|M_{\mathbf{a}, \cdot}\|_{\ell_2}$ and given a constant $r \in \mathbb{R}_0^+$ the ellipsoid $\ell_2^{\mathbf{a}, r} := \{b \in \ell_2^{\mathbf{a}} : \|b\|_{\mathbf{a}^{-1}} \leq r\} \subseteq \ell_2^{\mathbf{a}}$. Since $(\mathbf{a}\mathbf{v})_\bullet \in \ell_\infty$, and hence $(\mathbf{a}\mathbf{v})_{(m)} := \|(\mathbf{a}\mathbf{v}) \cdot \mathbf{1}_\bullet^{m \perp}\|_{\ell_\infty} \in \mathbb{R}^+$ for each $m \in \mathbb{N}$ we have $\ell_2^{\mathbf{a}} \subseteq \ell_2(\mathbf{v}^2)$ (**Property** §11.15), and $\|b \cdot \mathbf{1}_\bullet^{m \perp}\|_{\mathbf{v}} \leq r (\mathbf{a}\mathbf{v})_{(m)}$ for all $b \in \ell_2^{\mathbf{a}, r}$ (**Lemma** §11.17). □

§20.13 **Lemma.** Under Assumption §20.11 set

$$\mathbb{F}_2^{\mathbf{a}, r} := \{h \in \mathbb{L}_2(\lambda_{[0,1]}) : h_\bullet = U h \in \ell_2^{\mathbf{a}, r}\}. \tag{20.04}$$

Then we have $\sup \{\|h\|_{\mathbb{L}_\infty(\lambda_{[0,1]})} : h \in \mathbb{F}_2^{\mathbf{a}, r}\} \leq r \tau_{\mathbf{a}, \mathbf{u}}$.

§20.14 **Proof** of **Lemma** §20.13. is given in the lecture. □

§20.15 **Proposition (Upper bound).** Let Assumptions §19.02 and §20.11 be satisfied. For $n \in \mathbb{N}$ considering $m_n^* \in \mathbb{N}$ and $R_n^*(\mathbf{a}, \mathbf{v}) \in \mathbb{R}^+$ as in (12.06) (**Proposition** §12.21) we have

$$\sup \{\mathcal{U}_f^{\otimes n} (\|\widehat{f}_\bullet^{m_n^*} - f_\bullet\|_{\mathbf{v}}^2) : f \in \mathbb{F}_2^{\mathbf{a}, r}\} \leq C R_n^*(\mathbf{a}, \mathbf{v}).$$

with constant $C = \sigma_\xi^2 + r^2 \tau_{\mathbf{a}, \mathbf{u}}^2 + r^2$.

§20.16 **Proof** of **Proposition** §20.15. is given in the lecture. □

§20.17 **Illustration.** The *trigonometric basis* given for $x \in [0, 1]$ by

$$\mathbf{u}_1 := \mathbf{1}_{[0,1]}, \mathbf{u}_{2k}(x) := \sqrt{2} \cos(2\pi kt), \mathbf{u}_{2k+1}(x) := \sqrt{2} \sin(2\pi kt), k \in \mathbb{N},$$

is an orthonormal basis of $\mathbb{L}_2(\lambda_{[0,1]})$, hence it satisfies Assumption §20.11 **(os1)**. Keeping in mind that $\|\mathbf{u}_j^2\|_{\mathbb{L}_2(\lambda_{[0,1]})} \leq 2$ for all $j \in \mathbb{N}$ also the Assumption §20.11 **(os3)** is satisfied for all $\mathbf{a} \in \ell_2$

because $\tau_{a,u}^2 \leq 2\|\mathbf{a}\|_{\ell_2}^2$ (see also [Illustration §16.18](#)). In [Table 02 \[§12\]](#) ([Illustration §12.26](#)) the order of the rate $R_n^*(\mathbf{a}, \mathbf{v})$ is depicted for the two cases **(o)** and **(s)**. We note that we have $\mathbf{a} \in \ell_2$ in case **(o)** for $a > 1/2$ while in case **(s)** for $a \in \mathbb{R}_V^+$. \square

§20.18 **Remark.** In [Proposition §20.15](#) an upper bound is shown under [Assumption §19.02](#) which amongst others imposes that $\mathcal{U}_f = \mathcal{U}_{[0,1]} \odot \mathbb{P}_{f(X)}^\xi$ with $\mathbb{P}^\xi \in \mathbb{P}_{\{0\} \times \mathbb{R}_V^+}$. Recall that $\mathbb{P}_{\{0\} \times \mathbb{R}_V^+} \subseteq \mathcal{W}(\mathcal{B})$ denotes the subset of all probability distributions over $(\mathbb{R}, \mathcal{B})$ with finite second moment and mean zero. For $\sigma^2 \in \mathbb{R}_V^+$ let us further introduce $\mathbb{P}_{\{0\} \times (0, \sigma^2]}^+ \subseteq \mathbb{P}_{\{0\} \times \mathbb{R}_V^+}$ containing only probability distributions with second moment bounded by σ^2 . In what follows we treat the distribution \mathbb{P}^ξ of the error term as a nuisance parameter and consider the maximal risk over both $\mathbb{F}_2^{a,r}$ and $\mathbb{P}_{\{0\} \times (0, \sigma^2]}^+$ (see [Definition §13.07](#)). \square

§20.19 **Corollary (Upper bound).** *Let Assumptions §19.02 and §20.11 be satisfied. For $n \in \mathbb{N}$ considering $m_n^* \in \mathbb{N}$ and $R_n^*(\mathbf{a}, \mathbf{v}) \in \mathbb{R}^+$ as in (12.06) ([Proposition §12.21](#)) we have*

$$\sup \left\{ (\mathcal{U}_{[0,1]} \odot \mathbb{P}_{f(X)}^\xi)^{\otimes n} (\|\widehat{f}^{m_n^*} - f\|_{\mathbb{V}}^2) : f \in \mathbb{F}_2^{a,r}, \mathbb{P}^\xi \in \mathbb{P}_{\{0\} \times (0, \sigma^2]}^+ \right\} \leq C R_n^*(\mathbf{a}, \mathbf{v}).$$

with constant $C = \sigma^2 + r^2 \tau_{a,u}^2 + r^2$.

§20.20 **Proof of Corollary §20.19.** is given in the lecture. \square

§20|02 Local and maximal local ϕ -risk

We measure secondly the accuracy of the OPE $\widehat{f}^m = \widehat{f} \mathbb{1}^m$ of $f^m = f \mathbb{1}^m \in \ell_2 \mathbb{1}^m$ with $f = Uf \in \ell_2$ by a local mean- ϕ -error, i.e. ϕ -risk.

§20.21 **Reminder.** If $\phi \in (\mathbb{R}_V^+)^{\mathbb{N}}$ and $f \in \text{dom}(\phi_{\mathbb{N}}) := \{a \in \ell_2 : \phi a \in \ell_1\}$, then we have $f^m = f \mathbb{1}^m \in \text{dom}(\phi_{\mathbb{N}})$ too and $|\phi_{\mathbb{N}}(f - f^m)| = o(1)$ as $m \rightarrow \infty$ ([Property §11.22](#)). Moreover, $\boldsymbol{\varepsilon} \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B})^{\otimes n} \otimes 2^{\mathbb{N}}$ given in [Definition §19.08](#) satisfies $\boldsymbol{\varepsilon} \mathbb{1}^m \in \text{dom}(\phi_{\mathbb{N}})$ (note that $\phi \mathbb{1}^m, \boldsymbol{\varepsilon} \mathbb{1}^m \in \ell_2$) and thus also

$$n^{-1/2} \boldsymbol{\varepsilon} \mathbb{1}^m + f^m = \widehat{f}^m \in \text{dom}(\phi_{\mathbb{N}}). \quad (20.05)$$

Finally, under [Assumptions §19.02](#) and [§19.05](#) and $f \in \mathbb{L}_\infty(\lambda_{[0,1]})$ due to [Lemma §19.10 \(i\)](#) the process $\boldsymbol{\varepsilon} \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B})^{\otimes n} \otimes 2^{\mathbb{N}}$ admits a covariance operator $\Gamma_f \in \mathbb{L}(\ell_2)$, i.e. $\boldsymbol{\varepsilon} \sim \mathbb{P}_{(a, \Gamma)}$, satisfying $\|\Gamma_f\|_{\mathbb{L}(\ell_2)} \leq \sigma_\xi^2 + \|f\|_{\mathbb{L}_\infty(\lambda_{[0,1]})}^2$. \square

§20|02|01 Local ϕ -risk

§20.22 **Proposition (Upper bound).** *Let Assumptions §19.02 and §19.05, $\phi \in (\mathbb{R}_V^+)^{\mathbb{N}}$ and $f \in \text{dom}(\phi_{\mathbb{N}})$ be satisfied and for all $n, m \in \mathbb{N}$ set*

$$\begin{aligned} R_n^m(f, \phi) &:= |\phi_{\mathbb{N}}(f \mathbb{1}^{m|\perp})|^2 + n^{-1} \|\mathbb{1}^m\|_\phi^2, \quad m_n^\circ := \arg \min \{R_n^m(f, \phi) : m \in \mathbb{N}\} \\ \text{and} \quad R_n^\circ(f, \phi) &:= R_n^{m_n^\circ}(f, \phi) := \min \{R_n^m(f, \phi) : m \in \mathbb{N}\}. \end{aligned} \quad (20.06)$$

If $f \in \mathbb{L}_\infty(\lambda_{[0,1]})$ then we have $\mathcal{U}_f^{\otimes n} (|\phi_{\mathbb{N}}(\widehat{f}^{m_n^\circ} - f)|^2) \leq 1 \vee (\sigma_\xi^2 + \|f\|_{\mathbb{L}_\infty(\lambda_{[0,1]})}^2) R_n^\circ(f, \phi)$.

§20.23 **Proof of Proposition §20.22.** is given in the lecture. \square

§20.24 **Oracle inequality.** Under Assumptions §19.02 and §19.05 let $\phi \in (\mathbb{R}_{\setminus 0})^{\mathbb{N}}$ and $f \in \text{dom}(\phi_{\mathcal{N}})$. If in addition $\mathfrak{v}_f := \max(\sigma_{\xi}^{-2}, \sigma_{\xi}^2 + \|f\|_{\mathbb{L}_{\infty}(\lambda_{0,1})}^2) \in \mathbb{R}_{\setminus 0}^+$ then $\max(\|\Gamma_f\|_{\mathbb{L}(\ell_2)}, \|\Gamma_f^{-1}\|_{\mathbb{L}(\ell_2)}) \leq \mathfrak{v}_f$ (see Lemma §19.10), and hence Property §12.36 implies

$$\begin{aligned} \mathfrak{v}_f^{-1} \mathbb{R}_n^m(f, \phi) &\leq \mathcal{U}_f^{\otimes n} (|\phi_{\mathcal{N}}(\hat{f}^m - f)|^2) = n^{-1} \|\phi \mathbb{1}^m\|_{\Gamma}^2 + |\phi_{\mathcal{N}}(f \mathbb{1}^{m\perp})|^2 \\ &\leq \mathfrak{v}_f \mathbb{R}_n^m(f, \phi) \quad \text{for all } m, n \in \mathbb{N}. \end{aligned}$$

As a consequence we immediately obtain the following *oracle inequality* (see Definition §12.34)

$$\begin{aligned} \mathfrak{v}_f^{-1} \mathbb{R}_n^{\circ}(f, \phi) &\leq \inf_{m \in \mathbb{N}} \mathcal{U}_f^{\otimes n} (|\phi_{\mathcal{N}}(\hat{f}^m - f)|^2) \leq \mathcal{U}_f^{\otimes n} (|\phi_{\mathcal{N}}(\hat{f}^{m_n^{\circ}} - f)|^2) \\ &\leq \mathfrak{v}_f \mathbb{R}_n^{\circ}(f, \phi) \leq \mathfrak{v}_f^2 \inf_{m \in \mathbb{N}} \mathcal{U}_f^{\otimes n} (|\phi_{\mathcal{N}}(\hat{f}^m - f)|^2), \quad (20.07) \end{aligned}$$

and hence, $\mathbb{R}_n^{\circ}(f, \phi)$, m_n° and the statistic $\hat{f}^{m_n^{\circ}}$, respectively, is an *oracle bound*, an *oracle dimension* and *oracle optimal* (up to the constant \mathfrak{v}_f^2). We observe that $\mathbb{R}_n^{\circ}(f, \phi) = o(1)$ as $n \rightarrow \infty$ (Remark §12.37), and thus, $\mathbb{R}_n^{\circ}(f, \phi)$ is an *oracle rate*. However, note that the oracle dimension $m_n^{\circ} = m_n^{\circ}(f, \phi)$ depends on the unknown regression coefficients f , and thus also the oracle optimal statistic $\hat{f}^{m_n^{\circ}}$. In other words $\hat{f}^{m_n^{\circ}}$ is not a feasible estimator. \square

§20.25 **Illustration.** We illustrate the last results considering usual behaviour for both the variance and the bias term. Similar to the two cases (p) and (np) in Illustration §20.10 we distinguish here the following two cases

(p) $\phi \in \ell_2$ or there is $K \in \mathbb{N}$ with $\sup\{|\phi_{\mathcal{N}}(f \mathbb{1}^{m\perp})|^2 : m \in \mathbb{N} \cap [K, \infty)\} = 0$,

(np) $\phi \notin \ell_2$ and for all $m \in \mathbb{N}$ holds $\sup\{|\phi_{\mathcal{N}}(f \mathbb{1}^{m\perp})|^2 : m \in \mathbb{N} \cap [K, \infty)\} \in \mathbb{R}_{\setminus 0}^+$.

In case (p) the oracle bound is again parametric, i.e. $n \mathbb{R}_n^{\circ}(f, \phi) = O(1)$, while in case (np) the oracle bound is nonparametric, i.e. $\lim_{n \rightarrow \infty} n \mathbb{R}_n^{\circ}(f, \phi) = \infty$. In case (np) consider the following two specifications

Table 02 [§20]

Order of the oracle rate $\mathbb{R}_n^{\circ}(f, \phi)$ as $n \rightarrow \infty$

$(j \in \mathbb{N})$	$(a \in \mathbb{R}_{\setminus 0}^+)$	(squared bias)	(variance)		
$\phi_j = j^{v-1/2}$	f_j	$ \phi_{\mathcal{N}}(f \mathbb{1}^{m\perp}) ^2$	$\ \mathbb{1}^m\ _{\phi}^2$	m_n°	$\mathbb{R}_n^{\circ}(f, \phi)$
(o)	$v \in (0, a)$	$j^{-2(a-v)}$	m^{2v}	$n^{\frac{1}{2a}}$	$n^{-\frac{(a-v)}{a}}$
	$v = 0$	$j^{-a-1/2}$	$\log m$	$(\frac{n}{\log n})^{\frac{1}{2a}}$	$\frac{\log n}{n}$
(s)	$v \in \mathbb{R}_{\setminus 0}^+$	$e^{-j^{2a}}$	m^{2v}	$(\log n)^{\frac{1}{2a}}$	$\frac{(\log n)^{\frac{v}{a}}}{n}$
	$v = 0$	$e^{-j^{2a}}$	$\log m$	$(\log n)^{\frac{1}{2a}}$	$\frac{\log \log n}{n}$

We note that in Table 02 [§20] the order of the oracle rate $\mathbb{R}_n^{\circ}(f, \phi)$ is depict for $v \geq 0$ only. For $v < 0$ the oracle rate $\mathbb{R}_n^{\circ}(f, \phi)$ is parametric. \square

§20|02|02 Maximal local ϕ -risk

§20.26 **Assumption.** Consider $\phi, \mathfrak{a} \in (\mathbb{R}_{\setminus 0})^{\mathbb{N}}$ with $\mathfrak{a} \in \ell_{\infty}$ and $(\mathfrak{a}\phi)_j := (\mathfrak{a}_j \phi_j)_{j \in \mathbb{N}} = \mathfrak{a} \cdot \phi \in \ell_2$, and hence $\|\mathfrak{a} \mathbb{1}^{m\perp}\|_{\phi} = \|(\mathfrak{a}\phi) \mathbb{1}^{m\perp}\|_{\ell_2} = o(1)$ as $m \rightarrow \infty$. The orthonormal system $(u_j)_{j \in \mathbb{N}}$ in $\mathbb{L}_2(\lambda_{0,1})$ is (os1) complete, i.e an orthonormal basis in $\mathbb{L}_2(\lambda_{0,1})$ and as process $u_j^2 = (u_j^2)_{j \in \mathbb{N}}$ on $([0, 1], \mathcal{B}_{[0,1]})$ satisfies (os2) $\|\mathcal{U}_{\mathbb{N}}(\mathfrak{a}^2 u_j^2)\|_{\mathbb{L}_{\infty}(\lambda_{0,1})} \leq \tau_{\mathfrak{a}, u}^2$ for $\tau_{\mathfrak{a}, u} \in [1, \infty)$. \square

§20.27 **Reminder.** Under Assumption §20.26 we have $\ell_2^{\mathfrak{a}} = \text{dom}(M_{\mathfrak{a},\cdot}) = \ell_2 \mathfrak{a} \subseteq \ell_2$ and the three measures $\nu_{\mathbb{N}}$, $\mathfrak{a}^{-2} \nu_{\mathbb{N}}$ and $|\phi| \nu_{\mathbb{N}}$ dominate mutually each other, i.e. they share the same null sets (see **Property** §11.05). We consider $\ell_2^{\mathfrak{a}}$ endowed with $\|\cdot\|_{\mathfrak{a}^{-1}} = \|M_{\mathfrak{a},\cdot}\|_{\ell_2}$ and given a constant $r \in \mathbb{R}_{>0}^+$ the ellipsoid $\ell_2^{\mathfrak{a},r} := \{\mathfrak{a} \in \ell_2^{\mathfrak{a}} : \|\mathfrak{a}\|_{\mathfrak{a}^{-1}} \leq r\} \subseteq \ell_2^{\mathfrak{a}}$. Since $(\mathfrak{a}\phi)_{\cdot} \in \ell_2$, and hence $\|\mathfrak{a} \mathbb{1}_{\cdot}^{m\perp}\|_{\phi} = \|(\mathfrak{a}\phi)_{\cdot} \mathbb{1}_{\cdot}^{m\perp}\|_{\ell_2} \in \mathbb{R}^+$ for each $m \in \mathbb{N}$ ($\|\mathfrak{a} \mathbb{1}_{\cdot}^{m\perp}\|_{\phi} = o(1)$ as $m \rightarrow \infty$ by dominated convergence) we have $\ell_2^{\mathfrak{a}} \subseteq \text{dom}(\phi \nu_{\mathbb{N}})$ (**Property** §11.27), and $|\phi \nu_{\mathbb{N}}(f \mathbb{1}_{\cdot}^{m\perp})| \leq r \|\mathfrak{a} \mathbb{1}_{\cdot}^{m\perp}\|_{\phi}$ for all $f \in \ell_2^{\mathfrak{a},r}$ (**Lemma** §11.29). \square

§20.28 **Remark.** Under Assumption §20.26 considering the set $\mathbb{F}_2^{\mathfrak{a},r}$ of regression functions in $\mathbb{L}_2(\lambda_{[0,1]})$ defined in (20.04) we have $\|f\|_{\mathbb{L}_{\infty}(\lambda_{[0,1]})} \leq r \tau_{\mathfrak{a},u}$ for all $f \in \mathbb{F}_2^{\mathfrak{a},r}$ due to **Lemma** §20.13. \square

§20.29 **Proposition (Upper bound).** Let Assumptions §19.02 and §20.26 be satisfied. For $n \in \mathbb{N}$ considering $m_n^* \in \mathbb{N}$ and $R_n^*(\mathfrak{a}, \phi) \in \mathbb{R}^+$ as in (12.13) (**Proposition** §12.42) we have

$$\sup \{ \mathcal{U}_f^{\otimes n} (|\phi \nu_{\mathbb{N}}(\widehat{f}^m - f)|^2) : f \in \mathbb{F}_2^{\mathfrak{a},r} \} \leq C R_n^*(\mathfrak{a}, \phi).$$

with constant $C = \sigma_{\xi}^2 + r^2 \tau_{\mathfrak{a},u}^2$.

§20.30 **Proof of Proposition** §20.29. is given in the lecture. \square

§20.31 **Illustration.** Consider the *trigonometric basis* as in **Illustration** §20.17 which satisfies Assumption §20.26 for all $\mathfrak{a} \in \ell_2$. In Table 04 [§12] the order of the rate $R_n^*(\mathfrak{a}, \phi)$ is depict for the two cases **(o)** and **(s)** introduced in **Illustration** §12.47. We note that we have $\mathfrak{a} \in \ell_2$ in case **(o)** for $a > 1/2$ while in case **(s)** for $a \in \mathbb{R}_{>0}^+$. \square

§20.32 **Corollary (Upper bound).** Let Assumptions §19.02 and §20.26 be satisfied. For $n \in \mathbb{N}$ considering $m_n^* \in \mathbb{N}$ and $R_n^*(\mathfrak{a}, \phi) \in \mathbb{R}^+$ as in (12.13) (**Proposition** §12.42) we have

$$\sup \{ (\mathcal{U}_{[0,1]} \odot \mathbb{P}_{f(x)}^{\xi})^{\otimes n} (|\phi \nu_{\mathbb{N}}((\widehat{f}^m - f) \mathbb{1}_{\cdot}^m)|^2) : f \in \mathbb{F}_2^{\mathfrak{a},r}, \mathbb{P}^{\xi} \in \mathbb{P}_{\{0\} \times (0, \sigma_{\xi}^2)} \} \leq C R_n^*(\mathfrak{a}, \phi).$$

with constant $C = \sigma^2 + r^2 \tau_{\mathfrak{a},u}^2$.

§20.33 **Proof of Corollary** §20.32. is given in the lecture. \square

§21 Minimax optimal regression

§21|01 Maximal local ϕ -risk

§21.01 **Reminder (Maximal local ϕ -risk).** Under Assumptions §19.02 and §20.26 the observable noisy version $\widehat{f} = f + n^{-1/2} \varepsilon$ of the regression coefficients $f = Uf \in \ell_2$ take the form of a *statistical direct problem* (see **Definition** §10.19) where the stochastic processes $\varepsilon \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B})^{\otimes n} \otimes 2^{\mathbb{N}}$ is given in **Definition** §19.08. Under Assumptions §19.02 and §20.26 in **Proposition** §20.29 is shown an upper bound for a maximal local ϕ -risk of an OPE over the class $\mathbb{F}_2^{\mathfrak{a},r} \subseteq \mathbb{L}_2(\lambda_{[0,1]})$ of regression functions defined in (20.04). More precisely, assuming $\mathbb{P}^{\xi} \in \mathbb{P}_{\{0\} \times \mathbb{R}_{>0}^+}$ with $\sigma_{\xi}^2 = \mathbb{P}^{\xi}(\text{id}_{\mathbb{R}}^2) \in \mathbb{R}_{>0}^+$ and for $f \in \mathbb{F}_2^{\mathfrak{a},r}$ setting $\mathcal{U}_f := \mathcal{U}_{[0,1]} \odot \mathbb{P}_{f(x)}^{\xi}$ the performance of the OPE $\widehat{f}^m = \widehat{f} \mathbb{1}_{\cdot}^m \in \ell_2 \mathbb{1}_{\cdot}^m \subseteq \text{dom}(\phi \nu_{\mathbb{N}})$ with dimension $m \in \mathbb{N}$ is measured by its maximal local ϕ -risk, that is

$$\mathcal{R}_n^{\phi}[\widehat{f}^m | \mathbb{F}_2^{\mathfrak{a},r}, \{\mathbb{P}^{\xi}\}] := \sup \{ \mathcal{U}_f^{\otimes n} (|\phi \nu_{\mathbb{N}}(\widehat{f}^m - f)|^2) : \mathcal{U}_f := \mathcal{U}_{[0,1]} \odot \mathbb{P}_{f(x)}^{\xi}, f \in \mathbb{F}_2^{\mathfrak{a},r} \}.$$

indicating explicitly the dependence on the error distribution $\mathbb{P}^{\xi} \in \mathbb{P}_{\{0\} \times \mathbb{R}_{>0}^+}$. Let us recall (12.13) (**Proposition** §12.42) where for $n, m \in \mathbb{N}$ we have defined

$$R_n^m(\mathfrak{a}, \phi) := \|\mathfrak{a} \mathbb{1}_{\cdot}^{m\perp}\|_{\phi}^2 + n^{-1} \|\mathbb{1}_{\cdot}^m\|_{\phi}^2, \quad m_n^* := \arg \min \{ R_n^m(\mathfrak{a}, \phi) : m \in \mathbb{N} \}$$

and $R_n^*(\mathfrak{a}, \phi) := R_n^{m_n^*}(\mathfrak{a}, \phi) = \min \{ R_n^m(\mathfrak{a}, \phi) : m \in \mathbb{N} \}.$ (21.01)

By **Proposition** §20.29 under Assumptions §19.02 and §20.26 the maximal local ϕ -risk of an OPE $\widehat{f}_n^{m_n^*}$ with optimally chosen dimension m_n^* as in (21.01) satisfies

$$\mathcal{R}_n^\phi[\widehat{f}_n^{m_n^*} | \mathbb{F}_2^{a,r}, \{\mathbb{P}^\xi\}] \leq C R_n^*(\mathbf{a}, \phi)$$

with $C = \sigma_\xi^2 + r^2 \tau_{a,u}^2$. □

§21.02 **Lemma (Lower bound based on two hypotheses).** Given $\mathbb{P}^\xi \in \mathbb{P}_{\{0\} \times \mathbb{R}_0^+}$ if there are $f^0, f^1 \in \mathbb{F}_2^{a,r}$ with associated probability measures $\mathbb{P}_0 := \mathcal{U}_{[0,1]} \odot \mathbb{P}_{f^0(x)}^\xi$ and $\mathbb{P}_1 := \mathcal{U}_{[0,1]} \odot \mathbb{P}_{f^1(x)}^\xi$ such that $\text{KL}(\mathbb{P}_0 | \mathbb{P}_1) \leq 2n^{-1}$ then for all $n \geq 2$ we have

$$\inf_{\widetilde{f}} \mathcal{R}_n^\phi[\widetilde{f} | \mathbb{F}_2^{a,r}, \{\mathbb{P}^\xi\}] \geq \frac{1}{64} |\phi \nu_N(f^0 - f^1)|^2.$$

where the infimum is taken over all possible estimators.

