

Vladimir Spokoiny

Script to VL WS2014

February 16, 2015

Springer

Berlin Heidelberg New York

Hong Kong London

Milan Paris Tokyo

Contents

1	Quasi maximum likelihood estimation in linear models	7
1.1	Linear Modeling	7
1.1.1	Estimation under homogeneous noise assumption	9
1.1.2	Linear basis transformation	10
1.1.3	Orthogonal and orthonormal design	12
1.1.4	Spectral representation	13
1.2	Properties of the response estimate $\tilde{\mathbf{f}}$	14
1.2.1	Decomposition into a deterministic and a stochastic component	14
1.2.2	Properties of the operator Π	15
1.2.3	Quadratic loss and risk of the response estimation	16
1.2.4	Misspecified “colored noise”	17
1.3	Properties of the MLE $\tilde{\boldsymbol{\theta}}$	18
1.3.1	Properties of the stochastic component	18
1.3.2	Properties of the deterministic component	19
1.3.3	Risk of estimation. R-efficiency	21
1.3.4	The case of a misspecified noise	23
1.4	Linear models and quadratic log-likelihood	24
1.4.1	Inference based on the maximum likelihood	27
1.4.2	A misspecified LPA	30
1.4.3	A misspecified noise structure	30
1.5	Random design regression	33
1.5.1	Independent measurements	33
1.5.2	Aggregated random design	36
1.5.3	Application to instrumental regression	38
1.6	Matrix Bernstein inequality	40

2	Sieve model selection in linear models	41
2.1	Projection estimation. Loss and risk	41
2.1.1	A linear model	41
2.1.2	Linear decomposition	42
2.1.3	Inhomogeneous errors	42
2.1.4	Linear decomposition	43
2.1.5	Quadratic loss. Bias-variance decomposition	43
2.1.6	Projection estimation and the model choice problem	44
2.2	Unbiased risk estimation	46
2.2.1	AIC and pairwise comparison	48
2.2.2	Pairwise analysis	50
2.2.3	Uniform bounds and the zone of insensitivity	52
2.2.4	A bound on the excess	53
2.3	The approach based on multiple testing. “Smallest accepted” rule	55
2.3.1	A LR test	56
2.3.2	Multiplicity correction	57
2.3.3	Definition of the oracle and propagation property	58
2.3.4	A bound on the loss	59
2.3.5	Role of β	61
3	Linear smoothers	63
3.1	Regularization and ridge regression	64
3.2	Penalized likelihood. Bias and variance	64
3.3	Inference for the penalized MLE	67
3.4	Projection and shrinkage estimates	68
3.5	Smoothness constraints and roughness penalty approach	70
3.6	Shrinkage in a linear inverse problem	71
3.7	Spectral cut-off and spectral penalization. Diagonal estimates	71
4	Ordered model selection for linear smoothers	75
4.1	Model and problem	75
4.1.1	Loss and risk	75
4.1.2	Oracle choice. Bias-variance trade-off	77
4.2	Unbiased risk estimation	78
4.2.1	Zone of insensitivity	79
4.2.2	An oracle bound	82
4.3	Smallest accepted (SmA) method	83
4.3.1	Decomposition of the test statistic. Bias and variance	83

4.3.2	Multiplicity correction. A Bonferroni bound and Monte-Carlo method.....	84
4.3.3	Oracle choice.....	85
4.3.4	Data-driven choice and the oracle inequality.....	85
4.3.5	Analysis of the payment for adaptation $\bar{z}(m^*)$	88
4.3.6	Choice of β and \mathbf{x}	88
4.3.7	Power loss function.....	88
4.3.8	Penalized MLE.....	92
4.4	Lepski's method.....	96
4.5	Intersection of confidence sets (ICI) method.....	98
5	Unordered case. Anisotropic sets and subset selection.....	101
5.1	Subset selection procedure.....	101
5.1.1	SmA procedure and multilevel synchronization.....	102
5.1.2	Prediction loss.....	105
5.1.3	Estimation loss.....	105
5.1.4	Linear functional estimation.....	106
5.1.5	Subset selection problem.....	107
5.2	Anisotropic models.....	109
6	Fisher and Wilks expansion.....	113
6.1	Main results.....	114
6.2	Conditions.....	115
6.3	Properties of the MLE $\tilde{\theta}$	117
6.4	Critical dimension. Examples.....	119
6.4.1	Linear and generalized linear models.....	120
6.4.2	Generalized linear models (GLM).....	120
6.4.3	I.i.d. case.....	122
6.5	Some auxiliary results and proofs.....	123
6.5.1	Local linear approximation of the gradient of the log-likelihood.....	123
6.5.2	Local quadratic approximation of the log-likelihood.....	125
6.5.3	Proof of Theorem 6.3.1.....	125
6.5.4	Proof of Theorem 6.3.2.....	126
6.5.5	Proof of Theorem 6.3.3.....	127
6.5.6	An entropy bound for the maximum of a random process.....	127
6.5.7	A bound for the norm of a vector random process.....	129
6.5.8	A deviation bound for the quadratic form $\ \xi\ ^2$	130
6.6	Some results for the normal law.....	131

6	Contents	
	6.6.1	Deviation bounds 131
	6.6.2	Gaussian comparison via KL-divergence and Pinsker’s inequality . . 135
7	Sieve Model Selection 139
	7.1	Sieve SmA procedure 139
	7.2	Resampling methods for parameter tuning in generalized regression 141
	7.2.1	Generalized regression 141
	7.2.2	Multiplier bootstrap 142
	7.2.3	Numerical issues 144
	7.3	Why does it work? Linear Gaussian case 145
	7.3.1	Small modeling bias condition 145
	7.3.2	The “large bias” case 152
	7.3.3	Bootstrap and the SmA procedure 154
	7.4	Linear non-Gaussian case and GAR 157
	7.5	Sieve Generalized Linear regression 157
	7.5.1	Sieve MLE 159
	7.5.2	Bootstrap counterpart 159
	7.5.3	Bootstrap for the SmA procedure 160
8	SmA and parameter tuning in high dimensional regression 163
	8.1	SmA subset selection in high dimensional regression 164
	8.1.1	Norm comparison for a family of Gaussian vectors 165
	8.2	Gaussian comparison in high dimension 167
	8.2.1	Stein identity, Slepian bridge, and Gaussian comparison 167
	8.2.2	Comparing of the maximum of Gaussians 170
	8.2.3	Anti-concentration for Gaussian maxima 171
9	Penalized model selection 175
	9.1	Complexity penalization 175
	9.1.1	Bootstrap based tuning 178
	9.2	Sparse penalty 179
	9.2.1	Basic inequality 181
	9.2.2	Dual problem and Danzig selector 182
	9.2.3	Data-driven choice of λ 182
	References 183

Quasi maximum likelihood estimation in linear models

1.1 Linear Modeling

A linear model assumes that the observations Y_i follow the equation:

$$Y_i = \Psi_i^\top \boldsymbol{\theta}^* + \varepsilon_i \quad (1.1)$$

for $i = 1, \dots, n$, where $\boldsymbol{\theta}^* = (\theta_1^*, \dots, \theta_p^*)^\top \in \mathbb{R}^p$ is an unknown parameter vector, Ψ_i are given vectors in \mathbb{R}^p and the ε_i 's are individual errors with zero mean. A typical example is given by linear regression (see Section ??) when the vectors Ψ_i are the values of a set of functions (e.g polynomial, trigonometric) series at the design points X_i .

A linear Gaussian model assumes in addition that the vector of errors $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top$ is normally distributed with zero mean and a covariance matrix Σ :

$$\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \Sigma).$$

In this chapter we suppose that Σ is given in advance. We will distinguish between three cases:

1. the errors ε_i are i.i.d. $\mathcal{N}(0, \sigma^2)$, or equivalently, the matrix Σ is equal to $\sigma^2 \mathbf{I}_n$ with \mathbf{I}_n being the unit matrix in \mathbb{R}^n .
2. the errors are independent but not homogeneous, that is, $\mathbb{E}\varepsilon_i^2 = \sigma_i^2$. Then the matrix Σ is diagonal: $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$.
3. the errors ε_i are dependent with a covariance matrix Σ .

In practical applications one mostly starts with the white Gaussian noise assumption and more general cases 2 and 3 are only considered if there are clear indications of the noise inhomogeneity or correlation. The second situation is typical e.g. for the eigenvector decomposition in an inverse problem. The last case is the most general and includes the first two.

Denote by $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ (resp. $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top$) the vector of observations (resp. of errors) in \mathbb{R}^n and by Ψ the $p \times n$ matrix with columns Ψ_i . Let also Ψ^\top denote its transpose. Then the model equation can be rewritten as:

$$\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \Sigma).$$

An equivalent formulation is that $\Sigma^{-1/2}(\mathbf{Y} - \Psi^\top \boldsymbol{\theta})$ is a standard normal vector in \mathbb{R}^n . The log-density of the distribution of the vector $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ w.r.t. the Lebesgue measure in \mathbb{R}^n is therefore of the form

$$\begin{aligned} L(\boldsymbol{\theta}) &= -\frac{n}{2} \log(2\pi) - \frac{\log(\det \Sigma)}{2} - \frac{1}{2} \|\Sigma^{-1/2}(\mathbf{Y} - \Psi^\top \boldsymbol{\theta})\|^2 \\ &= -\frac{n}{2} \log(2\pi) - \frac{\log(\det \Sigma)}{2} - \frac{1}{2} (\mathbf{Y} - \Psi^\top \boldsymbol{\theta})^\top \Sigma^{-1} (\mathbf{Y} - \Psi^\top \boldsymbol{\theta}). \end{aligned}$$

In case 1 this expression can be rewritten as

$$L(\boldsymbol{\theta}) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \Psi_i^\top \boldsymbol{\theta})^2.$$

In case 2 the expression is similar:

$$L(\boldsymbol{\theta}) = -\sum_{i=1}^n \left\{ \frac{1}{2} \log(2\pi\sigma_i^2) + \frac{(Y_i - \Psi_i^\top \boldsymbol{\theta})^2}{2\sigma_i^2} \right\}.$$

The *maximum likelihood estimate* (MLE) $\tilde{\boldsymbol{\theta}}$ of $\boldsymbol{\theta}^*$ is defined by maximizing the log-likelihood $L(\boldsymbol{\theta})$:

$$\tilde{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \mathbb{R}^p}{\operatorname{argmax}} L(\boldsymbol{\theta}) = \underset{\boldsymbol{\theta} \in \mathbb{R}^p}{\operatorname{argmin}} (\mathbf{Y} - \Psi^\top \boldsymbol{\theta})^\top \Sigma^{-1} (\mathbf{Y} - \Psi^\top \boldsymbol{\theta}). \quad (1.2)$$

We omit the other terms in the expression of $L(\boldsymbol{\theta})$ because they do not depend on $\boldsymbol{\theta}$. This estimate is the *least squares estimate* (LSE) because it minimizes the sum of squared distances between the observations Y_i and the linear responses $\Psi_i^\top \boldsymbol{\theta}$. Note that (1.2) is a quadratic optimization problem which has a closed form solution. Differentiating the right hand-side of (1.2) w.r.t. $\boldsymbol{\theta}$ yields the *normal equation*

$$\Psi \Sigma^{-1} \Psi^\top \tilde{\boldsymbol{\theta}} = \Psi \Sigma^{-1} \mathbf{Y}.$$

If the $p \times p$ -matrix $\Psi \Sigma^{-1} \Psi^\top$ is non-degenerate then the normal equation has the unique solution

$$\tilde{\boldsymbol{\theta}} = (\Psi \Sigma^{-1} \Psi^\top)^{-1} \Psi \Sigma^{-1} \mathbf{Y} = \mathcal{S} \mathbf{Y}, \quad (1.3)$$

where

$$\mathcal{S} = (\Psi \Sigma^{-1} \Psi^\top)^{-1} \Psi \Sigma^{-1}$$

is a $p \times n$ matrix. We denote by $\tilde{\theta}_j$ the entries of the vector $\tilde{\boldsymbol{\theta}}$, $j = 1, \dots, p$.

If the matrix $\Psi \Sigma^{-1} \Psi^\top$ is degenerate, then the normal equation has infinitely many solutions. However, one can still apply the formula (1.3) where $(\Psi \Sigma^{-1} \Psi^\top)^{-1}$ is a pseudo-inverse of the matrix $\Psi \Sigma^{-1} \Psi^\top$.

The ML-approach leads to the *parameter estimate* $\tilde{\boldsymbol{\theta}}$. Note that due to the model (1.1), the product $\tilde{\mathbf{f}} = \Psi^\top \tilde{\boldsymbol{\theta}}$ is an estimate of the mean $\mathbf{f}^* \stackrel{\text{def}}{=} \mathbb{E}\mathbf{Y}$ of the vector of observations \mathbf{Y} :

$$\tilde{\mathbf{f}} = \Psi^\top \tilde{\boldsymbol{\theta}} = \Psi^\top (\Psi \Sigma^{-1} \Psi^\top)^{-1} \Psi \Sigma^{-1} \mathbf{Y} = \Pi \mathbf{Y},$$

where

$$\Pi = \Psi^\top (\Psi \Sigma^{-1} \Psi^\top)^{-1} \Psi \Sigma^{-1}$$

is an $n \times n$ matrix (linear operator) in \mathbb{R}^n . The vector $\tilde{\mathbf{f}}$ is called a *prediction* or *response* regression estimate.

Below we study the properties of the estimates $\tilde{\boldsymbol{\theta}}$ and $\tilde{\mathbf{f}}$. In this study we try to address both types of possible model misspecification: due to a wrong assumption about the error distribution and due to a possibly wrong linear parametric structure. Namely we consider the model

$$Y_i = f_i + \varepsilon_i, \quad \varepsilon \sim \mathcal{N}(0, \Sigma_0). \quad (1.4)$$

The response values f_i are usually treated as the value of the regression function $f(\cdot)$ at the design points X_i . The parametric model (1.1) can be viewed as an approximation of (1.4) while Σ is an approximation of the true covariance matrix Σ_0 . If \mathbf{f}^* is indeed equal to $\Psi^\top \boldsymbol{\theta}^*$ and $\Sigma = \Sigma_0$, then $\tilde{\boldsymbol{\theta}}$ and $\tilde{\mathbf{f}}$ are MLEs, otherwise quasi MLEs. In our study we mostly restrict ourselves to the case 1 assumption about the noise ε : $\varepsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$. The general case can be reduced to this one by a simple data transformation, namely, by multiplying the equation (1.4) $\mathbf{Y} = \mathbf{f}^* + \varepsilon$ with the matrix $\Sigma^{-1/2}$, see Section 1.4.1 for more detail.

1.1.1 Estimation under homogeneous noise assumption

If a homogeneous noise is assumed, that is, if $\Sigma = \sigma^2 \mathbf{I}_n$ and $\varepsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$, then the formulae for the MLEs $\tilde{\boldsymbol{\theta}}$, $\tilde{\mathbf{f}}$ slightly simplify. In particular, the variance σ^2 cancels and the resulting estimate is the *ordinary least squares* (oLSE):

$$\tilde{\boldsymbol{\theta}} = (\Psi \Psi^\top)^{-1} \Psi \mathbf{Y} = \mathcal{S} \mathbf{Y}$$

with $\mathcal{S} = (\Psi\Psi^\top)^{-1}\Psi$. Also

$$\tilde{\mathbf{f}} = \Psi^\top (\Psi\Psi^\top)^{-1}\Psi\mathbf{Y} = \Pi\mathbf{Y}$$

with $\Pi = \Psi^\top (\Psi\Psi^\top)^{-1}\Psi$.

Exercise 1.1.1. Derive the formulae for $\tilde{\boldsymbol{\theta}}, \tilde{\mathbf{f}}$ directly from the log-likelihood $L(\boldsymbol{\theta})$ for homogeneous noise.

If the assumption $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2\mathbf{I}_n)$ about the errors is not precisely fulfilled, then the oLSE can be viewed as a quasi MLE.

1.1.2 Linear basis transformation

Denote by $\boldsymbol{\psi}_1^\top, \dots, \boldsymbol{\psi}_p^\top$ the rows of the matrix Ψ . Then the $\boldsymbol{\psi}_i$'s are vectors in \mathbb{R}^n and we call them *the basis vectors*. In the linear regression case the $\boldsymbol{\psi}_i$'s are obtained as the values of the basis functions at the design points. Our linear parametric assumption simply means that the underlying vector \mathbf{f}^* can be represented as a linear combination of the vectors $\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_p$:

$$\mathbf{f}^* = \theta_1^*\boldsymbol{\psi}_1 + \dots + \theta_p^*\boldsymbol{\psi}_p.$$

In other words, \mathbf{f}^* belongs to the linear subspace in \mathbb{R}^n spanned by the vectors $\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_p$. It is clear that this assumption still holds if we select another basis in this subspace.

Let U be any linear orthogonal transformation in \mathbb{R}^p with $UU^\top = \mathbf{I}_p$. Then the linear relation $\mathbf{f}^* = \Psi^\top\boldsymbol{\theta}^*$ can be rewritten as

$$\mathbf{f}^* = \Psi^\top UU^\top\boldsymbol{\theta}^* = \check{\Psi}^\top\mathbf{u}^*$$

with $\check{\Psi} = U^\top\Psi$ and $\mathbf{u}^* = U^\top\boldsymbol{\theta}^*$. Here the columns of $\check{\Psi}$ mean the new basis vectors $\check{\boldsymbol{\psi}}_m$ in the same subspace while \mathbf{u}^* is the vector of coefficients describing the decomposition of the vector \mathbf{f}^* w.r.t. this new basis:

$$\mathbf{f}^* = u_1^*\check{\boldsymbol{\psi}}_1 + \dots + u_p^*\check{\boldsymbol{\psi}}_p.$$

The natural question is how the expression for the MLEs $\tilde{\boldsymbol{\theta}}$ and $\tilde{\mathbf{f}}$ change with the change of the basis. The answer is straightforward. For notational simplicity, we only consider the case with $\Sigma = \sigma^2\mathbf{I}_n$. The model can be rewritten as

$$\mathbf{Y} = \check{\Psi}^\top\mathbf{u}^* + \boldsymbol{\varepsilon}$$

yielding the solutions

$$\tilde{\mathbf{u}} = (\check{\Psi}\check{\Psi}^\top)^{-1}\check{\Psi}\mathbf{Y} = \check{\mathcal{S}}\mathbf{Y}, \quad \tilde{\mathbf{f}} = \check{\Psi}^\top(\check{\Psi}\check{\Psi}^\top)^{-1}\check{\Psi}\mathbf{Y} = \check{\mathcal{H}}\mathbf{Y},$$

where $\check{\Psi} = U^\top\Psi$ implies

$$\begin{aligned}\check{\mathcal{S}} &= (\check{\Psi}\check{\Psi}^\top)^{-1}\check{\Psi} = U^\top\mathcal{S}, \\ \check{\mathcal{H}} &= \check{\Psi}^\top(\check{\Psi}\check{\Psi}^\top)^{-1}\check{\Psi} = \mathcal{H}.\end{aligned}$$

This yields

$$\tilde{\mathbf{u}} = U^\top\tilde{\boldsymbol{\theta}}$$

and moreover, the estimate $\tilde{\mathbf{f}}$ is not changed for any linear transformation of the basis. The first statement can be expected in view of $\boldsymbol{\theta}^* = U\mathbf{u}^*$, while the second one will be explained in the next section: \mathcal{H} is the linear projector on the subspace spanned by the basis vectors and this projector is invariant w.r.t. basis transformations.

Exercise 1.1.2. Consider univariate polynomial regression of degree $p - 1$. This means that f is a polynomial function of degree $p - 1$ observed at the points X_i with errors ε_i that are assumed to be i.i.d. normal. The function f can be represented as

$$f(x) = \theta_1^* + \theta_2^*x + \dots + \theta_p^*x^{p-1}$$

using the basis functions $\psi_j(x) = x^{j-1}$ for $j = 0, \dots, p - 1$. At the same time, for any point x_0 , this function can also be written as

$$f(x) = u_1^* + u_2^*(x - x_0) + \dots + u_p^*(x - x_0)^{p-1}$$

using the basis functions $\check{\psi}_j = (x - x_0)^{j-1}$.

- Write the matrices Ψ and $\Psi\Psi^\top$ and similarly $\check{\Psi}$ and $\check{\Psi}\check{\Psi}^\top$.
- Describe the linear transformation A such that $\mathbf{u} = A\boldsymbol{\theta}$ for $p = 1$.
- Describe the transformation A such that $\mathbf{u} = A\boldsymbol{\theta}$ for $p > 1$.