§21.03 **Proof of Lemma** §21.02. is given in the lecture. □

§21.04 **Remark.** If we consider furthermore candidate regression functions $f^0 := f^*$ and $f^1 = -f^*$ for some $f^* \in \mathbb{F}_2^{a,r}$, and hence by definition $f^0, f^1 \in \mathbb{F}_2^{a,r}$, then trivially $|\phi \nu(f^0 - f^1)|^2 = 4|\phi \nu_N(f^*)|^2$. If the associated probability measures $\mathcal{U}_{f^0} = \mathcal{U}_{[0,1]} \odot \mathbb{P}_{f^0(x)}^\xi$ and $\mathcal{U}_{f^1} = \mathcal{U}_{[0,1]} \odot \mathbb{P}_{f^1(x)}^\xi$ satisfy $\text{KL}(\mathcal{U}_{f^0} | \mathcal{U}_{f^1}) \leq 2n^{-1}$ then due to **Lemma** §21.02 for all $n \geq 2$ we have

$$\inf_{\widetilde{f}} \mathcal{R}_n^\phi[\widetilde{f} | \mathbb{F}_2^{a,r}, \{\mathbb{P}^\xi\}] \geq \frac{1}{16} |\phi \nu_N(f^*)|^2. \quad (21.02)$$

We find a minimax-optimal lower bound by constructing a candidate $f^* = U^* f^* \in \mathbb{F}_2^{a,r}$ that has the largest possible $|\phi \nu_N(f^*)|^2$ -value but \mathcal{U}_{f^0} and \mathcal{U}_{f^1} are still statistically indistinguishable in the sense that $\text{KL}(\mathcal{U}_{f^0} | \mathcal{U}_{f^1}) \leq 2n^{-1}$. □

§21.05 **Assumption.** The distribution $\mathbb{P}^\xi \in \mathcal{W}(\mathcal{B})$ admits a Lebesgue-density $\mathbb{p}^\xi := d\mathbb{P}^\xi/d\lambda$ and $\xi + x \sim \mathbb{P}_x^\xi$ for all $x \in \mathbb{R}$. There exist constants $C_\xi, x_\xi \in \mathbb{R}_0^+$ such that

$$\forall x \in [-x_\xi, x_\xi] : \quad \text{KL}(\mathbb{P}^\xi | \mathbb{P}_x^\xi) = \int \mathbb{p}^\xi(u) \log \left(\frac{\mathbb{p}^\xi(u)}{\mathbb{p}^\xi(u-x)} \right) \lambda(du) \leq C_\xi x^2. \quad \square$$

§21.06 **Lemma.** Let $\mathbb{P}^\xi \in \mathcal{W}(\mathcal{B})$ satisfy Assumption §21.05 with constants $C_\xi, x_\xi \in \mathbb{R}_0^+$ and under Assumption §20.26 let $f^* \in \ell_2^{a,r}$ fulfill $\|f^*\|_{a^{-1}} \leq x_\xi / (2\tau_{a,u})$. Setting $f^0 := U^* f^*$ and $f^1 := -U^* f^*$ the distributions $\mathcal{U}_{f^r} := \mathcal{U}_{[0,1]} \odot \mathbb{P}_{f^r(x)}^\xi$, $\tau \in \{0, 1\}$ satisfy $\text{KL}(\mathcal{U}_{f^0} | \mathcal{U}_{f^1}) \leq 4C_\xi \|f^*\|_{\mathbb{L}_2(\lambda_{[0,1]})}^2$.

§21.07 **Proof of Lemma** §21.06. is given in the lecture. □

§21.08 **Reminder.** Under Assumption §20.26 let in addition $\alpha_2^* \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ (see **Notation** §13.23), then Assumption §13.24 is satisfied. If $\alpha_2^* > n^{-1}$ then exploiting the definition (21.01) of m_n^* we have $\alpha_{m_n^*}^2 > n^{-1} \geq \alpha_{m_n^*+1}^2$ (see **Comment** §13.25) which we use in the next proof. □

§21.09 **Proposition (Lower bound).** Let $\mathbb{P}^\xi \in \mathcal{W}(\mathcal{B})$ satisfy Assumption §21.05 with constants $C_\xi, x_\xi \in \mathbb{R}_0^+$ and let Assumptions §19.02 and §20.26 be fulfilled. If $\alpha_2^* \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ then for all $n \in \mathbb{N} \cap (1 \vee \alpha_2^{*-2}, \infty)$ we have

$$\inf_{\widetilde{f}} \mathcal{R}_n^\phi[\widetilde{f} | \mathbb{F}_2^{a,r}, \{\mathbb{P}^\xi\}] \geq C R_n^*(\mathbf{a}, \phi) \quad (21.03)$$

with constant $C := 16^{-1} (r^2 \wedge x_\xi^2 / (4\tau_{a,u}^2) \wedge 1 / (2C_\xi))$ and infimum taken over all estimators.

§21.10 **Proof** of **Proposition** §21.09. is given in the lecture. \square

§21.11 **Comment**. If ξ is normally distributed with mean zero and variance $\sigma_\xi^2 \in \mathbb{R}_0^+$, i.e. $\xi \sim N_{(0, \sigma_\xi^2)}$, then for all $x \in \mathbb{R}$ we have

$$\text{KL}(N_{(0, \sigma_\xi^2)} | N_{(x, \sigma_\xi^2)}) = N_{(0, \sigma_\xi^2)} \left(\log \frac{dN_{(0, \sigma_\xi^2)}}{dN_{(x, \sigma_\xi^2)}} \right) = \frac{x^2}{2\sigma_\xi^2}$$

and thus Assumption §21.05 holds with $C_\xi = 1/(2\sigma_\xi^2)$ and $x_\xi^2 = \infty$ (see **Proof** §13.15). Consequently, from **Proposition** §21.09 we obtain immediately,

$$\inf_{\tilde{f}} \mathcal{R}_n^\phi[\tilde{f} | \mathbb{F}_2^{\text{a.r.}}, \{N_{(0, \sigma_\xi^2)}\}] \geq C R_n^*(\mathbf{a}, \phi) \quad (21.04)$$

with constant $C := 16^{-1}(r^2 \wedge \sigma_\xi^2)$ and infimum taken over all estimators. \square

§21.12 **Corollary (Lower bound)**. Let Assumptions §19.02 and §20.26 be fulfilled and let $\sigma^2 \in \mathbb{R}_0^+$. If $\alpha_*^2 \in (\mathbb{R}_0^+)^{\mathbb{N}}$ then for all $n \in \mathbb{N} \cap (1 \vee \alpha_*^2, \infty)$ we have

$$\inf_{\tilde{f}} \mathcal{R}_n^\phi[\tilde{f} | \mathbb{F}_2^{\text{a.r.}}, P_{\{0\} \times (0, \sigma^2)}] \geq C R_n^*(\mathbf{a}, \phi) \quad (21.05)$$

with constant $C := 16^{-1}(r^2 \wedge \sigma^2)$ and infimum taken over all estimators.

§21.13 **Proof** of **Corollary** §21.12. is given in the lecture. \square

§21.14 **Illustration**. Consider the *trigonometric basis* as in **Illustration** §20.17 which satisfies Assumption §20.26 for all $\mathbf{a} \in \ell_2$ (see **Illustration** §20.31). In Table 04 [§12] the order of the rate $R_n^*(\mathbf{a}, \phi)$ is depicted for the two cases **(o)** and **(s)** introduced in **Illustration** §12.47. We note that we have $\mathbf{a} \in \ell_2$ in case **(o)** for $a > 1/2$ while in case **(s)** for $a \in \mathbb{R}_0^+$. In both cases the additional assumption $\alpha_*^2 \in (\mathbb{R}_0^+)^{\mathbb{N}}$ is satisfied. Consequently, due to **Proposition** §21.09 the Table 04 [§12] presents the order of the *minimax rate* $R_n^*(\mathbf{a}, \phi)$ which is attained by the *minimax-optimal* OPE $\hat{f}_*^{m_n^*} = \hat{f}_* \mathbb{1}_*^{m_n^*} \in \ell_2 \mathbb{1}_*^{m_n^*} \subseteq \text{dom}(\phi_{\mathbf{u}_n})$ with optimally selected dimension m_n^* (**Proposition** §20.29). We shall stress, that the order of m_n^* given in the Table 04 [§12] depends on the parameter $a \in \mathbb{R}_0^+$ characterising the (abstract) smoothness of the density of interest which is generally not known in advance. \square

§21|02 Maximal global v-risk

§21.15 **Reminder (Maximal global v-risk)**. Under Assumptions §19.02 and §20.11 the observable noisy version $\hat{f}_* = f_* + n^{-1/2} \boldsymbol{\varepsilon}_*$ of the regression coefficients $f_* = Uf \in \ell_2$ take the form of a *statistical direct problem* (see **Definition** §10.19) where the stochastic processes $\boldsymbol{\varepsilon}_* \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B})^{\otimes n} \otimes 2^{\mathbb{N}}$ is given in **Definition** §19.08. Under Assumptions §19.02 and §20.11 in **Proposition** §20.15 is shown an upper bound for a maximal global v-risk of an OPE over the class $\mathbb{F}_2^{\text{a.r.}} \subseteq \mathbb{L}_2(\lambda_{[0,1]})$ of regression functions defined in (20.04). More precisely, assuming $\mathbb{P}^\xi \in P_{\{0\} \times \mathbb{R}_0^+}$ with $\sigma_\xi^2 = \mathbb{P}^\xi(\text{id}_{\mathbb{R}}^2) \in \mathbb{R}_0^+$ and for $f \in \mathbb{F}_2^{\text{a.r.}}$ setting $\mathcal{U}_f := \mathcal{U}_{[0,1]} \odot \mathbb{P}_{f(x)}^\xi$ the performance of the OPE $\hat{f}_*^m = \hat{f}_* \mathbb{1}_*^m \in \ell_2 \mathbb{1}_*^m \subseteq \ell_2(\mathfrak{v}^2)$ with dimension $m \in \mathbb{N}$ is measured by its maximal global v-risk, that is

$$\mathcal{R}_n^{\mathfrak{v}}[\hat{f}_*^m | \mathbb{F}_2^{\text{a.r.}}, \{\mathbb{P}^\xi\}] := \sup \{ \mathcal{U}_f^{\otimes n} (\|\hat{f}_*^m - f\|_{\mathfrak{v}}^2) : \mathcal{U}_f := \mathcal{U}_{[0,1]} \odot \mathbb{P}_{f(x)}^\xi, f \in \mathbb{F}_2^{\text{a.r.}} \}.$$

Let us recall (12.06) (**Proposition** §12.21) where for $n, m \in \mathbb{N}$ we have defined $(\mathfrak{av})_{(m)}^2 = \|\mathfrak{av}\|_{\mathbb{1}_*^{m\perp}}^2$ and

$$R_n^m(\mathbf{a}, \mathfrak{v}) := [(\mathfrak{av})_{(m)}^2 \vee n^{-1} \|\mathbb{1}_*^m\|_{\mathfrak{v}}^2], \quad m_n^* := \arg \min \{ R_n^m(\mathbf{a}, \mathfrak{v}) : m \in \mathbb{N} \}$$

and $R_n^*(\mathbf{a}, \mathfrak{v}) := R_n^{m_n^*}(\mathbf{a}, \mathfrak{v}) = \min \{ R_n^m(\mathbf{a}, \mathfrak{v}) : m \in \mathbb{N} \}.$ (21.06)

By **Proposition** §20.15 under Assumptions §19.02 and §20.11 the maximal global \mathfrak{v} -risk of an OPE $\widehat{f}_n^{m_n^*}$ with optimally chosen dimension m_n^* as in (21.06) satisfies

$$\mathcal{R}_n^{\mathfrak{v}}[\widehat{f}_n^{m_n^*} | \mathbb{F}_2^{\mathfrak{a},r}, \{\mathbb{P}^\xi\}] \leq C R_n^*(\mathfrak{a}, \mathfrak{v})$$

with $C = \sigma_\xi^2 + r^2 \tau_{\mathfrak{a},u}^2 + r^2$. Furthermore, as in **Notation** §13.29 for $m \in \mathbb{N}$ we set $\mathcal{T}_m := \{-1, 1\}^m$ and for each $\tau := (\tau_j)_{j \in \llbracket m \rrbracket} \in \mathcal{T}_m$ and $j \in \llbracket m \rrbracket$ we introduce $\tau^{(j)} \in \mathcal{T}_m$ given by $\tau_j^{(j)} := -\tau_j$ and $\tau_l^{(j)} := \tau_l$ for $l \in \llbracket m \rrbracket \setminus \{j\}$. \square

§21.16 **Lemma (Assouad's cube technique)**. Given $\mathbb{P}^\xi \in \mathbb{P}_{\{0\} \times \mathbb{R}_0^+}$ if for each $\tau \in \mathcal{T}_m$ there is $f^\tau \in \mathbb{F}_2^{\mathfrak{a},r}$ with associated probability measure $\mathcal{U}_{f^\tau} := \mathcal{U}_{[0,1]} \odot \mathbb{P}_{f^\tau(x)}^\xi$ such that for all $\tau \in \mathcal{T}_m$ and $j \in \llbracket m \rrbracket$ we have $\text{KL}(\mathcal{U}_{f^\tau} | \mathcal{U}_{f^{(j)}}) \leq 2n^{-1}$ then for all $n \geq 2$

$$\inf_{\widehat{f}} \mathcal{R}_n^{\mathfrak{v}}[\widehat{f} | \mathbb{F}_2^{\mathfrak{a},r}, \{\mathbb{P}^\xi\}] \geq 2^{-m} \sum_{\tau \in \mathcal{T}_m} \frac{1}{64} \sum_{j \in \llbracket m \rrbracket} (\mathfrak{v}_j^2 |f_j^\tau - f_j^{\tau^{(j)}}|^2)$$

where the infimum is taken over all possible estimators.

§21.17 **Proof of Lemma** §21.16. is given in the lecture. \square

§21.18 **Remark**. Assume candidate regression functions $f^\tau := U^* f_\cdot^\tau$ with $f_\cdot^\tau := (\tau_j f_j^* \mathbf{1}_j^m)_{j \in \mathbb{N}}$, $\tau \in \mathcal{T}_m$, for some $f_\cdot^* \in \ell_2^{\mathfrak{a},r}$, where evidently $f_\cdot^\tau \in \ell_2^{\mathfrak{a},r}$ too, then trivially $\sum_{j \in \llbracket m \rrbracket} (\mathfrak{v}_j^2 |f_j^\tau - f_j^{\tau^{(j)}}|^2) = 4 \|f_\cdot^* \mathbf{1}_\cdot^m\|_{\mathfrak{v}}^2$. If for all $\tau \in \mathcal{T}_m$ and $j \in \llbracket m \rrbracket$ the associated probability measures $\mathcal{U}_{f^\tau} = \mathcal{U}_{[0,1]} \odot \mathbb{P}_{f^\tau(x)}^\xi$ and $\mathcal{U}_{f^{(j)}} := \mathcal{U}_{[0,1]} \odot \mathbb{P}_{f^{(j)}(x)}^\xi$ satisfy $\text{KL}(\mathcal{U}_{f^\tau} | \mathcal{U}_{f^{(j)}}) \leq 2n^{-1}$ then due to **Lemma** §21.16 for all $n \geq 2$ we have

$$\inf_{\widehat{f}} \mathcal{R}_n^{\mathfrak{v}}[\widehat{f} | \mathbb{F}_2^{\mathfrak{a},r}, \{\mathbb{P}^\xi\}] \geq 2^{-m} \sum_{\tau \in \mathcal{T}_m} \frac{1}{16} \|f_\cdot^* \mathbf{1}_\cdot^m\|_{\mathfrak{v}}^2 = \frac{1}{16} \|f_\cdot^* \mathbf{1}_\cdot^m\|_{\mathfrak{v}}^2. \quad (21.07)$$

We find a minimax-optimal lower bound by choosing the parameter m and the function f_\cdot^* that have the largest possible $\|f_\cdot^* \mathbf{1}_\cdot^m\|_{\mathfrak{v}}^2$ -value although that the associated \mathcal{U}_{f^τ} , $\tau \in \mathcal{T}_m$ are still statistically indistinguishable in the sense that $\text{KL}(\mathcal{U}_{f^\tau} | \mathcal{U}_{f^{(j)}}) \leq 2n^{-1}$ for all $j \in \llbracket m \rrbracket$ and $\tau \in \mathcal{T}_m$. \square

§21.19 **Lemma**. Let $\mathbb{P}^\xi \in \mathcal{W}(\mathcal{B})$ satisfy Assumption §21.05 with constants $C_\xi, x_\xi \in \mathbb{R}_{>0}^+$ and under Assumption §20.11 let $f_\cdot^* \in \ell_2^{\mathfrak{a},r}$ fulfill $\|f_\cdot^*\|_{\mathfrak{a}^{-1}} \leq x_\xi / (2\tau_{\mathfrak{a},u})$. For each $\tau \in \mathcal{T}_m$ introduce $f_\cdot^\tau := (\tau_j f_j^* \mathbf{1}_j^m)_{j \in \mathbb{N}} \in \ell_2^{\mathfrak{a},r}$ and $f^\tau := U^* f_\cdot^\tau \in \mathbb{F}_2^{\mathfrak{a},r}$ with associated probability measure $\mathcal{U}_{f^\tau} = \mathcal{U}_{[0,1]} \odot \mathbb{P}_{f^\tau(x)}^\xi$. Then for each $j \in \llbracket m \rrbracket$ we have $\text{KL}(\mathcal{U}_{f^\tau} | \mathcal{U}_{f^{(j)}}) \leq 4C_\xi \|f_\cdot^* \mathbf{1}_\cdot^m\|_{\ell_\infty}^2$.

§21.20 **Proof of Lemma** §21.19. is given in the lecture. \square

§21.21 **Reminder**. For $w_\cdot \in \ell_\infty$ we set $w_{(0)}^2 := \|w_\cdot\|_{\ell_\infty}^2$ and $w_{(j)}^2 = (w_{(j)}^2 := \|w_\cdot \mathbf{1}_\cdot^{j+1}\|_{\ell_\infty}^2)_{j \in \mathbb{N}}$ (**Notation** §13.34) where by construction $w_{(j)}^2 = \sup \{w_i^2 : i \in \mathbb{N} \cap [j+1, \infty)\}$, $j \in \mathbb{N}_0$ and $w_{(0)}^2 \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$. Under Assumption §20.11 let in addition $(\mathfrak{a}\mathfrak{v})_{(s)}^2 \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ and there exists $C_{(\mathfrak{a}\mathfrak{v})} \in (0, 1]$ such that $C_{(\mathfrak{a}\mathfrak{v})} \|(\mathfrak{a}\mathfrak{v})_{(s)}^{-2} \mathbf{1}_\cdot^m\|_{\ell_\infty} \leq (\mathfrak{a}\mathfrak{v})_{(m-1)}^{-2}$ or in equal

$$(\mathfrak{a}\mathfrak{v})_{(m-1)}^2 \geq \min \{(\mathfrak{a}\mathfrak{v})_j^2 : j \in \llbracket m \rrbracket\} \geq C_{(\mathfrak{a}\mathfrak{v})} (\mathfrak{a}\mathfrak{v})_{(m-1)}^2$$

for all $m \in \mathbb{N}$, then Assumption §13.35 is satisfied. For m_n^* and $R_n^* := R_n^{m_n^*}(\mathfrak{a}, \mathfrak{v})$ as in (21.06) we distinguish case i) : $R_n^* = n^{-1} \|\mathbf{1}_\cdot^{m_n^*}\|_{\mathfrak{v}}^2 > (\mathfrak{a}\mathfrak{v})_{(m_n^*)}^2$ and case ii) : $R_n^* = (\mathfrak{a}\mathfrak{v})_{(m_n^*)}^2 \geq n^{-1} \|\mathbf{1}_\cdot^{m_n^*}\|_{\mathfrak{v}}^2$. Due to **Comment** §13.36 if $(\mathfrak{a}\mathfrak{v})_{(1)}^2 > n^{-1} \mathfrak{v}_1^2$ then in case i) we obtain $(\mathfrak{a}\mathfrak{v})_{(m_n^*-1)}^2 \geq n^{-1} \|\mathbf{1}_\cdot^{m_n^*}\|_{\mathfrak{v}}^2$, while in case ii) setting (the defining set is not empty since $(\mathfrak{a}\mathfrak{v})_{(s)}^2 \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$)

$$m_n^\diamond := \min \{m \in \mathbb{N} \cap [m_n^* + 1, \infty) : n^{-1} \|\mathbf{1}_\cdot^m\|_{\mathfrak{v}}^2 \geq (\mathfrak{a}\mathfrak{v})_{(m)}^2\} \quad (21.08)$$

we have $(\mathfrak{a}\mathfrak{v})_{(m_n^\diamond)}^2 = (\mathfrak{a}\mathfrak{v})_{(m_n^\diamond-1)}^2 \leq n^{-1} \|\mathbf{1}_\cdot^{m_n^\diamond}\|_{\mathfrak{v}}^2$. We use those estimates in the next proof. \square

§21.22 **Proposition (Lower bound).** Let $\mathbb{P}^\xi \in \mathcal{W}(\mathcal{B})$ satisfy Assumption §21.05 with constants $C_\xi, x_\xi \in \mathbb{R}_0^+$ and let Assumptions §19.02 and §20.11 be fulfilled. If $(\mathbf{av})_{(\bullet)}^2 \in (\mathbb{R}_{(0)}^+)^{\mathbb{N}}$ and there exists $C_{(\mathbf{av})} \in (0, 1]$ such that $C_{(\mathbf{av})} \|(\mathbf{av})_{\bullet}^{-2} \mathbb{1}_{\bullet}^m\|_{\ell_\infty} \leq (\mathbf{av})_{(m-1)}^{-2}$ for all $m \in \mathbb{N}$, then for all $n \in \mathbb{N} \cap (1 \vee \mathbf{v}_1^2(\mathbf{av})_{(1)}^{-2}, \infty)$ we have

$$\inf_{\tilde{f}} \mathcal{R}_n^{\mathbf{v}}[\tilde{f} | \mathbb{F}_2^{\mathbf{a}, \mathbf{r}}, \{\mathbb{P}^\xi\}] \geq C R_n^*(\mathbf{a}, \mathbf{v}) \tag{21.09}$$

with constant $C := 16^{-1}(C_{(\mathbf{av})} x_\xi^2 / (4\tau_{\mathbf{a}, \mathbf{u}}^2) \wedge C_{(\mathbf{av})} \mathbf{r}^2 \wedge 1 / (2C_\xi))$ and infimum taken over all estimators.

§21.23 **Proof of Proposition §21.22.** is given in the lecture. □

§21.24 **Comment.** If $\xi \sim N_{(0, \sigma_\xi^2)}$ with $\sigma_\xi^2 \in \mathbb{R}_0^+$ then Assumption §21.05 holds with $C_\xi = 1 / (2\sigma_\xi^2)$ and $x_\xi^2 = \infty$ (see **Comment §21.11**). Consequently, from **Proposition §21.22** we obtain immediately

$$\inf_{\tilde{f}} \mathcal{R}_n^{\mathbf{v}}[\tilde{f} | \mathbb{F}_2^{\mathbf{a}, \mathbf{r}}, \{N_{(0, \sigma_\xi^2)}\}] \geq C R_n^*(\mathbf{a}, \mathbf{v}) \tag{21.10}$$

with constant $C := 16^{-1}(C_{(\mathbf{av})} \mathbf{r}^2 \wedge \sigma_\xi^2)$ and infimum taken over all estimators. □

§21.25 **Corollary (Lower bound).** Let Assumptions §19.02 and §20.11 be fulfilled and let $\sigma^2 \in \mathbb{R}_0^+$. If $(\mathbf{av})_{(\bullet)}^2 \in (\mathbb{R}_{(0)}^+)^{\mathbb{N}}$ and there exists $C_{(\mathbf{av})} \in (0, 1]$ such that $C_{(\mathbf{av})} \|(\mathbf{av})_{\bullet}^{-2} \mathbb{1}_{\bullet}^m\|_{\ell_\infty} \leq (\mathbf{av})_{(m-1)}^{-2}$ for all $m \in \mathbb{N}$, then for all $n \in \mathbb{N} \cap (1 \vee \mathbf{v}_1^2(\mathbf{av})_{(1)}^{-2}, \infty)$ we have

$$\inf_{\tilde{f}} \mathcal{R}_n^{\phi}[\tilde{f} | \mathbb{F}_2^{\mathbf{a}, \mathbf{r}}, P_{\{0\} \times (0, \sigma^2)}] \geq C R_n^*(\mathbf{a}, \phi) \tag{21.11}$$

with constant $C := 16^{-1}(C_{(\mathbf{av})} \mathbf{r}^2 \wedge \sigma^2)$ and infimum taken over all estimators.