Hint: use the formula

$$u_j^* = \frac{1}{(j-1)!}f^{(j-1)}(x_0), \quad j = 1, \dots, p$$

to identify the coefficient u_j^* via $\theta_j^*, \dots, \theta_p^*$.

1.1.3 Orthogonal and orthonormal design

Orthogonality of the design matrix Ψ means that the basis vectors ψ_1, \dots, ψ_p are orthonormal in the sense

$$\psi_j^\top \psi_{j'} = \sum_{i=1}^n \psi_{m,i} \psi_{m',i} = \begin{cases} 0 & \text{if } j \neq j', \\ \lambda_j & \text{if } j = j', \end{cases}$$

for some positive values $\lambda_1, \dots, \lambda_p$. Equivalently one can write

$$\Psi\Psi^\top = \Lambda = \text{diag}(\lambda_1, \dots, \lambda_p).$$

This feature of the design is very useful and it essentially simplifies the computation and analysis of the properties of $\tilde{\boldsymbol{\theta}}$. Indeed, $\Psi\Psi^\top = \Lambda$ implies

$$\tilde{\boldsymbol{\theta}} = \Lambda^{-1}\Psi\mathbf{Y}, \quad \tilde{\mathbf{f}} = \Psi^\top\tilde{\boldsymbol{\theta}} = \Psi^\top\Lambda^{-1}\Psi\mathbf{Y}$$

with $\Lambda^{-1} = \text{diag}(\lambda_1^{-1}, \dots, \lambda_p^{-1})$. In particular, the first relation means

$$\tilde{\theta}_j = \lambda_j^{-1} \sum_{i=1}^n Y_i \psi_{j,i},$$

that is, $\tilde{\theta}_j$ is the scalar product of the data and the basis vector ψ_j for $j = 1, \dots, p$. The estimate of the response \mathbf{f} reads as

$$\tilde{\mathbf{f}} = \tilde{\theta}_1\psi_1 + \dots + \tilde{\theta}_p\psi_p.$$

Theorem 1.1.1. *Consider the model $\mathbf{Y} = \Psi^\top\boldsymbol{\theta} + \boldsymbol{\varepsilon}$ with homogeneous errors $\boldsymbol{\varepsilon} : \mathbb{E}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top = \sigma^2\mathbf{I}_n$. If the design Ψ is orthogonal, that is, if $\Psi\Psi^\top = \Lambda$ for a diagonal matrix Λ , then the estimated coefficients $\tilde{\theta}_j$ are uncorrelated: $\text{Var}(\tilde{\boldsymbol{\theta}}) = \sigma^2\Lambda^{-1}$. Moreover, if $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2\mathbf{I}_n)$, then $\tilde{\boldsymbol{\theta}} \sim \mathcal{N}(\boldsymbol{\theta}^*, \sigma^2\Lambda^{-1})$.*

An important message of this result is that the orthogonal design allows for splitting the original multivariate problem into a collection of independent univariate problems: each coefficient θ_j^* is estimated by $\tilde{\theta}_j$ independently on the remaining coefficients.

The calculus can be further simplified in the case of an orthogonal design with $\Psi\Psi^\top = \mathbf{I}_p$. Then one speaks about an *orthonormal design*. This also implies that every basis function (vector) ψ_j is standardized: $\|\psi_j\|^2 = \sum_{i=1}^n \psi_{j,i}^2 = 1$. In the case of an orthonormal design, the estimate $\tilde{\boldsymbol{\theta}}$ is particularly simple: $\tilde{\boldsymbol{\theta}} = \Psi\mathbf{Y}$. Correspondingly, the target of estimation $\boldsymbol{\theta}^*$ satisfies $\boldsymbol{\theta}^* = \Psi\mathbf{f}^*$. In other words, the target is the collection (θ_j^*) of the Fourier coefficients of the underlying function (vector) \mathbf{f}^* w.r.t. the basis Ψ while the estimate $\tilde{\boldsymbol{\theta}}$ is the collection of empirical Fourier coefficients $\tilde{\theta}_j$:

$$\theta_j^* = \sum_{i=1}^n f_i \psi_{j,i}, \quad \tilde{\theta}_j = \sum_{i=1}^n Y_i \psi_{j,i}$$

An important feature of the orthonormal design is that it preserves the noise homogeneity:

$$\text{Var}(\tilde{\boldsymbol{\theta}}) = \sigma^2 I_p.$$

1.1.4 Spectral representation

Consider a linear model

$$\mathbf{Y} = \boldsymbol{\Psi}^\top \boldsymbol{\theta} + \boldsymbol{\varepsilon} \quad (1.5)$$

with homogeneous errors $\boldsymbol{\varepsilon}$: $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 I_n$. The rows of the matrix $\boldsymbol{\Psi}$ can be viewed as basis vectors in \mathbb{R}^n and the product $\boldsymbol{\Psi}^\top \boldsymbol{\theta}$ is a linear combinations of these vectors with the coefficients $(\theta_1, \dots, \theta_p)$. Effectively linear least squares estimation does a kind of projection of the data onto the subspace generated by the basis functions. This projection is of course invariant w.r.t. a basis transformation within this linear subspace. This fact can be used to reduce the model to the case of an orthogonal design considered in the previous section. Namely, one can always find a linear orthogonal transformation $U : \mathbb{R}^p \rightarrow \mathbb{R}^p$ ensuring the orthogonality of the transformed basis. This means that the rows of the matrix $\check{\boldsymbol{\Psi}} = U\boldsymbol{\Psi}$ are orthogonal and the matrix $\check{\boldsymbol{\Psi}}\check{\boldsymbol{\Psi}}^\top$ is diagonal:

$$\check{\boldsymbol{\Psi}}\check{\boldsymbol{\Psi}}^\top = U\boldsymbol{\Psi}\boldsymbol{\Psi}^\top U^\top = \Lambda = \text{diag}(\lambda_1, \dots, \lambda_p).$$

The original model reads after this transformation in the form

$$\mathbf{Y} = \check{\boldsymbol{\Psi}}^\top \mathbf{u} + \boldsymbol{\varepsilon}, \quad \check{\boldsymbol{\Psi}}\check{\boldsymbol{\Psi}}^\top = \Lambda,$$

where $\mathbf{u} = U\boldsymbol{\theta} \in \mathbb{R}^p$. Within this model, the transformed parameter \mathbf{u} can be estimated using the empirical Fourier coefficients $Z_j = \check{\boldsymbol{\psi}}_j^\top \mathbf{Y}$, where $\check{\boldsymbol{\psi}}_j$ is the j th row of $\check{\boldsymbol{\Psi}}$, $j = 1, \dots, p$. The original parameter vector $\boldsymbol{\theta}$ can be recovered via the equation $\boldsymbol{\theta} = U^\top \mathbf{u}$. This set of equations can be written in the form

$$\mathbf{Z} = \Lambda \mathbf{u} + \Lambda^{1/2} \boldsymbol{\xi} \quad (1.6)$$

where $\mathbf{Z} = \check{\boldsymbol{\Psi}}\mathbf{Y} = U\boldsymbol{\Psi}\mathbf{Y}$ is a vector in \mathbb{R}^p and $\boldsymbol{\xi} = \Lambda^{-1/2}\check{\boldsymbol{\Psi}}\boldsymbol{\varepsilon} = \Lambda^{-1/2}U\boldsymbol{\Psi}\boldsymbol{\varepsilon} \in \mathbb{R}^p$. The equation (1.6) is called the *spectral representation* of the linear model (1.5). The reason is that the basic transformation U can be built by a singular value decomposition of $\boldsymbol{\Psi}$. This representation is widely used in context of linear inverse problems; see Section 3.6.

Theorem 1.1.2. Consider the model (1.5) with homogeneous errors ε , that is, $\mathbb{E}\varepsilon\varepsilon^\top = \sigma^2 I_n$. Then there exists an orthogonal transform $U : \mathbb{R}^p \rightarrow \mathbb{R}^p$ leading to the spectral representation (1.6) with homogeneous uncorrelated errors $\xi : \mathbb{E}\xi\xi^\top = \sigma^2 I_p$. If $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$, then the vector ξ is normal as well: $\xi = \mathcal{N}(0, \sigma^2 I_p)$.

Exercise 1.1.3. Prove the result of Theorem 1.1.2.

Hint: select any U ensuring $U^\top \Psi \Psi^\top U = \Lambda$. Then

$$\mathbb{E}\xi\xi^\top = \Lambda^{-1/2} U \Psi \mathbb{E}\varepsilon\varepsilon^\top \Psi^\top U^\top \Lambda^{-1/2} = \sigma^2 \Lambda^{-1/2} U^\top \Psi \Psi^\top U \Lambda^{-1/2} = \sigma^2 I_p.$$

A special case of the spectral representation corresponds to the orthonormal design with $\Psi \Psi^\top = I_p$. In this situation, the spectral model reads as $\mathbf{Z} = \mathbf{u} + \xi$, that is, we simply observe the target \mathbf{u} corrupted with a homogeneous noise ξ . Such an equation is often called the *sequence space model* and it is intensively used in the literature for the theoretical study; cf. Section 3 below.

1.2 Properties of the response estimate $\tilde{\mathbf{f}}$

This section discusses some properties of the estimate $\tilde{\mathbf{f}} = \Psi^\top \tilde{\boldsymbol{\theta}} = \Pi \mathbf{Y}$ of the response vector \mathbf{f}^* . It is worth noting that the first and essential part of the analysis does not rely on the underlying model distribution, only on our parametric assumptions that $\mathbf{f} = \Psi^\top \boldsymbol{\theta}^*$ and $\text{Cov}(\varepsilon) = \Sigma = \sigma^2 I_n$. The real model only appears when studying the risk of estimation. We will comment on the cases of misspecified \mathbf{f} and Σ .

When $\Sigma = \sigma^2 I_n$, the operator Π in the representation $\tilde{\mathbf{f}} = \Pi \mathbf{Y}$ of the estimate $\tilde{\mathbf{f}}$ reads as

$$\Pi = \Psi^\top (\Psi \Psi^\top)^{-1} \Psi. \quad (1.7)$$

First we make use of the linear structure of the model (1.1) and of the estimate $\tilde{\mathbf{f}}$ to derive a number of its simple but important properties.

1.2.1 Decomposition into a deterministic and a stochastic component

The model equation $\mathbf{Y} = \mathbf{f}^* + \varepsilon$ yields

$$\tilde{\mathbf{f}} = \Pi \mathbf{Y} = \Pi(\mathbf{f}^* + \varepsilon) = \Pi \mathbf{f}^* + \Pi \varepsilon. \quad (1.8)$$

The first element of this sum, $\Pi \mathbf{f}^*$ is purely deterministic, but it depends on the unknown response vector \mathbf{f}^* . Moreover, it will be shown in the next lemma that $\Pi \mathbf{f}^* = \mathbf{f}^*$ if the parametric assumption holds and the vector \mathbf{f}^* indeed can be represented as $\Psi^\top \boldsymbol{\theta}^*$.

The second element is stochastic as a linear transformation of the stochastic vector $\boldsymbol{\varepsilon}$ but is independent of the model response \mathbf{f}^* . The properties of the estimate $\tilde{\mathbf{f}}$ heavily rely on the properties of the linear operator Π from (1.7) which we collect in the next section.

1.2.2 Properties of the operator Π

Let $\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_p$ be the columns of the matrix Ψ^\top . These are the vectors in \mathbb{R}^n also called *the basis vectors*.

Lemma 1.2.1. *Let the matrix $\Psi\Psi^\top$ be non-degenerate. Then the operator Π fulfills the following conditions:*

- (i) Π is symmetric (self-adjoint), that is, $\Pi^\top = \Pi$.
- (ii) Π is a projector in \mathbb{R}^n , i.e. $\Pi^\top \Pi = \Pi^2 = \Pi$ and $\Pi(\mathbf{1}_n - \Pi) = 0$, where $\mathbf{1}_n$ means the unity operator in \mathbb{R}^n .
- (iii) For an arbitrary vector v from \mathbb{R}^n , it holds $\|v\|^2 = \|\Pi v\|^2 + \|v - \Pi v\|^2$.
- (iv) The trace of Π is equal to the dimension of its image, $\text{tr } \Pi = p$.
- (v) Π projects the linear space \mathbb{R}^n on the linear subspace $L_p = \langle \boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_p \rangle$, which is spanned by the basis vectors $\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_p$, that is,

$$\|\mathbf{f}^* - \Pi \mathbf{f}^*\| = \inf_{\mathbf{g} \in L_p} \|\mathbf{f}^* - \mathbf{g}\|.$$

- (vi) The matrix Π can be represented in the form

$$\Pi = U^\top \Lambda_p U$$

where U is an orthonormal matrix and Λ_p is a diagonal matrix with the first p diagonal elements equal to 1 and the others equal to zero:

$$\Lambda_p = \text{diag}\{\underbrace{1, \dots, 1}_p, \underbrace{0, \dots, 0}_{n-p}\}.$$

Proof. It holds

$$\{\Psi^\top (\Psi\Psi^\top)^{-1} \Psi\}^\top = \Psi^\top (\Psi\Psi^\top)^{-1} \Psi$$

and

$$\Pi^2 = \Psi^\top (\Psi\Psi^\top)^{-1} \Psi \Psi^\top (\Psi\Psi^\top)^{-1} \Psi = \Psi^\top (\Psi\Psi^\top)^{-1} \Psi = \Pi,$$

which proves the first two statements of the lemma. The third one follows directly from the first two. Next,

$$\text{tr } \Pi = \text{tr } \Psi^\top (\Psi \Psi^\top)^{-1} \Psi = \text{tr } \Psi \Psi^\top (\Psi \Psi^\top)^{-1} = \text{tr } I_p = p.$$

The second property means that Π is a projector in \mathbb{R}^n and the fourth one means that the dimension of its image space is equal to p . The basis vectors ψ_1, \dots, ψ_p are the rows of the matrix Ψ . It is clear that

$$\Pi \Psi^\top = \Psi^\top (\Psi \Psi^\top)^{-1} \Psi \Psi^\top = \Psi^\top.$$

Therefore, the vectors ψ_j are invariants of the operator Π and in particular, all these vectors belong to the image space of this operator. If now \mathbf{g} is a vector in L_p , then it can be represented as $\mathbf{g} = c_1 \psi_1 + \dots + c_p \psi_p$ and therefore, $\Pi \mathbf{g} = \mathbf{g}$ and $\Pi L_p = L_p$. Finally, the non-singularity of the matrix $\Psi \Psi^\top$ means that the vectors ψ_1, \dots, ψ_p forming the rows of Ψ are linearly independent. Therefore, the space L_p spanned by the vectors ψ_1, \dots, ψ_p is of dimension p , and hence it coincides with the image space of the operation Π .

The last property is the usual diagonal decomposition of a projector.

Exercise 1.2.1. Consider the case of an orthogonal design with $\Psi \Psi^\top = I_p$. Specify the projector Π of Lemma 1.2.1 for this situation, particularly its decomposition from (vi).

1.2.3 Quadratic loss and risk of the response estimation

In this section we study the quadratic risk of estimating the response \mathbf{f}^* . The reason for studying the quadratic risk of estimating the response \mathbf{f}^* will be made clear when we discuss the properties of the fitted likelihood in the next section.

The loss $\varphi(\tilde{\mathbf{f}}, \mathbf{f}^*)$ of the estimate $\tilde{\mathbf{f}}$ can be naturally defined as the squared norm of the difference $\tilde{\mathbf{f}} - \mathbf{f}^*$:

$$\varphi(\tilde{\mathbf{f}}, \mathbf{f}^*) = \|\tilde{\mathbf{f}} - \mathbf{f}^*\|^2 = \sum_{i=1}^n |f_i - \tilde{f}_i|^2.$$

Correspondingly, the quadratic risk of the estimate $\tilde{\mathbf{f}}$ is the mean of this loss

$$\mathcal{R}(\tilde{\mathbf{f}}) = \mathbb{E} \varphi(\tilde{\mathbf{f}}, \mathbf{f}^*) = \mathbb{E} [(\tilde{\mathbf{f}} - \mathbf{f}^*)^\top (\tilde{\mathbf{f}} - \mathbf{f}^*)]. \quad (1.9)$$

The next result describes the loss and risk decomposition for two cases: when the parametric assumption $\mathbf{f}^* = \Psi^\top \boldsymbol{\theta}^*$ is correct and in the general case.

Theorem 1.2.1. *Suppose that the errors ε_i from (1.1) are independent with $\mathbb{E} \varepsilon_i = 0$ and $\mathbb{E} \varepsilon_i^2 = \sigma^2$, i.e. $\Sigma = \sigma^2 I_n$. Then the loss $\varphi(\tilde{\mathbf{f}}, \mathbf{f}^*) = \|\Pi \mathbf{Y} - \mathbf{f}^*\|^2$ and the risk $\mathcal{R}(\tilde{\mathbf{f}})$ of the LSE $\tilde{\mathbf{f}}$ fulfill*

$$\begin{aligned}\wp(\tilde{\mathbf{f}}, \mathbf{f}^*) &= \|\mathbf{f}^* - \Pi \mathbf{f}^*\|^2 + \|\Pi \boldsymbol{\varepsilon}\|^2, \\ \mathcal{R}(\tilde{\mathbf{f}}) &= \|\mathbf{f}^* - \Pi \mathbf{f}^*\|^2 + p\sigma^2.\end{aligned}$$

Moreover, if $\mathbf{f}^* = \Psi^\top \boldsymbol{\theta}^*$, then

$$\begin{aligned}\wp(\tilde{\mathbf{f}}, \mathbf{f}^*) &= \|\Pi \boldsymbol{\varepsilon}\|^2, \\ \mathcal{R}(\tilde{\mathbf{f}}) &= p\sigma^2.\end{aligned}$$

Proof. We apply (1.9) and the decomposition (1.8) of the estimate $\tilde{\mathbf{f}}$. It follows

$$\begin{aligned}\wp(\tilde{\mathbf{f}}, \mathbf{f}^*) &= \|\tilde{\mathbf{f}} - \mathbf{f}^*\|^2 = \|\mathbf{f}^* - \Pi \mathbf{f}^* - \Pi \boldsymbol{\varepsilon}\|^2 \\ &= \|\mathbf{f}^* - \Pi \mathbf{f}^*\|^2 + 2(\mathbf{f}^* - \Pi \mathbf{f}^*)^\top \Pi \boldsymbol{\varepsilon} + \|\Pi \boldsymbol{\varepsilon}\|^2.\end{aligned}$$

This implies the decomposition for the loss of $\tilde{\mathbf{f}}$ by Lemma 1.2.1, (ii). Next we compute the mean of $\|\Pi \boldsymbol{\varepsilon}\|^2$ applying again Lemma 1.2.1. Indeed

$$\begin{aligned}\mathbb{E}\|\Pi \boldsymbol{\varepsilon}\|^2 &= \mathbb{E}(\Pi \boldsymbol{\varepsilon})^\top \Pi \boldsymbol{\varepsilon} = \mathbb{E} \operatorname{tr}\{\Pi \boldsymbol{\varepsilon}(\Pi \boldsymbol{\varepsilon})^\top\} = \mathbb{E} \operatorname{tr}(\Pi \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top \Pi^\top) \\ &= \operatorname{tr}\{\Pi \mathbb{E}(\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top) \Pi\} = \sigma^2 \operatorname{tr}(\Pi^2) = p\sigma^2.\end{aligned}$$

Now consider the case when $\mathbf{f}^* = \Psi^\top \boldsymbol{\theta}^*$. By Lemma 1.2.1 $\mathbf{f}^* = \Pi \mathbf{f}^*$ and the last two statements of the theorem clearly follow.

1.2.4 Misspecified “colored noise”

Here we briefly comment on the case when $\boldsymbol{\varepsilon}$ is not a white noise. So, our assumption about the errors ε_i is that they are uncorrelated and homogeneous, that is, $\Sigma = \sigma^2 I_n$ while the true covariance matrix is given by Σ_0 . Many properties of the estimate $\tilde{\mathbf{f}} = \Pi \mathbf{Y}$ which are simply based on the linearity of the model (1.1) and of the estimate $\tilde{\mathbf{f}}$ itself continue to apply. In particular, the loss $\wp(\tilde{\mathbf{f}}, \mathbf{f}^*) = \|\tilde{\mathbf{f}} - \mathbf{f}^*\|^2$ can again be decomposed as

$$\|\tilde{\mathbf{f}} - \mathbf{f}^*\|^2 = \|\mathbf{f}^* - \Pi \mathbf{f}^*\|^2 + \|\Pi \boldsymbol{\varepsilon}\|^2.$$

Theorem 1.2.2. *Suppose that $\mathbb{E}\boldsymbol{\varepsilon} = 0$ and $\operatorname{Var}(\boldsymbol{\varepsilon}) = \Sigma_0$. Then the loss $\wp(\tilde{\mathbf{f}}, \mathbf{f}^*)$ and the risk $\mathcal{R}(\tilde{\mathbf{f}})$ of the LSE $\tilde{\mathbf{f}}$ fulfill*

$$\begin{aligned}\wp(\tilde{\mathbf{f}}, \mathbf{f}^*) &= \|\mathbf{f}^* - \Pi \mathbf{f}^*\|^2 + \|\Pi \boldsymbol{\varepsilon}\|^2, \\ \mathcal{R}(\tilde{\mathbf{f}}) &= \|\mathbf{f}^* - \Pi \mathbf{f}^*\|^2 + \operatorname{tr}(\Pi \Sigma_0 \Pi).\end{aligned}$$

Moreover, if $\mathbf{f}^* = \Psi^\top \boldsymbol{\theta}^*$, then

$$\begin{aligned}\wp(\tilde{\mathbf{f}}, \mathbf{f}^*) &= \|\Pi\boldsymbol{\varepsilon}\|^2, \\ \mathcal{R}(\tilde{\mathbf{f}}) &= \text{tr}(\Pi\Sigma_0\Pi).\end{aligned}$$

Proof. The decomposition of the loss from Theorem 1.2.1 only relies on the geometric properties of the projector Π and does not use the covariance structure of the noise. Hence, it only remains to check the expectation of $\|\Pi\boldsymbol{\varepsilon}\|^2$. Observe that

$$\mathbb{E}\|\Pi\boldsymbol{\varepsilon}\|^2 = \mathbb{E}\text{tr}[\Pi\boldsymbol{\varepsilon}(\Pi\boldsymbol{\varepsilon})^\top] = \text{tr}[\Pi\mathbb{E}(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top)\Pi] = \text{tr}(\Pi\Sigma_0\Pi)$$

as required.

1.3 Properties of the MLE $\tilde{\boldsymbol{\theta}}$

In this section we focus on the properties of the quasi MLE $\tilde{\boldsymbol{\theta}}$ built for the idealized linear Gaussian model $\mathbf{Y} = \Psi^\top\boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$ with $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2\mathbf{I}_n)$. As in the previous section, we do not assume the parametric structure of the underlying model and consider a more general model $\mathbf{Y} = \mathbf{f}^* + \boldsymbol{\varepsilon}$ with an unknown vector \mathbf{f}^* and errors $\boldsymbol{\varepsilon}$ with zero mean and covariance matrix Σ_0 . Due to (1.3), it holds $\tilde{\boldsymbol{\theta}} = \mathcal{S}\mathbf{Y}$ with $\mathcal{S} = (\Psi\Psi^\top)^{-1}\Psi$. An important feature of this estimate is its linear dependence on the data. The linear model equation $\mathbf{Y} = \mathbf{f}^* + \boldsymbol{\varepsilon}$ and linear structure of the estimate $\tilde{\boldsymbol{\theta}} = \mathcal{S}\mathbf{Y}$ allow us for decomposing the vector $\tilde{\boldsymbol{\theta}}$ into a deterministic and stochastic terms:

$$\tilde{\boldsymbol{\theta}} = \mathcal{S}\mathbf{Y} = \mathcal{S}(\mathbf{f}^* + \boldsymbol{\varepsilon}) = \mathcal{S}\mathbf{f}^* + \mathcal{S}\boldsymbol{\varepsilon}. \quad (1.10)$$

The first term $\mathcal{S}\mathbf{f}^*$ is deterministic but depends on the unknown vector \mathbf{f}^* while the second term $\mathcal{S}\boldsymbol{\varepsilon}$ is stochastic but it does not involve the model response \mathbf{f}^* . Below we study the properties of each component separately.

1.3.1 Properties of the stochastic component

The next result describes the distributional properties of the stochastic component $\boldsymbol{\delta} = \mathcal{S}\boldsymbol{\varepsilon}$ for $\mathcal{S} = (\Psi\Psi^\top)^{-1}\Psi$ and thus, of the estimate $\tilde{\boldsymbol{\theta}}$.

Theorem 1.3.1. *Assume $\mathbf{Y} = \mathbf{f}^* + \boldsymbol{\varepsilon}$ with $\mathbb{E}\boldsymbol{\varepsilon} = 0$ and $\text{Var}(\boldsymbol{\varepsilon}) = \Sigma_0$. The stochastic component $\boldsymbol{\delta} = \mathcal{S}\boldsymbol{\varepsilon}$ in (1.10) fulfills*

$$\mathbb{E}\boldsymbol{\delta} = 0, \quad W^2 \stackrel{\text{def}}{=} \text{Var}(\boldsymbol{\delta}) = \mathcal{S}\Sigma_0\mathcal{S}^\top, \quad \mathbb{E}\|\boldsymbol{\delta}\|^2 = \text{tr}W^2 = \text{tr}(\mathcal{S}\Sigma_0\mathcal{S}^\top).$$

Moreover, if $\Sigma = \Sigma_0 = \sigma^2\mathbf{I}_n$, then

$$W^2 = \sigma^2(\Psi\Psi^\top)^{-1}, \quad \mathbb{E}\|\boldsymbol{\delta}\|^2 = \text{tr}(W^2) = \sigma^2 \text{tr}[(\Psi\Psi^\top)^{-1}]. \quad (1.11)$$

Similarly for the estimate $\tilde{\boldsymbol{\theta}}$ it holds

$$\mathbb{E}\tilde{\boldsymbol{\theta}} = \mathcal{S}\mathbf{f}^*, \quad \text{Var}(\tilde{\boldsymbol{\theta}}) = W^2.$$

If the errors $\boldsymbol{\varepsilon}$ are Gaussian, then the both $\boldsymbol{\delta}$ and $\tilde{\boldsymbol{\theta}}$ are Gaussian as well:

$$\boldsymbol{\delta} \sim \mathcal{N}(0, W^2) \quad \tilde{\boldsymbol{\theta}} \sim \mathcal{N}(\mathcal{S}\mathbf{f}^*, W^2).$$

Proof. For the variance W^2 of $\boldsymbol{\delta}$ holds

$$\text{Var}(\boldsymbol{\delta}) = \mathbb{E}\boldsymbol{\delta}\boldsymbol{\delta}^\top = \mathbb{E}\mathcal{S}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top\mathcal{S}^\top = \mathcal{S}\Sigma_0\mathcal{S}^\top.$$

Next we use that $\mathbb{E}\|\boldsymbol{\delta}\|^2 = \mathbb{E}\boldsymbol{\delta}^\top\boldsymbol{\delta} = \mathbb{E}\text{tr}(\boldsymbol{\delta}\boldsymbol{\delta}^\top) = \text{tr}W^2$. If $\Sigma = \Sigma_0 = \sigma^2\mathbf{I}_n$, then (1.11) follows by simple algebra.

If $\boldsymbol{\varepsilon}$ is a Gaussian vector, then $\boldsymbol{\delta}$ as its linear transformation is Gaussian as well. The properties of $\tilde{\boldsymbol{\theta}}$ follow directly from the decomposition (1.10).

With $\Sigma_0 \neq \sigma^2\mathbf{I}_n$, the variance W^2 can be represented as

$$W^2 = (\Psi\Psi^\top)^{-1}\Psi\Sigma_0\Psi^\top(\Psi\Psi^\top)^{-1}.$$

Exercise 1.3.1. Let $\boldsymbol{\delta}$ be the stochastic component of $\tilde{\boldsymbol{\theta}}$ built for the misspecified linear model $\mathbf{Y} = \Psi^\top\boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$ with $\text{Var}(\boldsymbol{\varepsilon}) = \Sigma$. Let also the true noise variance is Σ_0 . Then $\text{Var}(\tilde{\boldsymbol{\theta}}) = W^2$ with

$$W^2 = (\Psi\Sigma^{-1}\Psi^\top)^{-1}\Psi\Sigma^{-1}\Sigma_0\Sigma^{-1}\Psi^\top(\Psi\Sigma^{-1}\Psi^\top)^{-1}.$$

The main finding in the presented study is that the stochastic part $\boldsymbol{\delta} = \mathcal{S}\boldsymbol{\varepsilon}$ of the estimate $\tilde{\boldsymbol{\theta}}$ is completely independent of the structure of the vector \mathbf{f}^* . In other words, the behavior of the stochastic component $\boldsymbol{\delta}$ does not change even if the linear parametric assumption is misspecified.

1.3.2 Properties of the deterministic component

Now we study the deterministic term starting with the parametric situation $\mathbf{f}^* = \Psi^\top\boldsymbol{\theta}^*$. Here we only specify the results for the case 1 with $\Sigma = \sigma^2\mathbf{I}_n$.

Theorem 1.3.2. Let $\mathbf{f}^* = \Psi^\top\boldsymbol{\theta}^*$. Then $\tilde{\boldsymbol{\theta}} = \mathcal{S}\mathbf{Y}$ with $\mathcal{S} = (\Psi\Psi^\top)^{-1}\Psi$ is unbiased, that is, $\mathbb{E}\tilde{\boldsymbol{\theta}} = \mathcal{S}\mathbf{f}^* = \boldsymbol{\theta}^*$.

Proof. For the proof, just observe that $\mathcal{S}\mathbf{f}^* = (\Psi\Psi^\top)^{-1}\Psi\Psi^\top\boldsymbol{\theta}^* = \boldsymbol{\theta}^*$.

Now we briefly discuss what happens when the linear parametric assumption is not fulfilled, that is, \mathbf{f}^* cannot be represented as $\Psi^\top \boldsymbol{\theta}^*$. In this case it is not yet clear what $\tilde{\boldsymbol{\theta}}$ really estimates. The answer is given in the context of the general theory of minimum contrast estimation. Namely, define $\boldsymbol{\theta}^*$ as the point which maximizes the expectation of the (quasi) log-likelihood $L(\boldsymbol{\theta})$:

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \mathbb{E}L(\boldsymbol{\theta}). \quad (1.12)$$

Theorem 1.3.3. *The solution $\boldsymbol{\theta}^*$ of the optimization problem (1.12) is given by*

$$\boldsymbol{\theta}^* = \mathcal{S}\mathbf{f}^* = (\Psi\Psi^\top)^{-1}\Psi\mathbf{f}^*.$$

Moreover,

$$\Psi^\top \boldsymbol{\theta}^* = \Pi\mathbf{f}^* = \Psi^\top (\Psi\Psi^\top)^{-1}\Psi\mathbf{f}^*.$$

In particular, if $\mathbf{f}^* = \Psi^\top \boldsymbol{\theta}^*$, then $\boldsymbol{\theta}^*$ follows (1.12).

Proof. The use of the model equation $\mathbf{Y} = \mathbf{f}^* + \boldsymbol{\varepsilon}$ and of the properties of the stochastic component $\boldsymbol{\delta}$ yield by simple algebra

$$\begin{aligned} \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \mathbb{E}L(\boldsymbol{\theta}) &= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \mathbb{E}(\mathbf{f}^* - \Psi^\top \boldsymbol{\theta} + \boldsymbol{\varepsilon})^\top (\mathbf{f}^* - \Psi^\top \boldsymbol{\theta} + \boldsymbol{\varepsilon}) \\ &= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \{(\mathbf{f}^* - \Psi^\top \boldsymbol{\theta})^\top (\mathbf{f}^* - \Psi^\top \boldsymbol{\theta}) + \mathbb{E}(\boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon})\} \\ &= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \{(\mathbf{f}^* - \Psi^\top \boldsymbol{\theta})^\top (\mathbf{f}^* - \Psi^\top \boldsymbol{\theta})\}. \end{aligned}$$

Differentiating w.r.t. $\boldsymbol{\theta}$ leads to the equation

$$\Psi(\mathbf{f}^* - \Psi^\top \boldsymbol{\theta}) = 0$$

and the solution $\boldsymbol{\theta}^* = (\Psi\Psi^\top)^{-1}\Psi\mathbf{f}^*$ which is exactly the expected value of $\tilde{\boldsymbol{\theta}}$ by Theorem 1.3.1.

Exercise 1.3.2. State the result of Theorems 1.3.2 and 1.3.3 for the MLE $\tilde{\boldsymbol{\theta}}$ built in the model $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$ with $\operatorname{Var}(\boldsymbol{\varepsilon}) = \Sigma$.

Hint: check that the statements continue to apply with $\mathcal{S} = (\Psi\Sigma^{-1}\Psi^\top)^{-1}\Psi\Sigma^{-1}$.

The last results and the decomposition (1.10) explain the behavior of the estimate $\tilde{\boldsymbol{\theta}}$ in a very general situation. The considered model is $\mathbf{Y} = \mathbf{f}^* + \boldsymbol{\varepsilon}$. We assume a linear parametric structure and independent homogeneous noise. The estimation procedure means in fact a kind of projection of the data \mathbf{Y} on a p -dimensional linear subspace in \mathbb{R}^n spanned by the given basis vectors $\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_p$. This projection, as a linear operator,

can be decomposed into a projection of the deterministic vector \mathbf{f}^* and a projection of the random noise ε . If the linear parametric assumption $\mathbf{f}^* \in \langle \boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_p \rangle$ is correct, that is, $\mathbf{f}^* = \theta_1^* \boldsymbol{\psi}_1 + \dots + \theta_p^* \boldsymbol{\psi}_p$, then this projection keeps \mathbf{f}^* unchanged and only the random noise is reduced via this projection. If \mathbf{f}^* cannot be exactly expanded using the basis $\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_p$, then the procedure recovers the projection of \mathbf{f}^* onto this subspace. The latter projection can be written as $\Psi^\top \boldsymbol{\theta}^*$ and the vector $\boldsymbol{\theta}^*$ can be viewed as the target of estimation.

1.3.3 Risk of estimation. R-efficiency

This section briefly discusses how the obtained properties of the estimate $\tilde{\theta}$ can be used to evaluate the risk of estimation. A particularly important question is the optimality of the MLE $\tilde{\theta}$. The main result of the section claims that $\tilde{\theta}$ is R-efficient if the model is correctly specified and is not if there is a misspecification.

We start with the case of a correct parametric specification $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \varepsilon$, that is, the linear parametric assumption $\mathbf{f}^* = \Psi^\top \boldsymbol{\theta}^*$ is exactly fulfilled and the noise ε is homogeneous: $\varepsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$. Later we extend the result to the case when the LPA $\mathbf{f}^* = \Psi^\top \boldsymbol{\theta}^*$ is not fulfilled and to the case when the noise is not homogeneous but still correctly specified. Finally we discuss the case when the noise structure is misspecified.

Under LPA $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \varepsilon$ with $\varepsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$, the estimate $\tilde{\theta}$ is also normal with mean $\boldsymbol{\theta}^*$ and the variance $W^2 = \sigma^2 \mathcal{S} \mathcal{S}^\top = \sigma^2 (\Psi \Psi^\top)^{-1}$. Define a $p \times p$ symmetric matrix D by the equation

$$D^2 = \frac{1}{\sigma^2} \sum_{i=1}^n \Psi_i \Psi_i^\top = \frac{1}{\sigma^2} \Psi \Psi^\top.$$

Clearly $W^2 = D^{-2}$.

Now we show that $\tilde{\theta}$ is R-efficient. Actually this fact can be derived from the Cramér-Rao Theorem because the Gaussian model is a special case of an exponential family. However, we check this statement directly by computing the Cramér-Rao efficiency bound. Recall that the Fisher information matrix $\mathbb{F}(\boldsymbol{\theta})$ for the log-likelihood $L(\boldsymbol{\theta})$ is defined as the variance of $\nabla L(\boldsymbol{\theta})$ under $\mathbb{P}_{\boldsymbol{\theta}}$.

Theorem 1.3.4 (Gauss-Markov). *Let $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \varepsilon$ with $\varepsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$. Then $\tilde{\theta}$ is R-efficient estimate of $\boldsymbol{\theta}^*$: $\mathbb{E} \tilde{\theta} = \boldsymbol{\theta}^*$,*

$$\mathbb{E}[(\tilde{\theta} - \boldsymbol{\theta}^*)(\tilde{\theta} - \boldsymbol{\theta}^*)^\top] = \text{Var}(\tilde{\theta}) = D^{-2},$$

and for any unbiased linear estimate $\hat{\theta}$ satisfying $\mathbb{E}_{\boldsymbol{\theta}} \hat{\theta} \equiv \boldsymbol{\theta}$, it holds

$$\text{Var}(\hat{\theta}) \geq \text{Var}(\tilde{\theta}) = D^{-2}.$$

Proof. Theorems 1.3.1 and 1.3.2 imply that $\tilde{\boldsymbol{\theta}} \sim \mathcal{N}(\boldsymbol{\theta}^*, W^2)$ with $W^2 = \sigma^2(\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)^{-1} = D^{-2}$. Next we show that for any $\boldsymbol{\theta}$

$$\text{Var}[\nabla L(\boldsymbol{\theta})] = D^2,$$

that is, the Fisher information does not depend on the model function \boldsymbol{f}^* . The log-likelihood $L(\boldsymbol{\theta})$ for the model $\boldsymbol{Y} \sim \mathcal{N}(\boldsymbol{\Psi}^\top \boldsymbol{\theta}^*, \sigma^2 I_n)$ reads as

$$L(\boldsymbol{\theta}) = -\frac{1}{2\sigma^2}(\boldsymbol{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta})^\top (\boldsymbol{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta}) - \frac{n}{2} \log(2\pi\sigma^2).$$

This yields for its gradient $\nabla L(\boldsymbol{\theta})$:

$$\nabla L(\boldsymbol{\theta}) = \sigma^{-2} \boldsymbol{\Psi}(\boldsymbol{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta})$$

and in view of $\text{Var}(\boldsymbol{Y}) = \Sigma = \sigma^2 I_n$, it holds

$$\text{Var}[\nabla L(\boldsymbol{\theta})] = \sigma^{-4} \boldsymbol{\Psi} \text{Var}(\boldsymbol{Y}) \boldsymbol{\Psi}^\top = \sigma^{-2} \boldsymbol{\Psi} \boldsymbol{\Psi}^\top$$

as required.

The R-efficiency $\tilde{\boldsymbol{\theta}}$ follows from the Cramér-Rao efficiency bound because $\{\text{Var}(\tilde{\boldsymbol{\theta}})\}^{-1} = \text{Var}\{\nabla L(\boldsymbol{\theta})\}$. However, we present an independent proof of this fact. Actually we prove a sharper result that the variance of a linear unbiased estimate $\hat{\boldsymbol{\theta}}$ coincides with the variance of $\tilde{\boldsymbol{\theta}}$ only if $\hat{\boldsymbol{\theta}}$ coincides almost surely with $\tilde{\boldsymbol{\theta}}$, otherwise it is larger. The idea of the proof is quite simple. Consider the difference $\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}$ and show that the condition $\mathbb{E}\hat{\boldsymbol{\theta}} = \mathbb{E}\tilde{\boldsymbol{\theta}} = \boldsymbol{\theta}^*$ implies orthogonality $\mathbb{E}\{\tilde{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}})^\top\} = 0$. This, in turns, implies $\text{Var}(\hat{\boldsymbol{\theta}}) = \text{Var}(\tilde{\boldsymbol{\theta}}) + \text{Var}(\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}) \geq \text{Var}(\tilde{\boldsymbol{\theta}})$. So, it remains to check the orthogonality of $\tilde{\boldsymbol{\theta}}$ and $\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}$. Let $\hat{\boldsymbol{\theta}} = A\boldsymbol{Y}$ for a $p \times n$ matrix A and $\mathbb{E}_\theta \hat{\boldsymbol{\theta}} \equiv \boldsymbol{\theta}$ and all $\boldsymbol{\theta}$. These two equalities and $\mathbb{E}\boldsymbol{Y} = \boldsymbol{\Psi}^\top \boldsymbol{\theta}^*$ imply that $A\boldsymbol{\Psi}^\top \boldsymbol{\theta}^* \equiv \boldsymbol{\theta}^*$, i.e. $A\boldsymbol{\Psi}^\top$ is the identity $p \times p$ matrix. The same is true for $\tilde{\boldsymbol{\theta}} = S\boldsymbol{Y}$ yielding $S\boldsymbol{\Psi}^\top = I_p$. Next, in view of $\mathbb{E}\hat{\boldsymbol{\theta}} = \mathbb{E}\tilde{\boldsymbol{\theta}} = \boldsymbol{\theta}^*$

$$\mathbb{E}\{(\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}})\tilde{\boldsymbol{\theta}}^\top\} = \mathbb{E}(A - S)\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top S^\top = \sigma^2(A - S)\boldsymbol{\Psi}^\top(\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)^{-1} = 0,$$

and the assertion follows.