§21.26 **Proof of Corollary §21.25.** is given in the lecture. □

§21.27 **Illustration.** Consider the *trigonometric basis* as in **Illustration §20.17** which satisfies Assumption §20.11 for all $\mathbf{a}_\bullet \in \ell_2$ (see **Illustration §20.17**). In Table 02 [§12] the order of the rate $R_n^*(\mathbf{a}, \mathbf{v})$ is depicted for the two cases **(o)** and **(s)** introduced in **Illustration §12.26**. We note that we have $\mathbf{a}_\bullet \in \ell_2$ in case **(o)** for $a > 1/2$ while in case **(s)** for $a \in \mathbb{R}_0^+$. In both cases the additional assumptions, $(\mathbf{av})_{(\bullet)}^2 \in (\mathbb{R}_{(0)}^+)^{\mathbb{N}}$ and there exists $C_{(\mathbf{av})} \in (0, 1]$ such that $C_{(\mathbf{av})} \|(\mathbf{av})_{\bullet}^{-2} \mathbb{1}_{\bullet}^m\|_{\ell_\infty} \leq (\mathbf{av})_{(m-1)}^{-2}$ for all $m \in \mathbb{N}$, are satisfied. Consequently, due to **Proposition §21.22** the Table 02 [§12] presents the order of the *minimax rate* $R_n^*(\mathbf{a}, \mathbf{v})$ which is attained by the *minimax-optimal* OPE $\hat{f}_n^{m_n^*} = \hat{f}_n \mathbb{1}_n^{m_n^*} \in \ell_2 \mathbb{1}_n^{m_n^*} \subseteq \ell_2(\mathbf{v}^2)$ with optimally selected dimension m_n^* (**Proposition §20.15**). We shall stress, that the order of m_n^* given in the Table 02 [§12] depends on the parameter $a \in \mathbb{R}_0^+$ characterising the (abstract) smoothness of the regression function of interest which is generally not known in advance. □

§22 Data-driven regression

§22|01 Data-driven global estimation by model selection

§22.01 **Reminder.** Talagrand's inequality stated in the form of **Lemma §18.01** provides again our key argument in order to control the deviations of the reminder term. Let us briefly recall how we intend to apply Talagrand's inequality (see **Remark §18.02** for a similar approach). Reconsider the stochastic process $\psi = (\psi_j(X, Y) := Y u_j(X))_{j \in \mathbb{N}} \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B}) \otimes 2^{\mathbb{N}}$ where $\psi_j \in \mathcal{L}_1(\mathcal{U}_j)$ for each $j \in \mathbb{N}$ and the OPE $\hat{f}^m = \hat{f} \mathbb{1}^m \in \ell_2 \mathbb{1}^m$ with dimension $m \in \mathbb{N}$ (**Definition §20.04**). $\hat{f}_\bullet = \hat{\mathbb{P}}_n \psi_\bullet = (\hat{\mathbb{P}}_n \psi_j)_{j \in \mathbb{N}}$ are noisy versions (**Definition §15.08**) of the regression coefficients $f_\bullet = Uf = \mathcal{U}_j(\psi_\bullet) = (\mathcal{U}_j \psi_j = \mathcal{U}_j(Y u_j(X)))_{j \in \mathbb{N}}$ (see **Notation §19.07**). For $m \in \mathbb{N}$ introduce

the unit ball $\mathbb{B}_m := \{a \in \ell_2(\mathfrak{v})\mathbb{1}^m : \|a\|_{\mathfrak{v}} \leq 1\}$ contained in the linear subspace $\ell_2(\mathfrak{v})\mathbb{1}^m$ spanned by $(\mathbb{1}^{(j)})_{j \in [m]}$. Clearly, for each $a \in \ell_2(\mathfrak{v})\mathbb{1}^m$ we have $r_a := \sum_{j \in [m]} \mathfrak{v}_j^2 a_j \psi_j = \nu_{\mathbb{N}}(\mathfrak{v}^2 a, \psi) \in \mathcal{B}_{[0,1]} \otimes \mathcal{B}$, i.e. it is a $\mathcal{B}_{[0,1]} \otimes \mathcal{B}$ - \mathcal{B} -measurable function, where $\widehat{\mathbb{P}}_n(r_a) = \nu_{\mathbb{N}}(\mathfrak{v}^2 a, \widehat{\mathbb{P}}_n \psi) = \nu_{\mathbb{N}}(\mathfrak{v}^2 a, \widehat{f})$, $\mathcal{U}_f(r_a) = \nu_{\mathbb{N}}(\mathfrak{v}^2 a, \mathcal{U}_f \psi) = \nu_{\mathbb{N}}(\mathfrak{v}^2 a, f)$ and hence $\bar{r}_a = \widehat{\mathbb{P}}_n(r_a) - \mathcal{U}_f(r_a) = \nu_{\mathbb{N}}(\mathfrak{v}^2 a, (\widehat{f} - f)) = \langle \widehat{f} - f, a \rangle_{\mathfrak{v}}$. Let \mathcal{B}_m be a countable dense subset of the unit ball \mathbb{B}_m (see Remark §18.02 for more details), then we obtain

$$\|\widehat{f}^m - f^m\|_{\mathfrak{v}}^2 = \sup \{ |\langle \widehat{f} - f, a \rangle_{\mathfrak{v}}|^2 : a \in \mathcal{B}_m \} = \sup \{ |\bar{r}_a|^2 : a \in \mathcal{B}_m \}.$$

The last identity provides the necessary argument to apply below Talagrand's inequality (§18.01) where we need to calculate the three constants h , H and v . We note that $\psi \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B}) \otimes 2^{\mathbb{N}}$ and thus $r_a = \nu_{\mathbb{N}}(\mathfrak{v}^2 a, \psi) \in \mathcal{B}_{[0,1]} \otimes \mathcal{B}$ is not bounded. Therefore, we decompose $\psi = \psi^b + \psi^u$ into two parts $\psi^b, \psi^u \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B}) \otimes 2^{\mathbb{N}}$, a bounded and a remaining unbounded one. To be more precise, for $a \in \ell_2(\mathfrak{v})\mathbb{1}^m$ setting $r_a^b := \nu_{\mathbb{N}}(\mathfrak{v}^2 a, \psi^b) \in \mathcal{B}_{[0,1]} \otimes \mathcal{B}$ and $r_a^u := \nu_{\mathbb{N}}(\mathfrak{v}^2 a, \psi^u) \in \mathcal{B}_{[0,1]} \otimes \mathcal{B}$ let $\sup \{ |r_a^b(x, y)| : a \in \mathcal{B}_d, x \in [0, 1], y \in \mathbb{R} \} \in \mathbb{R}$ be satisfied. Introducing further $\bar{r}_a^b := \widehat{\mathbb{P}}_n(r_a^b) - \mathcal{U}_f(r_a^b)$ and $\bar{r}_a^u := \widehat{\mathbb{P}}_n(r_a^u) - \mathcal{U}_f(r_a^u)$ we evidently have

$$\begin{aligned} \|\widehat{f}^m - f^m\|_{\mathfrak{v}}^2 &= \sup \{ |\bar{r}_a^b + \bar{r}_a^u|^2 : a \in \mathcal{B}_m \} \\ &\leq 2 \sup \{ |\bar{r}_a^b|^2 : a \in \mathcal{B}_m \} + 2 \sup \{ |\bar{r}_a^u|^2 : a \in \mathcal{B}_m \} \\ &= 2 \|(\widehat{\mathbb{P}}_n \psi^b - \mathcal{U}_f \psi^b)\mathbb{1}^m\|_{\mathfrak{v}}^2 + 2 \|(\widehat{\mathbb{P}}_n \psi^u - \mathcal{U}_f \psi^u)\mathbb{1}^m\|_{\mathfrak{v}}^2 \quad (22.01) \end{aligned}$$

Considering the first term on the right hand side provided that

$$\begin{aligned} \sup \{ \|\psi^b(x, y)\mathbb{1}^m\|_{\mathfrak{v}} : x \in [0, 1], y \in \mathbb{R} \} &= \sup \{ |r_a^b(x, y)| : a \in \mathcal{B}_d, x \in [0, 1], y \in \mathbb{R} \} \leq h_m, \\ \mathcal{U}_f^{\otimes n}(\|(\widehat{\mathbb{P}}_n \psi^b - \mathcal{U}_f \psi^b)\mathbb{1}^m\|_{\mathfrak{v}}^2) &= \mathcal{U}_f^{\otimes n}(\sup \{ |\bar{r}_a^b|^2 : a \in \mathcal{B}_d \}) \leq H_m^2, \\ \sup \{ \mathcal{U}_f(\|\nu_{\mathbb{N}}(\mathfrak{v}^2 a, (\psi^b - \mathcal{U}_f \psi^b))\|_{\mathfrak{v}}^2) : a \in \mathcal{B}_d \} &= \sup \{ n \mathcal{U}_f^{\otimes n}(|\bar{r}_a^b|^2) : a \in \mathcal{B}_d \} \leq v_m \quad (22.02) \end{aligned}$$

we eventually apply Talagrand's inequality (§18.01) and we obtain

$$\mathcal{U}_f^{\otimes n}(\|(\widehat{\mathbb{P}}_n \psi^b - \mathcal{U}_f \psi^b)\mathbb{1}^m\|_{\mathfrak{v}}^2 - 6H_m^2)_+ \leq C_{\text{tal}} \left\{ \frac{v_m}{n} \exp\left(\frac{-nH_m^2}{6v_m}\right) + \frac{h_m^2}{n^2} \exp\left(\frac{-nH_m}{100h_m}\right) \right\} \quad (22.03)$$

for some universal numerical constant $C_{\text{tal}} \in [1, \infty)$. \square

§22|01|01 Global \mathfrak{v} -risk

§22.02 **Assumption.** The weights $\mathfrak{v} \in (\mathbb{R}_{>0})^{\mathbb{N}}$ satisfy

$$\forall x \in \mathbb{R}_0^+ : \sum_{m \in \mathbb{N}} \{x \|\mathfrak{v}^2 \mathbb{1}^m\|_{\ell_\infty} \exp(-\|\mathfrak{v} \mathbb{1}^m\|_{\ell_2}^2 / (x \|\mathfrak{v}^2 \mathbb{1}^m\|_{\ell_\infty}))\} =: C_{\mathfrak{v}}(x) \in \mathbb{R}^+. \quad (22.04)$$

The orthonormal system $(u_j)_{j \in \mathbb{N}}$ in $\mathbb{L}_2(\lambda_{[0,1]})$ is **(os1) complete**, i.e. an *orthonormal basis* in $\mathbb{L}_2(\lambda_{[0,1]})$ and satisfies as process $u^2 = (u_j^2)_{j \in \mathbb{N}}$ on $([0, 1], \mathcal{B}_{[0,1]})$ for $\tau_{\mathfrak{v}, u} \in [1, \infty)$ and for all $m \in \mathbb{N}$ **(os3)** $\sup \{ \|u(x)\mathbb{1}^m\|_{\mathfrak{v}}^2 : x \in [0, 1] \} \leq \tau_{\mathfrak{v}, u}^2 \|\mathbb{1}^m\|_{\mathfrak{v}}^2 \in \mathbb{R}^+$. \square

§22.03 **Remark.** Under Assumption §22.02 (18.05) we have $\|\mathfrak{v} \mathbb{1}^m\|_{\ell_2}^{-2} = o(1)$ as $m \rightarrow \infty$ (**Comment** §14.22), see also **Illustration** §14.23 for an example when (18.05) is not satisfied. \square

§22.04 **Reminder (Global oracle \mathfrak{v} -risk).** Given Assumptions §19.02 and §22.02 we consider an OPE as in **Definition** §20.04. Here the observable noisy version $\widehat{f} = f + n^{-1/2}\varepsilon$ of the regression

coefficients $f_{\bullet} = Uf \in \ell_2$ take the form of a *statistical direct problem* (see Definition §10.19) where the stochastic processes $\varepsilon_{\bullet} \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B})^{\otimes n} \otimes 2^{\mathbb{N}}$ is given in Definition §19.08. Under Assumptions §19.02 and §22.02, (and hence Assumption §19.05 and $\mathbf{v} \in (\mathbb{R}_{>0})^{\mathbb{N}}$) and $f_{\bullet} \in \ell_2(\mathbf{v}^2)$ in §20.09 an *oracle inequality* for the global \mathbf{v} -risk of the OPE's is shown. More precisely, as in (20.02) (Proposition §20.07) for all $n, m \in \mathbb{N}$ setting

$$\begin{aligned} R_n^m(f, \mathbf{v}) &:= \|f_{\bullet} \mathbb{1}^{m \perp}\|_{\mathbf{v}}^2 + n^{-1} \|\mathbb{1}^m\|_{\mathbf{v}}^2, \quad m_n^{\circ} := \arg \min \{R_n^m(f, \mathbf{v}) : m \in \mathbb{N}\} \\ &\text{and } R_n^{\circ}(f, \mathbf{v}) := R_n^{m_n^{\circ}}(f, \mathbf{v}) = \min \{R_n^m(f, \mathbf{v}) : m \in \mathbb{N}\}. \end{aligned} \quad (22.05)$$

and assuming $\mathbf{v}_f := \max(\sigma_{\xi}^{-2}, \sigma_{\xi}^2 + \|f\|_{\mathbb{L}_{\infty}(\lambda_{0,1})}^2) \in \mathbb{R}_{>0}^+$ due to Property §20.09 the (infeasible) OPE $\widehat{f}_{\bullet}^{m_n^{\circ}} = \widehat{f}_{\bullet} \mathbb{1}^{m_n^{\circ}} \in \ell_2(\mathbf{v}^2) \mathbb{1}^{m_n^{\circ}} \subseteq \ell_2(\mathbf{v}^2)$ with oracle dimension m_n° as in (22.05) satisfies

$$\begin{aligned} \mathbf{v}_f^{-1} R_n^{\circ}(f, \mathbf{v}) &\leq \inf_{m \in \mathbb{N}} \mathcal{U}_f^{\otimes n}(\|\widehat{f}_{\bullet}^m - f_{\bullet}\|_{\mathbf{v}}^2) \leq \mathcal{U}_f^{\otimes n}(\|\widehat{f}_{\bullet}^{m_n^{\circ}} - f_{\bullet}\|_{\mathbf{v}}^2) \\ &\leq \mathbf{v}_f R_n^{\circ}(f, \mathbf{v}) \leq \mathbf{v}_f^2 \inf_{m \in \mathbb{N}} \mathcal{U}_f^{\otimes n}(\|\widehat{f}_{\bullet}^m - f_{\bullet}\|_{\mathbf{v}}^2), \end{aligned}$$

and hence it is *oracle optimal* (with constant \mathbf{v}_f^2). □

Partially known penalty sequence

§22.05 **Notation.** Consider a sequence of penalties $\text{pen}_{\bullet}^{f, \mathbf{v}} = (\text{pen}_m^{f, \mathbf{v}})_{m \in \mathbb{N}} \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ given by

$$\text{pen}_m^{f, \mathbf{v}} := 48\mathbf{v}_f^2 \tau_{\mathbf{v}, u}^2 n^{-1} \|\mathbb{1}^m\|_{\mathbf{v}}^2, \quad \text{for each } m \in \mathbb{N} \text{ with } \mathbf{v}_f^2 := 1 + \mathcal{U}_f(Y^2) \quad (22.06)$$

which is obviously only *partially known* in advance, and the in advance known upper bound (where the defining set is not empty)

$$\mathbb{M}^{\mathbf{v}} := \max \{m \in \mathbb{N} : \|\mathbb{1}^m\|_{\mathbf{v}}^2 \leq n\mathbf{v}_f^2, m \leq n^{-2/3} \exp(\frac{n^{1/6}}{100})\}. \quad (22.07)$$

Considering the partially data-driven OSE $\widehat{f}_{\bullet}^{\widehat{m}} = \widehat{f}_{\bullet} \mathbb{1}^{\widehat{m}}$ with dimension parameter

$$\widehat{m} := \arg \min \{-\|\widehat{f}_{\bullet}^m\|_{\mathbf{v}} + \text{pen}_m^{f, \mathbf{v}} : m \in \llbracket \mathbb{M}^{\mathbf{v}} \rrbracket\} \quad (22.08)$$

we derive below an upper bound for its global \mathbf{v} -risk, $\mathcal{U}_f^{\otimes n}(\|\widehat{f}_{\bullet}^{\widehat{m}} - f_{\bullet}\|_{\mathbf{v}}^2)$. □

§22.06 **Lemma.** Under Assumptions §19.02 and §22.02, $Y \in \mathcal{L}_5(\mathcal{U}_f)$ and $f \in \mathbb{L}_{\infty}(\lambda_{0,1})$ for $\text{pen}_{\bullet}^{f, \mathbf{v}} \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ as in (22.06) and $\mathbb{M}^{\mathbf{v}} \in \mathbb{N}$ as in (22.07) we have

$$\begin{aligned} \mathcal{U}_f^{\otimes n}(\max \{(\|\widehat{f}_{\bullet}^{\widehat{m}} - f_{\bullet}\|_{\mathbf{v}}^2 - \text{pen}_{\widehat{m}}^{f, \mathbf{v}}/4) : m \in \llbracket \mathbb{M}^{\mathbf{v}} \rrbracket\}) \\ \leq 2C_{\text{tal}} \tau_{\mathbf{v}, u}^2 (C_{\mathbf{v}}(x_{\xi, f}) + \mathbf{v}_1^2)(1 + \mathcal{U}_f(Y^2) + \mathcal{U}_f(|Y|^5)) n^{-1} \end{aligned} \quad (22.09)$$

for some universal numerical constant $C_{\text{tal}} \in [1, \infty)$ and $x_{\xi, f} = 6(\sigma_{\xi}^2 + \|f\|_{\mathbb{L}_{\infty}(\lambda_{0,1})}^2)/(\mathbf{v}_f^2 \tau_{\mathbf{v}, u}^2) \in \mathbb{R}^+$.

§22.07 **Proof** of Lemma §22.06. is given in the lecture. □

§22.08 **Proposition (Upper bound).** Under Assumptions §19.02 and §22.02, $Y \in \mathcal{L}_5(\mathcal{U}_f)$ and $f \in \mathbb{L}_{\infty}(\lambda_{0,1})$ for $\mathbb{M}^{\mathbf{v}} \in \mathbb{N}$ as in (22.07) and $\text{pen}_{\bullet}^{f, \mathbf{v}} \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ as in (22.06) the partially data-driven OPE $\widehat{f}_{\bullet}^{\widehat{m}} = \widehat{f}_{\bullet} \mathbb{1}^{\widehat{m}} \in \ell_2(\mathbf{v}^2) \mathbb{1}^{\widehat{m}} \subseteq \ell_2(\mathbf{v}^2)$ of $f_{\bullet} \in \ell_2(\mathbf{v}^2)$ with data-driven dimension $\widehat{m} \in \llbracket \mathbb{M}^{\mathbf{v}} \rrbracket$ as in (22.08) satisfies

$$\begin{aligned} \mathcal{U}_f^{\otimes n}(\|\widehat{f}_{\bullet}^{\widehat{m}} - f_{\bullet}\|_{\mathbf{v}}^2) &\leq 192(1 + \mathcal{U}_f(Y^2)) \tau_{\mathbf{v}, u}^2 \min \{R_n^m(f, \mathbf{v}) : m \in \llbracket \mathbb{M}^{\mathbf{v}} \rrbracket\} \\ &\quad + C \tau_{\mathbf{v}, u}^2 (C_{\mathbf{v}}(x_{\xi, f}) + \mathbf{v}_1^2)(1 + \mathcal{U}_f(Y^2) + \mathcal{U}_f(|Y|^5)) n^{-1} \end{aligned} \quad (22.10)$$

for some universal numerical constant $C = 16C_{\text{tal}} \in [1, \infty)$ and $x_{\xi, f} = 6(\sigma_{\xi}^2 + \|f\|_{\mathbb{L}_{\infty}(\lambda_{0,1})}^2) \mathbf{v}_f^{-2} \tau_{\mathbf{v}, u}^{-2} \in \mathbb{R}^+$.

§22.09 **Proof** of **Proposition** §22.08. is given in the lecture. \square

§22.10 **Comment.** The oracle bound $R_n^\circ(f, \mathbf{v}) = R_n^{m_n^\circ}(f, \mathbf{v}) = \min \{R_n^m(f, \mathbf{v}): m \in \mathbb{N}\}$ (for details see **Reminder** §22.04) satisfies $nR_n^\circ(f, \mathbf{v}) \geq \|\mathbb{1}_v^{m_n^\circ}\|_v^2 \geq \mathbf{v}_1^2$. Consequently, the last upper bound in (22.10) and the oracle bound $R_n^\circ(f, \mathbf{v})$ coincide up to a constant $(192(1 + \mathcal{U}_f(Y^2))\tau_{v,u}^2 + C\tau_{v,u}^2(C_v(x_{\varepsilon,f})\mathbf{v}_1^{-2} + 1)(1 + \mathcal{U}_f(Y^2) + \mathcal{U}_f(|Y|^5)))$ provided the oracle dimension fulfils $m_n^\circ \in \llbracket M^p \rrbracket$. Therefore, we wish the upper bound M^p to be as large as possible. The next assertion shows that M^p as in (22.07) is a suitable choice for the upper bound. \square

§22.11 **Corollary.** Under the assumptions of **Proposition** §22.08 for each $n \in \mathbb{N}$ such that $R_n^\circ(f, \mathbf{v}) \leq \mathbf{v}_1^2$ and $m_n^\circ \leq n^{-2/3} \exp(\frac{n^{1/6}}{100})$ we have

$$\mathcal{U}_f^{\otimes n}(\|\widehat{f}_\bullet^{\widehat{m}} - f_\bullet\|_v^2) \leq KR_n^\circ(f, \mathbf{v})$$

and, hence up to the constant $K = 32(C_{\text{tal}} + 12)\tau_{v,u}^2(C_v(x_{\varepsilon,f})\mathbf{v}_1^{-2} + 1)(1 + \mathcal{U}_f(|Y|^5))$ the infeasible partially data-driven estimator $\widehat{f}_\bullet^{\widehat{m}}$ is **oracle optimal**.

§22.12 **Proof** of **Corollary** §22.11. is given in the lecture. \square

§22.13 **Illustration.** Consider the *trigonometric basis* as in **Illustration** §20.17 which satisfies Assumption §22.02 (os1), (os3) for all $\mathbf{a}_\bullet \in \ell_2$. In Table 01 [§12] (**Illustration** §12.19) the order of the rate $R_n^\circ(f, \mathbf{v})$ is depict for the two specifications (o) and (s). We note that we have $\mathbf{a}_\bullet \in \ell_2$ in case (o) for $a > 1/2$ while in case (s) for $a \in \mathbb{R}_v^+$. The sequence \mathbf{v} satisfies Assumption §22.02, i.e. (22.04), for $v \geq -1/2$. Moreover, the optimal dimension m_n° given in Table 01 [§12] satisfies $m_n^\circ \leq n^{-2/3} \exp(\frac{n^{1/6}}{100})$, and thus (under the above restrictions) the partially data-driven (hence not feasible) density estimator attains the oracle rate $R_n^\circ(f, \mathbf{v})$ up to the constant given in **Corollary** §22.11. \square

Estimated penalty sequence

§22.14 **Notation.** The penalty sequence $\text{pen}_\bullet^{f,v} \in (\mathbb{R}_v^+)^{\mathbb{N}}$ given in (22.06) still depends amongst others on characteristics of the unknown regression function f . More precisely, for $m \in \mathbb{N}$ the term $\text{pen}_m^{f,v}$ involves the quantity $\mathbf{v}_f^2 = 1 + \mathcal{U}_f(Y^2)$ which we eventually estimate without bias by $\widehat{\mathbf{v}}^2 := 1 + \widehat{\mathbb{P}}_n(Y^2)$ (keeping in mind that we identify Y and the coordinate map $\Pi_{\mathbb{R}}$). Therewith, let us introduce a fully data-driven sequence of penalties $\widehat{\text{pen}}_\bullet^v = (\widehat{\text{pen}}_m^v)_{m \in \mathbb{N}} \in (\mathbb{R}_v^+)^{\mathbb{N}}$ given by

$$\widehat{\text{pen}}_m^v := 2 \times 48 \widehat{\mathbf{v}}^2 \tau_{v,u}^2 n^{-1} \|\mathbb{1}_v^m\|_v^2 \quad \text{for each } m \in \mathbb{N} \quad \text{with } \widehat{\mathbf{v}}^2 := 1 + \widehat{\mathbb{P}}_n(Y^2) \quad (22.11)$$

and the upper bound $M^p \in \mathbb{N}$ given in (22.07) which are both *fully known* in advance. Considering the data-driven OSE $\widehat{f}_\bullet^{\widehat{m}} = \widehat{f}_\bullet \mathbb{1}_v^{\widehat{m}}$ with dimension parameter selected by

$$\widehat{m} := \arg \min \{ -\|\widehat{f}_\bullet^{\widehat{m}}\|_v + \widehat{\text{pen}}_m^v : m \in \llbracket M^p \rrbracket \} \quad (22.12)$$

we derive below an upper bound for its global \mathbf{v} -risk, $\mathcal{U}_f^{\otimes n}(\|\widehat{f}_\bullet^{\widehat{m}} - f_\bullet\|_v^2)$. \square

§22.15 **Proposition (Upper bound).** Under Assumptions §19.02 and §22.02, $Y \in \mathcal{L}_5(\mathcal{U}_f)$ and $f \in \mathbb{L}_\infty(\lambda_{0,1})$ for $M^p \in \mathbb{N}$ as in (22.07) and $\widehat{\text{pen}}_m^v \in (\mathbb{R}_v^+)^{\mathbb{N}}$ as in (22.11) the fully data-driven OPE $\widehat{f}_\bullet^{\widehat{m}} = \widehat{f}_\bullet \mathbb{1}_v^{\widehat{m}} \in \ell_2(\mathbf{v}) \mathbb{1}_v^{\widehat{m}} \subseteq \ell_2(\mathbf{v}^2)$ of $f_\bullet \in \ell_2(\mathbf{v}^2)$ with data-driven dimension $\widehat{m} \in \llbracket M^p \rrbracket$ as in (22.08) satisfies

$$\begin{aligned} \mathcal{U}_f^{\otimes n}(\|\widehat{f}_\bullet^{\widehat{m}} - f_\bullet\|_v^2) &\leq 288 \tau_{v,u}^2 (1 + \mathcal{U}_f(Y^2)) \min \{R_n^m(f, \mathbf{v}): m \in \llbracket M^p \rrbracket\} \\ &\quad + C \tau_{v,u}^2 (C_v(x_{\varepsilon,f}) + \mathbf{v}_1^2) (1 + \mathcal{U}_f(|Y|^5)) n^{-1} \end{aligned} \quad (22.13)$$

for $x_{\xi,f} = 6(\sigma_\xi^2 + \|f\|_{\mathbb{L}_\infty(\lambda_{0,1})}^2) \mathbf{v}_f^{-2} \tau_{\mathbf{v},u}^{-2} \in \mathbb{R}^+$ and some universal numerical constant $C = 3(16C_{\text{tal}} + 384) \in [1, \infty)$.

§22.16 **Proof** of **Proposition** §22.15. is given in the lecture. \square

§22.17 **Comment**. We shall stress that the last upper bound (22.13) in **Proposition** §22.15 (for the fully data-driven procedure) and the upper bound (22.10) in **Proposition** §22.08 (for the partially data-driven procedure) differ only in the constants. Thus, **Comment** §22.10 still applies here and the proof of the next results follows line by line their counterparts above. \square

§22.18 **Corollary**. Under the assumptions of **Proposition** §22.15 for each $n \in \mathbb{N}$ such that $R_n^\circ(f, \mathbf{v}) \leq \mathbf{v}_1^2$ and $m_n^\circ \leq n^{-2/3} \exp(\frac{n^{1/6}}{100})$ we have

$$\mathcal{U}_f^{\otimes n} (\|\widehat{f}_\cdot^m - f_\cdot\|_{\mathbf{v}}^2) \leq \text{KR}_n^\circ(f, \mathbf{v})$$

and, hence up to the constant $K = 5(16C_{\text{tal}} + 384) \tau_{\mathbf{v},u}^2 (C_{\mathbf{v}}(x_{\xi,f}) \mathbf{v}_1^{-2} + 1)(1 + \mathcal{U}_f(|Y|^5))$ the feasible fully data-driven estimator \widehat{f}_\cdot^m is *oracle optimal*.