Exercise 1.3.3. Check the details of the proof of the theorem. Show that the statement $\text{Var}(\hat{\boldsymbol{\theta}}) \geq \text{Var}(\tilde{\boldsymbol{\theta}})$ only uses that $\hat{\boldsymbol{\theta}}$ is unbiased and that $\mathbb{E}\boldsymbol{Y} = \boldsymbol{\Psi}^\top \boldsymbol{\theta}^*$ and $\text{Var}(\boldsymbol{Y}) = \sigma^2 I_n$.

Exercise 1.3.4. Compute $\nabla^2 L(\boldsymbol{\theta})$. Check that it is non-random, does not depend on $\boldsymbol{\theta}$, and fulfills for every $\boldsymbol{\theta}$ the identity

$$\nabla^2 L(\boldsymbol{\theta}) \equiv -\text{Var}[\nabla L(\boldsymbol{\theta})] = -D^2.$$

A colored noise

The majority of the presented results continue to apply in the case of heterogeneous and even dependent noise with $\text{Var}(\boldsymbol{\varepsilon}) = \Sigma_0$. The key facts behind this extension are the decomposition (1.10) and the properties of the stochastic component $\boldsymbol{\delta}$ from Section 1.3.1: $\boldsymbol{\delta} \sim \mathcal{N}(0, W^2)$. In the case of a colored noise, the definition of W and D is changed for

$$D^2 \stackrel{\text{def}}{=} W^{-2} = \Psi \Sigma_0^{-1} \Psi^\top.$$

Exercise 1.3.5. State and prove the analog of Theorem 1.3.4 for the colored noise $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \Sigma_0)$.

A misspecified LPA

An interesting feature of our results so far is that they equally apply for the correct linear specification $\boldsymbol{f}^* = \Psi^\top \boldsymbol{\theta}^*$ and for the case when the identity $\boldsymbol{f}^* = \Psi^\top \boldsymbol{\theta}$ is not precisely fulfilled whatever $\boldsymbol{\theta}$ is taken. In this situation the target of analysis is the vector $\boldsymbol{\theta}^*$ describing the best linear approximation of \boldsymbol{f}^* by $\Psi^\top \boldsymbol{\theta}$. We already know from the results of Section 1.3.1 and 1.3.2 that the estimate $\tilde{\boldsymbol{\theta}}$ is also normal with mean $\boldsymbol{\theta}^* = \mathcal{S} \boldsymbol{f}^* = (\Psi \Psi^\top)^{-1} \Psi \boldsymbol{f}^*$ and the variance $W^2 = \sigma^2 \mathcal{S} \mathcal{S}^\top = \sigma^2 (\Psi \Psi^\top)^{-1}$.

Theorem 1.3.5. Assume $\boldsymbol{Y} = \boldsymbol{f}^* + \boldsymbol{\varepsilon}$ with $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 I_n)$. Let $\boldsymbol{\theta}^* = \mathcal{S} \boldsymbol{f}^*$. Then $\tilde{\boldsymbol{\theta}}$ is R -efficient estimate of $\boldsymbol{\theta}^*$: $\mathbb{E} \tilde{\boldsymbol{\theta}} = \boldsymbol{\theta}^*$,

$$\mathbb{E}[(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)^\top] = \text{Var}(\tilde{\boldsymbol{\theta}}) = D^{-2},$$

and for any unbiased linear estimate $\hat{\boldsymbol{\theta}}$ satisfying $\mathbb{E}_{\boldsymbol{\theta}} \hat{\boldsymbol{\theta}} \equiv \boldsymbol{\theta}$, it holds

$$\text{Var}(\hat{\boldsymbol{\theta}}) \geq \text{Var}(\tilde{\boldsymbol{\theta}}) = D^{-2}.$$

Proof. The proofs only utilize that $\tilde{\boldsymbol{\theta}} \sim \mathcal{N}(\boldsymbol{\theta}^*, W^2)$ with $W^2 = D^{-2}$. The only small remark concerns the equality $\text{Var}[\nabla L(\boldsymbol{\theta})] = D^2$ from Theorem 1.3.4.

Exercise 1.3.6. Check the identity $\text{Var}[\nabla L(\boldsymbol{\theta})] = D^2$ from Theorem 1.3.4 for $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \Sigma_0)$.

1.3.4 The case of a misspecified noise

Here we again consider the linear parametric assumption $\boldsymbol{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$. However, contrary to the previous section, we admit that the noise $\boldsymbol{\varepsilon}$ is not homogeneous normal:

$\varepsilon \sim \mathcal{N}(0, \Sigma_0)$ while our estimation procedure is the quasi MLE based on the assumption of noise homogeneity $\varepsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$. We already know that the estimate $\tilde{\boldsymbol{\theta}}$ is unbiased with mean $\boldsymbol{\theta}^*$ and variance $W^2 = \mathcal{S}\Sigma_0\mathcal{S}^\top$, where $\mathcal{S} = (\Psi\Psi^\top)^{-1}\Psi$. This gives

$$W^2 = (\Psi\Psi^\top)^{-1}\Psi\Sigma_0\Psi^\top(\Psi\Psi^\top)^{-1}.$$

The question is whether the estimate $\tilde{\boldsymbol{\theta}}$ based on the misspecified distributional assumption is efficient. The Cramér-Rao result delivers the lower bound for the quadratic risk in form of $\text{Var}(\tilde{\boldsymbol{\theta}}) \geq [\text{Var}(\nabla L(\boldsymbol{\theta}))]^{-1}$. We already know that the use of the correctly specified covariance matrix of the errors leads to an R-efficient estimate $\tilde{\boldsymbol{\theta}}$. The next result show that the use of a misspecified matrix Σ results in an estimate which is unbiased but not R-efficient, that is, the best estimation risk is achieved if we apply the correct model assumptions.

Theorem 1.3.6. *Let $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \varepsilon$ with $\varepsilon \sim \mathcal{N}(0, \Sigma_0)$. Then*

$$\text{Var}[\nabla L(\boldsymbol{\theta})] = \Psi\Sigma_0^{-1}\Psi^\top.$$

The estimate $\tilde{\boldsymbol{\theta}} = (\Psi\Psi^\top)^{-1}\Psi\mathbf{Y}$ is unbiased, that is, $\mathbb{E}\tilde{\boldsymbol{\theta}} = \boldsymbol{\theta}^$, but it is not R-efficient unless $\Sigma_0 = \Sigma$.*

Proof. Let $\tilde{\boldsymbol{\theta}}_0$ be the MLE for the correct model specification with the noise $\varepsilon \sim \mathcal{N}(0, \Sigma_0)$. As $\tilde{\boldsymbol{\theta}}$ is unbiased, the difference $\tilde{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}_0$ is orthogonal to $\tilde{\boldsymbol{\theta}}_0$ and it holds for the variance of $\tilde{\boldsymbol{\theta}}$

$$\text{Var}(\tilde{\boldsymbol{\theta}}) = \text{Var}(\tilde{\boldsymbol{\theta}}_0) + \text{Var}(\tilde{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}_0);$$

cf. with the proof of Gauss-Markov-Theorem 1.3.4.

Exercise 1.3.7. Compare directly the variances of $\tilde{\boldsymbol{\theta}}$ and of $\tilde{\boldsymbol{\theta}}_0$.

1.4 Linear models and quadratic log-likelihood

Linear Gaussian modeling leads to a specific log-likelihood structure; see Section 1. Namely, the log-likelihood function $L(\boldsymbol{\theta})$ is quadratic in $\boldsymbol{\theta}$, the coefficients of the quadratic terms are deterministic and the cross term is linear both in $\boldsymbol{\theta}$ and in the observations Y_i . Here we show that this geometric structure of the log-likelihood characterizes linear models. We say that $L(\boldsymbol{\theta})$ is *quadratic* if it is a quadratic function of $\boldsymbol{\theta}$ and there is a deterministic symmetric matrix D^2 such that for any $\boldsymbol{\theta}^\circ, \boldsymbol{\theta}$

$$L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^\circ) = (\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top \nabla L(\boldsymbol{\theta}^\circ) - (\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top D^2 (\boldsymbol{\theta} - \boldsymbol{\theta}^\circ) / 2. \quad (1.13)$$

Here $\nabla L(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \frac{dL(\boldsymbol{\theta})}{d\boldsymbol{\theta}}$. As usual we define

$$\begin{aligned}\tilde{\boldsymbol{\theta}} &\stackrel{\text{def}}{=} \underset{\boldsymbol{\theta}}{\operatorname{argmax}} L(\boldsymbol{\theta}), \\ \boldsymbol{\theta}^* &= \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \mathbb{E}L(\boldsymbol{\theta}).\end{aligned}$$

The next result describes some properties of the estimate $\tilde{\boldsymbol{\theta}}$ which are entirely based on the geometric (quadratic) structure of the function $L(\boldsymbol{\theta})$. All the results are stated by using the matrix D^2 and the vector $\boldsymbol{\zeta} = \nabla L(\boldsymbol{\theta}^*)$.

Theorem 1.4.1. *Let $L(\boldsymbol{\theta})$ be quadratic for a matrix $D^2 > 0$. Then for any $\boldsymbol{\theta}^\circ$*

$$\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^\circ = D^{-2}\nabla L(\boldsymbol{\theta}^\circ). \quad (1.14)$$

In particular, with $\boldsymbol{\theta}^\circ = 0$, it holds

$$\tilde{\boldsymbol{\theta}} = D^{-2}\nabla L(0).$$

Taking $\boldsymbol{\theta}^\circ = \boldsymbol{\theta}^$ yields*

$$\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^* = D^{-2}\boldsymbol{\zeta} \quad (1.15)$$

with $\boldsymbol{\zeta} \stackrel{\text{def}}{=} \nabla L(\boldsymbol{\theta}^)$. Moreover, $\mathbb{E}\boldsymbol{\zeta} = 0$, and it holds with $V^2 = \operatorname{Var}(\boldsymbol{\zeta}) = \mathbb{E}\boldsymbol{\zeta}\boldsymbol{\zeta}^\top$*

$$\begin{aligned}\mathbb{E}\tilde{\boldsymbol{\theta}} &= \boldsymbol{\theta}^* \\ \operatorname{Var}(\tilde{\boldsymbol{\theta}}) &= D^{-2}V^2D^{-2}.\end{aligned}$$

Further, for any $\boldsymbol{\theta}$,

$$L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}) = (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta})^\top D^2(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta})/2 = \|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta})\|^2/2. \quad (1.16)$$

Finally, it holds for the excess $L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^) \stackrel{\text{def}}{=} L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*)$*

$$2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) = (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)^\top D^2(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) = \boldsymbol{\zeta}^\top D^{-2}\boldsymbol{\zeta} = \|\boldsymbol{\xi}\|^2 \quad (1.17)$$

with $\boldsymbol{\xi} = D^{-1}\boldsymbol{\zeta}$.

Proof. The extremal point equation $\nabla L(\boldsymbol{\theta}) = 0$ for the quadratic function $L(\boldsymbol{\theta})$ from (1.13) yields (1.14). The equation (1.13) with $\boldsymbol{\theta}^\circ = \boldsymbol{\theta}^*$ implies for any $\boldsymbol{\theta}$

$$\nabla L(\boldsymbol{\theta}) = \nabla L(\boldsymbol{\theta}^\circ) - D^2(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ) = \boldsymbol{\zeta} - D^2(\boldsymbol{\theta} - \boldsymbol{\theta}^*). \quad (1.18)$$

Therefore, it holds for the expectation $\mathbb{E}L(\boldsymbol{\theta})$

$$\nabla \mathbb{E}L(\boldsymbol{\theta}) = \mathbb{E}\boldsymbol{\zeta} - D^2(\boldsymbol{\theta} - \boldsymbol{\theta}^*),$$

and the equation $\nabla \mathbb{E}L(\boldsymbol{\theta}^*) = 0$ implies $\mathbb{E}\boldsymbol{\zeta} = 0$.

To show (1.16), apply again the property (1.13) with $\boldsymbol{\theta}^\circ = \tilde{\boldsymbol{\theta}}$:

$$\begin{aligned} L(\boldsymbol{\theta}) - L(\tilde{\boldsymbol{\theta}}) &= (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})^\top \nabla L(\tilde{\boldsymbol{\theta}}) - (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})^\top D^2(\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})/2 \\ &= -(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta})^\top D^2(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta})/2. \end{aligned}$$

Here we used that $\nabla L(\tilde{\boldsymbol{\theta}}) = 0$ because $\tilde{\boldsymbol{\theta}}$ is an extreme point of $L(\boldsymbol{\theta})$. The last result (1.17) is a special case with $\boldsymbol{\theta} = \boldsymbol{\theta}^*$ in view of (1.15).

This theorem delivers an important message: the main properties of the MLE $\tilde{\boldsymbol{\theta}}$ can be explained via the geometric (quadratic) structure of the log-likelihood. An interesting question to clarify is whether a quadratic log-likelihood structure is specific for linear Gaussian model. The answer is positive: there is one-to-one correspondence between linear Gaussian models and quadratic log-likelihood functions. Indeed, the identity (1.18) with $\boldsymbol{\theta}^\circ = \boldsymbol{\theta}^*$ can be rewritten as

$$\nabla L(\boldsymbol{\theta}) + D^2\boldsymbol{\theta} \equiv \boldsymbol{\zeta} + D^2\boldsymbol{\theta}^*.$$

If we fix any $\boldsymbol{\theta}$ and define $\mathbf{Y} = \nabla L(\boldsymbol{\theta}) + D^2\boldsymbol{\theta}$, this yields

$$\mathbf{Y} = D^2\boldsymbol{\theta}^* + \boldsymbol{\zeta}.$$

Similarly, $\mathbf{Y} \stackrel{\text{def}}{=} D^{-1}\{\nabla L(\boldsymbol{\theta}) + D^2\boldsymbol{\theta}\}$ yields the equation

$$\mathbf{Y} = D\boldsymbol{\theta}^* + \boldsymbol{\xi}, \tag{1.19}$$

where $\boldsymbol{\xi} = D^{-1}\boldsymbol{\zeta}$. We can summarize as follows.

Theorem 1.4.2. *Let $L(\boldsymbol{\theta})$ be quadratic with a non-degenerated matrix D^2 . Then $\mathbf{Y} \stackrel{\text{def}}{=} D^{-1}\{\nabla L(\boldsymbol{\theta}) + D^2\boldsymbol{\theta}\}$ does not depend on $\boldsymbol{\theta}$ and $L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^*)$ is the quasi log-likelihood ratio for the linear Gaussian model (1.19) with $\boldsymbol{\xi}$ standard normal. It is the true log-likelihood if and only if $\boldsymbol{\zeta} \sim \mathcal{N}(0, D^2)$.*

Proof. The model (1.19) with $\boldsymbol{\xi} \sim \mathcal{N}(0, I_p)$ leads to the log-likelihood ratio

$$(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top D(\mathbf{Y} - D\boldsymbol{\theta}^*) - \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2 = (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \boldsymbol{\zeta} - \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2$$

in view of the definition of \mathbf{Y} . The definition (1.13) implies

$$L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^*) = (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla L(\boldsymbol{\theta}^*) - \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2.$$

As these two expressions coincide, it follows that $L(\boldsymbol{\theta})$ is the true log-likelihood if and only if $\boldsymbol{\xi} = D^{-1}\boldsymbol{\zeta}$ is standard normal.

1.4.1 Inference based on the maximum likelihood

All the results presented above for linear models were based on the explicit representation of the (quasi) MLE $\tilde{\boldsymbol{\theta}}$. Here we present the approach based on the analysis of the maximum likelihood. This approach does not require to fix any analytic expression for the point of maximum of the (quasi) likelihood process $L(\boldsymbol{\theta})$. Instead we work directly with the maximum of this process. We establish exponential inequalities for the *excess* or the *maximum likelihood* $L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)$. We also show how these results can be used to study the accuracy of the MLE $\tilde{\boldsymbol{\theta}}$, in particular, for building confidence sets.

One more benefit of the ML-based approach is that it equally applies to a homogeneous and to a heterogeneous noise provided that the noise structure is not misspecified. The celebrated chi-squared result about the maximum likelihood $L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)$ claims that the distribution of $2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)$ is chi-squared with p degrees of freedom χ_p^2 and it does not depend on the noise covariance; see Section 1.4.1.

Now we specify the setup. The starting point of the ML-approach is the linear Gaussian model assumption $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$ with $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \Sigma)$. The corresponding log-likelihood ratio $L(\boldsymbol{\theta})$ can be written as

$$L(\boldsymbol{\theta}) = -\frac{1}{2}(\mathbf{Y} - \Psi^\top \boldsymbol{\theta})^\top \Sigma^{-1}(\mathbf{Y} - \Psi^\top \boldsymbol{\theta}) + R, \quad (1.20)$$

where the remainder term R does not depend on $\boldsymbol{\theta}$. Now one can see that $L(\boldsymbol{\theta})$ is a quadratic function of $\boldsymbol{\theta}$. Moreover, $\nabla^2 L(\boldsymbol{\theta}) = \Psi \Sigma^{-1} \Psi^\top$, so that $L(\boldsymbol{\theta})$ is quadratic with $D^2 = \Psi \Sigma^{-1} \Psi^\top$. This enables us to apply the general results of Section 1.4 which are only based on the geometric (quadratic) structure of the log-likelihood $L(\boldsymbol{\theta})$: the true data distribution can be arbitrary.

Theorem 1.4.3. *Consider $L(\boldsymbol{\theta})$ from (1.20). For any $\boldsymbol{\theta}$, it holds with $D^2 = \Psi \Sigma^{-1} \Psi^\top$*

$$L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}) = (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta})^\top D^2 (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}) / 2. \quad (1.21)$$

In particular, if $\Sigma = \sigma^2 \mathbf{I}_n$ then the fitted log-likelihood is proportional to the quadratic loss $\|\tilde{\mathbf{f}} - \mathbf{f}_\theta\|^2$ for $\tilde{\mathbf{f}} = \Psi^\top \tilde{\boldsymbol{\theta}}$ and $\mathbf{f}_\theta = \Psi^\top \boldsymbol{\theta}$:

$$L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}) = \frac{1}{2\sigma^2} \|\Psi^\top (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta})\|^2 = \frac{1}{2\sigma^2} \|\tilde{\mathbf{f}} - \mathbf{f}_\theta\|^2.$$

If $\boldsymbol{\theta}^ \stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta}} \mathbb{E}L(\boldsymbol{\theta}) = D^{-2} \Psi \Sigma^{-1} \mathbf{f}^*$ for $\mathbf{f}^* = \mathbb{E}\mathbf{Y}$, then*

$$2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) = \boldsymbol{\zeta}^\top D^{-2} \boldsymbol{\zeta} = \|\boldsymbol{\xi}\|^2 \quad (1.22)$$

with $\boldsymbol{\zeta} = \nabla L(\boldsymbol{\theta}^)$ and $\boldsymbol{\xi} \stackrel{\text{def}}{=} D^{-1} \boldsymbol{\zeta}$.*

Proof. The results (1.21) and (1.22) follow from Theorem 1.4.1; see (1.16) and (1.17).

If the model assumptions are not misspecified one can establish the remarkable χ^2 result.

Theorem 1.4.4. *Let $L(\boldsymbol{\theta})$ from (1.20) be the log-likelihood for the model $\mathbf{Y} = \boldsymbol{\Psi}^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$ with $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \Sigma)$. Then $\boldsymbol{\xi} = D^{-1}\boldsymbol{\zeta} \sim \mathcal{N}(0, \mathbf{I}_p)$ and $2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \sim \chi_p^2$ is chi-squared with p degrees of freedom.*

Proof. By direct calculus

$$\boldsymbol{\zeta} = \nabla L(\boldsymbol{\theta}^*) = \boldsymbol{\Psi} \Sigma^{-1} (\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta}^*) = \boldsymbol{\Psi} \Sigma^{-1} \boldsymbol{\varepsilon}.$$

So, $\boldsymbol{\zeta}$ is a linear transformation of a Gaussian vector \mathbf{Y} and thus it is Gaussian as well. By Theorem 1.4.1, $\mathbb{E}\boldsymbol{\zeta} = 0$. Moreover, $\text{Var}(\boldsymbol{\varepsilon}) = \Sigma$ implies

$$\text{Var}(\boldsymbol{\zeta}) = \mathbb{E} \boldsymbol{\Psi}^\top \Sigma^{-1} \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top \Sigma^{-1} \boldsymbol{\Psi} = \boldsymbol{\Psi} \Sigma^{-1} \boldsymbol{\Psi}^\top = D^2$$

yielding that $\boldsymbol{\xi} = D^{-1}\boldsymbol{\zeta}$ is standard normal.

The last result $2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \sim \chi_p^2$ is sometimes called the “chi-squared phenomenon”: the distribution of the maximum likelihood only depends on the number of parameters to be estimated and is independent of the design $\boldsymbol{\Psi}$, of the noise covariance matrix Σ , etc. This particularly explains the use of word “phenomenon” in the name of the result.

Exercise 1.4.1. Check that the linear transformation $\check{\mathbf{Y}} = \Sigma^{-1/2} \mathbf{Y}$ of the data does not change the value of the log-likelihood ratio $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ and hence, of the maximum likelihood $L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)$.

Hint: use the representation

$$\begin{aligned} L(\boldsymbol{\theta}) &= \frac{1}{2} (\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta})^\top \Sigma^{-1} (\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta}) + R \\ &= \frac{1}{2} (\check{\mathbf{Y}} - \check{\boldsymbol{\Psi}}^\top \boldsymbol{\theta})^\top (\check{\mathbf{Y}} - \check{\boldsymbol{\Psi}}^\top \boldsymbol{\theta}) + R \end{aligned}$$

and check that the transformed data $\check{\mathbf{Y}}$ is described by the model $\check{\mathbf{Y}} = \check{\boldsymbol{\Psi}}^\top \boldsymbol{\theta}^* + \check{\boldsymbol{\varepsilon}}$ with $\check{\boldsymbol{\Psi}} = \boldsymbol{\Psi} \Sigma^{-1/2}$ and $\check{\boldsymbol{\varepsilon}} = \Sigma^{-1/2} \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I}_n)$ yielding the same log-likelihood ratio as in the original model.

Exercise 1.4.2. Assume homogeneous noise in (1.20) with $\Sigma = \sigma^2 \mathbf{I}_n$. Then it holds

$$2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) = \sigma^{-2} \|\Pi \boldsymbol{\varepsilon}\|^2$$

where $\Pi = \boldsymbol{\Psi}^\top (\boldsymbol{\Psi} \boldsymbol{\Psi}^\top)^{-1} \boldsymbol{\Psi}$ is the projector in \mathbb{R}^n on the subspace spanned by the vectors $\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_p$.

Hint: use that $\zeta = \sigma^{-2}\Psi\varepsilon$, $D^2 = \sigma^{-2}\Psi\Psi^\top$, and

$$\sigma^{-2}\|H\varepsilon\|^2 = \sigma^{-2}\varepsilon^\top H^\top H\varepsilon = \sigma^{-2}\varepsilon^\top H\varepsilon = \zeta^\top D^{-2}\zeta.$$

We write the result of Theorem 1.4.3 in the form $2L(\tilde{\theta}, \theta^*) \sim \chi_p^2$, where χ_p^2 stands for the chi-squared distribution with p degrees of freedom. This result can be used to build likelihood-based confidence ellipsoids for the parameter θ^* . Given $\mathfrak{z} > 0$, define

$$\mathcal{E}(\mathfrak{z}) = \{\theta : L(\tilde{\theta}, \theta) \leq \mathfrak{z}\} = \left\{ \theta : \sup_{\theta'} L(\theta') - L(\theta) \leq \mathfrak{z} \right\}. \quad (1.23)$$

Theorem 1.4.5. Assume $\mathbf{Y} = \Psi^\top \theta^* + \varepsilon$ with $\varepsilon \sim \mathcal{N}(0, \Sigma)$ and consider the MLE $\tilde{\theta}$. Define \mathfrak{z}_α by $P(\chi_p^2 > 2\mathfrak{z}_\alpha) = \alpha$. Then $\mathcal{E}(\mathfrak{z}_\alpha)$ from (1.23) is an α -confidence set for θ^* .

Exercise 1.4.3. Let $D^2 = \Psi\Sigma^{-1}\Psi^\top$. Check that the likelihood-based CS $\mathcal{E}(\mathfrak{z}_\alpha)$ and estimate-based CS $E(z_\alpha) = \{\theta : \|D(\tilde{\theta} - \theta)\| \leq z_\alpha\}$, $z_\alpha^2 = 2\mathfrak{z}_\alpha$, coincide in the case of the linear modeling:

$$\mathcal{E}(\mathfrak{z}_\alpha) = \{\theta : \|D(\tilde{\theta} - \theta)\|^2 \leq 2\mathfrak{z}_\alpha\}.$$

Another corollary of the chi-squared result is a concentration bound for the maximum likelihood. A similar result was stated for the univariate exponential family model: the value $L(\tilde{\theta}, \theta^*)$ is stochastically bounded with exponential moments, and the bound does not depend on the particular family, parameter value, sample size, etc. Now we can extend this result to the case of a linear Gaussian model. Indeed, Theorem 1.4.3 states that the distribution of $2L(\tilde{\theta}, \theta^*)$ is chi-squared and only depends on the number of parameters to be estimated. The latter distribution concentrates on the ball of radius of order $p^{1/2}$ and the deviation probability is exponentially small.

Theorem 1.4.6. Assume $\mathbf{Y} = \Psi^\top \theta^* + \varepsilon$ with $\varepsilon \sim \mathcal{N}(0, \Sigma)$. Then for every $\mathbf{x} > 0$, it holds with $\varkappa \geq 6.6$

$$\begin{aligned} & \mathbb{P}(2L(\tilde{\theta}, \theta^*) > p + \sqrt{\varkappa xp} \vee (\varkappa x)) \\ &= \mathbb{P}(\|D(\tilde{\theta} - \theta^*)\|^2 > p + \sqrt{\varkappa xp} \vee (\varkappa x)) \leq \exp(-\mathbf{x}). \end{aligned} \quad (1.24)$$

Proof. Define $\xi \stackrel{\text{def}}{=} D(\tilde{\theta} - \theta^*)$. By Theorem 1.3.4 ξ is standard normal vector in \mathbb{R}^p and by Theorem 1.4.3 $2L(\tilde{\theta}, \theta^*) = \|\xi\|^2$. Now the statement (1.24) follows from the general deviation bound for the Gaussian quadratic forms; see Theorem ??.

The main message of this result can be explained as follows: the deviation probability that the estimate $\tilde{\boldsymbol{\theta}}$ does not belong to the elliptic set $E(z) = \{\boldsymbol{\theta} : \|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta})\| \leq z\}$ starts to vanish when z^2 exceeds the dimensionality p of the parameter space. Similarly, the coverage probability that the true parameter $\boldsymbol{\theta}^*$ is not covered by the confidence set $\mathcal{E}(\mathfrak{z})$ starts to vanish when $2\mathfrak{z}$ exceeds p .

Corollary 1.4.1. *Assume $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$ with $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \Sigma)$. Then for every $\mathbf{x} > 0$, it holds with $2\mathfrak{z} = p + \sqrt{\varkappa \mathbf{x} p} \vee (\varkappa \mathbf{x})$ for $\varkappa \geq 6.6$*

$$\mathbb{P}(\mathcal{E}(\mathfrak{z}) \not\ni \boldsymbol{\theta}^*) \leq \exp(-\mathbf{x}).$$

Exercise 1.4.4. Compute \mathfrak{z} ensuring the covering of 95% in the dimension $p = 1, 2, 10, 20$.

1.4.2 A misspecified LPA

Now we discuss the behavior of the fitted log-likelihood for the misspecified linear parametric assumption $\mathbb{E}\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^*$. Let the response function \mathbf{f}^* not be linearly expandable as $\mathbf{f}^* = \Psi^\top \boldsymbol{\theta}^*$. Following to Theorem 1.3.3, define $\boldsymbol{\theta}^* = \mathcal{S}\mathbf{f}^*$ with $\mathcal{S} = (\Psi \Sigma^{-1} \Psi^\top)^{-1} \Psi \Sigma^{-1}$. This point provides the best approximation of the nonlinear response \mathbf{f}^* by a linear parametric fit $\Psi^\top \boldsymbol{\theta}$.

Theorem 1.4.7. *Assume $\mathbf{Y} = \mathbf{f}^* + \boldsymbol{\varepsilon}$ with $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \Sigma)$. Let $\boldsymbol{\theta}^* = \mathcal{S}\mathbf{f}^*$. Then $\tilde{\boldsymbol{\theta}}$ is an R-efficient estimate of $\boldsymbol{\theta}^*$ and*

$$2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) = \boldsymbol{\zeta}^\top D^{-2} \boldsymbol{\zeta} = \|\boldsymbol{\xi}\|^2 \sim \chi_p^2,$$

where $D^2 = \Psi \Sigma^{-1} \Psi^\top$, $\boldsymbol{\zeta} = \nabla L(\boldsymbol{\theta}^*) = \Psi \Sigma^{-1} \boldsymbol{\varepsilon}$, $\boldsymbol{\xi} = D^{-1} \boldsymbol{\zeta}$ is standard normal vector in \mathbb{R}^p and χ_p^2 is a chi-squared random variable with p degrees of freedom. In particular, $\mathcal{E}(\mathfrak{z}_\alpha)$ is an α -CS for the vector $\boldsymbol{\theta}^*$ and the bound of Corollary 1.4.1 applies.

Exercise 1.4.5. Prove the result of Theorem 1.4.7.

1.4.3 A misspecified noise structure

This section addresses the question about the features of the maximum likelihood in the case when the likelihood is built under a wrong assumption about the noise structure. As one can expect, the chi-squared result is not valid anymore in this situation and the distribution of the maximum likelihood depends on the true noise covariance. However, the nice geometric structure of the maximum likelihood manifested by Theorems 1.4.3

and 1.4.5 does not rely on the true data distribution and it is only based on our structural assumptions on the considered model. This helps to get rigorous results about the behaviors of the maximum likelihood and particularly about its concentration properties.

Theorem 1.4.8. *Let $\tilde{\boldsymbol{\theta}}$ be built for the model $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$ with $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \Sigma)$, while the true noise covariance is Σ_0 : $\mathbb{E}\boldsymbol{\varepsilon} = 0$ and $\text{Var}(\boldsymbol{\varepsilon}) = \Sigma_0$. Then*

$$\begin{aligned}\mathbb{E}\tilde{\boldsymbol{\theta}} &= \boldsymbol{\theta}^*, \\ \text{Var}(\tilde{\boldsymbol{\theta}}) &= D^{-2}W^2D^{-2},\end{aligned}$$

where

$$\begin{aligned}D^2 &= \Psi \Sigma^{-1} \Psi^\top, \\ W^2 &= \Psi \Sigma^{-1} \Sigma_0 \Sigma^{-1} \Psi^\top.\end{aligned}$$

Further,

$$2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) = \|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2 = \|\boldsymbol{\xi}\|^2, \quad (1.25)$$

where $\boldsymbol{\xi}$ is a random vector in \mathbb{R}^p with $\mathbb{E}\boldsymbol{\xi} = 0$ and

$$\text{Var}(\boldsymbol{\xi}) = B \stackrel{\text{def}}{=} D^{-1}W^2D^{-1}.$$

Moreover, if $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \Sigma_0)$, then $\tilde{\boldsymbol{\theta}} \sim \mathcal{N}(\boldsymbol{\theta}^*, D^{-2}W^2D^{-2})$ and $\boldsymbol{\xi} \sim \mathcal{N}(0, B)$.

Proof. The moments of $\tilde{\boldsymbol{\theta}}$ have been computed in Theorem 1.4.1 while the equality $2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) = \|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2 = \|\boldsymbol{\xi}\|^2$ is given in Theorem 1.4.3. Next, $\boldsymbol{\zeta} = \nabla L(\boldsymbol{\theta}^*) = \Psi \Sigma^{-1} \boldsymbol{\varepsilon}$ and

$$W^2 \stackrel{\text{def}}{=} \text{Var}(\boldsymbol{\zeta}) = \Psi \Sigma^{-1} \text{Var}(\boldsymbol{\varepsilon}) \Sigma^{-1} \Psi^\top = \Psi \Sigma^{-1} \Sigma_0 \Sigma^{-1} \Psi^\top.$$

This implies that

$$\text{Var}(\boldsymbol{\xi}) = \mathbb{E}\boldsymbol{\xi}\boldsymbol{\xi}^\top = D^{-1} \text{Var}(\boldsymbol{\zeta}) D^{-1} = D^{-1}W^2D^{-1}.$$

It remains to note that if $\boldsymbol{\varepsilon}$ is a Gaussian vector, then $\boldsymbol{\zeta} = \Psi \Sigma^{-1} \boldsymbol{\varepsilon}$, $\boldsymbol{\xi} = D^{-1} \boldsymbol{\zeta}$, and $\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^* = D^{-2} \boldsymbol{\zeta}$ are Gaussian as well.

Exercise 1.4.6. Check that $\Sigma_0 = \Sigma$ leads back to the χ^2 -result.

One can see that the chi-squared result is not valid any more if the noise structure is misspecified. An interesting question is whether the CS $\mathcal{E}(\mathfrak{z})$ can be applied in the

case of a misspecified noise under some proper adjustment of the value \mathfrak{z} . Surprisingly, the answer is not entirely negative. The reason is that the vector $\boldsymbol{\xi}$ from (1.25) is zero mean and its norm has a similar behavior as in the case of the correct noise specification: the probability $\mathbb{P}(\|\boldsymbol{\xi}\| > z)$ starts to degenerate when z^2 exceeds $\mathbb{E}\|\boldsymbol{\xi}\|^2$. A general bound from Theorem 6.6.1 in Section 6.6 implies the following bound for the coverage probability.

Corollary 1.4.2. *Under the conditions of Theorem 1.4.8, for every $\mathbf{x} > 0$, it holds with $\mathfrak{p} = \text{tr}(B)$, $v^2 = 2 \text{tr}(B^2)$, and $a^* = \|B\|_\infty$*

$$\mathbb{P}(2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) > \mathfrak{p} + (2v\mathbf{x}^{1/2}) \vee (6a^*\mathbf{x})) \leq \exp(-\mathbf{x}).$$

Exercise 1.4.7. Show that an overestimation of the noise in the sense $\Sigma \geq \Sigma_0$ preserves the coverage probability for the CS $\mathcal{E}(\mathfrak{z}_\alpha)$, that is, if $2\mathfrak{z}_\alpha$ is the $1 - \alpha$ quantile of χ_p^2 , then $\mathbb{P}(\mathcal{E}(\mathfrak{z}_\alpha) \not\ni \boldsymbol{\theta}^*) \leq \alpha$.

1.5 Random design regression

Consider the linear regression equation

$$\mathbf{Y} = \mathbf{\Psi}^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}, \quad (1.26)$$

where $\boldsymbol{\theta}^* \in \mathbb{R}^p$ is the target parameter, \mathbf{Y} is the n -vector of responses, $\boldsymbol{\varepsilon}$ is the n -vector of errors, and $\mathbf{\Psi} = (\Psi_1, \dots, \Psi_n)$ is a $p \times n$ design matrix with columns $\Psi_i \in \mathbb{R}^p$. We suppose that the error vector $\boldsymbol{\varepsilon}$ satisfies

$$\mathbb{E}\boldsymbol{\varepsilon} = 0, \quad \text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}_n. \quad (1.27)$$

The corresponding Gaussian log-likelihood looks as

$$L(\boldsymbol{\theta}) = -\frac{1}{2\sigma^2} \|\mathbf{Y} - \mathbf{\Psi}^\top \boldsymbol{\theta}\|^2$$

yielding the MLE $\tilde{\boldsymbol{\theta}} = (\mathbf{\Psi}\mathbf{\Psi}^\top)^{-1}\mathbf{\Psi}\mathbf{Y}$. Below we discuss the case of a random design $\mathbf{\Psi}$. The problem for the analysis comes from the inversion of the random matrix $\mathbf{\Psi}\mathbf{\Psi}^\top$. The results below present some conditions under which this random matrix can be replaced by its expectation. We distinguish between two cases. The classical random design assumes that the columns Ψ_i of $\mathbf{\Psi}$ are independent. The other case is when $\mathbf{\Psi}$ is a sum of independent random matrices.

1.5.1 Independent measurements

Let Ψ_1, \dots, Ψ_n be independent. We assume that the design is non degenerate, so that the matrix $\mathbf{M}^2 \stackrel{\text{def}}{=} \mathbb{E}(\mathbf{\Psi}\mathbf{\Psi}^\top)$ is positive. Consider

$$\mathbf{M}^{-1}\{\mathbf{\Psi}\mathbf{\Psi}^\top - \mathbb{E}(\mathbf{\Psi}\mathbf{\Psi}^\top)\}\mathbf{M}^{-1} = \mathbf{A}_1 + \dots + \mathbf{A}_n,$$

where

$$\mathbf{A}_i = \mathbf{M}^{-1}\{\Psi_i\Psi_i^\top - \mathbb{E}(\Psi_i\Psi_i^\top)\}\mathbf{M}^{-1} \quad (1.28)$$

is a symmetric $p \times p$ random matrix with $\mathbb{E}\mathbf{A}_i = 0$. Also define the variance parameter

$$S_n^2 \stackrel{\text{def}}{=} \|\mathbb{E}(\mathbf{A}_1^2 + \dots + \mathbf{A}_n^2)\|. \quad (1.29)$$

We also assume that all design vectors Ψ_i are uniformly bounded with probability one. This implies a uniform bound

$$\|\mathbf{A}_i\| \leq u_n \quad a.s. \quad (1.30)$$

for a small constant u_n . In the case of an i.i.d. design, define

$$\begin{aligned} M_1^2 &\stackrel{\text{def}}{=} \mathbb{E}(\Psi_1 \Psi_1^\top), \\ \sigma_1^2 &\stackrel{\text{def}}{=} \mathbb{E}(M_1^{-1} \Psi_1 \Psi_1^\top M_1^{-1} - I_p)^2. \end{aligned}$$

Also suppose that with probability one

$$\|M_1^{-1} \Psi_1 \Psi_1^\top M_1^{-1} - I_p\| \leq u^*.$$

Then it holds

$$\begin{aligned} M^2 &= n M_1^2, \\ S_n^2 &= n^{-1} \sigma_1^2, \\ u_n &\leq n^{-1} u^* \end{aligned} \tag{1.31}$$

The matrix Bernstein inequality; see Theorem 1.6.1, yields:

Theorem 1.5.1. *Suppose that Ψ_i are independent and A_i from (1.28) fulfill (1.30). Then with $M^2 = \mathbb{E}(\Psi \Psi^\top)$ and S_n^2 defined by (1.29), it holds for all $z > 0$*

$$\mathbb{P}(\|M^{-1} \Psi \Psi^\top M^{-1} - I_p\| > z) \leq 2p \exp\left\{-\frac{z^2}{2S_n^2 + 2u_n z/3}\right\}.$$

If Ψ_i are i.i.d. and (1.31) holds then

$$\mathbb{P}(n^{1/2} \|M^{-1} \Psi \Psi^\top M^{-1} - I_p\| > z) \leq 2p \exp\left\{-\frac{z^2}{2\sigma_1^2 + 2n^{-1/2} u^* z/3}\right\}.$$

Proof. (please check)

For any fixed \mathbf{x} and $\delta > 0$, one can fix any n satisfying

$$n \geq (2\sigma_1^2 \delta^{-2} + 2u\delta^{-1}/3) \{\mathbf{x} + \log(2p)\} \tag{1.32}$$

to ensure

$$\mathbb{P}(\|M^{-1} \Psi \Psi^\top M^{-1} - I_p\| > \delta) \leq e^{-\mathbf{x}}. \tag{1.33}$$

If n and \mathbf{x} are fixed, then one can claim (1.33) for

$$\delta = \sqrt{\frac{2\sigma_1^2}{n} (\mathbf{x} + \log(2p))} + \frac{2u}{3n} (\mathbf{x} + \log(2p)) \tag{1.34}$$

Corollary 1.5.1. *Suppose Ψ_i are i.i.d. and (1.31) holds. If n fulfills (1.32) for some fixed δ and \mathbf{x} , then (1.33) holds true. Similarly, if n and \mathbf{x} are fixed, then δ from (1.34) ensures (1.33).*

Proof. (please check).

The result (1.33) guarantees that on a set $\Omega(\mathbf{x})$ with $\mathbb{P}(\Omega(\mathbf{x})) \leq e^{-x}$, for any $\gamma \in \mathbb{R}^p$

$$(1 - \delta)\gamma^\top \mathbf{M}^2 \gamma \leq \gamma^\top \boldsymbol{\Psi} \boldsymbol{\Psi}^\top \gamma \leq (1 + \delta)\gamma^\top \mathbf{M}^2 \gamma. \quad (1.35)$$

(please check).