§22.19 **Proof** of **Corollary** §22.18. is given in the lecture. \square

§22|01|02 Maximal global \mathbf{v} -risk

§22.20 **Assumption**. Consider weights $\mathbf{a}_\cdot, \mathbf{v}_\cdot \in (\mathbb{R}_{\setminus 0})^{\mathbb{N}}$ with $\mathbf{a}_\cdot \in \ell_\infty$ and $(\mathbf{a}\mathbf{v})_\cdot := (\mathbf{a}_j \mathbf{v}_j)_{j \in \mathbb{N}} = \mathbf{a}_\cdot \mathbf{v}_\cdot \in \ell_\infty$. We write $(\mathbf{a}\mathbf{v})_{(m)} := \|(\mathbf{a}\mathbf{v})_\cdot \mathbb{1}_\cdot^{m \perp}\|_{\ell_\infty} \in \mathbb{R}^+$ for each $m \in \mathbb{N}$. The weights $\mathbf{v}_\cdot \in (\mathbb{R}_{\setminus 0})^{\mathbb{N}}$ satisfy (18.05). The orthonormal system $(\mathbf{u}_j)_{j \in \mathbb{N}}$ in $\mathbb{L}_2(\lambda_{0,1})$ is **(os1) complete**, i.e an *orthonormal basis* in $\mathbb{L}_2(\lambda_{0,1})$ and as process $\mathbf{u}_\cdot^2 = (\mathbf{u}_j^2)_{j \in \mathbb{N}}$ on $([0, 1], \mathcal{B}_{[0,1]})$ satisfies **(os2)** $\|\mathcal{U}_N(\mathbf{a}_\cdot^2 \mathbf{u}_\cdot^2)\|_{\mathbb{L}_2(\lambda_{0,1})} \leq \tau_{\mathbf{a},u}^2$ and for all $m \in \mathbb{N}$, **(os3)** $\sup \{\|\mathbf{u}_\cdot(x) \mathbb{1}_\cdot^m\|_{\mathbf{v}}^2 : x \in [0, 1]\} \leq \tau_{\mathbf{v},u}^2 \|\mathbb{1}_\cdot^m\|_{\mathbf{v}}^2 \in \mathbb{R}^+$. for $\tau_{\mathbf{a},u}, \tau_{\mathbf{v},u} \in [1, \infty)$. \square

§22.21 **Reminder** (*Maximal global \mathbf{v} -risk*). Given Assumptions §19.02 and §22.20 we consider an OPE as in **Definition** §20.04. Here the observable noisy version $\widehat{f}_\cdot = f_\cdot + n^{-1/2} \boldsymbol{\varepsilon}_\cdot$ of the regression coefficients $f_\cdot = \mathbf{U}f \in \ell_2$ take the form of a *statistical direct problem* (see **Definition** §10.19) where the stochastic processes $\boldsymbol{\varepsilon}_\cdot \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B})^{\otimes n} \otimes 2^{\mathbb{N}}$ is given in **Definition** §19.08. Under Assumptions §19.02 and §22.20 in **Proposition** §16.16 an upper bound for a maximal global \mathbf{v} -risk of an OPE is shown over the set $\mathbb{F}_2^{\mathbf{a},r}$ given in (20.04) (**Lemma** §20.13). More precisely, the performance of the OPE $\widehat{f}_\cdot^m = \widehat{f}_\cdot \mathbb{1}_\cdot^m \in \ell_2(\mathbf{v}^2) \mathbb{1}_\cdot^m \subseteq \ell_2(\mathbf{v}^2)$ with dimension $m \in \mathbb{N}$ is measured by its maximal global \mathbf{v} -risk over the ellipsoid $\mathbb{F}_2^{\mathbf{a},r}$, that is

$$\mathcal{R}_n^{\mathbf{v}}[\widehat{f}_\cdot^m | \mathbb{F}_2^{\mathbf{a},r}] := \sup \{ \mathcal{U}_f^{\otimes n} (\|\widehat{f}_\cdot^m - f_\cdot\|_{\mathbf{v}}^2) : f \in \mathbb{F}_2^{\mathbf{a},r} \}.$$

As in (12.06) (**Proposition** §12.21) for $n, m \in \mathbb{N}$ setting $(\mathbf{a}\mathbf{v})_{(m)}^2 := \|(\mathbf{a}\mathbf{v})_\cdot \mathbb{1}_\cdot^{m \perp}\|_{\ell_\infty}^2$ and

$$\begin{aligned} R_n^m(\mathbf{a}_\cdot, \mathbf{v}_\cdot) &:= (\mathbf{a}\mathbf{v})_{(m)}^2 \vee n^{-1} \|\mathbb{1}_\cdot^m\|_{\mathbf{v}}^2, & m_n^* &:= \arg \min \{ R_n^m(\mathbf{a}_\cdot, \mathbf{v}_\cdot) : m \in \mathbb{N} \} \\ &\text{and } R_n^*(\mathbf{a}_\cdot, \mathbf{v}_\cdot) &:= R_n^{m_n^*}(\mathbf{a}_\cdot, \mathbf{v}_\cdot) = \min \{ R_n^m(\mathbf{a}_\cdot, \mathbf{v}_\cdot) : m \in \mathbb{N} \} \end{aligned} \quad (22.14)$$

by **Proposition** §20.15 under Assumptions §19.02 and §22.02 the maximal global \mathbf{v} -risk of an OPE $\widehat{f}_\cdot^{m_n^*}$ with optimally chosen dimension m_n^* as in (22.14) satisfies

$$\mathcal{R}_n^{\mathbf{v}}[\widehat{f}_\cdot^{m_n^*} | \mathbb{F}_2^{\mathbf{a},r}] \leq C R_n^*(\mathbf{a}_\cdot, \mathbf{v}_\cdot)$$

with $C = \sigma_\xi^2 + r^2 \tau_{\mathbf{a},u}^2 + r^2$. Moreover, due to **Proposition** §21.22 $R_n^*(\mathbf{a}_\cdot, \mathbf{v}_\cdot)$ provides (up to a constant) also a lower bound of the maximal global \mathbf{v} -risk over the ellipsoid $\mathbb{F}_2^{\mathbf{a},r}$ for any estimator. Consequently, (up to a constant) $R_n^*(\mathbf{a}_\cdot, \mathbf{v}_\cdot)$ is a minimax bound and $\widehat{f}_\cdot^{m_n^*}$ is minimax optimal. However, the optimal dimension m_n^* depends on $\mathbf{a}_\cdot \in (\mathbb{R}_{\setminus 0}^+)^{\mathbb{N}}$ characterising the ellipsoid $\mathbb{F}_2^{\mathbf{a},r}$. \square

§22.22 **Proposition (Upper bound).** Under Assumptions §19.02 and §22.20 and $\mathbb{P}^\xi \in \mathbb{P}_{\{0\} \times \mathbb{R}_v^+}$ with $\sigma_\xi^2 := \mathbb{P}^\xi(\text{id}_{\mathbb{R}}^2) \in \mathbb{R}_v^+$ and $\kappa_\xi^5 := \mathbb{P}^\xi(|\text{id}_{\mathbb{R}}|^5) \in \mathbb{R}_v^+$ for $M^p \in \mathbb{N}$ as in (22.07) and $\widehat{\text{pen}}_v^p \in (\mathbb{R}_v^+)^{\mathbb{N}}$ as in (22.11) the fully data-driven OPE $\widehat{f}_\cdot^{\widehat{m}} = \widehat{f}_\cdot \mathbb{1}_\cdot^{\widehat{m}} \in \ell_2(v_\cdot^2) \mathbb{1}_\cdot^{\widehat{m}} \subseteq \ell_2(v_\cdot^2)$ of $f_\cdot \in \ell_2(v_\cdot^2)$ with fully data-driven dimension $\widehat{m} \in \llbracket M^p \rrbracket$ as in (22.12) satisfies

$$\mathcal{R}_n^v[\widehat{f}_\cdot^{\widehat{m}} | \mathbb{F}_2^{\text{a,r}}] \leq (3r^2 + 288\tau_{v,u}^2(1 + \sigma_\xi^2 + r^2\tau_{a,u}^2)) \min \{R_n^m(f_\cdot, v_\cdot) : m \in \llbracket M^p \rrbracket\} \\ + C\tau_{v,u}^2(C_v(x_\xi) + v_1^2)(1 + \kappa_\xi^5 + r^5\tau_{a,u}^5)n^{-1} \quad (22.15)$$

for $x_\xi := 6(\sigma_\xi^2 + r^2\tau_{a,u}^2)\tau_{v,u}^{-2} \in \mathbb{R}^+$ and some universal numerical constant $C = 96(16C_{\text{tal}} + 384) \in [1, \infty)$.

§22.23 **Proof of Proposition §22.22.** is given in the lecture. \square

§22.24 **Comment.** The minimax bound $R_n^*(\mathbf{a}_\cdot, v_\cdot) = R_n^{m_n^*}(\mathbf{a}_\cdot, v_\cdot) = \min \{R_n^m(\mathbf{a}_\cdot, v_\cdot) : m \in \mathbb{N}\}$ (for details see **Reminder §18.12**) satisfies $nR_n^*(\mathbf{a}_\cdot, v_\cdot) \geq \|\mathbb{1}_\cdot^{m_n^*}\|_v^2 \geq v_1^2$. Consequently, the last upper bound in (22.15) and the minimax bound $R_n^*(\mathbf{a}_\cdot, v_\cdot)$ coincide up to a constant $3r^2 + 288\tau_{v,u}^2(1 + \sigma_\xi^2 + r^2\tau_{a,u}^2) + C\tau_{v,u}^2(C_v(x_\xi)v_1^{-2} + 1)(1 + \kappa_\xi^5 + r^5\tau_{a,u}^5)$ provided the minimax dimension fulfils $m_n^* \in \llbracket M_n \rrbracket$. Therefore, we wish the upper bound M^p to be as large as possible. The next assertion shows that M^p as in (22.07) is a suitable choice for the upper bound. \square

§22.25 **Corollary.** Under the assumptions of **Proposition §22.22** for each $n \in \mathbb{N}$ such that $R_n^*(\mathbf{a}_\cdot, v_\cdot) \leq v_1^2$ and $m_n^* \leq n^{-2/3} \exp(\frac{n^{1/6}}{100})$ we have

$$\mathcal{R}_n^v[\widehat{f}_\cdot^{\widehat{m}} | \mathbb{F}_2^{\text{a,r}}] \leq (3r^2 + 288\tau_{v,u}^2(1 + \sigma_\xi^2 + r^2\tau_{a,u}^2)) \min \{R_n^m(\mathbf{a}_\cdot, v_\cdot) : m \in \llbracket M^p \rrbracket\} \\ + 96(16C_{\text{tal}} + 384)\tau_{v,u}^2(C_v(x_\xi) + v_1^2)(1 + \kappa_\xi^5 + r^5\tau_{a,u}^5)n^{-1} \\ \leq KR_n^*(\mathbf{a}_\cdot, v_\cdot) \quad (22.16)$$

and, hence up to the constant $K := 3r^2 + C\tau_{v,u}^2(C_v(x_\xi)v_1^{-2} + 1)(1 + \kappa_\xi^5 + r^5\tau_{a,u}^5)$ with universal numerical constant $C = 99(16C_{\text{tal}} + 384) \in [1, \infty)$ the feasible data-driven estimator $\widehat{f}_\cdot^{\widehat{m}}$ is **minimax optimal**.

§22.26 **Proof of Corollary §22.25.** is given in the lecture. \square

§22.27 **Illustration.** Consider the **trigonometric basis** as in **Illustration §20.17** which satisfies Assumption §20.11 for all $\mathbf{a}_\cdot \in \ell_2$ (see **Illustration §20.17**). In Table 02 [§12] the order of the rate $R_n^*(\mathbf{a}_\cdot, v_\cdot)$ is depicted for the two cases **(o)** and **(s)** introduced in **Illustration §12.26**. We note that we have $\mathbf{a}_\cdot \in \ell_2$ in case **(o)** for $a > 1/2$ while in case **(s)** for $a \in \mathbb{R}_v^+$. The sequence v_\cdot satisfies Assumption §22.20, i.e. (22.04), for $v \geq -1/2$. Moreover, the optimal dimension m_n^* given in Table 02 [§12] satisfies $m_n^* \leq n^{-2/3} \exp(\frac{n^{1/6}}{100})$, and thus (under the above restrictions) the adaptive density estimator attains the minimax optimal rate $R_n^*(\mathbf{a}_\cdot, v_\cdot)$ up to the constant given in **Corollary §22.25**. \square

§22|02 Data-driven local estimation by Goldenshluger and Lepskij's method

§22.28 **Reminder.** The Bernstein inequality stated in the form of **Lemma §18.19** provides again our key argument in order to control the deviations of the reminder term. Let us briefly recall how we intend to apply the Bernstein inequality (see **Remark §18.21** for a similar approach). Reconsider the stochastic process $\psi_\cdot = (\psi_j(X, Y) := Y u_j(X))_{j \in \mathbb{N}} \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B}) \otimes 2^{\mathbb{N}}$ where $\psi_j \in \mathcal{L}_1(\mathcal{U}_j)$ for each $j \in \mathbb{N}$ and the OPE $\widehat{f}_\cdot^{\widehat{m}} = \widehat{f}_\cdot \mathbb{1}_\cdot^{\widehat{m}} \in \ell_2 \mathbb{1}_\cdot^{\widehat{m}}$ with dimension $m \in \mathbb{N}$ (**Definition §20.04**). $\widehat{f}_\cdot = \widehat{\mathbb{P}}_n \psi_\cdot = (\widehat{\mathbb{P}}_n \psi_j)_{j \in \mathbb{N}}$ are noisy versions (**Definition §15.08**) of the regression coefficients $f_\cdot = Uf = \mathcal{U}_j(\psi_\cdot) = (\mathcal{U}_j \psi_j = \mathcal{U}_j(Y u_j(X)))_{j \in \mathbb{N}}$ (see **Notation §19.07**). Clearly, $r_m :=$

$\phi_{\mathcal{U}_N}(\psi \mathbf{1}^m)$ is a $\mathcal{B}_{[0,1]} \otimes \mathcal{B}$ - \mathcal{B} -measurable function, where $\widehat{\mathbb{P}}_n(r_m) = \phi_{\mathcal{U}_N}(\widehat{\mathbb{P}}_n(\psi) \mathbf{1}^m) = \phi_{\mathcal{U}_N}(\widehat{f} \mathbf{1}^m)$ and $\mathcal{U}_f(r_m) = \phi_{\mathcal{U}_N}(\mathcal{U}_f(\psi) \mathbf{1}^m) = \phi_{\mathcal{U}_N}(f \mathbf{1}^m)$, and thus $\bar{r}_m = \widehat{\mathbb{P}}_n(r_m) - \mathcal{U}_f(r_m) = \phi_{\mathcal{U}_N}(\widehat{f}^m - f^m)$. We note that ψ and thus $r_m = \phi_{\mathcal{U}_N}(\psi \mathbf{1}^m)$ is not bounded. Therefore, we decompose $\psi = \psi^b + \psi^u$ into two parts $\psi^b, \psi^u \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B}) \otimes 2^{\mathbb{N}}$, a bounded and a remaining unbounded one. To be more precise, setting $r_m^b := \phi_{\mathcal{U}_N}(\psi^b \mathbf{1}^m) \in \mathcal{B}_{[0,1]} \otimes \mathcal{B}$ and $r_m^u := \phi_{\mathcal{U}_N}(\psi^u \mathbf{1}^m) \in \mathcal{B}_{[0,1]} \otimes \mathcal{B}$ let $\sup \{ |r_m^b(x, y)| : x \in [0, 1], y \in \mathbb{R} \} \in \mathbb{R}$ be satisfied. Introducing further $\bar{r}_m^b := \widehat{\mathbb{P}}_n(r_m^b) - \mathcal{U}_f(r_m^b)$ and $\bar{r}_m^u := \widehat{\mathbb{P}}_n(r_m^u) - \mathcal{U}_f(r_m^u)$ we evidently have $\bar{r}_m = \bar{r}_m^b + \bar{r}_m^u$ and hence

$$\begin{aligned} |\phi_{\mathcal{U}_N}(\widehat{f}^m - f^m)| &= |\bar{r}_m^b + \bar{r}_m^u|^2 \leq 2|\bar{r}_m^b|^2 + 2|\bar{r}_m^u|^2 \\ &= 2|\widehat{\mathbb{P}}_n(r_m^b) - \mathcal{U}_f(r_m^b)|^2 + 2|\widehat{\mathbb{P}}_n(r_m^u) - \mathcal{U}_f(r_m^u)|^2. \end{aligned} \quad (22.17)$$

Considering the first term on the right hand side provided that

$$\begin{aligned} \mathcal{U}_f(|r_m^b - \mathcal{U}_f(r_m^b)|^2) &= \mathcal{U}_f(|\phi_{\mathcal{U}_N}(\psi^b \mathbf{1}^m) - \mathcal{U}_f(\phi_{\mathcal{U}_N}(\psi^b \mathbf{1}^m))|^2) \leq v_{f,m}^2 \in \mathbb{R}^+, \\ \sup \{ |r_m^b(x, y)| : x \in [0, 1], y \in \mathbb{R} \} &\leq b_m \in \mathbb{R}^+, \text{ and hence } |r_m^b - \mathcal{U}_f(r_m^b)| \leq 2b_m, \end{aligned} \quad (22.18)$$

due to the Bernstein inequality (**Lemma** §18.19 (18.16)) we have

$$\mathcal{U}_f^{\otimes n} \left(\left(|n^{1/2} \bar{r}_m^b|^2 - (4v_{f,m}^2 + 32b_m^2 (\log K)n^{-1}) \log K \right)_+ \right) \leq 8K^{-1} \{v_{f,m}^2 + 16b_m^2 n^{-1}\}. \quad (22.19)$$

for any $K \in [1, \infty)$. □

§22|02|01 Local ϕ -risk

§22.29 **Assumption.** Let $\phi \in (\mathbb{R}_{>0})^{\mathbb{N}}$ and the orthonormal system $(u_j)_{j \in \mathbb{N}}$ in $\mathbb{L}_2(\lambda_{[0,1]})$ is **(os1)** complete and satisfies as process $\mathbf{u}_\cdot = (u_j)_{j \in \mathbb{N}}$ on $([0, 1], \mathcal{B}_{[0,1]})$ for $\tau_u \in [1, \infty)$ and for all $m \in \mathbb{N}$ satisfies **(os3)** $\sup \{ \|\mathbf{u}_\cdot(x) \mathbf{1}^m\|_{\ell_2}^2 : x \in [0, 1] \} \leq \tau_u^2 m \in \mathbb{R}^+$. □

§22.30 **Remark.** Keeping **Reminder** §22.28 in mind we define $\psi^b, \psi^u \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B}) \otimes 2^{\mathbb{N}}$ as

$$\begin{aligned} \psi^b(x, y) &= (\psi_j^b(x, y) := y \mathbf{1}_{[0, n^{1/q}]}(|y|) u_j(x))_{j \in \mathbb{N}} \text{ and} \\ \psi^u(x, y) &= (\psi_j^u(x, y) := y \mathbf{1}_{(n^{1/q}, \infty)}(|y|) u_j(x))_{j \in \mathbb{N}}, \quad x \in [0, 1], y \in \mathbb{R} \end{aligned}$$

where evidently $\psi^b + \psi^u = \psi$ and

$$|\bar{r}_m^b| = |\phi_{\mathcal{U}_N}(\psi^b \mathbf{1}^m)| = |Y \mathbf{1}_{[0, n^{1/q}]}(|Y|) \phi_{\mathcal{U}_N}(\mathbf{u}_\cdot(X) \mathbf{1}^m)| \leq |Y \phi_{\mathcal{U}_N}(\mathbf{u}_\cdot(X) \mathbf{1}^m)| = |\phi_{\mathcal{U}_N}(\psi \mathbf{1}^m)| = |r_m|.$$

We use in the sequel that under **Assumption** §22.29 **(os3)** for each $m \in \mathbb{N}$

$$\begin{aligned} \sup \{ |r_m^b(x, y)|^2 : x \in [0, 1], y \in \mathbb{R} \} &= \sup \{ |y \mathbf{1}_{[0, n^{1/q}]}(|y|) \phi_{\mathcal{U}_N}(\mathbf{u}_\cdot(x) \mathbf{1}^m)|^2 : x \in [0, 1], y \in \mathbb{R} \} \\ &\leq n^{1/3} \|\mathbf{1}^m\|_{\phi}^2 \sup \{ \|\mathbf{u}_\cdot(x) \mathbf{1}^m\|_{\ell_2}^2 : x \in [0, 1] \} \leq n^{1/3} \tau_u^2 m \|\mathbf{1}^m\|_{\phi}^2 =: b_m^2 \end{aligned} \quad (22.20)$$

by applying the Cauchy Schwarz inequality and moreover (see **Proof** §15.11)

$$\begin{aligned} \mathcal{U}_f(|r_m^b - \mathcal{U}_f(r_m^b)|^2) &\leq \mathcal{U}_f(|r_m^b|^2) = \mathcal{U}_f(|\phi_{\mathcal{U}_N}(\psi^b \mathbf{1}^m)|^2) =: v_{f,m}^2 \\ &\leq \mathcal{U}_f(|\phi_{\mathcal{U}_N}(\psi \mathbf{1}^m)|^2) \leq (\sigma_{\xi}^2 + \|f\|_{\mathbb{L}_{\infty}(\lambda_{[0,1]})}^2) \|\mathbf{1}^m\|_{\phi}^2 \in \mathbb{R}^+. \end{aligned} \quad (22.21)$$

exploiting (??) in **Proof** §19.11. Combining (18.19), (18.20) and (18.18) (**Remark** §18.21) we obtain

$$\begin{aligned} \mathcal{U}_f^{\otimes n} \left(\left(|n^{1/2} \bar{r}_m^b|^2 - (4v_{f,m}^2 + 32b_m^2 (\log K)n^{-1}) \log K \right)_+ \right) \\ \leq 8K^{-1} \{ \sigma_{\xi}^2 + \|f\|_{\mathbb{L}_{\infty}(\lambda_{[0,1]})}^2 + 16\tau_u^2 m n^{-2/3} \} \|\mathbf{1}^m\|_{\phi}^2 \end{aligned} \quad (22.22)$$

for any $m \in \mathbb{N}$ and $K \in [1, \infty)$. □

§22.31 **Reminder (Local oracle ϕ -risk).** Given Assumptions §19.02 and §22.29 we consider an OPE as in Definition §20.04. Here the observable noisy version $\widehat{f}_\bullet = f_\bullet + n^{-1/2}\boldsymbol{\varepsilon}_\bullet$ of the regression coefficients $f_\bullet = Uf \in \ell_2$ take the form of a *statistical direct problem* (see Definition §10.19) where the stochastic processes $\boldsymbol{\varepsilon}_\bullet \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B})^{\otimes n} \otimes 2^{\mathbb{N}}$ is given in Definition §19.08. Under Assumptions §19.02 and §22.29 Assumptions §15.02 and §18.22, (and hence Assumption §19.05 and $\phi_\bullet \in (\mathbb{R}_{\setminus 0})^{\mathbb{N}}$) and $f_\bullet \in \text{dom}(\phi_{\nu_\bullet})$ in §20.24 an *oracle inequality* for the local ϕ -risk of the OPE's is shown. More precisely, as in (20.06) (Proposition §20.22) for all $n, m \in \mathbb{N}$ setting

$$\begin{aligned} R_n^m(f, \phi) &:= |\phi_{\nu_\bullet}(f_\bullet \mathbb{1}_\bullet^{m\perp})|^2 + n^{-1} \|\mathbb{1}_\bullet^m\|_\phi^2, \quad m_n^\circ := \arg \min \{R_n^m(f, \phi) : m \in \mathbb{N}\} \\ \text{and } R_n^\circ(f, \phi) &:= R_n^{m_n^\circ}(f, \phi) = \min \{R_n^m(f, \phi) : m \in \mathbb{N}\}. \end{aligned} \quad (22.23)$$

and assuming $v_f := \max(\sigma_\xi^{-2}, \sigma_\xi^2 + \|f\|_{\mathbb{L}_\infty(\lambda_{0,1})}^2) \in \mathbb{R}_{\setminus 0}^+$ due to Property §20.09 the (infeasible) OPE $\widehat{f}_\bullet^{m_n^\circ} = \widehat{f}_\bullet \mathbb{1}_\bullet^{m_n^\circ} \in \ell_2 \mathbb{1}_\bullet^{m_n^\circ} \subseteq \text{dom}(\phi_{\nu_\bullet})$ with oracle dimension m_n° as in (22.23) satisfies

$$\begin{aligned} v_f^{-1} R_n^\circ(f, \phi) &\leq \inf_{m \in \mathbb{N}} \mathcal{U}_f^{\otimes n}(|\phi_{\nu_\bullet}(\widehat{f}_\bullet^m - f_\bullet)|^2) \leq \mathcal{U}_f^{\otimes n}(|\phi_{\nu_\bullet}(\widehat{f}_\bullet^{m_n^\circ} - f_\bullet)|^2) \\ &\leq v_f R_n^\circ(f, \phi) \leq v_f^2 \inf_{m \in \mathbb{N}} \mathcal{U}_f^{\otimes n}(|\phi_{\nu_\bullet}(\widehat{f}_\bullet^m - f_\bullet)|^2), \end{aligned}$$

and hence it is *oracle optimal* (with constant v_f^2). \square

Partially known penalty sequence

§22.32 **Notation.** Consider first a sequence of penalties $\text{pen}_m^{f, \phi} = (\text{pen}_m^{f, \phi})_{m \in \mathbb{N}} \in (\mathbb{R}_{\setminus 0}^+)^{\mathbb{N}}$ given by

$$\begin{aligned} \text{pen}_m^{f, \phi} &:= 24n^{-1}(v_{f,m}^2 + 8b_m^2(\log K_m)n^{-1})(\log K_m) \quad \text{with } v_{f,m}^2 := \mathcal{U}_f(|\phi_{\nu_\bullet}(\psi_\bullet^b \mathbb{1}_\bullet^m)|^2), \\ b_m^2 &:= \tau_u^2 n^{1/3} m \|\mathbb{1}_\bullet^m\|_\phi^2, \quad \text{and } K_m := (1 \vee \|\mathbb{1}_\bullet^m\|_\phi^2) m^3 \in [1, \infty) \quad \text{for each } m \in \mathbb{N}, \end{aligned} \quad (22.24)$$

which is obviously only *partially known* in advance, and fully known upper bound

$$M^\phi := \max \{m \in \mathbb{N} : m \|\mathbb{1}_\bullet^m\|_\phi^2 \leq n^2 \phi_1^2\} \in \mathbb{N} \quad (22.25)$$

where the defining set is not empty and finite (i.e. $M^\phi \leq n^2$). Considering the data-driven OSE $\widehat{f}_\bullet^m = \widehat{f}_\bullet \mathbb{1}_\bullet^m$ with dimension parameter selected by Goldenshluger and Lepskij's method

$$\begin{aligned} \widehat{m} &:= \arg \min \{\text{contr}_m^{f, \phi} + \text{pen}_m^{f, \phi} : m \in \llbracket M^\phi \rrbracket\} \quad \text{and} \\ \text{contr}_m^{f, \phi} &:= \max \{(|\phi_{\nu_\bullet}(\widehat{f}_\bullet^j - \widehat{f}_\bullet^m)|^2 - \text{pen}_j^{f, \phi} - \text{pen}_m^{f, \phi})_+ : j \in \llbracket m, M^\phi \rrbracket\}, \quad m \in \llbracket M^\phi \rrbracket. \end{aligned} \quad (22.26)$$

Moreover, studying a ϕ -error the bias term introduced in (14.31) becomes

$$\text{bias}_m(f, \phi) = \sup \{|\phi_{\nu_\bullet}(f_\bullet^j - f_\bullet^m)| = |\phi_{\nu_\bullet}(f_\bullet \mathbb{1}_\bullet^{[m, j]})| : j \in \llbracket m, \infty \rrbracket\} \quad \forall m \in \mathbb{N}.$$