Now we apply this result to the MLE in the regression model (1.26). The moments of $\boldsymbol{\varepsilon}$ in (1.27) are understood conditionally on the design $\boldsymbol{\Psi}$. In addition, we assume that the errors $\boldsymbol{\varepsilon}$ are conditionally normal given the design $\boldsymbol{\Psi}$. The log-likelihood ratio can be written in the form

$$L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \frac{1}{\sigma^2} (\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta}^*)^\top \boldsymbol{\Psi}^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) - \frac{1}{2\sigma^2} (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \boldsymbol{\Psi} \boldsymbol{\Psi}^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*).$$

Introduce also an approximating quadratic log-likelihood defined by

$$\begin{aligned} \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) &= \frac{1}{\sigma^2} (\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta}^*)^\top \boldsymbol{\Psi}^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) - \frac{1}{2\sigma^2} (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \mathbf{M}^2 (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \\ &= \nabla^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) - \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \mathbf{D}^2 (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \end{aligned}$$

with the symmetric $p \times p$ matrix $\mathbf{D}^2 = \sigma^2 \mathbf{M}^2 = \sigma^{-2} \mathbb{E}(\boldsymbol{\Psi} \boldsymbol{\Psi}^\top)$ and the random p vector $\nabla = \sigma^{-2} \boldsymbol{\Psi} (\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta}^*)$ satisfying

$$\mathbb{E}(\nabla \mid \boldsymbol{\Psi}) = 0, \quad \text{Var}(\nabla \mid \boldsymbol{\Psi}) = \sigma^{-2} \boldsymbol{\Psi} \boldsymbol{\Psi}^\top.$$

Note that two expressions $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$, $\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ differ only in the quadratic term. We also know that $\boldsymbol{\Psi} \boldsymbol{\Psi}^\top$ is close to its expectation $\mathbf{M}^2 = \mathbb{E} \boldsymbol{\Psi} \boldsymbol{\Psi}^\top$; see (1.35). Moreover, the difference between them can be evaluated using the bound (1.33):

$$\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \frac{1}{2\sigma^2} (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top (\boldsymbol{\Psi} \boldsymbol{\Psi}^\top - \mathbf{M}^2) (\boldsymbol{\theta} - \boldsymbol{\theta}^*).$$

This helps to study the properties of the MLE $\tilde{\boldsymbol{\theta}}$ and of the maximum log-likelihood $L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}) \stackrel{\text{def}}{=} L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta})$. It holds

$$\tilde{\boldsymbol{\theta}} = (\boldsymbol{\Psi} \boldsymbol{\Psi}^\top)^{-1} \boldsymbol{\Psi} \mathbf{Y}, \quad L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}) = \frac{1}{2\sigma^2} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta})^\top \boldsymbol{\Psi} \boldsymbol{\Psi}^\top (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}).$$

Define

$$\boldsymbol{\xi} \stackrel{\text{def}}{=} \mathbf{D}^{-1} \boldsymbol{\Psi} (\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta}^*).$$

Conditioned on the design, $\boldsymbol{\xi}$ behaves as a nearly standard Gaussian vector:

$$\mathbb{E}(\boldsymbol{\xi} | \boldsymbol{\Psi}) = 0, \quad \text{Var}(\boldsymbol{\xi} | \boldsymbol{\Psi}) = D^{-1}(\sigma^{-2}\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)D^{-1}$$

and $\mathbb{E}\{D^{-1}(\sigma^{-2}\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)D^{-1}\} = I_p$. Moreover, restricting to $\Omega(\mathbf{x})$ helps to bound the moment generating function of $\boldsymbol{\xi}$: if $\boldsymbol{\varepsilon}$ is Gaussian conditioned on $\boldsymbol{\Psi}$, then

$$\log \mathbb{E} \left\{ \exp \left(\lambda \frac{\boldsymbol{\gamma}^\top \boldsymbol{\xi}}{\|D\boldsymbol{\gamma}\|} \right) \mathbb{I}(\Omega(\mathbf{x})) \right\} \leq (1 + \delta)^2 \lambda^2 / 2.$$

(please check. this place is nontrivial). Under the PA (1.26) this implies

$$D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \boldsymbol{\xi} = [D(\sigma^{-2}\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)^{-1}D - I_p]\boldsymbol{\xi}$$

yielding by (1.33) on $\Omega(\mathbf{x})$

$$\|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \boldsymbol{\xi}\| \leq \frac{\delta}{1 - \delta} \|\boldsymbol{\xi}\|. \quad (1.36)$$

(please check).

Similarly

$$\begin{aligned} 2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \|\boldsymbol{\xi}\|^2 &= (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)^\top \sigma^{-2}\boldsymbol{\Psi}\boldsymbol{\Psi}^\top (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \|\boldsymbol{\xi}\|^2 \\ &= \boldsymbol{\xi}^\top [D^{-1}(\sigma^{-2}\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)D^{-1} - I_p]\boldsymbol{\xi} \end{aligned}$$

yielding on $\Omega(\mathbf{x})$

$$|2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \|\boldsymbol{\xi}\|^2| \leq \delta \|\boldsymbol{\xi}\|^2. \quad (1.37)$$

These arguments yield the following statement:

Corollary 1.5.2. *Consider the model (1.26) and suppose $\|M^{-1}\boldsymbol{\Psi}\boldsymbol{\Psi}^\top M^{-1} - I_p\| \leq \delta$ for $M^2 = \mathbb{E}(\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)$ and some $\delta < 1/2$ on a dominating set $\Omega(\mathbf{x})$. Then the MLE $\tilde{\boldsymbol{\theta}}$ fulfills (1.36) and (1.37) on this set $\Omega(\mathbf{x})$ with $\boldsymbol{\xi} = D^{-1}\boldsymbol{\Psi}(\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta}^*)$.*

1.5.2 Aggregated random design

Here we discuss another random design setup for the regression model (1.26). Namely, assume that the design matrix $\boldsymbol{\Psi}$ can be represented as a sum of independent matrices $\boldsymbol{\Psi}_1, \dots, \boldsymbol{\Psi}_n$:

$$\boldsymbol{\Psi} = \boldsymbol{\Psi}_1 + \dots + \boldsymbol{\Psi}_n. \quad (1.38)$$

It is natural to expect that this model is close to the usual regression model in which the random matrix $\boldsymbol{\Psi}$ is replaced by its expectation $\mathbb{E}\boldsymbol{\Psi}$. Consider the product $\boldsymbol{\Psi}\boldsymbol{\Psi}^\top$ which

has to be close to the corresponding product $\mathbf{M}^2 \stackrel{\text{def}}{=} \mathbf{E}\Psi \mathbf{E}(\Psi^\top)$. We aim at bounding the normalized difference $\Psi - \mathbf{E}\Psi$ in the operator norm. Define for $i = 1, \dots, n$

$$V_i^2 \stackrel{\text{def}}{=} \mathbf{E}(\Psi_i \Psi_i^\top) - \mathbf{E}\Psi_i \mathbf{E}\Psi_i^\top,$$

and

$$S_n^2 \stackrel{\text{def}}{=} \|\mathbf{M}^{-1}(V_1^2 + \dots + V_n^2)\mathbf{M}^{-1}\|.$$

Typically S_n^2 is inversely proportional to n . We again assume that all design vectors Ψ_i are uniformly bounded with probability one. This implies a uniform bound

$$\|\mathbf{M}^{-1}(\Psi_i - \mathbf{E}\Psi_i)\| \leq u_n \quad a.s.$$

for a small constant u_n . Now consider

$$\mathbf{M}^{-1}(\Psi - \mathbf{E}\Psi) = \sum_{i=1}^n \mathbf{M}^{-1}(\Psi_i - \mathbf{E}\Psi_i).$$

The matrix Bernstein inequality; see Theorem 1.6.1, yields for any $z \geq 0$

$$\mathbb{P}(\|\mathbf{M}^{-1}(\Psi - \mathbf{E}\Psi)\| \geq z) \leq 2(p+q) \exp\left\{-\frac{z^2}{2S_n^2 + 2u_n z/3}\right\} \quad (1.39)$$

In the case with i.i.d. Ψ_i , define

$$\begin{aligned} M_1^2 &\stackrel{\text{def}}{=} \mathbf{E}\Psi_1 \mathbf{E}\Psi_1^\top, \\ \sigma_1^2 &\stackrel{\text{def}}{=} M_1^{-1} \mathbf{E}(\Psi_1 \Psi_1^\top) M_1^{-1} - I_p, \end{aligned}$$

and suppose that

$$\|M_1^{-1}(\Psi_1 - \mathbf{E}\Psi_1)\| \leq u_1.$$

Then it holds

$$\begin{aligned} \mathbf{M}^2 &= n^2 M_1, \\ u_n &\leq n^{-1} u_1, \end{aligned}$$

and

$$\mathbb{P}\left(n^{1/2} \|\mathbf{M}^{-1}(\Psi - \mathbf{E}\Psi)\| \geq z\right) \leq 2(p+q) \exp\left\{-\frac{z^2}{2\sigma_1^2 + 2u_1 z/(3n^{1/2})}\right\}.$$

The result (1.39) implies that Ψ is close to $\mathbf{E}\Psi$. The MLE $\tilde{\theta}$ also involves the product $\Psi\Psi^\top$ which has to be close to the corresponding product $\mathbf{M}^2 = \mathbf{E}\Psi \mathbf{E}(\Psi^\top)$. Below we assume that \mathbf{M} is sufficiently large and bound the difference $\mathbf{M}^{-1}\Psi\Psi^\top\mathbf{M}^{-1} - I_p$.

Proposition 1.5.1. *Let $\|M^{-1}(\Psi - E\Psi)\| \leq \delta$ for some $\delta > 0$. Then*

$$\|M^{-1}\Psi\Psi^\top M^{-1} - I_p\| \leq \delta^2 + 2\delta.$$

Proof. One can bound

$$M^{-1}\Psi\Psi^\top M^{-1} - I_p = M^{-1}(\Psi - E\Psi)(\Psi - E\Psi)^\top M^{-1} + 2M^{-1}(\Psi - E\Psi)E(\Psi^\top)M^{-1}.$$

For any unit vector $\gamma \in \mathbb{R}^p$, the definition of M implies

$$\|E(\Psi^\top)M^{-1}\gamma\|^2 = \gamma^\top M^{-1}E(\Psi)E(\Psi^\top)M^{-1}\gamma = \|\gamma\|^2 = 1.$$

Therefore,

$$\|M^{-1}(\Psi - E\Psi)E(\Psi^\top)M^{-1}\gamma\| \leq \|M^{-1}(\Psi - E\Psi)\|$$

thus

$$\|M^{-1}\Psi\Psi^\top M^{-1} - I_p\| \leq \|M^{-1}(\Psi - E\Psi)\|^2 + 2\|M^{-1}(\Psi - E\Psi)\|,$$

and the result follows.

Corollary 1.5.2 applies in this situation without any change.

1.5.3 Application to instrumental regression

Observed: a sample from (Y, X, W) . Model

$$Y = f(X) + U, \quad E[U | W] = 0.$$

where Y , an explained variable, X , an explanatory variable, W , an instrument. The target is the regression function $f(\cdot)$.

Let $\psi_1(x), \dots, \psi_j(x), \dots$ be a functional basis. Consider a finite approximation

$$f(x) = \theta_1\psi_1(x) + \dots + \theta_p\psi_p(x)$$

or in vector form

$$f(x) = \psi(x)^\top \theta$$

with $\psi(x) = (\psi_1(x), \dots, \psi_p(x))^\top \in \mathbb{R}^p$ and $\theta = (\theta_1, \dots, \theta_p)^\top \in \mathbb{R}^p$. This leads to an approximating model

$$Y = \psi(X)^\top \theta^* + U, \quad E[U | W] = 0.$$

The constraint $\mathbb{E}[U | W] = 0$ means that for any function $\phi(W)$

$$\mathbb{E}[Y\phi(W)] = \mathbb{E}[\phi(W)\psi(X)^\top] \boldsymbol{\theta}^*.$$

We apply a *discretization* or *finite dimensional approximation*: for a finite collection of functions $\boldsymbol{\phi}(w) = (\phi_1(w), \dots, \phi_q(w))^\top$, it holds

$$\mathbb{E}[Y\boldsymbol{\phi}(W)] = \mathbb{E}[\boldsymbol{\psi}(X)\boldsymbol{\phi}(W)^\top]^\top \boldsymbol{\theta}^* = \mathbf{T}^\top \boldsymbol{\theta}^*$$

with

$$\mathbf{T} = \mathbb{E}[\boldsymbol{\psi}(X)\boldsymbol{\phi}(W)^\top] \in \mathbb{R}^{p \times q}.$$

Define

$$\begin{aligned} \mathbf{Z} &= \mathbb{E}_n[\mathbf{Y} \boldsymbol{\phi}(W)] &= n^{-1} \sum_i Y_i \boldsymbol{\phi}(W_i) &\in \mathbb{R}^q, \\ \mathbf{T}_n &= \mathbb{E}_n[\boldsymbol{\psi}(X) \boldsymbol{\phi}(W)^\top] &= n^{-1} \sum_i \boldsymbol{\psi}(X_i) \boldsymbol{\phi}(W_i)^\top &\in \mathbb{R}^{q \times p}, \\ \boldsymbol{\varepsilon} &= \mathbb{E}_n[\boldsymbol{\phi}(W) \mathbf{U}] &= n^{-1} \sum_i \boldsymbol{\phi}(W_i) U_i &\in \mathbb{R}^q. \end{aligned} \tag{1.40}$$

The original problems reduces to

$$\mathbf{Z} = \mathbf{T}^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon},$$

where $\boldsymbol{\varepsilon}$ is the error q -vector, $\mathbf{T} = \mathbb{E}[\boldsymbol{\psi}(X) \boldsymbol{\phi}(W)^\top]$ is an unknown $p \times q$ matrix and only its empirical counterpart \mathbf{T}_n is available. In such cases one speaks of *an inverse problem with error in operator*. The main problem for the analysis in this model is that \mathbf{T}_n is random and correlated with \mathbf{Z} and $\boldsymbol{\varepsilon}$. The goal is to build an estimator $\tilde{\boldsymbol{\theta}}$ of the vector $\boldsymbol{\theta}^*$ leading to the estimator $\tilde{f}(x) = \boldsymbol{\psi}(x)^\top \tilde{\boldsymbol{\theta}}$ of the response.

The natural plug-in approach suggests to replace the unknown operator \mathbf{T} by its empirical counterpart \mathbf{T}_n leading to the approximating linear model

$$\mathbf{Z} = \mathbf{T}_n^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}.$$

with the random design $\mathbf{T}_n = n^{-1} \sum_i \boldsymbol{\psi}(X_i) \boldsymbol{\phi}(W_i)^\top$ so that the setup (1.38) applies. The corresponding least square estimator of $\boldsymbol{\theta}^*$ reads as

$$\tilde{\boldsymbol{\theta}} = (\mathbf{T}_n \mathbf{T}_n^\top)^{-1} \mathbf{T}_n \mathbf{Z}. \tag{1.41}$$

The results of Proposition 1.5.1 justify that the random matrix $\mathbf{T}_n \mathbf{T}_n^\top$ is very close to the product $\mathbb{E}(\mathbf{T}_n) \mathbb{E}(\mathbf{T}_n^\top) = \mathbf{T} \mathbf{T}^\top$ and the theoretical study of the properties of the estimator $\tilde{\boldsymbol{\theta}}$ can be done with $\mathbf{T} \mathbf{T}^\top$ in place of $\mathbf{T}_n \mathbf{T}_n^\top$ in (1.41). Similarly one can justify that the product $\mathbf{T}_n \mathbf{Z}$ behaves nearly as $\mathbf{T} \mathbf{Z}$.

Below we assume for simplicity that all triples (Y_i, X_i, W_i) are i.i.d. so that $T_i = \psi(X_i) \phi(W_i)^\top$ are also i.i.d. Define $\mathbf{M}^2 = \mathbf{T}\mathbf{T}^\top$ and

$$\sigma_1^2 = \|\mathbf{M}^{-1} \mathbb{E}(T_i T_i^\top) \mathbf{M}^{-1} - I_p\|$$

and suppose that it holds almost surely

$$\|\mathbf{M}^{-1}(T_i - \mathbf{T})\| \leq u. \quad (1.42)$$

Theorem 1.5.2. *Let (Y_i, X_i, W_i) be i.i.d. and T_i from (1.40) fulfill (1.42). Then for any $z > 0$*

$$\mathbb{P}\left(\sqrt{n}\|\mathbf{M}^{-1}(\mathbf{T}_n - \mathbf{T})\| > z\right) \leq 2(p+q) \exp\left\{-\frac{z^2}{2\sigma_1^2 + 2uz/(3n^{1/2})}\right\}.$$

Moreover, if $\|\mathbf{M}^{-1}(\mathbf{T}_n - \mathbf{T})\| \leq \delta$, then

$$\|\mathbf{M}^{-1} \mathbf{T}_n \mathbf{T}_n^\top \mathbf{M}^{-1} - I_p\| \leq \delta^2 + 2\delta.$$

1.6 Matrix Bernstein inequality

This section collects some useful facts about deviation of stochastic matrices from their mean. We mainly follows ?.

Theorem 1.6.1 (Matrix Bernstein). *Let $\mathbf{S}_1, \dots, \mathbf{S}_n$ be independent random matrices with common dimension $d_1 \times d_2$. Assume that each matrix has bounded deviation from its mean:*

$$\|\mathbf{S}_i - \mathbb{E}\mathbf{S}_i\| \leq R, \quad i = 1, \dots, n.$$

Form the sum $\mathbf{Z} = \sum_{i=1}^n \mathbf{S}_i$, and introduce a variance parameter

$$\sigma^2 \stackrel{\text{def}}{=} \max\left\{\|\mathbb{E}[(\mathbf{Z} - \mathbb{E}\mathbf{Z})(\mathbf{Z} - \mathbb{E}\mathbf{Z})^\top]\|, \|\mathbb{E}[(\mathbf{Z} - \mathbb{E}\mathbf{Z})^\top(\mathbf{Z} - \mathbb{E}\mathbf{Z})]\|\right\}.$$

Then

$$\mathbb{P}(\|\mathbf{Z} - \mathbb{E}\mathbf{Z}\| > z) \leq (d_1 + d_2) \exp\left\{-\frac{z^2/2}{\sigma^2 + Rz/3}\right\}.$$

Furthermore,

$$\mathbb{E}\|\mathbf{Z} - \mathbb{E}\mathbf{Z}\| \leq \sqrt{2\sigma^2 \log(d_1 + d_2)} + \frac{1}{3}R \log(d_1 + d_2).$$

Sieve model selection in linear models

Here we consider the problem of sieve model selection in linear regression model. A high dimensional linear model is approximated by its projection, the main issue is a proper choice of the cut-off parameter.

2.1 Projection estimation. Loss and risk

This section presents the main definitions and properties of the projection estimate.

2.1.1 A linear model

The linear parametric assumption can be stated in the following form: the observed vector \mathbf{Y} is supposed to follow the linear regression model with a homogeneous Gaussian error vector $\boldsymbol{\varepsilon}$:

$$\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n).$$

Here Ψ is the design matrix formed by a collection of basis vectors $\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_p \in \mathbb{R}^n$, and p is the parameter dimension assumed to be large or even infinity. The parametric assumption is only an idealization, for the true data distribution \mathbb{P} of the vector $\mathbf{Y} \in \mathbb{R}^n$ we assume that the observations Y_i are independent and the errors $\varepsilon_i = Y_i - \mathbb{E}Y_i$ have some moments. We also assume that the response vector $\mathbf{f} = \mathbb{E}\mathbf{Y}$ can be well approximated by $\Psi^\top \boldsymbol{\theta}$ for a proper choice of $\boldsymbol{\theta}$, or, equivalently,

$$\mathbf{f} \approx \theta_1^* \boldsymbol{\psi}_1 + \dots + \theta_p^* \boldsymbol{\psi}_p.$$

The MLE or oLSE of the parameter vector $\boldsymbol{\theta}^*$ for this model reads as

$$\tilde{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \|\mathbf{Y} - \Psi^\top \boldsymbol{\theta}\|^2 = (\Psi\Psi^\top)^{-1} \Psi\mathbf{Y} = \mathcal{S}\mathbf{Y}$$

with $\mathcal{S} = (\Psi\Psi^\top)^{-1} \Psi$.