If $f_\bullet \in \text{dom}(\phi_{\nu_\bullet})$ and hence $\nu_\bullet(|\phi_\bullet f_\bullet|) \in \mathbb{R}$ then $\text{bias}_m(f, \phi) \leq \nu_\bullet(|\phi_\bullet f_\bullet| \mathbb{1}_\bullet^{m\perp}) = o(1)$ as $m \rightarrow \infty$ by dominated convergence. Considering the data-driven OSE $\widehat{f}_\bullet^m = \widehat{f}_\bullet \mathbb{1}_\bullet^m$ with dimension parameter \widehat{m} selected as in (22.26) with penalty sequence $\text{pen}_m^{f, \phi}$ given in (22.24) and upper bound $M^\phi \in \mathbb{N}$ as in (22.25) we derive below an upper bound for its local ϕ -risk, $\mathcal{U}_f^{\otimes n}(|\phi_{\nu_\bullet}(\widehat{f}_\bullet^{\widehat{m}} - f_\bullet)|^2)$. \square

§22.33 **Lemma.** Under Assumptions §19.02 and §22.29, $f \in \mathbb{L}_\infty(\lambda_{0,1})$ and $Y \in \mathcal{L}_{14}(\mathcal{U}_f)$ for $\text{pen}_m^{f, \phi} \in (\mathbb{R}_{\setminus 0}^+)^{\mathbb{N}}$ as in (22.24) and $M^\phi \in \mathbb{N}$ as in (22.25) we have

$$\begin{aligned} \mathcal{U}_f^{\otimes n} \left(\max \{(|\phi_{\nu_\bullet}(\widehat{f}_\bullet^{\widehat{m}} - f_\bullet^m)|^2 - \text{pen}_m^{f, \phi}/3) : m \in \llbracket M^\phi \rrbracket\} \right) \\ \leq \{28\sigma_\xi^2 + 28\|f\|_{\mathbb{L}_\infty(\lambda_{0,1})}^2 + 448\tau_u^2 n^{-2/3} + 2\phi_1^2 \tau_u^2 \mathcal{U}_f(Y^{14})\} n^{-1} \end{aligned} \quad (22.27)$$

§22.34 **Proof of Lemma** §22.33. is given in the lecture. \square

§22.35 **Proposition (Upper bound).** Under Assumptions §19.02 and §22.29, $f \in \mathbb{L}_{\infty}(\lambda_{0,1})$ and $Y \in \mathcal{L}_{14}(\mathcal{U}_f)$ for $\text{pen}^{f,\phi} \in (\mathbb{R}_0^+)^{\mathbb{N}}$ as in (22.24) and $M^\phi \in \mathbb{N}$ as in (22.25) the OPE $\widehat{f}^{\widehat{m}} = \widehat{f} \mathbf{1}^{\widehat{m}} \in \ell_2 \mathbf{1}^{\widehat{m}} \subseteq \text{dom}(\phi_{\mathbb{N}})$ of $f_* \in \text{dom}(\phi_{\mathbb{N}})$ with partially data-driven dimension $\widehat{m} \in \llbracket M^\phi \rrbracket$ as in (22.26) satisfies for all $n \in \mathbb{N}$

$$\begin{aligned} \mathcal{U}_f^{\otimes n} (|\phi_{\mathbb{N}}(\widehat{f}^{\widehat{m}} - f)|^2) &\leq 128(\sigma_\xi^2 + \|f\|_{\mathbb{L}_{\infty}(\lambda_{0,1})}^2 + 8\tau_u^2) \\ &\quad \times \min \left\{ \text{bias}_m^2(\mathbb{P}, \phi) + n^{-1} \|\mathbf{1}^m\|_\phi^2 (\log K_m) (1 \vee (\log K_m) m n^{-2/3}); m \in \llbracket M \rrbracket \right\} \\ &\quad + 56(14\sigma_\xi^2 + 14\|f\|_{\mathbb{L}_{\infty}(\lambda_{0,1})}^2 + 224\tau_u^2 n^{-2/3} + \phi_1^2 \tau_u^2 \mathcal{U}_f(Y^{14})) n^{-1}. \end{aligned} \quad (22.28)$$

§22.36 **Proof of Proposition** §22.35. is given in the lecture. \square

§22.37 **Comment.** Let us compare the dominating part of the upper bound given in (22.28), that is

$$\min \left\{ \text{bias}_m^2(f, \phi) + n^{-1} \|\mathbf{1}^m\|_\phi^2 (\log K_m) (1 \vee (\log K_m) m n^{-2/3}); m \in \llbracket M^\phi \rrbracket \right\} \quad (22.29)$$

with the oracle bound $R_n^\circ(f, \phi) = \min \left\{ |\phi_{\mathbb{N}}(f^m - f)|^2 + n^{-1} \|\mathbf{1}^m\|_\phi^2 : m \in \mathbb{N} \right\}$ (for details see **Reminder** §22.31). In (22.29) we face eventually a deterioration by three sources. First, we generally have $\text{bias}_m(f, \phi) \geq |\phi_{\mathbb{N}}(f^m - f)|$, but note that for $f_* \phi \in (\mathbb{R}^+)^{\mathbb{N}}$ equality holds, that is

$$\text{bias}_m(f, \phi) = \sup \left\{ \nu_{\mathbb{N}}(\phi_* f \mathbf{1}^{\lfloor m, \infty \rfloor}) : j \in \llbracket m, \infty \rrbracket \right\} = \nu_{\mathbb{N}}(\phi_* f \mathbf{1}^{m \perp}) = |\phi_{\mathbb{N}}(f^m - f)|$$

for all $m \in \mathbb{N}$. Secondly, the variance term features an additional factor $(\log K_m)(1 \vee (\log K_m) m n^{-2/3})$, and finally the upper bound M^ϕ might impose an additional deterioration. The next assertion shows that M^ϕ is a suitable choice for the upper bound. Moreover, we estimate the bias term by $\text{bias}_m(f, \phi) \leq \nu(|\phi_* f \mathbf{1}^{m \perp}|)$ where equality holds whenever $f_* \phi \in (\mathbb{R}^+)^{\mathbb{N}}$. \square

§22.38 **Corollary.** For $n, m \in \mathbb{N}$ we set

$$\begin{aligned} R_n^m(f, \phi) &:= (\nu_{\mathbb{N}}(|\phi_* f \mathbf{1}^{m \perp}|))^2 \\ &\quad + (1 + (\log \|\mathbf{1}^m\|_\phi^2)_+ + \log m) (1 + ((\log \|\mathbf{1}^m\|_\phi^2)_+ + \log m) m n^{-2/3}) n^{-1} \|\mathbf{1}^m\|_\phi^2, \\ m^\diamond &:= \arg \min \left\{ R_n^m(f, \phi) : m \in \mathbb{N} \right\} \quad \text{and} \\ R_n^\diamond(f, \phi) &:= R_n^{m^\diamond}(f, \phi) = \min \left\{ R_n^m(f, \phi) : m \in \mathbb{N} \right\}. \end{aligned} \quad (22.30)$$

Under the assumptions of **Proposition** §22.35 for each $n \in \mathbb{N}$ such that $m^\diamond \in \llbracket M^\phi \rrbracket$ we have

$$\begin{aligned} \mathcal{U}_f^{\otimes n} (|\phi_{\mathbb{N}}(\widehat{f}^{\widehat{m}} - f)|^2) &\leq 1152(\sigma_\xi^2 + \|f\|_{\mathbb{L}_{\infty}(\lambda_{0,1})}^2 + 8\tau_u^2) R_n^\diamond(f, \phi) \\ &\quad + 56(14\sigma_\xi^2 + 14\|f\|_{\mathbb{L}_{\infty}(\lambda_{0,1})}^2 + 224\tau_u^2 n^{-2/3} + \phi_1^2 \tau_u^2 \mathcal{U}_f(Y^{14})) n^{-1} \\ &\leq 9216((\sigma_\xi^2 + \|f\|_{\mathbb{L}_{\infty}(\lambda_{0,1})}^2)(1 + \phi_1^{-2}) + \tau_u^2(1 + 2\phi_1^{-2} n^{-2/3} + \mathcal{U}_f(Y^{14}))) R_n^\diamond(f, \phi). \end{aligned} \quad (22.31)$$

§22.39 **Proof of Proof** §22.39. is given in the lecture. \square

§22.40 **Comment.** The data-driven bound $R_n^\diamond(f, \phi)$ compared to the oracle bound $R_n^\circ(f, \phi)$ features a deterioration of the variance term at least by a logarithmic factor. The appearance of the logarithmic factor within the bound is a known fact in the context of local estimation (cf. Laurent et al. [2008] who consider model selection given direct Gaussian observations). Brown and Low [1996] show that it is unavoidable in the context of nonparametric Gaussian regression and hence it is widely considered as an acceptable price for adaptation. \square

§22.41 **Illustration.** We illustrate the last results considering the two specifications **(o)** and **(s)** given in Table 03 [§12] (**Illustration** §12.40). We restrict ourselves to the case $\phi \notin \ell_2$ only.

Table 01 [§22]

Order of the oracle rate $R_n^\circ(f, \phi)$ and the data-driven rate $R_n^\circ(f, \phi)$ as $n \rightarrow \infty$

$(j \in \mathbb{N})$ $\phi = j^{v-1/2}$	$(a \in \mathbb{R}_v^+)$ f_j	(squared bias) $(\mathcal{U}_N(\phi f \cdot \mathbb{1}^m ^2))^2$	(variance) $\ \mathbb{1}^m\ _\phi^2$	M^ϕ	m°	$R_n^\circ(f, \phi)$	$R_n^\circ(f, \phi)$		
(o)	$v \in (0, a)$	$m^{-2(a-v)}$	m^{2v}	$n^{\frac{2}{2v+1}}$	$(\frac{n}{\log n})^{\frac{1}{2a}}$	$n^{-\frac{(a-v)}{a}}$	$(\frac{\log n}{n})^{\frac{(a-v)}{a}}$		
	$a \in (3/4, \infty)$							$(\frac{n^{5/6}}{\log n})^{\frac{1}{a+1/2}}$	$(\frac{\log n}{n^{5/6}})^{\frac{2(a-v)}{a+1/2}}$
	$a \in (0, 3/4]$	m^{-2a}	$\log m$	$\frac{n^2}{\log n}$	$(\frac{n}{(\log n)^2})^{\frac{1}{2a}}$	$\frac{\log n}{n}$	$\frac{(\log n)^2}{n}$		
	$v = 0$							$(\frac{n^{2/3}}{(\log n)^3})^{\frac{1}{2a+1}}$	$(\frac{(\log n)^3}{n^{5/3}})^{\frac{a}{a+1/2}}$
	$a \in (3/4, \infty)$								
$a \in (0, 3/4]$									
(s)	$v \in \mathbb{R}_v^+$	$m^{(1-2(a-v))_+} e^{-2m^{2a}}$	m^{2v}	$n^{\frac{2}{2v+1}}$	$(\log n)^{\frac{1}{2a}}$	$\frac{(\log n)^{\frac{v}{a}}}{n}$	$\frac{(\log n)^{\frac{v}{a}} (\log \log n)}{n}$		
	$v = 0$	$m^{(1-2a)_+} e^{-2m^{2a}}$	$\log m$	$\frac{n^2}{\log n}$	$(\log n)^{\frac{1}{2a}}$	$\frac{\log \log n}{n}$	$\frac{(\log \log n)^2}{n}$		

We note that in Table 01 [§22] the order of the oracle rate $R_n^\circ(f, \phi)$ and the data-driven rate $R_n^\circ(f, \phi)$ is depict for $v \geq 0$ only. In case $v < 0$ we have $\phi \in \ell_2$. Moreover, in case **(s)** for $a \in \mathbb{R}_v^+$ and **(o)** for $a \in (3/4, \infty)$ the rate $R_n^\circ(f, \phi)$ features only an additional logarithmic factor compared with the oracle rate $R_n^\circ(f, \phi)$. \square

Estimated penalty sequence

§22.42 **Notation.** The penalty sequence $\text{pen}_m^{f,v} \in (\mathbb{R}_v^+)^{\mathbb{N}}$ given in (22.24) still depends on characteristics of the unknown regression function f . More precisely, for $m \in \mathbb{N}$ the term $\text{pen}_m^{f,v}$ involves the quantity $\mathcal{V}_{f,m}^2 = \mathcal{U}_f(|\phi \mathcal{U}_N(\psi^b \mathbb{1}^m)|^2)$ which we eventually estimate without bias by $\widehat{\mathcal{V}}_m^2 := \widehat{\mathbb{P}}_n(|\phi \mathcal{U}_N(\psi^b \mathbb{1}^m)|^2)$. Based on this estimator let us introduce a fully data-driven sequence of penalties $\widehat{\text{pen}}_\bullet^\phi = (\widehat{\text{pen}}_m^\phi)_{m \in \mathbb{N}} \in (\mathbb{R}_v^+)^{\mathbb{N}}$ given by

$$\begin{aligned} \widehat{\text{pen}}_m^\phi &:= 24n^{-1} (2\widehat{\mathcal{V}}_m^2 + 3 \times 8b_m^2 (\log K_m) n^{-1}) (\log K_m) \quad \text{with} \quad \widehat{\mathcal{V}}_m^2 := \widehat{\mathbb{P}}_n(|\phi \mathcal{U}_N(\psi^b \mathbb{1}^m)|^2), \\ b_m^2 &:= \tau_u^2 n^{1/3} m \|\mathbb{1}^m\|_\phi^2, \quad \text{and} \quad K_m := (1 \vee \|\mathbb{1}^m\|_\phi^2) m^3 \in [1, \infty) \quad \text{for each } m \in \mathbb{N}, \end{aligned} \quad (22.32)$$

which is now *fully known* in advance, and fully known upper bound $M^\phi \in \mathbb{N}$ defined in (22.25). Considering the data-driven OSE $\widehat{f}^{\widehat{m}} = \widehat{f} \mathbb{1}^{\widehat{m}}$ with dimension parameter selected by Goldenshluger and Lepskij’s method

$$\begin{aligned} \widehat{m} &:= \arg \min \{ \widehat{\text{contr}}_m^\phi + \widehat{\text{pen}}_m^\phi : m \in \llbracket M^\phi \rrbracket \} \quad \text{and} \\ \widehat{\text{contr}}_m^\phi &:= \max \{ (|\phi \mathcal{U}_N(\widehat{\mathbb{P}}^j - \widehat{\mathbb{P}}^m)|^2 - \widehat{\text{pen}}_j^\phi - \widehat{\text{pen}}_m^\phi)_+ : j \in \llbracket m, M^\phi \rrbracket \}, \quad m \in \llbracket M^\phi \rrbracket \end{aligned} \quad (22.33)$$

we derive below an upper bound for its local ϕ -risk, $\mathcal{U}_f^{\otimes n}(|\phi \mathcal{U}_N(\widehat{f}^{\widehat{m}} - f)|^2)$. \square

§22.43 **Lemma.** Under Assumptions §19.02 and §22.29 and $f \in \mathbb{L}_\infty(\lambda_{(0,1)})$ for $\text{pen}_\bullet^{f,\phi}, \widehat{\text{pen}}_\bullet^\phi \in (\mathbb{R}_v^+)^{\mathbb{N}}$ as in (22.24) and (22.32), respectively, and for any $M \in \mathbb{N}$ we have

$$\mathcal{U}_f^{\otimes n} \left(\max \{ (\text{pen}_j^{f,\phi} - \widehat{\text{pen}}_j^\phi)_+ : j \in \llbracket M \rrbracket \} \right) \leq 80 \{ \sigma_\xi^2 + \|f\|_{\mathbb{L}_\infty(\lambda_{(0,1)})}^2 + 6\tau_u^2 n^{-2/3} \} n^{-1}. \quad (22.34)$$

§22.44 **Proof of Lemma** §22.43. is given in the lecture. \square

§22.45 **Proposition (Upper bound)**. Under Assumptions §19.02 and §22.29, $Y \in \mathcal{L}_{14}(\mathcal{U}_f)$ and $f \in \mathbb{L}_\infty(\lambda_{0,1})$ for $\widehat{\text{pen}}_\phi \in (\mathbb{R}_0^+)^{\mathbb{N}}$ as in (22.32) and for $M^\phi \in \mathbb{N}$ as in (22.25) the OPE $\widehat{f}_\cdot^{\widehat{m}} = \widehat{f}_\cdot \mathbf{1}_\cdot^{\widehat{m}} \in \ell_2 \mathbf{1}_\cdot^{\widehat{m}} \subseteq \text{dom}(\phi_{\nu_n})$ of $f_\cdot \in \text{dom}(\phi_{\nu_n})$ with fully data-driven dimension $\widehat{m} \in \llbracket M^\phi \rrbracket$ as in (22.33) satisfies for all $n \in \mathbb{N}$

$$\begin{aligned} \mathcal{U}_f^{\otimes n} (|\phi_{\nu_n}(\widehat{f}_\cdot^{\widehat{m}} - f_\cdot)|^2) &\leq 224(\sigma_\xi^2 + \|f\|_{\mathbb{L}_\infty(\lambda_{0,1})}^2 + 12\tau_u^2) \\ &\quad \times \min \{ \text{bias}_m^2(f, \phi) + n^{-1} \|\mathbf{1}_\cdot^m\|_\phi^2 (\log K_m) (1 \vee (\log K_m) m n^{-2/3}) : m \in \llbracket M^\phi \rrbracket \} \\ &\quad + 72(40\sigma_\xi^2 + 40\|f\|_{\mathbb{L}_\infty(\lambda_{0,1})}^2 + 240\tau_u^2 n^{-2/3} + \phi_1^2 \tau_u^2 \mathcal{U}_f(Y^{14})) n^{-1}. \end{aligned} \quad (22.35)$$

§22.46 **Proof of Proposition** §22.45. is given in the lecture. \square

§22.47 **Comment**. We shall stress that the last upper bound (22.35) in **Proposition** §22.45 (for the fully data-driven procedure) and the upper bound (22.28) in **Proposition** §22.35 (for the partially data-driven procedure) differ only in the numerical constants. Thus, thus the proof of the next result follows line by line their counterparts above. \square

§22.48 **Corollary**. Under the assumptions of **Proposition** §22.45 for each $n \in \mathbb{N}$ such that $m^\circ \in \llbracket M^\phi \rrbracket$ we have

$$\begin{aligned} \mathcal{U}_f^{\otimes n} (|\phi_{\nu_n}(\widehat{f}_\cdot^{\widehat{m}} - f_\cdot)|^2) &\leq 2016(\sigma_\xi^2 + \|f\|_{\mathbb{L}_\infty(\lambda_{0,1})}^2 + 12\tau_u^2) R_n^\circ(f, \phi) \\ &\quad + 72(40\sigma_\xi^2 + 40\|f\|_{\mathbb{L}_\infty(\lambda_{0,1})}^2 + 240\tau_u^2 n^{-2/3} + \phi_1^2 \tau_u^2 \mathcal{U}_f(Y^{14})) n^{-1} \\ &\leq 2016((\sigma_\xi^2 + \|f\|_{\mathbb{L}_\infty(\lambda_{0,1})}^2)(1 + \phi_1^{-2}) + 12\tau_u^2(1 + \phi_1^{-2} n^{-2/3} + \mathcal{U}_f(Y^{14}))) R_n^\circ(f, \phi). \end{aligned} \quad (22.36)$$

§22.49 **Proof of Proof** §22.49. is given in the lecture. \square

§22.50 **Comment**. The fullay data-driven bound $R_n^\circ(f, \phi)$ equals up to the numerical constants the bound in the partially known case. Therefore, the **Comment** §22.40 and the **Illustration** §22.41 apply here equally. \square

§22|02|02 Maximal local ϕ -risk

§22.51 **Assumption**. Consider $\phi_\cdot, \mathbf{a}_\cdot \in (\mathbb{R}_{0,1})^{\mathbb{N}}$ with $\mathbf{a}_\cdot \in \ell_\infty$ and $(\mathbf{a}\phi)_\cdot := (\mathbf{a}_j \phi_j)_{j \in \mathbb{N}} = \mathbf{a}_\cdot \phi_\cdot \in \ell_2$, and hence $\|\mathbf{a}_\cdot \mathbf{1}_\cdot^{m \perp}\|_\phi = \|(\mathbf{a}\phi)_\cdot \mathbf{1}_\cdot^{m \perp}\|_{\ell_2} = o(1)$ as $m \rightarrow \infty$. The orthonormal system $(\mathbf{u}_j)_{j \in \mathbb{N}}$ in $\mathbb{L}_2(\lambda_{0,1})$ is **(os1) complete**, i.e an orthonormal basis in $\mathbb{L}_2(\lambda_{0,1})$ and as process $\mathbf{u}_\cdot^2 = (\mathbf{u}_j^2)_{j \in \mathbb{N}}$ on $([0, 1], \mathcal{B}_{[0,1]})$ satisfies **(os2)** $\|\nu_n(\mathbf{a}_\cdot^2 \mathbf{u}_\cdot^2)\|_{\mathbb{L}_2(\lambda_{0,1})} \leq \tau_{\mathbf{a}, \mathbf{u}}^2$ and for all $m \in \mathbb{N}$, **(os3)** $\sup \{ \|\mathbf{u}_\cdot(x) \mathbf{1}_\cdot^m\|_{\ell_2}^2 : x \in [0, 1] \} \leq \tau_u^2 \|\mathbf{1}_\cdot^m\|_{\mathbb{U}}^2 \in \mathbb{R}^+$. for $\tau_{\mathbf{a}, \mathbf{u}}, \tau_u \in [1, \infty)$. \square

§22.52 **Remark**. Under Assumption §22.51 considering the set \mathbb{F}_2^{ar} of regression functions in $\mathbb{L}_2(\lambda_{0,1})$ defined in (20.04) we have $\|f\|_{\mathbb{L}_\infty(\lambda_{0,1})} \leq r \tau_{\mathbf{a}, \mathbf{u}}$ for all $f \in \mathbb{F}_2^{\text{ar}}$ due to **Lemma** §20.13. Consequently, given in addition Assumption §19.02 all assumptions of **Proposition** §22.45 are satisfied. \square

§22.53 **Reminder (Maximal local ϕ -risk)**. Given Assumptions §19.02 and §22.51 we consider an OPE as in **Definition** §20.04. Here the observable noisy version $\widehat{f}_\cdot = f_\cdot + n^{-1/2} \boldsymbol{\varepsilon}_\cdot$ of the regression coefficients $f_\cdot = \mathbf{U} f \in \ell_2$ take the form of a *statistical direct problem* (see **Definition** §10.19) where the stochastic processes $\boldsymbol{\varepsilon}_\cdot \in (\mathcal{B}_{[0,1]} \otimes \mathcal{B})^{\otimes n} \otimes 2^{\mathbb{N}}$ is given in **Definition** §19.08. Under Assumptions §19.02 and §22.51 (and hence Assumption §20.26) in **Proposition** §20.29 an upper bound for a maximal local ϕ -risk of an OPE is shown over the set \mathbb{F}_2^{ar} given in (20.04) (**Lemma** §20.13)

More precisely, the performance of the OPE $\widehat{f}_\cdot^m = \widehat{f}_\cdot \mathbf{1}_\cdot^m \in \ell_2 \mathbf{1}_\cdot^m \subseteq \text{dom}(\phi_{\nu_K})$ with dimension $m \in \mathbb{N}$ is measured by its maximal global ϕ -risk over the ellipsoid $\mathbb{F}_2^{\text{a,r}}$, that is

$$\mathcal{R}_n^\phi[\widehat{f}_\cdot^m | \mathbb{F}_2^{\text{a,r}}] := \sup \{ \mathcal{U}_f^{\otimes n}(|\phi_{\nu_K}(\widehat{f}_\cdot^m - f_\cdot)|^2) : f \in \mathbb{F}_2^{\text{a,r}} \}.$$

As in (12.13) (**Proposition** §12.42) for all $n, m \in \mathbb{N}$ setting

$$\begin{aligned} \mathbb{R}_n^m(\mathbf{a}, \phi) &:= \|\mathbf{a} \cdot \mathbf{1}_\cdot^{m \perp}\|_\phi^2 + n^{-1} \|\mathbf{1}_\cdot^m\|_\phi^2, \quad m_n^* := \arg \min \{ \mathbb{R}_n^m(\mathbf{a}, \phi) : m \in \mathbb{N} \} \\ \text{and } \mathbb{R}_n^*(\mathbf{a}, \phi) &:= \mathbb{R}_n^{m_n^*}(\mathbf{a}, \phi) = \min \{ \mathbb{R}_n^m(\mathbf{a}, \phi) : m \in \mathbb{N} \}. \end{aligned} \quad (22.37)$$

by **Proposition** §20.29 under Assumptions §19.02 and §22.51 the maximal local ϕ -risk of an OPE $\widehat{f}_\cdot^{m_n^*}$ with optimally chosen dimension m_n^* as in (22.37) satisfies

$$\mathcal{R}_n^\phi[\widehat{f}_\cdot^{m_n^*} | \mathbb{F}_2^{\text{a,r}}] \leq C \mathbb{R}_n^*(\mathbf{a}, \phi)$$

with $C = \sigma_\xi^2 + r^2 \tau_{\text{a,u}}^2$. Moreover, due to **Proposition** §21.09 $\mathbb{R}_n^*(\mathbf{a}, \phi)$ provides (up to a constant) also a lower bound of the maximal global ϕ -risk over the ellipsoid $\mathbb{F}_2^{\text{a,r}}$ for any estimator. Consequently, (up to a constant) $\mathbb{R}_n^*(\mathbf{a}, \phi)$ is a minimax bound and $\widehat{f}_\cdot^{m_n^*}$ is minimax optimal. However, the optimal dimension m_n^* depends on $\mathbf{a} \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ characterising the ellipsoid $\mathbb{F}_2^{\text{a,r}}$. \square

§22.54 **Proposition (Upper bound)**. Under Assumptions §19.02 and §22.51 and $\mathbb{P}^\xi \in \mathbb{P}_{\{0\} \times \mathbb{R}_0^+}$ with $\sigma_\xi^2 := \mathbb{P}^\xi(\text{id}_{\mathbb{R}}^2) \in \mathbb{R}_{>0}^+$ and $\kappa_\xi^{14} := \mathbb{P}^\xi(\text{id}_{\mathbb{R}}^{14}) \in \mathbb{R}_{>0}^+$ for $M^\phi \in \mathbb{N}$ as in (22.25) and and $\widehat{\text{pen}}_\cdot^\phi \in (\mathbb{R}_{>0}^+)^{\mathbb{N}}$ as in (22.32) the OPE $\widehat{f}_\cdot^{\widehat{m}} = \widehat{f}_\cdot \mathbf{1}_\cdot^{\widehat{m}} \in \ell_2 \mathbf{1}_\cdot^{\widehat{m}} \subseteq \text{dom}(\phi_{\nu_K})$ with fully data-driven dimension $\widehat{m} \in \llbracket M^\phi \rrbracket$ as in (22.33) satisfies for all $n \in \mathbb{N}$

$$\begin{aligned} \mathcal{R}_n^\phi[\widehat{f}_\cdot^{\widehat{m}} | \mathbb{F}_2^{\text{a,r}}] &\leq 224(\sigma_\xi^2 + r^2 \tau_{\text{a,u}}^2 + 12\tau_u^2) \\ &\quad \times \min \left\{ \|\mathbf{a} \cdot \mathbf{1}_\cdot^{m \perp}\|_\phi^2 + n^{-1} \|\mathbf{1}_\cdot^m\|_\phi^2 (\log K_m) (1 \vee (\log K_m) m n^{-2/3}) : m \in \llbracket M^\phi \rrbracket \right\} \\ &\quad + 576(5\sigma_\xi^2 + 5r^2 \tau_{\text{a,u}}^2 + 30\tau_u^2 n^{-2/3} + 2^{11} \phi_1^2 \tau_u^2 (\kappa_\xi^{14} + r^{14} \tau_{\text{a,u}}^{14})) n^{-1}. \end{aligned} \quad (22.38)$$

§22.55 **Proof of Proposition** §22.54. is given in the lecture. \square

§22.56 **Corollary**. Under the assumptions of **Proposition** §22.54 for $n, m \in \mathbb{N}$ we set

$$\begin{aligned} \mathbb{R}_n^m(\mathbf{a}, \phi) &:= \|\mathbf{a} \cdot \mathbf{1}_\cdot^{m \perp}\|_\phi^2 \\ &\quad + (1 + (\log \|\mathbf{1}_\cdot^m\|_\phi^2)_+ + \log m) (1 + ((\log \|\mathbf{1}_\cdot^m\|_\phi^2)_+ + \log m) m n^{-2/3}) n^{-1} \|\mathbf{1}_\cdot^m\|_\phi^2, \\ m^\diamond &:= \arg \min \{ \mathbb{R}_n^m(\mathbf{a}, \phi) : m \in \mathbb{N} \} \quad \text{and} \\ \mathbb{R}_n^\diamond(\mathbf{a}, \phi) &:= \mathbb{R}_n^{m^\diamond}(\mathbf{a}, \phi) = \min \{ \mathbb{R}_n^m(\mathbf{a}, \phi) : m \in \mathbb{N} \}. \end{aligned} \quad (22.39)$$

For each $n \in \mathbb{N}$ such that $m^\diamond \in \llbracket M^\phi \rrbracket$ we have

$$\begin{aligned} \mathcal{R}_n^\phi[\widehat{f}_\cdot^{\widehat{m}} | \mathbb{F}_2^{\text{a,r}}] &\leq 2016(\sigma_\xi^2 + r^2 \tau_{\text{a,u}}^2 + 12\tau_u^2) \mathbb{R}_n^\diamond(\mathbf{a}, \phi) \\ &\quad + 576(5\sigma_\xi^2 + 5r^2 \tau_{\text{a,u}}^2 + 30\tau_u^2 n^{-2/3} + 2^{11} \phi_1^2 \tau_u^2 (\kappa_\xi^{14} + r^{14} \tau_{\text{a,u}}^{14})) n^{-1} \\ &\leq 576((4 + 5\phi_1^{-2})(\sigma_\xi^2 + r^2 \tau_{\text{a,u}}^2 + 12\tau_u^2) + 2^{11} \tau_u^2 (\kappa_\xi^{14} + r^{14} \tau_{\text{a,u}}^{14})) \mathbb{R}_n^\diamond(\mathbf{a}, \phi). \end{aligned} \quad (22.40)$$

§22.57 **Proof of Corollary** §22.56. is given in the lecture. \square

§22.58 **Comment**. The data-driven bound $\mathbb{R}_n^\diamond(\mathbf{a}, \phi)$ compared to the minimax bound $\mathbb{R}_n^*(\mathbf{a}, \phi)$ features a deterioration of the variance term at least by a logarithmic factor. The appearance of the logarithmic factor within the bound is a known fact in the context of local estimation (cf. Laurent et al. [2008] who consider model selection given direct Gaussian observations). Brown and Low [1996] show that it is unavoidable in the context of nonparametric Gaussian regression and hence it is widely considered as an acceptable price for adaptation. \square

§22.59 **Illustration.** We illustrate the last results considering the two specifications (o) and (o) given in Table 04 [§12] (Illustration §12.47). We restrict ourselves again to the case $\phi \notin \ell_2$ only.

Table 02 [§22]

Order of the minimax rate $R_n^*(\mathbf{a}, \phi)$ and the data-driven rate $R_n^\circ(\mathbf{a}, \phi)$ as $n \rightarrow \infty$

$(j \in \mathbb{N})$ $\phi = j^{v-1/2}$	$(a \in \mathbb{R}_{\setminus 0}^+)$ \mathbf{a}_j^2	(squared bias) $\ \mathbf{a} \cdot \mathbf{1}^{m \perp}\ _\phi^2$	(variance) $\ \mathbf{1}^m\ _\phi^2$	M^ϕ	m°	$R_n^*(\mathbf{a}, \phi)$	$R_n^\circ(\mathbf{a}, \phi)$	
(o)	$v \in (0, a)$	j^{-a}	$m^{-2(a-v)}$	m^{2v}	$n^{\frac{2}{2v+1}}$	$n^{-\frac{(a-v)}{a}}$	$\left(\frac{\log n}{n}\right)^{\frac{(a-v)}{a}}$	
	$a \in (3/4, \infty)$							$\left(\frac{n}{\log n}\right)^{\frac{1}{2a}}$
	$a \in (0, 3/4]$	j^{-a}	m^{-2a}	$\log m$	$\frac{n^2}{\log n}$	$\frac{\log n}{n}$	$\left(\frac{\log n}{n^{5/6}}\right)^{\frac{1}{a+1/2}}$	
	$v = 0$						$\left(\frac{n}{(\log n)^2}\right)^{\frac{1}{2a}}$	$\frac{(\log n)^2}{n}$
	$a \in (0, 3/4]$					$\left(\frac{n^{5/3}}{(\log n)^3}\right)^{\frac{1}{2a+1}}$	$\left(\frac{(\log n)^3}{n^{5/3}}\right)^{\frac{a}{a+1/2}}$	
(s)	$v \in \mathbb{R}_{\setminus 0}^+$	$e^{-j^{2a}}$	$m^{2(v-a)+} e^{-m^{2a}}$	m^{2v}	$n^{\frac{1}{2v}}$	$(\log n)^{\frac{1}{2a}}$	$\frac{(\log n)^{\frac{v}{a}}}{n}$	$\frac{(\log n)^{\frac{v}{a}} (\log \log n)}{n}$
	$v = 0$	$e^{-j^{2a}}$	$e^{-m^{2a}}$	$\log m$	e^n	$(\log n)^{\frac{1}{2a}}$	$\frac{\log \log n}{n}$	$\frac{(\log \log n)^2}{n}$

We note that in Table 02 [§22] the order of the minimax rate $R_n^*(\mathbf{a}, \phi)$ and the data-driven rate $R_n^\circ(\mathbf{a}, \phi)$ is depicted for $v \geq 0$ only. In case $v < 0$ we have $\phi \in \ell_2$. Moreover, in case (s) for $a \in \mathbb{R}_{\setminus 0}^+$ and (o) for $a \in (3/4, \infty)$ the rate $R_n^\circ(\mathbf{a}, \phi)$ features only an additional logarithmic factor compared with the minimax rate $R_n^*(\mathbf{a}, \phi)$. □

Appendix A

Probability theory

Elements of the PROBABILITY THEORY are recalled along the lines of the text book Klenke [2008] where a detailed exposition with many examples can be found.

§19 Fundamentals

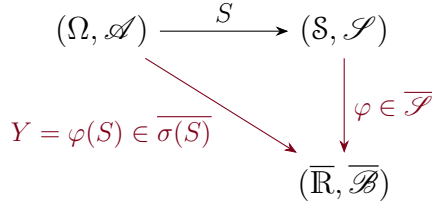
§19.01 **Notation.** For $x, y \in \mathbb{R}$ we agree on the following notations $\lfloor x \rfloor := \max\{k \in \mathbb{Z} : k \leq x\}$ (integer part), $x \vee y = \max(x, y)$ (maximum), $x \wedge y = \min(x, y)$ (minimum), $x^+ = \max(x, 0)$ (positive part), $x^- = \max(-x, 0)$ (negative part) and $|x| = x^- + x^+$ (modulus).

- (i) We set $\mathbb{R}^+ := [0, \infty)$, $\mathbb{R}_0^+ := (0, \infty)$, $\mathbb{R}_0 := \mathbb{R} \setminus \{0\}$, $\overline{\mathbb{R}} := [-\infty, \infty]$, $\overline{\mathbb{R}}^+ := [0, \infty]$.
- (ii) For $a, b \in \mathbb{R}$ with $a < b$ we write $\llbracket a, b \rrbracket := [a, b] \cap \mathbb{Z}$, $\llbracket a, b \llbracket := [a, b) \cap \mathbb{Z}$ and $\rrbracket a, b \rrbracket := (a, b] \cap \mathbb{Z}$. Moreover, let $\llbracket n \rrbracket := \llbracket 1, n \rrbracket$ and $\llbracket n \llbracket := \llbracket 1, n \llbracket$ for $n \in \mathbb{N}$.
- (iii) For $a^n = (a_i)_{i \in \llbracket n \rrbracket}$, $b^n = (b_i)_{i \in \llbracket n \rrbracket} \in \overline{\mathbb{R}}^n$ we write $a^n < b^n$, if $a_i < b_i$ for all $i \in \llbracket n \rrbracket$. For $a^n < b^n$, define the open *rectangle* as the Cartesian product $(a^n, b^n) := \prod_{i=1}^n (a_i, b_i) := (a_1, b_1) \times (a_2, b_2) \times \cdots \times (a_n, b_n)$. Analogously, we define $[a^n, b^n]$, $(a^n, b^n]$ and $[a^n, b^n)$.
- (iv) We call $\overline{\mathcal{B}} := \mathcal{B}_{\overline{\mathbb{R}}}$ the Borel- σ -field over the compactified real line $\overline{\mathbb{R}}$, where the sets $\{-\infty\}$, $\{\infty\}$ and \mathbb{R} are in $\overline{\mathbb{R}}$ closed and open, respectively, and hence Borel-measurable. In particular, the trace $\mathcal{B} := \mathcal{B}_{\mathbb{R}} = \overline{\mathcal{B}} \cap \mathbb{R}$ of $\overline{\mathcal{B}}$ over \mathbb{R} is the Borel- σ -field over \mathbb{R} . Furthermore, we write $\overline{\mathcal{B}}^+ := \overline{\mathcal{B}} \cap \overline{\mathbb{R}}^+$, $\mathcal{B}^+ := \mathcal{B} \cap \mathbb{R}^+$ and $\mathcal{B}_0^+ := \mathcal{B} \cap \mathbb{R}_0^+$.
- (v) Given a measurable space (Ω, \mathcal{A}) a Borel-measurable function $g : \Omega \rightarrow \mathbb{R}$ and $f : \Omega \rightarrow \overline{\mathbb{R}}$ is called *real* and *numerical*, respectively, and we write $g \in \mathcal{A}$ and $f \in \overline{\mathcal{A}}$ for short. g respectively f is called positive if $g(\Omega) \in \mathbb{R}^+$ respectively $f(\Omega) \in \overline{\mathbb{R}}^+$, then we write $g \in \mathcal{A}^+$ and $f \in \overline{\mathcal{A}}^+$ - We call a Borel-measurable function $f^k = (f_i)_{i \in \llbracket k \rrbracket} : \Omega \rightarrow \overline{\mathbb{R}}^k$, that is $f_i \in \overline{\mathcal{A}}$ for each $i \in \llbracket k \rrbracket$, and $g^k = (g_i)_{i \in \llbracket k \rrbracket} : \Omega \rightarrow \mathbb{R}^k$, *numerical* and *real*, respectively and we write $f^k \in \overline{\mathcal{A}}^k$ and $g^k \in \mathcal{A}^k$ for short. □

§19.02 **Property.**

- (i) For $X, Y \in \overline{\mathcal{A}}$ and $a \in \mathbb{R}$ holds: $aX \in \overline{\mathcal{A}}$ (with convention $0 \times \infty = 0$); $X \vee Y := \max(X, Y)$, $X \wedge Y := \min(X, Y) \in \overline{\mathcal{A}}$ and particularly $X^+ := X \vee 0$, $X^- := (-X)^+ \in \overline{\mathcal{A}}^+$, $|X| \in \overline{\mathcal{A}}^+$, $\{X < Y\}$, $\{X \leq Y\}$, $\{X = Y\} \in \mathcal{A}$, and $\lfloor X \rfloor \in \overline{\mathcal{A}}^+$.
- (ii) For $X^n = (X_i)_{i \in \llbracket n \rrbracket} \in \mathcal{A}^n$, i.e., $X_i \in \mathcal{A}$, $i \in \llbracket n \rrbracket$, and Borel-measurable $h : \mathbb{R}^n \rightarrow \mathbb{R}^m$ holds $h(X^n) \in \mathcal{A}^m$, and in particular $X_1 + X_2$, $X_1 - X_2$, $X_1 X_2 \in \mathcal{A}$, and $X_1 / X_2 \in \overline{\mathcal{A}}$.
- (iii) Let $(X_n)_{n \in \mathbb{N}}$ be a sequence in $\overline{\mathcal{A}}$. Then $\sup_{n \in \mathbb{N}} X_n \in \overline{\mathcal{A}}$, $\inf_{n \in \mathbb{N}} X_n \in \overline{\mathcal{A}}$, $X_* = \liminf_{n \rightarrow \infty} X_n \in \overline{\mathcal{A}}$ and $X^* = \limsup_{n \rightarrow \infty} X_n \in \overline{\mathcal{A}}$. If $X := \lim_{n \rightarrow \infty} X_n$ exists, then $X \in \overline{\mathcal{A}}$.
- (iv) Let $S : (\Omega, \mathcal{A}) \rightarrow (\mathcal{S}, \mathcal{S})$ be measurable, $\sigma(S) := S^{-1}(\mathcal{S}) \subseteq \mathcal{A}$ the sub- σ -field generated by S and $Y : \Omega \rightarrow \overline{\mathbb{R}}$. Then the following conditions are *equivalent*: (a) Y is $\sigma(S)$ -measurable, symbolically $Y \in \overline{\sigma(S)}$; (b) There exists a measurable $\varphi : (\mathcal{S}, \mathcal{S}) \rightarrow (\overline{\mathbb{R}}, \overline{\mathcal{B}})$, in short $\varphi \in \overline{\mathcal{S}}$, with $Y = \varphi(S)$. If Y is real, bounded or positive, then φ has each of those

properties too.



The function φ is uniquely determined by Y on $S(\Omega)$, and for all $s \notin S(\Omega)$ it can be arbitrarily be extended.

- (v) For every $X \in \overline{\mathcal{A}}^+$ the sequence of simple random variables $(X_n)_{n \in \mathbb{N}}$ in $\overline{\mathcal{A}}^+$ given by $X_n := (2^{-n} \lfloor 2^n X \rfloor) \wedge n$ satisfies (a) $X_n \uparrow X$; (b) $X_n \leq X \wedge n$; (c) For each $c \in \mathbb{R}^+$ holds $\lim_{n \rightarrow \infty} X_n = X$ uniformly on $\{X \leq c\}$. □

§19.03 **Notation.** For a measure μ on (Ω, \mathcal{A}) we denote the integral of $f \in \overline{\mathcal{A}}$ with respect to μ by $\mu f := \int f d\mu$, if it exists. For $s \in \mathbb{R}_0^+$ define $\|f\|_{\mathcal{L}_s(\mu)} := (\mu|f|^s)^{1/s}$, and $\|f\|_{\mathcal{L}_\infty(\mu)} := \inf\{c \in \mathbb{R}^+ : \mu(|f| > c) = 0\}$. For $s \in \overline{\mathbb{R}}_0^+ := (0, \infty]$ a function $f \in \overline{\mathcal{A}}$ is called $\mathcal{L}_s(\mu)$ -integrable, if $\|f\|_{\mathcal{L}_s(\mu)} < \infty$. We denote the set of all $\mathcal{L}_s(\mu)$ -integrable functions by $\mathcal{L}_s(\mu) := \mathcal{L}_s(\mathcal{A}, \mu) := \{f \in \overline{\mathcal{A}} : \|f\|_{\mathcal{L}_s(\mu)} < \infty\}$. Note that $\|\cdot\|_{\mathcal{L}_s(\mu)}$ is a seminorm on $\mathcal{L}_s(\mu)$ for each $s \in [1, \infty]$. Given a metric space (X, d) equipped with its Borel- σ -field \mathcal{B}_X we denote by $\mathcal{C}_b := \mathcal{C}_b(X)$ the set of all bounded and continuous functions mapping X into \mathbb{R} . For any finite measure μ on (X, \mathcal{B}_X) we have $\|h\|_{\mathcal{L}_\infty(\mu)} < \infty$ for all $h \in \mathcal{C}_b$ and thus $\mathcal{C}_b \subseteq \mathcal{L}_\infty(\mathcal{B}_X, \mu)$ in equal. We denote by λ the Lebesgue measure on $(\mathbb{R}, \mathcal{B})$ and write shortly $\mathcal{L}_s := \mathcal{L}_s(\mathcal{B}) := \mathcal{L}_s(\mathcal{B}, \lambda)$. □

§19.04 **Notation.** We understand a vector $a^k = (a_i)_{i \in \llbracket k \rrbracket}$ as a column vector, i.e., $a^k = (a_1 \cdots a_k)^t \in \overline{\mathbb{R}}^k$ and hence we identify $\overline{\mathbb{R}}^k$ and $\overline{\mathbb{R}}^{(k,1)}$. We denote by $\|\cdot\|$ and $\langle \cdot, \cdot \rangle$ the Euclidean norm and inner product on \mathbb{R}^k , respectively, i.e., $\|a^k\| = (\sum_{i \in \llbracket k \rrbracket} |a_i|^2)^{1/2}$ and $\langle a^k, b^k \rangle = \sum_{i \in \llbracket k \rrbracket} a_i b_i = (b^k)^t a^k$ for all $a^k, b^k \in \overline{\mathbb{R}}^k$. For $s \in \mathbb{R}_0^+$ we define $\|a^k\|_s := (\sum_{i \in \llbracket k \rrbracket} |a_i|^s)^{1/s}$ and $\|a^k\|_\infty := \max_{i \in \llbracket k \rrbracket} |a_i|$. Note that $f^k \in \overline{\mathcal{A}}^k$ and $g^k \in \mathcal{A}^k$ imply $\|f^k\|_s \in \overline{\mathcal{A}}$ and $\|g^k\|_s \in \mathcal{A}$ for any $s \in \overline{\mathbb{R}}_0^+$. We call $f^k = (f_i)_{i \in \llbracket k \rrbracket}$ $\mathcal{L}_s^k(\mu)$ -integrable if $\|f^k\|_s \in \mathcal{L}_s(\mu)$ or equivalently $f_i \in \mathcal{L}_s(\mu)$ for each $i \in \llbracket k \rrbracket$. We define $\|f^k\|_{\mathcal{L}_s^k(\mu)} := \|\|f^k\|_p\|_{\mathcal{L}_s(\mu)}$ and $\mathcal{L}_s^k(\mu) := \mathcal{L}_s^k(\mathcal{A}, \mu) := \{f^k \in \overline{\mathcal{A}}^k : \|f^k\|_{\mathcal{L}_s^k(\mu)} < \infty\}$ with a slight abuse of notation. □

§19.05 **Notation.** Let X be a random variable, i.e. a measurable function, defined on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ with values in a measurable space (X, \mathcal{X}) . The probability measure on (X, \mathcal{X}) induced by X is denoted by $\mathbb{P}^X := \mathbb{P} \circ X^{-1}$ and we write $X \sim \mathbb{P}^X$ for short. For $f \in \overline{\mathcal{X}}$ the expectation of f with respect to \mathbb{P}^X or equivalently of $f(X)$ with respect to \mathbb{P} (if it exists) is denoted by $\mathbb{E}_X^f := \mathbb{P}^X f = \mathbb{P} f(X) =: \mathbb{E} f(X)$ for short. For example, when applied to the empirical measure $\widehat{\mathbb{P}}_n$ given by $\widehat{\mathbb{P}}_n(x^n) := \frac{1}{n} \sum_{i \in \llbracket n \rrbracket} \delta_{x_i}$ for $x^n = (x_i)_{i \in \llbracket n \rrbracket} \in \mathcal{X}^n$ this yields $\widehat{\mathbb{P}}_n f \in \overline{\mathcal{X}}$ with $x^n \mapsto (\widehat{\mathbb{P}}_n f)(x^n) := \frac{1}{n} \sum_{i \in \llbracket n \rrbracket} f(x_i)$. In other words, for each $x^n \in \mathcal{X}^n$, $(\widehat{\mathbb{P}}_n f)(x^n)$ is an abbreviation for the average $\frac{1}{n} \sum_{i \in \llbracket n \rrbracket} f(x_i)$. We denote by $\mathcal{W}(\mathcal{X})$ the set of all probability measures on (X, \mathcal{X}) and for \mathbb{R}^n equipped with its Borel- σ -field $\mathcal{B}^n := \mathcal{B}_{\mathbb{R}^n}$ by $\mathcal{W}_s(\mathcal{B}^n) \subseteq \mathcal{W}(\mathcal{B}^n)$ the subset of all probability measures on $(\mathbb{R}^n, \mathcal{B}^n)$ with finite $s \in \mathbb{R}^+$ absolute mean, that is, for all $\mathbb{P} \in \mathcal{W}_s(\mathcal{B}^n)$ the identity mapping $\text{id}_n : \mathbb{R}^n \rightarrow \mathbb{R}^n$ belongs to $\mathcal{L}_s^n(\mathbb{P})$. Furthermore, for $Y \sim \mathbb{P}$ we write $\mathbb{E}(Y) = \mathbb{P}(Y) := \mathbb{P}(\text{id}_n) = (\mathbb{P}(\Pi_i))_{i \in \llbracket n \rrbracket}$ using for $i \in \llbracket n \rrbracket$ the coordinate map $\Pi_i : \mathbb{R}^n \rightarrow \mathbb{R}$ with $x^n = (x_i)_{i \in \llbracket n \rrbracket} \mapsto \Pi_i(x^n) := x_i$. □

§19.06 **Property.** Let $X \in \mathcal{L}_2^k(\mathbb{P})$, i.e. $\|X\|_{\mathcal{L}_2^k(\mathbb{P})}^2 = \mathbb{P}(\|X\|^2) < \infty$. For each $b \in \mathbb{R}^n$ and $A \in \mathbb{R}^{(n,k)}$ we have $Y := AX + b \in \mathcal{L}_2^n(\mathbb{P})$. If we further denote by $\mu := \mathbb{P}X \in \mathbb{R}^k$ and $\Sigma := \text{Cov}(X) = \mathbb{P}(X - \mu)(X - \mu)^t = \mathbb{P}(XX^t) - \mu\mu^t \in \mathbb{R}^{(k,k)}$ expectation vector and covariance matrix of X , respectively, then $\mathbb{P}(Y) = A\mu + b \in \mathbb{R}^n$ and $\text{Cov}(Y) = A\Sigma A^t \in \mathbb{R}^{(n,n)}$. \square

§19.07 **Definition.** A $\mathcal{L}_2^k(\mathbb{P})$ -random vector X with $\mu := \mathbb{P}(X)$ and $\Sigma := \text{Cov}(X)$ is *multivariate normally distributed*, $X \sim N_{(\mu, \Sigma)}$ for short, if for each $c \in \mathbb{R}^k$ the real random variable $\langle X, c \rangle$ is normally distributed with mean $\langle \mu, c \rangle$ and variance $\langle \Sigma c, c \rangle$, i.e., $\langle X, c \rangle \sim N_{(\langle \mu, c \rangle, \langle \Sigma c, c \rangle)}$. If Id_k denotes the k -dimensional identity matrix, then $X \sim N_{(0, \text{Id}_k)}$ is called a *standard normal random vector*. \square

§19.08 **Property.** A random vector $X = (X_i)_{i \in \llbracket k \rrbracket}$ is *standard normal*, i.e., $X \sim N_{(0, \text{Id}_k)}$ if and only if its components $\{X_i : i \in \llbracket k \rrbracket\} \in \mathbb{K}$ are independent and identically $N_{(0,1)}$ -distributed. \square

§19.09 **Remark.** In other words, a multivariate $N_{(0, \text{Id}_k)}$ -distribution equals the product of its marginal $N_{(0,1)}$ -distributions, or $N_{(0, \text{Id}_k)} = N_{(0,1)}^{\otimes k} := \bigotimes_{i \in \llbracket k \rrbracket} N_{(0,1)}$ for short. \square

§20 Convergence of random variables

Here and subsequently, a metric space is equipped with its Borel- σ -field.