2.1.2 Linear decomposition

The estimate $\tilde{\boldsymbol{\theta}} = \mathcal{S}\mathbf{Y} = (\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)^{-1}\boldsymbol{\Psi}\mathbf{Y}$ is linear in \mathbf{Y} . This implies by the model equation $\mathbf{Y} = \mathbf{f}^* + \boldsymbol{\varepsilon}$ the decomposition

$$\tilde{\boldsymbol{\theta}} = \mathcal{S}\mathbf{Y} = \mathcal{S}\mathbf{f}^* + \mathcal{S}\boldsymbol{\varepsilon}.$$

We already know that $\tilde{\boldsymbol{\theta}}$ is unbiased, that is,

$$\mathbb{E}\tilde{\boldsymbol{\theta}} = \boldsymbol{\theta}^*,$$

where the target $\boldsymbol{\theta}^*$ can be defined as the vector of coefficients of the best linear fit:

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \|\mathbf{f} - \boldsymbol{\Psi}^\top\boldsymbol{\theta}\|^2 = (\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)^{-1}\boldsymbol{\Psi}\mathbf{f}.$$

If the error homogeneity assumption is correct, that is, if $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2\mathbf{I}_n)$, then the variance of $\tilde{\boldsymbol{\theta}}$ follows

$$\operatorname{Var}(\tilde{\boldsymbol{\theta}}) = \mathbb{E}\{(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)^\top\} = (\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)^{-1}\boldsymbol{\Psi}\mathbb{E}(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top)\boldsymbol{\Psi}^\top(\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)^{-1} = \sigma^2(\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)^{-1}.$$

Moreover, for any $q \times p$ matrix W , it holds in the same way

$$\mathbb{E}[W(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)]^2 = \operatorname{tr}[W(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)^\top W^\top] = \sigma^2 \operatorname{tr}[W(\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)^{-1}W^\top].$$

2.1.3 Inhomogeneous errors

In the general case of $\operatorname{Var}(\boldsymbol{\varepsilon}) = \boldsymbol{\Sigma}$,

$$\begin{aligned} \operatorname{Var}(\tilde{\boldsymbol{\theta}}) &= \mathbb{E}\mathbb{E}\{(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)^\top\} = (\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)^{-1}\boldsymbol{\Psi}\mathbb{E}(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top)\boldsymbol{\Psi}^\top(\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)^{-1} \\ &= (\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)^{-1}\boldsymbol{\Psi}\boldsymbol{\Sigma}\boldsymbol{\Psi}^\top(\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)^{-1}. \end{aligned}$$

Exercise 2.1.1. Consider the regression model

$$Y_i = \theta_1^* \psi_1(X_i) + \dots + \theta_p^* \psi_p(X_i) + \varepsilon_i \quad (2.1)$$

with independent heterogeneous errors $\operatorname{Var}(\varepsilon_i) = \sigma_i^2$. Consider the MLE $\tilde{\boldsymbol{\theta}}$ and the LSE $\tilde{\boldsymbol{\theta}}_{LSE} = (\boldsymbol{\Psi}\boldsymbol{\Psi}^\top)^{-1}\boldsymbol{\Psi}\mathbf{Y}$ and corresponding to homogeneous errors.

- Compute $\tilde{\boldsymbol{\theta}}$
- Show that $\mathbb{E}\tilde{\boldsymbol{\theta}} = \mathbb{E}\tilde{\boldsymbol{\theta}}_{LSE} = \boldsymbol{\theta}^*$.
- Compute the variance $\operatorname{Var}(\tilde{\boldsymbol{\theta}})$ and the variance $\operatorname{Var}(\tilde{\boldsymbol{\theta}}_{LSE})$;
- show that $\operatorname{Var}(\tilde{\boldsymbol{\theta}}_{LSE}) \geq \operatorname{Var}(\tilde{\boldsymbol{\theta}})$;

- check that $\text{Var}(\tilde{\boldsymbol{\theta}}_{LSE}) = \text{Var}(\tilde{\boldsymbol{\theta}})$ iff all the σ_i are equal to each other.

To illustrate the performance of the MLE $\tilde{\boldsymbol{\theta}}$, consider the case of the orthonormal design $\Psi\Psi^\top = I_p$ and homogeneous errors $\boldsymbol{\varepsilon}$. Then

$$\tilde{\boldsymbol{\theta}} = \Psi\mathbf{Y}, \quad \text{Var}(\tilde{\boldsymbol{\theta}}) = \sigma^2 I_p.$$

and for any symmetric positive matrix W , it holds

$$\mathbb{E}\{W(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\}^2 = \sigma^2 \text{tr}(WW^\top)$$

Now we consider the prediction $\tilde{\mathbf{f}} = \Psi^\top \tilde{\boldsymbol{\theta}}$.

2.1.4 Linear decomposition

It follows

$$\tilde{\mathbf{f}} = \Psi^\top \tilde{\boldsymbol{\theta}} = \Psi^\top (\Psi\Psi^\top)^{-1} \Psi\mathbf{Y} = \Pi\mathbf{Y}.$$

Here Π is the projector in the space \mathbb{R}^n on the linear subspace spanned by the basis vectors $\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_p$. The model equation $\mathbf{Y} = \mathbf{f}^* + \boldsymbol{\varepsilon}$ implies the decomposition

$$\tilde{\mathbf{f}} = \Pi\mathbf{f} + \Pi\boldsymbol{\varepsilon}.$$

It implies

$$\mathbb{E}\tilde{\mathbf{f}} = \Pi\mathbf{f}.$$

Moreover, under the noise homogeneity $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 I_n$, it holds

$$\text{Var}(\tilde{\mathbf{f}}) = \mathbb{E}(\Pi\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top \Pi) = \sigma^2 \Pi.$$

2.1.5 Quadratic loss. Bias-variance decomposition

For the quadratic loss function $\varrho(\tilde{\mathbf{f}}, \mathbf{f}^*) = \|\tilde{\mathbf{f}} - \mathbf{f}^*\|^2$, it follows

$$\begin{aligned} \varrho(\tilde{\mathbf{f}}, \mathbf{f}^*) &= \|\tilde{\mathbf{f}} - \mathbf{f}^*\|^2 = \|\Pi\mathbf{f} - \mathbf{f}^* + \Pi\boldsymbol{\varepsilon}\|^2 \\ &= \|(I_p - \Pi)\mathbf{f}^*\|^2 + \|\Pi\boldsymbol{\varepsilon}\|^2 + 2(\Pi\boldsymbol{\varepsilon})^\top (I_p - \Pi)\mathbf{f}^* \\ &= \|(I_p - \Pi)\mathbf{f}^*\|^2 + \|\Pi\boldsymbol{\varepsilon}\|^2. \end{aligned} \tag{2.2}$$

Here we have used that Π is a projector in \mathbb{R}^n which implies $\Pi^\top(I_p - \Pi) = 0$. Therefore, the quadratic risk $\mathcal{R}(\tilde{\mathbf{f}}, \mathbf{f}) = \mathbb{E}\varrho(\tilde{\mathbf{f}}, \mathbf{f})$ satisfies under homogeneous errors $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 I_p$

$$\mathcal{R}(\tilde{\mathbf{f}}, \mathbf{f}^*) = \|(I_p - \Pi)\mathbf{f}^*\|^2 + \|\Pi\boldsymbol{\varepsilon}\|^2 = \|(I_p - \Pi)\mathbf{f}^*\|^2 + \sigma p$$

as

$$\mathbb{E}\|\Pi\boldsymbol{\varepsilon}\|^2 = \mathbb{E}(\Pi\boldsymbol{\varepsilon})^\top(\Pi\boldsymbol{\varepsilon}) = \mathbb{E}\text{tr}(\Pi\boldsymbol{\varepsilon})(\Pi\boldsymbol{\varepsilon})^\top = \text{tr}(\Pi\mathbb{E}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top\Pi^\top) = \sigma^2\text{tr}(\Pi\Pi^\top) = \sigma^2 p.$$

The term $\|(I_p - \Pi)\mathbf{f}^*\|^2$ is usually called *the squared bias* and it describes the accuracy of approximation of \mathbf{f}^* by its projection $\Pi\mathbf{f}^*$ on the space generated by the basis functions ψ_1, \dots, ψ_p . The term $\sigma^2 p$ called *the variance* measures the statistical error related to estimation of p unknown coefficients in the decomposition of $\Pi\mathbf{f}^* = \theta_1^*\psi_1 + \dots + \theta_p^*\psi_p$.

Exercise 2.1.2. Consider the nonparametric model

$$Y_i = f(X_i) + \varepsilon_i, \quad \text{Var}(\varepsilon_i) = \sigma_i^2$$

and the parametric approximation (2.1). Consider the qMLE (LSE) $\tilde{\mathbf{f}}_{LSE} = \Pi\mathbf{Y}$ for $\Pi = \Psi^\top(\Psi\Psi^\top)^{-1}\Psi$.

- Derive the bias-variance decomposition for the quadratic losses $\|\tilde{\mathbf{f}}_{LSE} - \mathbf{f}\|^2$ and of the risk $\mathbb{E}\|\tilde{\mathbf{f}}_{LSE} - \mathbf{f}\|^2$.
- Compute the variance term of $\tilde{\mathbf{f}}_{LSE}$ and of $\tilde{\mathbf{f}} = \Psi^\top\tilde{\boldsymbol{\theta}}$ for the MLE $\tilde{\boldsymbol{\theta}}$.

2.1.6 Projection estimation and the model choice problem

In this section we consider the linear model $\mathbf{Y} = \Psi^\top\boldsymbol{\theta} + \boldsymbol{\varepsilon}$ with a p -dimensional parameter p which is large or even infinity. The full dimensional estimation of the parameter $\boldsymbol{\theta}$ can be highly inefficient. Here we consider the simplest method of complexity reduction called *projection*. The idea is to use just a submodes corresponding to the reduced subset of parameters.

We associate the rows of the design matrix Ψ with basis vectors in \mathbb{R}^n . By Ψ_m we denote a sub matrix of Ψ composed of the first m rows ψ_1, \dots, ψ_m . It corresponds to the reduced regression model

$$\mathbf{Y} = \Psi_m^\top\boldsymbol{\theta}_m + \boldsymbol{\varepsilon}$$

with the parameter $\boldsymbol{\theta}_m$ from \mathbb{R}^m . The corresponding estimate $\tilde{\boldsymbol{\theta}}_m$ and the predictor $\tilde{\mathbf{f}}_m$ read as

$$\begin{aligned}\tilde{\boldsymbol{\theta}}_m &= (\Psi_m \Psi_m^\top)^{-1} \Psi_m \mathbf{Y}, \\ \tilde{\mathbf{f}}_m &= \Psi_m (\Psi_m \Psi_m^\top)^{-1} \Psi_m \mathbf{Y} = \Pi_m \mathbf{Y}\end{aligned}$$

where Π_m is a projector in \mathbb{R}^n on the subspace spanned by the basis vectors $\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_m$.

In the case of an orthonormal design, one just considers the first m empirical coefficients z_1, \dots, z_m and drop the others. The corresponding parameter estimate $\tilde{\boldsymbol{\theta}}_m$ reads as

$$\tilde{\theta}_{m,j} = \begin{cases} z_j & \text{if } j \leq m, \\ 0 & \text{otherwise} \end{cases}$$

with $z_j = \boldsymbol{\psi}_j^\top \mathbf{Y}$. The response vector $\mathbf{f}^* = \mathbb{E}\mathbf{Y}$ is estimated by $\Psi^\top \tilde{\boldsymbol{\theta}}_m$ leading to the representation

$$\tilde{\mathbf{f}}_m = z_1 \boldsymbol{\psi}_1 + \dots + z_m \boldsymbol{\psi}_m.$$

In other words, $\tilde{\mathbf{f}}_m$ is just a projection of the observed vector \mathbf{Y} onto the subspace L_m spanned by the first m basis vectors $\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_m$: $L_m = \langle \boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_m \rangle$. This explains the name of the method. Clearly one can study the properties of $\tilde{\boldsymbol{\theta}}_m$ or $\tilde{\mathbf{f}}_m$ using the methods of previous sections. However, one more question for this approach is still open: a proper choice of m . The standard way of accessing this issue is based on the analysis of the quadratic risk.

Consider first the prediction risk defined as $\mathcal{R}(\tilde{\mathbf{f}}_m, \mathbf{f}^*) = \mathbb{E} \|\tilde{\mathbf{f}}_m - \mathbf{f}^*\|^2$. Below we focus on the case of a homogeneous noise with $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}_p$. An extension to the colored noise is possible. Recall that $\tilde{\mathbf{f}}_m$ effectively estimates the vector $\mathbf{f}_m = \Pi_m \mathbf{f}^*$, where Π_m is the projector on L_m ; see Section 1.2.3. Moreover, the quadratic risk $\mathcal{R}(\tilde{\mathbf{f}}_m, \mathbf{f}^*)$ can be decomposed as

$$\mathcal{R}(\tilde{\mathbf{f}}_m, \mathbf{f}^*) = \|\mathbf{f}^* - \Pi_m \mathbf{f}^*\|^2 + \sigma^2 m = \sigma^2 m + \sum_{j=m+1}^p \theta_j^{*2}. \quad (2.3)$$

Obviously the squared bias $\|\mathbf{f}^* - \Pi_m \mathbf{f}^*\|^2$ decreases with m while the variance $\sigma^2 m$ linearly grows with m . Risk minimization leads to the so called *bias-variance trade-off*: one selects m which minimizes the risk $\mathcal{R}(\tilde{\mathbf{f}}_m, \mathbf{f}^*)$ over all possible m :

$$m^* \stackrel{\text{def}}{=} \underset{m}{\operatorname{argmin}} \mathcal{R}(\tilde{\mathbf{f}}_m, \mathbf{f}^*) = \underset{m}{\operatorname{argmin}} \{ \|\mathbf{f}^* - \Pi_m \mathbf{f}^*\|^2 + \sigma^2 m \}. \quad (2.4)$$

Unfortunately this choice requires some information about the bias $\|\mathbf{f}^* - \Pi_m \mathbf{f}^*\|$ which depends on the unknown vector \mathbf{f}^* . As this information is not available in typical situation, the value m^* is also called an *oracle* choice. A data-driven choice of m is one of the central issue in the nonparametric statistics.

The situation is not changed if we consider the estimation risk $\mathbb{E}\|\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*\|^2$. Indeed, the basis orthogonality $\Psi\Psi^\top = I_p$ implies for $\mathbf{f}^* = \Psi^\top \boldsymbol{\theta}^*$

$$\|\tilde{\mathbf{f}}_m - \mathbf{f}^*\|^2 = \|\Psi^\top \tilde{\boldsymbol{\theta}}_m - \Psi^\top \boldsymbol{\theta}^*\|^2 = \|\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*\|^2$$

and minimization of the estimation risk coincides with minimization of the prediction risk.

The problem of selecting the model m^* can be stated in different ways depending on what is the target and objective of the method. Usually the problem is formulated as the problem of *adaptive estimation* and the one aims at constructing an estimate $\hat{\boldsymbol{\theta}}$ whose risk is close to the risk of the oracle $\tilde{\boldsymbol{\theta}}_{m^*}$. The problem of *model selection* mainly focuses choosing a proper model \hat{m} on the base of available data. The latter problem is appealing if one is concerned with inference, prediction, or some other model-based question. To understand the difference between two possible setups, consider the ideal situation when the risk is completely flat: $\mathcal{R}_m \equiv \mathbf{C}$. Then any model choice yields the same risk and one can free to take any model. In terms of building a confidence statement, for prediction or testing, the model choice matters a lot and a smaller model (in term of complexity) will be much more useful. In some sense, two mentioned objectives are contradictory: a flat risk is very good for estimation and enables us to apply a simple rule-of-thumb for choosing the parameter m . However, identification of a good model is very hard for models with a flat risk function. At the same time, the case of a profiled risk makes the choice of the model crucial but it can be identified much easier. Below we try to address both issues: estimation of the parameter $\boldsymbol{\theta}^*$ and of the oracle model m^* .

2.2 Unbiased risk estimation

The “oracle” choice m^* cannot be implemented because the bias term $\|\mathbf{f}^* - \Pi_m \mathbf{f}^*\|^2$ depends on the target object \mathbf{f}^* . Now we want to develop a data-driven rule which attempts to reproduce (mimic) the oracle. The first naive idea is to look at the empirical risk (data fit) $\|\mathbf{Y} - \tilde{\mathbf{f}}_m\|^2$ in which we replace the function \mathbf{f} by the data \mathbf{Y} . Unfortunately, this rule leads to the trivial solution

$$\hat{m} = \underset{m}{\operatorname{argmin}} \|\mathbf{Y} - \tilde{\mathbf{f}}_m\|^2 = p.$$

Indeed, the value $\|\mathbf{Y} - \tilde{\mathbf{f}}_m\|^2$ monotonously decreases with m as follows from the next lemma.

Lemma 2.2.1. *Consider the projection estimator $\tilde{\mathbf{f}}_m = \Pi_m \mathbf{Y}$. For two different values $m' > m$, the following statements hold:*

- $\Pi_{m',m} \stackrel{\text{def}}{=} \Pi_{m'} - \Pi_m$ is a projector in \mathbb{R}^n .
- If Ψ is orthogonal then $\Pi_{m',m}$ projects onto subspace generated by $\psi_{m+1}, \dots, \psi_{m'}$.
- The next identity is fulfilled:

$$\|\mathbf{Y} - \Pi_m \mathbf{Y}\|^2 - \|\mathbf{Y} - \Pi_{m'} \mathbf{Y}\|^2 = \|\Pi_{m',m} \mathbf{Y}\|^2 = \|\Pi_{m',m} \mathbf{f} + \Pi_{m',m} \boldsymbol{\varepsilon}\|^2 \geq 0.$$

Proof. It obviously holds

$$\begin{aligned} \|\mathbf{Y} - \Pi_m \mathbf{Y}\|^2 - \|\mathbf{Y} - \Pi_{m'} \mathbf{Y}\|^2 &= \mathbf{Y}^\top (I - \Pi_m) \mathbf{Y} - \mathbf{Y}^\top (I - \Pi_{m'}) \mathbf{Y} \\ &= \mathbf{Y}^\top (\Pi_{m'} - \Pi_m) \mathbf{Y} = \|\Pi_{m',m} \mathbf{Y}\|^2. \end{aligned}$$

Therefore, empirical risk minimization always tries to select the largest possible model which provides the best data fit. In the extreme case of $m = n$, we obtain the perfect fit $\tilde{\mathbf{f}}_m = \mathbf{Y}$, that is, the estimate coincides with the data. This is formally correct but the corresponding squared risk is equal to $\sigma^2 p$ which can be a very large number. So, the empirical risk minimization does not do the required job, it does not mimic the “oracle” risk minimizer. Now we try to look more attentively at the empirical risk to understand the origin of the problem. First compute its expectation.

Lemma 2.2.2. *It holds under homogeneous errors $\boldsymbol{\varepsilon}$:*

- For each m

$$\mathbb{E} \|\mathbf{Y} - \Pi_m \mathbf{Y}\|^2 = \|\mathbf{f}^* - \Pi_m \mathbf{f}^*\|^2 + \sigma^2(n - m). \quad (2.5)$$

- For any $m' > m$

$$\mathbb{E} \|\mathbf{Y} - \Pi_m \mathbf{Y}\|^2 - \mathbb{E} \|\mathbf{Y} - \Pi_{m'} \mathbf{Y}\|^2 = \|\Pi_{m',m} \mathbf{f}^*\|^2 + \sigma^2(m' - m)$$

Proof. Obvious.

The first term in the statement (2.5) is exactly the squared bias which is a good news: the empirical risk contains the same term which we need in the squared risk evaluation. Unfortunately, the second term $\sigma^2(n - m)$ behaves differently than the similar variance term $\sigma^2 m$. Another good news is that both variance terms are known to us. Therefore, one can easily make a correction of the empirical risk which delivers an unbiased risk estimate: just add $\sigma^2(2m - n)$. Define

$$\tilde{\mathcal{R}}_m = \|\mathbf{Y} - \Pi_m \mathbf{Y}\|^2 + 2\sigma^2 m.$$

Then it holds

$$\mathbb{E}\tilde{\mathcal{R}}_m = \mathcal{R}(\tilde{\mathbf{f}}_m, \mathbf{f}^*) + \sigma^2 n.$$

In words, the expectation of $\tilde{\mathcal{R}}_m$ is equal to the risk $\mathcal{R}(\tilde{\mathbf{f}}_m, \mathbf{f}^*)$ up to the fixed term $\sigma^2 n$ which does not affect the model choice. This suggests to define

$$\hat{m} \stackrel{\text{def}}{=} \underset{m}{\operatorname{argmin}} \tilde{\mathcal{R}}_m = \underset{m}{\operatorname{argmin}} (\|\mathbf{Y} - \Pi_m \mathbf{Y}\|^2 + 2\sigma^2 m). \quad (2.6)$$

This rule is known as Akaike information criteria (AIC) and it is very popular in practical applications. It suggests to balance the data fit measured by $\|\mathbf{Y} - \Pi_m \mathbf{Y}\|^2$ and the model complexity $2\sigma^2 m$. One can say that this rule selects a model with a possibly small complexity $\sigma^2 m$ still providing a reasonable data fit $\|\mathbf{Y} - \Pi_m \mathbf{Y}\|^2$.