§20.01 **Definition.** Let X and X_n , $n \in \mathbb{N}$, be random variables on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ with values in a metric space (\mathcal{X}, d) . The sequence $(X_n)_{n \in \mathbb{N}}$ *converges to* X :

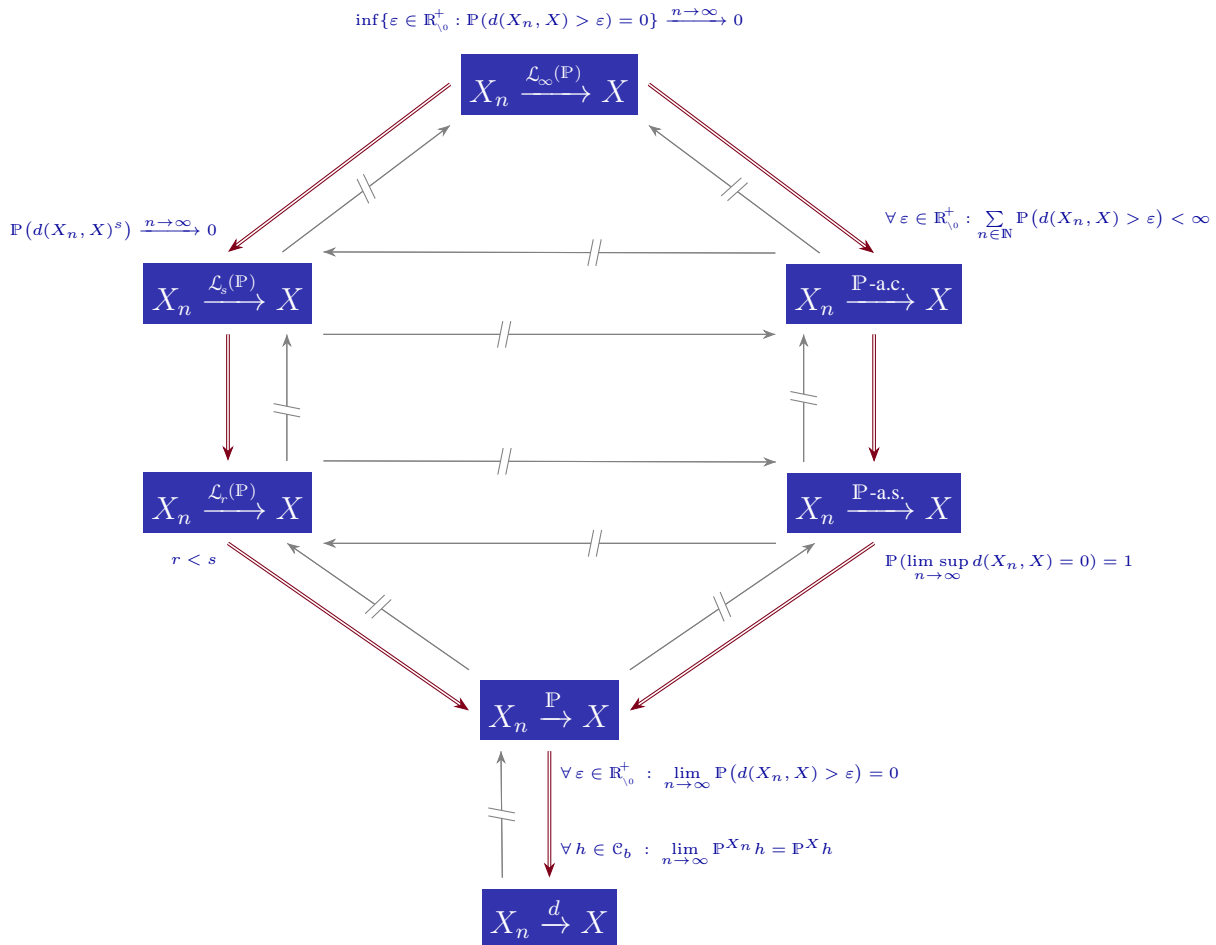
- (a) *almost surely* (\mathbb{P} -a.s.), if $\mathbb{P}(\lim_{n \rightarrow \infty} d(X_n, X) = 0) = 1$. We write $X_n \xrightarrow{n \rightarrow \infty} X$ \mathbb{P} -a.s., or briefly, $X_n \xrightarrow{\mathbb{P}\text{-a.s.}} X$.
- (b) *almost completely* (\mathbb{P} -a.c.), if $\sum_{n \in \mathbb{N}} \mathbb{P}(d(X_n, X) > \varepsilon) < \infty$ for all $\varepsilon \in \mathbb{R}_0^+$. We write $X_n \xrightarrow{n \rightarrow \infty} X$ \mathbb{P} -a.c., or briefly, $X_n \xrightarrow{\mathbb{P}\text{-a.c.}} X$.
- (c) *in probability*, if $\lim_{n \rightarrow \infty} \mathbb{P}(d(X_n, X) \geq \varepsilon) = 0$ for all $\varepsilon \in \mathbb{R}_0^+$. We write $X_n \xrightarrow{n \rightarrow \infty} X$ in \mathbb{P} , or briefly, $X_n \xrightarrow{\mathbb{P}} X$.
- (d) *in distribution*, if $\lim_{n \rightarrow \infty} \mathbb{P}^{X_n} f = \mathbb{P}^X f$ for any $f \in \mathcal{C}_b(\mathcal{X})$. We write $X_n \xrightarrow{n \rightarrow \infty} X$ in distribution, or briefly, $X_n \xrightarrow{d} X$ and with a slight abuse of notation also $X_n \xrightarrow{d} \mathbb{P}^X$.
- (e) *in $\mathcal{L}_s(\mathbb{P})$ or s -th mean*, if $\lim_{n \rightarrow \infty} \mathbb{P}(d(X_n, X)^s) = 0$. We write $X_n \xrightarrow{n \rightarrow \infty} X$ in $\mathcal{L}_s(\mathbb{P})$, or briefly, $X_n \xrightarrow{\mathcal{L}_s(\mathbb{P})} X$. \square

§20.02 **Remark.** Let X and X_n , $n \in \mathbb{N}$, be random vectors in \mathbb{R}^k , i.e., $(\mathbb{R}^k, \mathcal{B}^k)$ -valued random variables, and $\|\cdot\|_s$ as in **Notation** §19.04. Convergence of $(X_n)_{n \in \mathbb{N}}$ to X in s -th mean, that is, $\mathbb{P}\|X_n - X\|_s^s = \mathbb{P}\|X_n - X\|_{\mathcal{L}_s^k(\mathbb{P})}^s \xrightarrow{n \rightarrow \infty} 0$, equals the component-wise convergence of $(X_n^i)_{n \in \mathbb{N}}$ to X^i in $\mathcal{L}_s(\mathbb{P})$, i.e., $\mathbb{P}|X_n^i - X^i|^s = \|X_n^i - X^i\|_{\mathcal{L}_s(\mathbb{P})}^s \xrightarrow{n \rightarrow \infty} 0$ for each $i \in \llbracket k \rrbracket$. \square

§20.03 **Property.** Let X and X_n , $n \in \mathbb{N}$, be random variables on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ with values in a metric space (\mathcal{X}, d) .

- (i) *The following statements are equivalent:* (a) $X_n \xrightarrow{\mathbb{P}\text{-a.s.}} X$; (b) $\sup_{m \geq n} d(X_m, X_n) \xrightarrow{\mathbb{P}} 0$; (c) $\forall \varepsilon, \delta \in \mathbb{R}_0^+ : \exists N \in \mathbb{N} : \forall n \geq N : \mathbb{P}(\bigcap_{j \geq n} \{d(X_j, X) \leq \varepsilon\}) \geq 1 - \delta$ and (d) $\sup_{m \geq n} d(X_m, X) \xrightarrow{\mathbb{P}} 0$.

- (ii) *(Continuous mapping theorem)* Let $g : \mathcal{X} \rightarrow \mathbb{R}$ be continuous and let $(X_n)_{n \in \mathbb{N}}$ converge to X \mathbb{P} -a.s. (respectively, in probability or in distribution). Then $(g(X_n))_{n \in \mathbb{N}}$ converges to $g(X)$ \mathbb{P} -a.s. (respectively, in probability or in distribution).
- (iii) Counter examples show, that the converse (in gray) of the following direct implications (in red) do not hold. □



§20.04 **Definition.** A family of $\{X_{n,j} : j \in \llbracket k_n \rrbracket, n \in \mathbb{N} \in \mathbb{K}\}$ of real \mathcal{L}_2 -random variables is called a standardised array, if for every $n \in \mathbb{N}$ the family $\{X_{n,j} : j \in \llbracket k_n \rrbracket \in \mathbb{K}\}$ is independent, centred and normed, i.e., $\mathbb{E}(X_{n,j}) = 0, j \in \llbracket k_n \rrbracket$ and $\sum_{j \in \llbracket k_n \rrbracket} \text{var}_m(X_{n,j}) = 1$. A standardised array $\{X_{n,j} : j \in \llbracket k_n \rrbracket, n \in \mathbb{N} \in \mathbb{K}\}$ is said to satisfy

- (a) the *Lindeberg condition*, if $\lim_{n \rightarrow \infty} \sum_{j \in \llbracket k_n \rrbracket} \mathbb{E}(X_{n,j}^2 \mathbb{1}_{\{|X_{n,j}| \geq \delta\}}) = 0$ for every $\delta \in \mathbb{R}_0^+$;
- (b) the *Lyapunov condition*, if there is $\delta \in \mathbb{R}_0^+$ such that $\lim_{n \rightarrow \infty} \sum_{j \in \llbracket k_n \rrbracket} \mathbb{E}|X_{n,j}|^{2+\delta} = 0$. □

§20.05 **Property.** Let $(X_n)_{n \in \mathbb{N}}$ be a sequence of independent real random variables.

- (i) (Law of Large Numbers) Let $X_n, n \in \mathbb{N}$, be identically distributed. Then $X_1 \in \mathcal{L}_1(\mathbb{P})$ if and only if $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i \in \llbracket n \rrbracket} X_i = \mathbb{P}(X_1)$ \mathbb{P} -a.s. (and then also in $\mathcal{L}_1(\mathbb{P})$).
- (ii) (*Lévy's equivalence theorem*) For partial sums $(S_n := \sum_{i \in \llbracket n \rrbracket} X_i)_{n \in \mathbb{N}}$ \mathbb{P} -a.s. convergence is equivalent to convergence in probability. Otherwise, they diverge with probability one.
(Kolmogorov's three-series theorem) $(S_n)_{n \in \mathbb{N}}$ converges \mathbb{P} -a.s. if and only if there is $\varepsilon \in \mathbb{R}_0^+$ such that each of the following three conditions holds: (a) $\sum_{n \in \mathbb{N}} \mathbb{P}(|X_n| > \varepsilon) < \infty$;

- (b) $\sum_{n \in \mathbb{N}} \mathbb{E}(X_n \mathbb{1}_{\{|X_n| \leq \varepsilon\}})$ converges; and (c) $\sum_{n \in \mathbb{N}} \text{var}_m(X_n \mathbb{1}_{\{|X_n| \leq \varepsilon\}}) < \infty$.

Let $\{X_{n,j} : j \in \llbracket k_n \rrbracket, n \in \mathbb{N} \in \mathbb{K}\}$ be a standardised array.

(iii) *The Lyapunov condition implies the Lindeberg condition.*

(iv) (Central Limit Theorem of Lindeberg (1922)) *If the Lindeberg condition hold, then (for the row sum) $S_n^* = \sum_{j \in \llbracket k_n \rrbracket} X_{nj} \xrightarrow{d} N_{(0,1)}$.* \square

§20.06 **Remark** (Law of Large Numbers). Let $X_n^k, n \in \mathbb{N}$, be i.i.d. random vector in \mathbb{R}^k . Then $\|X_1^k\|_{\mathcal{L}_1^k(\mathbb{P})} = \mathbb{P}\|X_1^k\|_1 < \infty$ if and only if $\frac{1}{n} \sum_{i \in \llbracket n \rrbracket} X_i^k \xrightarrow{\mathbb{P}\text{-a.s.}} \mathbb{E}(X_1^k)$ (then also in $\mathcal{L}_1^k(\mathbb{P})$). \square

§20.07 **Property** (Portemanteau). Let X and $X_n, n \in \mathbb{N}$, be random variables on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ with values in a metric space (\mathcal{X}, d) . The following statements are equivalent:

- (i) $X_n \xrightarrow{d} X$;
- (ii) $\liminf_{n \rightarrow \infty} \mathbb{P}(X_n \in U) \geq \mathbb{P}(X \in U)$ for all open $U \subseteq \mathcal{X}$;
- (iii) $\limsup_{n \rightarrow \infty} \mathbb{P}(X_n \in F) \leq \mathbb{P}(X \in F)$ for all closed $F \subseteq \mathcal{X}$;
- (iv) $\lim_{n \rightarrow \infty} \mathbb{P}(X_n \in B) = \mathbb{P}(X \in B)$ for all measurable B with $\mathbb{P}(X \in \partial B) = 0$ where \bar{B} , B and $\partial B = \bar{B} \setminus B$ is the closure, interior and the boundary of B , respectively. \square

§20.08 **Property** (Helly-Bray). Let X and $X_n, n \in \mathbb{N}$, be random vectors in \mathbb{R}^k with cumulative distribution function (c.d.f.) for each $x \in \mathbb{R}^k$ given by $\mathbb{F}(x) := \mathbb{P}(X \leq x)$ and $\mathbb{F}_n(x) := \mathbb{P}(X_n \leq x)$. Then the following statements are equivalent: (i) $X_n \xrightarrow{d} X$ and (ii) $\lim_{n \rightarrow \infty} \mathbb{F}_n(x) = \mathbb{F}(x)$ for all points of continuity x of \mathbb{F} . \square

§20.09 **Property** (Continuous mapping theorem). Let (\mathcal{X}_1, d_1) and (\mathcal{X}_2, d_2) be metric spaces and let $\varphi : \mathcal{X}_1 \rightarrow \mathcal{X}_2$ be measurable. Denote by U_φ the set of points of discontinuity of φ . If X and $X_n, n \in \mathbb{N}$, are \mathcal{X}_1 -valued random variables with $\mathbb{P}(X \in U_\varphi) = 0$ and $X_n \xrightarrow{d} X$, then $\varphi(X_n) \xrightarrow{d} \varphi(X)$. \square

§20.10 **Property** (Slutzky's lemma). Let X and $X_n, Y_n, n \in \mathbb{N}$, be random variables taking values in a common metric space (\mathcal{X}, d) and satisfying $X_n \xrightarrow{d} X$ and $d(X_n, Y_n) \xrightarrow{\mathbb{P}} 0$. Then $Y_n \xrightarrow{d} X$. \square

§20.11 **Example**. Let X and $X_n, n \in \mathbb{N}$, be a random vector in \mathbb{R}^k satisfying $X_n \xrightarrow{d} X$.

- (a) If $Y_n, n \in \mathbb{N}$, are random vector in \mathbb{R}^k and $c \in \mathbb{R}^k$ such that $Y_n \xrightarrow{d} c$, then $X_n + Y_n \xrightarrow{d} X + c$.
- (b) If $\Sigma_n, n \in \mathbb{N}$ are random matrices in $\mathbb{R}^{(k,k)}$ and Σ is a matrix in $\mathbb{R}^{(k,k)}$ such that $\Sigma_n \xrightarrow{d} \Sigma$, then $\Sigma_n X_n \xrightarrow{d} \Sigma X$. If in addition Σ is strictly positive definite, and thus invertible, then $\Sigma_n^{-1} X_n \xrightarrow{d} \Sigma^{-1} X$ and $\Sigma_n^{-1/2} X_n \xrightarrow{d} \Sigma^{-1/2} X$. \square

§20.12 **Property** (Cramér-Wold device). Let $X_n, n \in \mathbb{N}$, be random vectors in \mathbb{R}^k . Then, the following are equivalent: (a) There is a random vector X with $X_n \xrightarrow{d} X$. (b) For any $v \in \mathbb{R}^k$, there is a real X^v with $\langle v, X_n \rangle \xrightarrow{d} X^v$. If (a) and (b) hold, then X^v and $\langle v, X \rangle$ are identically distributed (i.d.), $X^v \stackrel{d}{=} \langle v, X \rangle$ for short, for all $v \in \mathbb{R}^k$. \square

§20.13 **Property** (Lindeberg-Feller CLT). For each $n \in \mathbb{N}$ let $\{Y_{n,j} : j \in \llbracket k_n \rrbracket \in \mathbb{K}\}$ be independent and centred \mathcal{L}_2^p -random vectors such that (i) $\sum_{j \in \llbracket k_n \rrbracket} \mathbb{E}\|Y_{n,j}\|^2 \mathbf{1}_{\{\|Y_{n,j}\| > \varepsilon\}} \xrightarrow{n \rightarrow \infty} 0$ for any $\varepsilon \in \mathbb{R}_{>0}^+$ and (ii) $\sum_{j \in \llbracket k_n \rrbracket} \mathbb{E}(Y_{n,j} Y_{n,j}^t) \xrightarrow{n \rightarrow \infty} \Sigma$. Then $\sum_{j \in \llbracket k_n \rrbracket} Y_{n,j} \xrightarrow{d} N_{(0,\Sigma)}$. \square

§20.14 **Example**. Let X and $X_n, n \in \mathbb{N}$, be i.i.d. $\mathcal{L}_2^k(\mathbb{P})$ -random vectors with $\mu = \mathbb{P}(X)$ and strictly positive definite $\Sigma = \text{Cov}(X)$.

- (a) (CLT) $\frac{1}{\sqrt{n}} \sum_{i \in [n]} (X_i - \mu) \xrightarrow{d} N_{(0, \Sigma)}$,
- (b) (LLN) $\bar{X}_n := \frac{1}{n} \sum_{i \in [n]} X_i \xrightarrow{\mathbb{P}} \mu$,
- (c) (LLN) $\frac{1}{n} \sum_{i \in [n]} X_i X_i^t \xrightarrow{\mathbb{P}} \mathbb{E}(X X^t)$,
- (d) $\hat{\Sigma}_n := \frac{1}{n} \sum_{i \in [n]} (X_i - \bar{X}_n)(X_i - \bar{X}_n)^t = \frac{1}{n} \sum_{i \in [n]} X_i X_i^t - \bar{X}_n \bar{X}_n^t \xrightarrow{\mathbb{P}} \mathbb{E}(X X^t) - \mu \mu^t = \text{Cov}(X) = \Sigma$ (using (b) and (c) and continuous mapping theorem §20.03)
- (e) $\sqrt{n} \Sigma_n^{-1/2} (\bar{X} - \mu) \xrightarrow{d} N_{(0, \text{Id})}$ (using (a), (d) and Slutsky's lemma §20.10 as in the Example §20.11 (b)) □

§20.15 **Remark.** A map $\phi : \mathbb{R}^k \rightarrow \mathbb{R}^m$, that is defined at least in a neighbourhood of θ_o , is called differentiable at θ_o , if there exists a linear map (matrix) $\dot{\phi}_{\theta_o} : \mathbb{R}^k \rightarrow \mathbb{R}^m$ such that

$$\lim_{\theta \rightarrow \theta_o} \frac{\|\phi(\theta) - \phi(\theta_o) - \dot{\phi}_{\theta_o}(\theta - \theta_o)\|}{\|\theta - \theta_o\|} = 0.$$

The linear map $x \mapsto \dot{\phi}_{\theta_o}(x)$ is called *(total) derivative* as opposed to partial derivatives. A sufficient condition for ϕ to be (totally) differentiable is that all partial derivatives $\partial \phi_j(\theta) / \partial \theta_l$ exist for θ in a neighbourhood of θ_o and are continuous at θ_o . □

§20.16 **Property (Delta method).** Let $\phi : \mathbb{R}^k \supseteq \mathcal{D}_\phi \rightarrow \mathbb{R}^m$ be a map defined on a subset \mathcal{D}_ϕ of \mathbb{R}^k and differentiable at θ_o . Let T and T_n , $n \in \mathbb{N}$ be random variables taking their values in the domain \mathcal{D}_ϕ of ϕ . If $r_n(T_n - \theta_o) \xrightarrow{d} T$ for numbers $r_n \rightarrow \infty$, then $r_n(\phi(T_n) - \phi(\theta_o)) \xrightarrow{d} \dot{\phi}_{\theta_o}(T)$. Moreover, the difference between $r_n(\phi(T_n) - \phi(\theta_o))$ and $\dot{\phi}_{\theta_o}(r_n(T_n - \theta_o))$ converges to zero in probability. □

§20.17 **Remark.** Commonly, $\sqrt{n}(T_n - \theta_o) \xrightarrow{d} N_{(\mu, \Sigma)}$. Then applying the delta method it follows that $\sqrt{n}(\phi(T_n) - \phi(\theta_o)) \xrightarrow{d} N_{(\dot{\phi}_{\theta_o}\mu, \dot{\phi}_{\theta_o}\Sigma\dot{\phi}_{\theta_o}^t)}$. □

§20.18 **Property (Markov's inequality).** If X is a $\mathcal{L}_s^k(\mathbb{P})$ -random vector for some $s \geq 1$, then $\mathbb{P}(\|X\|_s > c) \leq c^{-s} \mathbb{P}(\|X\|_s^s) = c^{-s} \|X\|_{\mathcal{L}_s^k(\mathbb{P})}^s$. □

§20.19 **Property (Monotone convergence).** Let $(X_n)_{n \in \mathbb{N}}$ be a sequence of monotonically increasing real $\mathcal{L}_1(\mathbb{P})$ -random variables converging \mathbb{P} -a.s. to a numerical random variable X , for short $X_n \uparrow X$ \mathbb{P} -a.s.. Then $\mathbb{P}X = \lim_{n \rightarrow \infty} \mathbb{P}X_n$. □

§20.20 **Property (Dominated convergence).** Let $(X_n)_{n \in \mathbb{N}}$ be a sequence of real $\mathcal{L}_1(\mathbb{P})$ -random variables converging \mathbb{P} -a.s. to a numerical random variable X , i.e., $X_n \xrightarrow{\mathbb{P}\text{-a.s.}} X$. If there is a real $\mathcal{L}_1(\mathbb{P})$ random variable Y with $\sup_{n \in \mathbb{N}} |X_n| \leq Y$ \mathbb{P} -a.s. (and thus $\sup_{n \in \mathbb{N}} |X_n| \in \mathcal{L}_1(\mathbb{P})$), then $X \in \mathcal{L}_1(\mathbb{P})$ and $X_n \xrightarrow{\mathcal{L}_1(\mathbb{P})} X$. □

§20.21 **Definition.** A sequence of random variables $(X_n)_{n \in \mathbb{N}}$ with values in a metric space (\mathcal{X}, d) is called *(uniformly) tight* (straff) or *bounded in probability*, if, for any $\varepsilon \in \mathbb{R}_0^+$, there exists a compact set $K_\varepsilon \subseteq \mathcal{X}$ such that $\mathbb{P}(X_n \in K_\varepsilon) \geq 1 - \varepsilon$ for all $n \in \mathbb{N}$. □

§20.22 **Remark.** If (\mathcal{X}, d) is Polish, i.e., separable and complete, then every \mathcal{X} -valued random variable is bounded in probability and thus so is every finite family. □

§20.23 **Example.** A sequence $(X_n)_{n \in \mathbb{N}}$ of random vectors in \mathbb{R}^k is bounded in probability, if for any $\varepsilon > 0$, there exists a constant K_ε such that $\mathbb{P}(\|X_n\| > K_\varepsilon) \leq \varepsilon$ for all $n \in \mathbb{N}$. □

§20.24 **Property (Prohorov's theorem).** Let X and X_n , $n \in \mathbb{N}$, be random variables with values in a Polish space.

- (i) If $X_n \xrightarrow{d} X$, then $(X_n)_{n \in \mathbb{N}}$ is bounded in probability.
- (ii) If $(X_n)_{n \in \mathbb{N}}$ is bounded in probability, then there exists a sub-sequence $(X_{n_k})_{k \in \mathbb{N}}$ which converges in distribution. \square

§20.25 **Landau notation.** Let X_n , $n \in \mathbb{N}$, be random variables on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ with values in a metric space (\mathcal{X}, d) and let x_n , $n \in \mathbb{N}$, belong to \mathcal{X} .

- (i) We write (a) $x_n = o(1)$, if $d(x_n, 0) \xrightarrow{n \rightarrow \infty} 0$, and (b) $x_n = O(1)$, if $\sup_{n \in \mathbb{N}} d(x_n, 0) < \infty$, and analogously (a) $X_n = o_{\mathbb{P}}(1)$, if $X_n \xrightarrow{\mathbb{P}} 0$, and (b) $X_n = O_{\mathbb{P}}(1)$, if $(X_n)_{n \in \mathbb{N}}$ is bounded in probability
- (ii) Let a_n , $n \in \mathbb{N}$, be strictly positive numbers. We write (a) $x_n = o(a_n)$, if $d(x_n, 0)/a_n = o(1)$, and that (b) $x_n = O(a_n)$, if $d(x_n, 0)/a_n = O(1)$, and analogously (a) $X_n = o_{\mathbb{P}}(a_n)$, if $d(X_n, 0)/a_n = o_{\mathbb{P}}(1)$, and (b) $X_n = O_{\mathbb{P}}(a_n)$, if $d(X_n, 0)/a_n = O_{\mathbb{P}}(1)$.
- (iii) Let A_n , $n \in \mathbb{N}$, be strictly positive random variables on $(\Omega, \mathcal{A}, \mathbb{P})$. We write (a) $X_n = o_{\mathbb{P}}(A_n)$, if $d(X_n, 0)/A_n = o_{\mathbb{P}}(1)$, and (b) $X_n = O_{\mathbb{P}}(A_n)$, if $d(X_n, 0)/A_n = O_{\mathbb{P}}(1)$. \square

§20.26 **Property (Exercise).** For real random variables the following properties hold:

- (i) $o_{\mathbb{P}}(1) + o_{\mathbb{P}}(1) = o_{\mathbb{P}}(1)$ meaning if $X_n = o_{\mathbb{P}}(1)$ and $Y_n = o_{\mathbb{P}}(1)$ then $X_n + Y_n = o_{\mathbb{P}}(1)$;
- (ii) $O_{\mathbb{P}}(1) + o_{\mathbb{P}}(1) = O_{\mathbb{P}}(1)$;
- (iii) $O_{\mathbb{P}}(1) \cdot o_{\mathbb{P}}(1) = o_{\mathbb{P}}(1)$;
- (iv) $(1 + o_{\mathbb{P}}(1))^{-1} = O_{\mathbb{P}}(1)$;
- (v) $o_{\mathbb{P}}(O_{\mathbb{P}}(1)) = o_{\mathbb{P}}(1)$ meaning if $X_n = O_{\mathbb{P}}(1)$ and $Y_n = o_{\mathbb{P}}(X_n)$ then $Y_n = o_{\mathbb{P}}(1)$. \square

§21 Conditional expectation

In the reminder of this section let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space, \mathbb{E} be the expectation with respect to \mathbb{P} and $\mathcal{F} \subseteq \mathcal{A}$ be a sub- σ -field of \mathcal{A} .

§21.01 **Notation.** We write shortly $X \in \overline{\mathcal{A}}^+$, if X is a positive numerical random variable on (Ω, \mathcal{A}) , i.e., $X : \Omega \rightarrow \overline{\mathbb{R}}^+$ is a \mathcal{A} - $\overline{\mathcal{B}}^+$ -measurable function. In particular, we have $\overline{\mathcal{F}}^+ \subseteq \overline{\mathcal{A}}^+$ and for $Y \in \overline{\mathcal{F}}^+$ its expectation $\mathbb{E}(Y)$ is well-defined. \square

§21.02 **Property.** For every $X \in \overline{\mathcal{A}}^+$ exists $Y \in \overline{\mathcal{F}}^+$ with $\mathbb{E}(\mathbf{1}_F Y) = \mathbb{E}(\mathbf{1}_F X)$ for all $F \in \mathcal{F}$, where Y is unique up to \mathbb{P} -a.s. equality. \square

§21.03 **Definition.** A map $Y : \Omega \rightarrow \overline{\mathbb{R}}^+$ is called a (version of the) *conditional expectation* of $X \in \overline{\mathcal{A}}^+$ given \mathcal{F} , symbolically $\mathbb{E}(X | \mathcal{F}) := Y$, if

- (CE1) Y is \mathcal{F} - $\overline{\mathcal{B}}^+$ -measurable, hence $Y \in \overline{\mathcal{F}}^+$ and
- (CE2) $\mathbb{E}(\mathbf{1}_F Y) = \mathbb{E}(\mathbf{1}_F X)$ for any $F \in \mathcal{F}$.

Any map $\mathbb{E}(\bullet | \mathcal{F}) : \overline{\mathcal{A}}^+ \rightarrow \overline{\mathcal{F}}^+$ with $X \mapsto \mathbb{E}(X | \mathcal{F})$ is called (version of the) *conditional expectation* with respect to \mathbb{P} given \mathcal{F} . It implies a map $\mathbb{P}(\bullet | \mathcal{F}) : \mathcal{A} \rightarrow \overline{\mathcal{F}}^+$ with $A \mapsto \mathbb{P}(A | \mathcal{F}) := \mathbb{E}(\mathbf{1}_A | \mathcal{F})$ called (version of the) *conditional distribution* of \mathbb{P} given \mathcal{F} . Exploiting (CE2) every version satisfies $\mathbb{E}(\mathbf{1}_F \mathbb{P}(A | \mathcal{F})) = \int_F \mathbb{P}(A | \mathcal{F}) d\mathbb{P} = \mathbb{P}(F \cap A)$ for all $F \in \mathcal{F}$ and $A \in \mathcal{A}$. \square

§21.04 **Reminder.** Let $X \in \overline{\mathcal{A}}$ be a numerical random variable. Considering the decomposition $X = X^+ - X^-$ with $X^+, X^- \in \overline{\mathcal{A}}^+$ we define for X with $\mathbb{P}(|X|) < \infty$, hence $\mathbb{E}(X^+) < \infty$ and $\mathbb{E}(X^-) < \infty$, the expectation $\mathbb{E}(X) := \mathbb{E}(X^+) - \mathbb{E}(X^-)$. Keep in mind that $\mathcal{L}_1(\mathcal{A}, \mathbb{P}) := \{X \in \overline{\mathcal{A}} : \mathbb{E}(|X|) < \infty\}$ and $\mathbb{E} : \mathcal{L}_1(\mathcal{A}, \mathbb{P}) \rightarrow \mathbb{R}$ denotes the uniquely determined expectation with respect to \mathbb{P} . Note that $\overline{\mathcal{F}} \subseteq \overline{\mathcal{A}}$ implies $\mathcal{L}_1(\mathcal{F}, \mathbb{P}) \subseteq \mathcal{L}_1(\mathcal{A}, \mathbb{P})$. Let $X \in \mathcal{L}_1(\mathcal{A}, \mathbb{P})$, and hence $\mathbb{E}(X^+) < \infty$ and for any version $\mathbb{E}(X^+|\mathcal{F})$ holds (CE1), $\mathbb{E}(X^+|\mathcal{F}) \in \overline{\mathcal{F}}^+$ and (CE2), $\mathbb{E}(\mathbb{1}_F \mathbb{E}(X^+|\mathcal{F})) = \mathbb{E}(\mathbb{1}_F X^+)$ for all $F \in \mathcal{F}$, in particular with $F = \Omega$ also $\mathbb{E}(\mathbb{E}(X^+|\mathcal{F})) = \mathbb{E}(X^+) < \infty$. Therewith, $\mathbb{E}(X^+|\mathcal{F}) \in \mathcal{L}_1(\mathcal{F}, \mathbb{P})$ and analogously also for any version $\mathbb{E}(X^-|\mathcal{F}) \in \mathcal{L}_1(\mathcal{F}, \mathbb{P})$. Consequently, $\mathbb{E}(X^+|\mathcal{F}) - \mathbb{E}(X^-|\mathcal{F}) \in \mathcal{L}_1(\mathcal{F}, \mathbb{P})$ satisfies (CE2) too. \square

§21.05 **Definition.** For $X \in \mathcal{L}_1(\mathcal{A}, \mathbb{P})$ and each version $\mathbb{E}(X^+|\mathcal{F}), \mathbb{E}(X^-|\mathcal{F}) \in \mathcal{L}_1(\mathcal{F}, \mathbb{P})$ we call $\mathbb{E}(X|\mathcal{F}) := \mathbb{E}(X^+|\mathcal{F}) - \mathbb{E}(X^-|\mathcal{F}) \in \mathcal{L}_1(\mathcal{F}, \mathbb{P})$ a (version of the) *conditional expectation* of X given \mathcal{F} . Any map

$$\mathbb{E}(\bullet|\mathcal{F}) : \mathcal{L}_1(\mathcal{A}, \mathbb{P}) \rightarrow \mathcal{L}_1(\mathcal{F}, \mathbb{P}) \text{ with } X \mapsto \mathbb{E}(X|\mathcal{F}) := \mathbb{E}(X^+|\mathcal{F}) - \mathbb{E}(X^-|\mathcal{F})$$

is called a (version of the) *conditional expectation* with respect to \mathbb{P} given \mathcal{F} . \square

§21.06 **Remark.** Due to **Property** §21.02 versions of the conditional expectation of $X \in \overline{\mathcal{A}}^+$ or $X \in \mathcal{L}_1(\mathcal{A}, \mathbb{P})$ given \mathcal{F} differ only on null sets. This property does in generally not extend to the version of the conditional expectation with respect to \mathbb{P} given \mathcal{F} , since for each X we obtain a null set, and their union in general is not a null set. \square

§21.07 **Definition.** Let $(\Omega_1, \mathcal{A}_1), (\Omega_2, \mathcal{A}_2)$ be measurable spaces. A map $\kappa : \Omega_1 \times \mathcal{A}_2 \rightarrow \mathbb{R}^+$ is called *Markov kernel* (from $(\Omega_1, \mathcal{A}_1)$ to $(\Omega_2, \mathcal{A}_2)$), if

(MK1) $A_2 \mapsto \kappa(\omega_1, A_2)$ is for all $\omega_1 \in \Omega_1$ a probability measure on $(\Omega_2, \mathcal{A}_2)$, symbolically $\kappa(\omega_1, \bullet) \in \mathcal{W}(\mathcal{A}_2)$;

(MK2) $\omega_1 \mapsto \kappa(\omega_1, A_2)$ is \mathcal{A}_1 - \mathcal{B} -measurable for all $A_2 \in \mathcal{A}_2$, symbolically $\kappa(\bullet, A_2) \in \mathcal{A}_1^+$. \square

§21.08 **Notation.** Consider a probability space $(\Omega_1, \mathcal{A}_1, \mathbb{P})$, a measurable space $(\Omega_2, \mathcal{A}_2)$ and a Markov kernel κ (from $(\Omega_1, \mathcal{A}_1)$ to $(\Omega_2, \mathcal{A}_2)$). Then there exists an unique probability measure $\kappa \odot \mathbb{P}$ on $(\Omega_2 \times \Omega_1, \mathcal{A}_2 \otimes \mathcal{A}_1)$ determined by

$$\kappa \odot \mathbb{P}(A_2 \times A_1) = \int_{A_1} \kappa(\omega_1, A_2) \mathbb{P}(d\omega_1), \quad \text{for all } A_1 \in \mathcal{A}_1, A_2 \in \mathcal{A}_2.$$

If $f \in \overline{\mathcal{A}_2 \otimes \mathcal{A}_1}^+$ or $f \in \mathcal{L}_1(\kappa \odot \mathbb{P})$ then

$$\kappa \odot \mathbb{P} f = \int_{\Omega_2 \times \Omega_1} f(\omega_2, \omega_1) \kappa \odot \mathbb{P}(d\omega_2, d\omega_1) = \int_{\Omega_1} \int_{\Omega_2} f(\omega_2, \omega_1) \kappa(\omega_1, d\omega_2) \mathbb{P}(d\omega_1).$$

Furthermore, we denote by $\kappa \mathbb{P}$ the marginal distribution on $(\Omega_2, \mathcal{A}_2)$ induced by $\kappa \odot \mathbb{P}$, i.e. $\kappa \mathbb{P}(A_2) = \kappa \odot \mathbb{P}(A_2 \times \Omega_1) = \int_{\Omega_1} \kappa(\omega_1, A_2) \mathbb{P}(d\omega_1)$ for all $A_2 \in \mathcal{A}_2$. \square

§21.09 **Definition.**

(a) $\mathbb{P}(\bullet|\mathcal{F})$ is called *regular* (version of the) conditional distribution of \mathbb{P} given \mathcal{F} , if $(\omega, A) \mapsto \mathbb{P}(A|\mathcal{F})(\omega)$ satisfies the conditions (MK1) and (MK2), i.e. $\mathbb{P}(\bullet|\mathcal{F})$ is a Markov kernel (from (Ω, \mathcal{F}) to (Ω, \mathcal{A})).

(b) $\mathbb{E}(\bullet|\mathcal{F})$ is called *regular* (version of the) conditional expectation with respect to \mathbb{P} given \mathcal{F} , if the implied conditional distribution $\mathbb{P}(\bullet|\mathcal{F})$ of \mathbb{P} given \mathcal{F} is regular, and for each $\omega \in \Omega$ is $X \mapsto \mathbb{E}(X|\mathcal{F})(\omega)$ the expectation with respect to $\mathbb{P}(\bullet|\mathcal{F})(\omega)$. \square

§21.10 **Property.**

- (i) Each regular conditional distribution of \mathbb{P} given \mathcal{F} is implied by a regular conditional expectation with respect to \mathbb{P} given \mathcal{F} .
- (ii) For any probability measure \mathbb{P} on a polish space (Ω, d) endowed with its Borel- σ -algebra \mathcal{B}_Ω and sub- σ -field $\mathcal{F} \subseteq \mathcal{B}_\Omega$ exists a regular conditional distribution of \mathbb{P} given \mathcal{F} . \square

§21.11 **Notation.**

- (i) Let X be a random variable on $(\Omega, \mathcal{A}, \mathbb{P})$ with values in a measurable space $(\mathcal{X}, \mathcal{X})$. For $h \in \mathcal{L}_1(\mathcal{X}, \mathbb{P}^X)$ denotes $\mathbb{E}_x(h|\mathcal{F}) := \mathbb{E}(h(X)|\mathcal{F}) \in \mathcal{L}_1(\mathcal{F}, \mathbb{P})$ a conditional expectation of $h(X)$ given \mathcal{F} and $\mathbb{E}_x(\bullet|\mathcal{F}) : \mathcal{L}_1(\mathcal{X}, \mathbb{P}^X) \rightarrow \mathcal{L}_1(\mathcal{F}, \mathbb{P})$ with $h \mapsto \mathbb{E}_x(h|\mathcal{F})$ a (*regular*) (version of the) *conditional expectation* with respect to \mathbb{P}^X given \mathcal{F} .
- (ii) Let S be a random variable on $(\Omega, \mathcal{A}, \mathbb{P})$ with values in a measurable space $(\mathcal{S}, \mathcal{S})$. For $h \in \mathcal{L}_1(\overline{\mathcal{A}}, \mathbb{P})$ we call $\mathbb{E}(h|\sigma(S)) \in \mathcal{L}_1(\sigma(S), \mathbb{P})$ be a conditional expectation of h given $\mathcal{F} = \sigma(S)$. Keeping $\mathbb{E}(h|\sigma(S)) \in \overline{\sigma(S)}$ in mind and applying **Property** §19.02 (iv) there is $\varphi \in \overline{\mathcal{S}}$ with $\mathbb{E}(h|\sigma(S)) = \varphi(S)$, that is, $\mathbb{E}(h|\sigma(S))(\omega) = \varphi(S(\omega))$, $\omega \in \Omega$. Then $\mathbb{E}(h|S) := \varphi \in \mathcal{L}_1(\mathcal{S}, \mathbb{P}^S)$ and $\mathbb{E}(h|S = s) := \varphi(s) \in \overline{\mathbb{R}}$ is called a (version of the) *conditional expectation* of h given S respectively $S = s$, and $\mathbb{E}(\bullet|S) : \mathcal{L}_1(\mathcal{A}, \mathbb{P}) \rightarrow \mathcal{L}_1(\mathcal{S}, \mathbb{P}^S)$ with $X \mapsto \mathbb{E}(X|S)$ a (*regular*) (version of the) *conditional expectation* with respect to \mathbb{P} given S .
- (iii) Let $(X, S) : (\Omega, \mathcal{A}) \rightarrow (\mathcal{X} \times \mathcal{S}, \mathcal{X} \otimes \mathcal{S})$ with joint distribution $\mathbb{P}^{(X,S)}$. We denote by $\Pi_x : \mathcal{X} \times \mathcal{S} \rightarrow \mathcal{X}$ and $\Pi_s : \mathcal{X} \times \mathcal{S} \rightarrow \mathcal{S}$ with $(x, s) \mapsto \Pi_x(x, s) := x$ and $(x, s) \mapsto \Pi_s(x, s) := s$, respectively, the corresponding coordinate maps. The marginal distribution of X respectively S is given by $\mathbb{P}^X = \mathbb{P} \circ X^{-1} = \mathbb{P} \circ \Pi_x^{-1}(X, S) = \mathbb{P}^{(X,S)} \circ \Pi_x^{-1}$ respectively $\mathbb{P}^S = \mathbb{P}^{(X,S)} \circ \Pi_s^{-1}$. For each version $\mathbb{P}^{(X,S)}(\bullet|\sigma(\Pi_s))$ of the conditional distribution with respect to $\mathbb{P}^{(X,S)}$ given $\sigma(\Pi_s)$, the map

$$\mathbb{P}^X(\bullet|S) : \mathcal{X} \rightarrow \overline{\mathcal{S}} \text{ with } B \mapsto \mathbb{P}^X(B|S) := \varphi \text{ determined by}$$

$$\mathbb{P}^X(B|\sigma(\Pi_s)) = \mathbb{P}^{(X,S)}(\Pi_x^{-1}(B)|\sigma(\Pi_s)) = \varphi(\Pi_s)$$

and analogously $\mathbb{P}^X(\bullet|S = s)$ is called (version of the) *conditional distribution* of X given S respectively $S = s$. We call a version *regular*, if $(s, B) \mapsto \mathbb{P}^X(B|S = s)$ is a Markov kernel (from $(\mathcal{S}, \mathcal{S})$ to $(\mathcal{X}, \mathcal{X})$), where due to **Definition** §21.03 (CE2) $\mathbb{P}^X(\bullet|S) \odot \mathbb{P}^S = \mathbb{P}^{(X,S)}$ (see **Notation** §21.08). Analogously, for $h \in \mathcal{L}_1(\mathcal{X}, \mathbb{P}^X)$ we define a (regular) version $\mathbb{E}_x(h|S) \in \mathcal{L}_1(\mathcal{S}, \mathbb{P}^S)$ and $\mathbb{E}_x(h|S = s) \in \overline{\mathbb{R}}$ of the conditional expectation of h given S respectively $S = s$. If $\mathbb{P}^X(\bullet|S)$ is a regular conditional distribution of X given S and for $s \in \mathcal{S}$ the probability measure $\mathbb{P}^X(\bullet|S = s)$ has for example a finite first absolute moment, i.e., $\mathbb{P}^X(\bullet|S = s) \in \mathcal{W}_1(\mathcal{B}^n)$ (see **Notation** §19.05) then $\mathbb{E}(X|S = s) = \mathbb{E}_x(\text{id}_X|S = s) = \int_{\mathcal{X}} x \mathbb{P}^X(dx|S = s)$.

- (iv) Suppose the joint distribution $\mathbb{P}^{(X,S)}$ is dominated by a product measure $\mu \otimes \nu$ where μ and ν is a σ -finite measure on \mathcal{X} and \mathcal{S} , respectively, $\mu \in \mathcal{M}_\sigma(\mathcal{X})$ and $\nu \in \mathcal{M}_\sigma(\mathcal{S})$ for short. Let $f^{(X,S)}$ denote a $(\mu \otimes \nu)$ -density of $\mathbb{P}^{(X,S)}$. A μ - and ν -density of the marginal distribution \mathbb{P}^X and \mathbb{P}^S is given by $f^X : x \mapsto \int_{\mathcal{S}} f^{(X,S)}(x, s) \nu(ds)$ and $f^S : s \mapsto \int_{\mathcal{X}} f^{(X,S)}(x, s) \mu(dx)$, respectively. The $f^{X|S} : \mathcal{S} \times \mathcal{X} \rightarrow \overline{\mathbb{R}}^+$ with

$$(s, x) \mapsto f^{X|S=s}(x) = \frac{f^{(X,S)}(x, s)}{f^S(s)} \mathbb{1}_{\{f^S(s) > 0\}} + f^X(x) \mathbb{1}_{\{f^S(s) = 0\}}$$

belongs to $\overline{\mathcal{S} \otimes \mathcal{X}^+}$ and it is a μ -density of the Markov kernel $\mathbb{P}^{X|S}$ from $(\mathcal{S}, \mathcal{S})$ to $(\mathcal{X}, \mathcal{X})$ defined by $(s, B) \mapsto \mathbb{P}^{X|S=s}(B) := \int_B \mathfrak{f}^{X|S=s}(x) \mu(dx)$. We call $\mathfrak{f}^{X|S=s}$ conditional density of X given $S = s$.

- (v) As an example let $(X, S) \in \mathcal{B}^{k+l}$ be multivariate normally distributed with $\text{Cov}(X, S) = \Sigma_{XS}$ and marginal distributions $X \sim N_{(\mu_X, \Sigma_X)}$ and $S \sim N_{(\mu_S, \Sigma_S)}$, i.e.,

$$\begin{pmatrix} X \\ S \end{pmatrix} \sim N_{(\mu, \Sigma)} \text{ with } \mu = \begin{pmatrix} \mu_X \\ \mu_S \end{pmatrix} \in \mathbb{R}^{k+l} \text{ and } \Sigma = \begin{pmatrix} \Sigma_X & \Sigma_{XS} \\ \Sigma_{XS}^t & \Sigma_S \end{pmatrix}.$$

Assuming $\Sigma > 0$ the joint distribution $\mathbb{P}^{(X,S)}$ admits a density with respect to the Lebesgue measure λ^{k+l} on $(\mathbb{R}^{k+l}, \mathcal{B}^{k+l})$. For each $s \in \mathbb{R}^l$ the conditional density $\mathfrak{f}^{X|S=s}$ as in (iv) is a density of the multivariate normal distribution $N_{(\mu_{X|S=s}, \Sigma_{X|S=s})}$ -distribution with

$$\mu_{X|S=s} := \mu_X + \Sigma_{XS} \Sigma_S^{-1} (s - \mu_S) \in \mathbb{R}^k \text{ und } \Sigma_{X|S=s} := \Sigma_X - \Sigma_{XS} \Sigma_S^{-1} \Sigma_{SX} > 0$$

which is thus a regular conditional distribution of X given $S = s$. □

§21.12 **Property.** Let $X, Y \in \mathcal{L}_1(\mathcal{A}, \mathbb{P})$ and $\mathcal{F} \subseteq \mathcal{A}$ be a sub- σ -field. Any version of the conditional expectation satisfies the following properties \mathbb{P} -a.s.:

- (i) For all $a, b \in \mathbb{R}$ holds $\mathbb{E}(aX + bY | \mathcal{F}) = a\mathbb{E}(X | \mathcal{F}) + b\mathbb{E}(Y | \mathcal{F})$; (linear)
- (ii) For $X \leq Y$ holds $\mathbb{E}(X | \mathcal{F}) \leq \mathbb{E}(Y | \mathcal{F})$; (monotone)
- (iii) $|\mathbb{E}(X | \mathcal{F})| \leq \mathbb{E}(|X| | \mathcal{F})$; (triangular inequality)
- (iv) For $S \in \overline{\mathcal{A}}$ with $\mathbb{E}(|S| | \mathcal{F}) < \infty$ holds $\mathbb{P}(|S| < \infty) = 1$. (finite)
- (v) For $\phi : \mathbb{R} \rightarrow \mathbb{R}$ convex with $\phi(X) \in \mathcal{L}_1(\overline{\mathcal{A}}, \mathbb{P})$ holds $\phi(\mathbb{E}(X | \mathcal{F})) \leq \mathbb{E}(\phi(X) | \mathcal{F})$. (Jensen's inequality)
- (vi) For $X_n \uparrow X$ \mathbb{P} -a.s. holds $\sup_{n \in \mathbb{N}} \mathbb{E}(X_n | \mathcal{F}) = \mathbb{E}(X | \mathcal{F})$. (monotone convergence)
- (vii) For $X_n \rightarrow X$ \mathbb{P} -a.s. with $|X_n| \leq Y$, $n \in \mathbb{N}$, holds $\lim_{n \rightarrow \infty} \mathbb{E}(X_n | \mathcal{F}) = \mathbb{E}(X | \mathcal{F})$ \mathbb{P} -a.s. and in $\mathcal{L}_1(\mathcal{A}, \mathbb{P})$. (dominated convergence)

If the version is regular, i.e., $\mathbb{E}(\bullet | \mathcal{F})(\omega)$ is an expectation for all $\omega \in \Omega$, then the statements (i)-(vii) holds for all $\omega \in \Omega$. □

§21.13 **Property.** Let $X, Y \in \mathcal{L}_1(\mathcal{A}, \mathbb{P})$ and $\mathcal{G} \subseteq \mathcal{F} \subseteq \mathcal{A}$ sub- σ -fields. Any version of the conditional expectation satisfies the following properties \mathbb{P} -a.s.:

- (i) For $\mathbb{E}(|XY|) < \infty$ and $Y \in \mathcal{F}$ holds
$$\mathbb{E}(XY | \mathcal{F}) = Y\mathbb{E}(X | \mathcal{F}) \text{ and } \mathbb{E}(Y | \mathcal{F}) = \mathbb{E}(Y | \sigma(Y)) = Y;$$
- (ii) $\mathbb{E}(\mathbb{E}(X | \mathcal{F}) | \mathcal{G}) = \mathbb{E}(\mathbb{E}(X | \mathcal{G}) | \mathcal{F}) = \mathbb{E}(X | \mathcal{G})$; (tower property)
- (iii) If $\sigma(X)$ and \mathcal{F} are independent, then $\mathbb{E}(X | \mathcal{F}) = \mathbb{E}(X)$; (independence)
- (iv) $\mathbb{E}(\mathbb{E}(X | \mathcal{F})) = \mathbb{E}(X)$. (total probability)
- (v) For $\overline{\mathcal{F}} := \{A \in \mathcal{A} | \mathbb{P}(A) \in \{0, 1\}\}$ holds $\mathbb{E}(X | \overline{\mathcal{F}}) = \mathbb{E}(X)$. □

§21.14 **Property.** Let $\mathcal{F} \subseteq \mathcal{A}$ be a sub- σ -field and $\mathbb{E}(\bullet | \mathcal{F})$ be a conditional expectation.

- (i) $\mathbb{E}(\bullet | \mathcal{F}) : \mathcal{L}_2(\mathcal{A}, \mathbb{P}) \rightarrow \mathcal{L}_2(\mathcal{F}, \mathbb{P})$ is an orthogonal projection, that is, for all $X \in \mathcal{L}_2(\mathcal{A}, \mathbb{P})$ and $Y \in \mathcal{L}_2(\mathcal{F}, \mathbb{P})$ holds

$$\|X - Y\|_{\mathcal{L}_2(\mathbb{P})}^2 = \mathbb{E}(|X - Y|^2) \geq \mathbb{E}(|X - \mathbb{E}(X | \mathcal{F})|^2) = \|X - \mathbb{E}(X | \mathcal{F})\|_{\mathcal{L}_2(\mathbb{P})}^2,$$

where equality holds if and only if $Y = \mathbb{E}(X | \mathcal{F})$ \mathbb{P} -a.s..

- (ii) $\mathbb{E}(\bullet | \mathcal{F}) : \mathcal{L}_s(\mathcal{A}, \mathbb{P}) \rightarrow \mathcal{L}_s(\mathcal{F}, \mathbb{P})$ is a contraction for $s \in [1, \infty]$, i.e., $\|\mathbb{E}(X | \mathcal{F})\|_{\mathcal{L}_s(\mathbb{P})} \leq \|X\|_{\mathcal{L}_s(\mathbb{P})}$, and thus bounded and continuous. If $(X_n)_{n \in \mathbb{N}}$ converges in $\mathcal{L}_s(\mathcal{A}, \mathbb{P})$, then $(\mathbb{E}(X_n | \mathcal{F}))_{n \in \mathbb{N}}$ converges in $\mathcal{L}_s(\mathcal{F}, \mathbb{P})$. \square

Bibliography

- A. Barron, L. Birgé, and P. Massart. Risk bounds for model selection via penalization. *Probability Theory and Related Fields*, 113(3):301–413, 1999.
- L. Birgé. An alternative point of view on Lepskij’s method. In Monogr., editor, *State of the art in probability and statistics*, volume 36 of *IMS Lecture Notes*, pages 113–133. (Leiden 1999), 2001.
- L. D. Brown and M. G. Low. A constrained risk inequality with applications to nonparametric functional estimation. *The Annals of Statistics*, 24(6):2524–2535, 1996.
- N. L. Carr. Kinetics of catalytic isomerization of n-pentane. *Industrial and Engineering Chemistry*, 52:391–396, 1960.
- F. Comte. *Estimation non-paramétrique*. Spartacus-idh, Paris, 2015.
- A. Goldenshluger and O. Lepskij. Bandwidth selection in kernel density estimation: Oracle inequalities and adaptive minimax optimality. *The Annals of Statistics*, 39:1608–1632, 2011.
- T. Kawata. *Fourier analysis in probability theory*. Academic Press, New York, 1972.
- T. Klein and E. Rio. Concentration around the mean for maxima of empirical processes. *The Annals of Probability*, 33(3):1060–1077, 2005.
- A. Klenke. *Probability theory. A comprehensive course*. London: Springer, 2008.
- A. Klenke. *Wahrscheinlichkeitstheorie*. Springer Spektrum, 3., überarbeitete und ergänzte Auflage, 2012.
- B. Laurent and P. Massart. Adaptive estimation of a quadratic functional by model selection. *The Annals of Statistics*, 28(5):1302–1338, 2000.
- B. Laurent, C. Ludena, and C. Prieur. Adaptive estimation of linear functionals by model selection. *Electronic Journal of Statistics*, 2(993-1020), 2008.
- L. Le Cam. Convergence of Estimates Under Dimensionality Restrictions. *The Annals of Statistics*, 1(1):38–53, 1973.
- O. Lepskij. On a problem of adaptive estimation in Gaussian white noise. *Theory of Probability and its Applications*, 35(3):454–466, 1990.
- O. Lepskij. Asymptotically minimax adaptive estimation. I: Upper bounds. Optimally adaptive estimates. *Theory of Probability and its Applications*, 36(4):682–697, 1991.
- O. Lepskij. Asymptotically minimax adaptive estimation. II: Schemes without optimal adaption: Adaptive estimators. *Theory of Probability and its Applications*, 37(3):433–448, 1992a.
- O. Lepskij. On problems of adaptive estimation in white Gaussian noise. In *Topics in non-parametric estimation*, volume 12 of *Pap. Semin. Math. Stat.*, pages 87–106, Moscow/USSR, 1992b. Adv. Sov. Math.

- H. B. Mann and D. R. Whitney. On a test of whether one of two random variables is stochastically larger than the other. *Annals of Mathematical Statistics*, 18(1):50–60, 1947.
- P. Massart. *Concentration inequalities and model selection. Ecole d'Été de Probabilités de Saint-Flour XXXIII – 2003*. Lecture Notes in Mathematics 1896. Berlin: Springer., 2007.
- M. Talagrand. New concentration inequalities in product spaces. *Inventiones mathematicae*, 126:505–563, 1996.
- A. B. Tsybakov. *Introduction to nonparametric estimation*. Springer Series in Statistics. Springer, New York, 2009.
- A. W. van der Vaart. *Asymptotic statistics*. Cambridge University Press, 1998.
- F. Wilcoxon. Individual comparisons by ranking methods. *Biometrics Bulletin*, 1(6):80–83, 1945.
- H. Witting. *Mathematische Statistik I: Parametrische Verfahren bei festem Stichprobenumfang*. Stuttgart: B. G. Teubner, 1985.
- H. Witting and U. Müller-Funk. *Mathematische Statistik II. Asymptotische Statistik: Parametrische Modelle und nichtparametrische Funktionale*. Stuttgart: B. G. Teubner, 1995.