Exercise 2.2.1. Consider the projection estimator $\tilde{\mathbf{f}}_m = \Pi_m \mathbf{Y}$ for the model (2.1) with $\Pi_m = \Psi_m^\top (\Psi_m \Psi_m^\top)^{-1} \Psi_m$. For two different values $m' > m$:

- Check that $\Pi_{m',m} \stackrel{\text{def}}{=} \Pi_{m'} - \Pi_m$ is a projector in \mathbb{R}^n . Describe its image in the orthogonal case when $\Psi \Psi^\top$ is a diagonal matrix.
- Check the identities

$$\begin{aligned} \|\mathbf{Y} - \Pi_m \mathbf{Y}\|^2 - \|\mathbf{Y} - \Pi_{m'} \mathbf{Y}\|^2 &= \|\tilde{\mathbf{f}}_{m'} - \tilde{\mathbf{f}}_m\|^2 = \|\Pi_{m',m} \mathbf{f} + \Pi_{m',m} \boldsymbol{\varepsilon}\|^2, \\ \|\tilde{\mathbf{f}}_{m'} - \mathbf{f}\|^2 - \|\tilde{\mathbf{f}}_m - \mathbf{f}\|^2 &= -\|\Pi_{m',m} \mathbf{f}\|^2 + \|\Pi_{m',m} \boldsymbol{\varepsilon}\|^2. \end{aligned}$$

- compute $\mathbb{E}\|\tilde{\mathbf{f}}_{m'} - \tilde{\mathbf{f}}_m\|^2$ and $\mathbb{E}[\|\tilde{\mathbf{f}}_{m'} - \mathbf{f}\|^2 - \|\tilde{\mathbf{f}}_m - \mathbf{f}\|^2]$.

2.2.1 AIC and pairwise comparison

Here we try to understand whether the AIC rule does a good job in model selection. In particular, whether it mimics the oracle. Our study will be based on pairwise comparison. More precisely, we check two situations: when the data-driven choice \hat{m} is larger than the oracle and the inverse case. The most important problem is to bound the probability and the risk associated with the event $\{\hat{m} > m^*\}$.

The definition of m^* (2.4) implies for $m > m^*$ with $\mathcal{R}_m \stackrel{\text{def}}{=} \mathcal{R}(\tilde{\mathbf{f}}_m, \mathbf{f}^*)$

$$\begin{aligned} \mathcal{R}_m - \mathcal{R}_{m^*} &= \|\mathbf{f}^* - \Pi_m \mathbf{f}^*\|^2 - \|\mathbf{f}^* - \Pi_{m^*} \mathbf{f}^*\|^2 + \sigma^2(m - m^*) \\ &= -\|\mathbf{b}_{m,m^*}\|^2 + \sigma^2(m - m^*) \geq 0, \end{aligned} \quad (2.7)$$

where $\mathbf{b}_{m,m^*} \stackrel{\text{def}}{=} \Pi_{m,m^*} \mathbf{f}^*$.

Exercise 2.2.2. Consider the model $\mathbf{Y} = \mathbf{f}^* + \boldsymbol{\varepsilon}$ with homogeneous errors $\sigma_i \equiv \sigma$. Let m^* be the oracle choice from (2.4).

- check (2.7);
- check that for $m < m^*$, it holds

$$\|\mathbf{b}_{m^*,m}\|^2 = \|\Pi_{m^*,m}\mathbf{f}\|^2 \geq \sigma^2(m^* - m). \quad (2.8)$$

- check that for $m > m^*$, it holds

$$\|\mathbf{b}_{m,m^*}\|^2 = \|\Pi_{m,m^*}\mathbf{f}\|^2 \leq \sigma^2(m - m^*) \quad (2.9)$$

In words, due to (2.8), it is reasonable to increase the model complexity towards m^* , the gain in the quality of approximation is larger than the additional complexity. However, (2.9) shows that if we increase the complexity of the model over the oracle m^* , then our additional loss due to increased complexity exceeds the gain due to bias reduction.

The next question is whether the data-driven choice \hat{m} reproduces this situation. The selected model \hat{m} is a winner in a pairwise competition with all other models, in particular, in competition with the “oracle” choice m^* . This means that $\tilde{\mathcal{R}}_{\hat{m}} \leq \tilde{\mathcal{R}}_{m^*}$. If the value $\tilde{\mathcal{R}}_m$ is close to its expectation \mathcal{R}_m and if \mathcal{R}_m is significantly larger than the oracle risk \mathcal{R}_{m^*} then the probability of the event $\tilde{\mathcal{R}}_m \leq \tilde{\mathcal{R}}_{m^*}$ is very small. So, one can expect that the selected model \hat{m} is mainly located on the set where the risk \mathcal{R}_m does not deviate much from \mathcal{R}_{m^*} . The next result quantifies this relation. We use the decomposition

$$\begin{aligned} \tilde{\mathcal{R}}_m - \tilde{\mathcal{R}}_{m^*} &= \|\mathbf{Y} - \Pi_m \mathbf{Y}\|^2 + 2\sigma^2 m - \|\mathbf{Y} - \Pi_{m^*} \mathbf{Y}\|^2 - 2\sigma^2 m^* \\ &= -\|\Pi_{m,m^*} \mathbf{Y}\|^2 + 2\sigma^2(m - m^*) \\ &= -\|\Pi_{m,m^*} \boldsymbol{\varepsilon} + \mathbf{b}_{m,m^*}\|^2 + 2\sigma^2(m - m^*) \\ &= \mathcal{R}_m - \mathcal{R}_{m^*} - \{\|\Pi_{m,m^*} \boldsymbol{\varepsilon}\|^2 - \sigma^2(m - m^*)\} - 2\mathbf{b}_{m,m^*}^\top \Pi_{m,m^*} \boldsymbol{\varepsilon}. \end{aligned} \quad (2.10)$$

The first stochastic term $\|\Pi_{m,m^*} \boldsymbol{\varepsilon}\|^2 - \sigma^2(m - m^*)$ of this difference is a centered quadratic form of the errors $\boldsymbol{\varepsilon}$ and the unknown regression function f does not show up there. The second one $2\mathbf{b}_{m,m^*}^\top \Pi_{m,m^*} \boldsymbol{\varepsilon}$ involves the bias \mathbf{b}_{m,m^*} but it is linear in $\boldsymbol{\varepsilon}$. Both terms can be easily bounded for the Gaussian errors $\boldsymbol{\varepsilon}$.

Lemma 2.2.3. *Let $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 I_n)$. Then it holds for $\mathbf{b}_{m,m^*} = \Pi_{m,m^*} \mathbf{f}^*$*

$$\begin{aligned} \mathbb{P}\left(\sigma^{-2}\|\Pi_{m,m^*} \boldsymbol{\varepsilon}\|^2 > \mathfrak{z}^+(m - m^*, \mathbf{x})\right) &\leq \frac{1}{2}e^{-\mathbf{x}}, \\ \mathbb{P}\left(\sigma^{-2}\|\Pi_{m,m^*} \boldsymbol{\varepsilon}\|^2 < \mathfrak{z}^-(m - m^*, \mathbf{x})\right) &\leq \frac{1}{2}e^{-\mathbf{x}}, \\ \mathbb{P}\left(\sigma^{-2}|\mathbf{b}_{m,m^*}^\top \Pi_{m,m^*} \boldsymbol{\varepsilon}| > \sigma^{-1}\|\mathbf{b}_{m,m^*}\|z_1(\mathbf{x})\right) &\leq e^{-\mathbf{x}}. \end{aligned}$$

where $\mathfrak{z}^+(k, \mathbf{x})$ is the upper $1 - 0.5e^{-\mathbf{x}}$ quantile of χ_k^2 , $\mathfrak{z}^-(k, \mathbf{x})$ is its lower $0.5e^{-\mathbf{x}}$ quantile, $z_1(\mathbf{x})$ is the quantile of ξ for a standard normal r.v. $\xi \sim \mathcal{N}(0, 1)$: $\mathbb{P}(|\xi| > z_1(\mathbf{x})) \leq e^{-\mathbf{x}}$.

It holds for any $k \geq 1$ and $\mathbf{x} > 0$ with $\mathbf{x}_1 = \mathbf{x} + \log(2)$

$$\begin{aligned}\mathfrak{z}^+(k, \mathbf{x}) &\leq k + \sqrt{6.6k \mathbf{x}_1} \vee (6.6\mathbf{x}_1), \\ \mathfrak{z}^-(k, \mathbf{x}) &\geq k - \sqrt{2k \mathbf{x}_1}.\end{aligned}\tag{2.11}$$

The proof only uses that $\sigma^{-1}\boldsymbol{\varepsilon}$ is a standard Gaussian vector in \mathbb{R}^n and thus, $\sigma^{-1}\Pi_{m,m^*}\boldsymbol{\varepsilon}$ is standard normal in \mathbb{R}^{m-m^*} , while $(\sigma\|\mathbf{b}_{m,m^*}\|)^{-1}\mathbf{b}_{m,m^*}^\top\Pi_{m,m^*}\boldsymbol{\varepsilon}$ is a standard normal r.v. if the bias \mathbf{b}_{m,m^*} does not vanish.

The presented bounds show that for moderate values of \mathbf{x}

$$z^\pm(k, \mathbf{x}) \stackrel{\text{def}}{=} \mathfrak{z}^+(k, \mathbf{x}) - \mathfrak{z}^-(k, \mathbf{x}) \leq \mathbf{C}\sqrt{k \mathbf{x}}.$$

for a fixed constant \mathbf{C} . Therefore, for large k , the interquartile range $z^\pm(k, \mathbf{x}) = \mathfrak{z}^+(k, \mathbf{x}) - \mathfrak{z}^-(k, \mathbf{x})$ is small relative to k which is the expectation of $\sigma^{-2}\mathbb{E}\|\Pi_k\boldsymbol{\varepsilon}\|^2$. This effect is called *concentration* and it explains why the AIC rule works: the difference between empirical risk and its population counterpart is small relatively to the risk itself.

2.2.2 Pairwise analysis

Now we make a more precise analysis of the term $\tilde{\mathfrak{R}}_m - \tilde{\mathfrak{R}}_{m^*}$ in (2.10). It is based on the following general property of the Gaussian distribution.

Lemma 2.2.4. *Let $\boldsymbol{\xi}$ be standard Gaussian vector in \mathbb{R}^k and $\boldsymbol{\delta}$ be a deterministic vector in \mathbb{R}^k with $\|\boldsymbol{\delta}\|^2 = \Delta$. Then*

- the distribution of $\|\boldsymbol{\xi} + \boldsymbol{\delta}\|^2$ only depends on k and Δ ;
- let, for a given \mathbf{x} , the quantiles $\mathfrak{z}^+(k, \Delta; \mathbf{x})$ and $\mathfrak{z}^-(k, \Delta; \mathbf{x})$ be defined as

$$\mathbb{P}(\|\boldsymbol{\xi} + \boldsymbol{\delta}\|^2 \geq \mathfrak{z}^+(k, \Delta; \mathbf{x})) = e^{-\mathbf{x}},\tag{2.12}$$

$$\mathbb{P}(\|\boldsymbol{\xi} + \boldsymbol{\delta}\|^2 \leq \mathfrak{z}^-(k, \Delta; \mathbf{x})) = e^{-\mathbf{x}}.$$

Then

$$\begin{aligned}\mathfrak{z}^+(k, \Delta; \mathbf{x}) &\leq \Delta + \mathfrak{z}^+(k, \mathbf{x}) + 2\Delta^{1/2}z_1(\mathbf{x}), \\ \mathfrak{z}^-(k, \Delta; \mathbf{x}) &\geq \Delta + \mathfrak{z}^-(k, \mathbf{x}) - 2\Delta^{1/2}z_1(\mathbf{x}),\end{aligned}\tag{2.13}$$

Proof. Use the decomposition

$$\|\boldsymbol{\xi} + \boldsymbol{\delta}\|^2 = \Delta + \|\boldsymbol{\xi}\|^2 + 2\boldsymbol{\xi}^\top\boldsymbol{\delta}$$

and Lemma 2.2.3.

We apply this result to $\pm\sigma^{-2}\|II_{m,m^*}\mathbf{Y}\|^2$ entering in the difference $\tilde{\mathcal{R}}_m - \tilde{\mathcal{R}}_{m^*}$. The bound $\tilde{\mathcal{R}}_m - \tilde{\mathcal{R}}_{m^*}$ can be rewritten for $m > m^*$ as $\sigma^{-2}\|II_{m,m^*}\mathbf{Y}\|^2 > 2(m - m^*)$ which is directly related to the upper quantile of non-central chi-squared. Define the value of non-centrality parameter Δ to have $\mathfrak{z}^\pm(k, \Delta^\pm; \mathbf{x})$ exactly equal to $2k$:

$$\mathfrak{z}^+(k, \Delta^+(k, \mathbf{x}); \mathbf{x}) = 2k, \quad \mathfrak{z}^-(k, \Delta^-(k, \mathbf{x}); \mathbf{x}) = 2k. \quad (2.14)$$

This definition can be rewritten as follows:

$$\mathcal{P}(\|\boldsymbol{\xi} + \boldsymbol{\delta}^+\|^2 > 2k) \leq e^{-\mathbf{x}}, \quad \text{if } \|\boldsymbol{\delta}^+\|^2 \leq \Delta^+(k, \mathbf{x}), \quad (2.15)$$

$$\mathcal{P}(\|\boldsymbol{\xi} + \boldsymbol{\delta}^-\|^2 < 2k) \leq e^{-\mathbf{x}}, \quad \text{if } \|\boldsymbol{\delta}^-\|^2 \geq \Delta^-(k, \mathbf{x}). \quad (2.16)$$

Exercise 2.2.3. For the quantities $\Delta^+(k, \mathbf{x}), \Delta^-(k, \mathbf{x})$ from (2.14) and $\mathbf{x} \geq 1$

- show that $\Delta^+(k, \mathbf{x}) < k$, $\Delta^-(k, \mathbf{x}) > k$;
- check that Lemma 2.2.4 implies

$$\Delta^+(k, \mathbf{x}) \geq \mathfrak{z}^-(k, \mathbf{x}) - 2k^{1/2}z_1(\mathbf{x}), \quad (2.17)$$

$$\Delta^-(k, \mathbf{x}) \leq \mathfrak{z}^+(k, \mathbf{x}) + 2k^{1/2}z_1(\mathbf{x}), \quad (2.18)$$

We conclude with the following statement.

Proposition 2.2.1. *Let the errors $\boldsymbol{\varepsilon}$ be normal and homogeneous: $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$. Then the inequalities*

$$\sigma^{-2}\|\mathbf{b}_{m,m^*}\|^2 \leq \Delta^+(m - m^*, \mathbf{x}), \quad m > m^* \quad (2.19)$$

$$\sigma^{-2}\|\mathbf{b}_{m^*,m}\|^2 \geq \Delta^-(m^* - m, \mathbf{x}), \quad m < m^*$$

ensures

$$\mathcal{P}(\tilde{\mathcal{R}}_m \leq \tilde{\mathcal{R}}_{m^*}) \leq e^{-\mathbf{x}}.$$

In particular, if the bias term $\|\mathbf{b}_m\|$ is uniformly bounded for all $m \geq m^*$ by a fixed constant $\mathbf{C}(\mathcal{F})$, then $\|\mathbf{b}_{m,m^*}\| \leq \mathbf{C}(\mathcal{F})$ and the inequality (2.19) is fulfilled if

$$m - m^* > (2\sigma^{-1}\mathbf{C}(\mathcal{F}) + \mathbf{C}\mathbf{x})^2 \quad (2.20)$$

for $\mathbf{C} \geq 3$.

Proof. Consider the case $m > m^*$. We apply the decomposition

$$\sigma^{-2}(\tilde{\mathcal{R}}_m - \tilde{\mathcal{R}}_{m^*}) = -\|\sigma^{-1}\Pi_{m,m^*}\boldsymbol{\varepsilon} + \sigma^{-1}\mathbf{b}_{m,m^*}\|^2 + 2(m - m^*).$$

Further, $\boldsymbol{\xi} \stackrel{\text{def}}{=} \sigma^{-1}\Pi_{m,m^*}\boldsymbol{\varepsilon}$ is standard normal in \mathbb{R}^k for $k = m - m^*$. The condition (2.15) with $\boldsymbol{\delta}^+ = \sigma^{-1}\mathbf{b}_{m,m^*}$ implies

$$\mathbb{P}(\tilde{\mathcal{R}}_m \leq \tilde{\mathcal{R}}_{m^*}) = \mathbb{P}(\|\boldsymbol{\xi} + \boldsymbol{\delta}^+\|^2 > 2k) \leq e^{-x}.$$

The case $m < m^*$ can be done in a similar way using (2.16) in place of (2.15).

The inequality $\|\mathbf{b}_m\| \leq \mathfrak{C}(\mathcal{F})$ implies $\|\mathbf{b}_{m,m^*}\| \leq \mathfrak{C}(\mathcal{F})$ for all $m > m^*$; see (2.22). Now it remains to check by (2.17) and (2.11) that (2.20) implies (2.19).

To be done: complete the proof

Exercise 2.2.4. Show that the value $\|\mathbf{b}_m\| = \|\mathbf{f}^* - \Pi_m \mathbf{f}^*\|$ monotonously decreases with m . Moreover, for any $m' > m$, the relative bias $\mathbf{b}_{m',m} = \Pi_{m',m} \mathbf{f}^*$ satisfies

$$\|\mathbf{b}_{m',m}\| \leq \|\mathbf{b}_m\|, \quad m' > m. \quad (2.21)$$

Check whether also holds

$$\|\mathbf{b}_{m',m}\| \leq \|\mathbf{b}_{m'}\|, \quad m' > m. \quad (2.22)$$

2.2.3 Uniform bounds and the zone of insensitivity

This section introduces the *set of insensitivity* $\mathcal{M}^\circ(\mathbf{x})$ which describes the quality of model selection. Namely, we aim at describing the set $\mathcal{M}^\circ(\mathbf{x})$ which contains all possible values of \hat{m} with a high probability. The ideal situation would be $\mathcal{M}^\circ(\mathbf{x}) = \{m^*\}$, but it is rare the case. Usually $\mathcal{M}^\circ(\mathbf{x})$ is a larger set containing m^* . Below we specify this set in terms of the difference $\mathcal{R}_m - \mathcal{R}_{m^*}$ and its decomposition (2.7).

The study of the previous section quantifies the pairwise relation $\tilde{\mathcal{R}}_m - \tilde{\mathcal{R}}_{m^*} \leq 0$: under $\sigma^{-2}\|\mathbf{b}_{m,m^*}\|^2 \leq \Delta^+(m - m^*, \mathbf{x})$; see (2.19), it holds

$$\mathbb{P}(\tilde{\mathcal{R}}_m \leq \tilde{\mathcal{R}}_{m^*}) \leq e^{-x}. \quad (2.23)$$

Now we need its uniform version over the complement of $\mathcal{M}^\circ(\mathbf{x})$:

$$\mathbb{P}\left(\max_{m \notin \mathcal{M}^\circ(\mathbf{x})} \{\tilde{\mathcal{R}}_m - \tilde{\mathcal{R}}_{m^*}\} \geq 0\right) \leq e^{-x}.$$

One can use a uniform adjustment in each bound (2.23) by increasing the value \mathbf{x} to another slightly larger level \mathbf{x}_s . A simple way is based on the so called Bonferroni correction: $\mathbf{x}_s \equiv \mathbf{x} + \log(|\mathcal{M}|)$.

Proposition 2.2.2. Let $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$. If $\mathcal{M}^\circ(\mathbf{x})$ is the set of indices m such that

$$\mathcal{M}^\circ(\mathbf{x}) \stackrel{\text{def}}{=} \begin{cases} \sigma^{-2} \|\mathbf{b}_{m,m^*}\|^2 \geq \Delta^+(m - m^*, \mathbf{x}_s), & m > m^*, \\ \sigma^{-2} \|\mathbf{b}_{m^*,m}\|^2 \leq \Delta^-(m^* - m, \mathbf{x}_s), & m < m^*, \end{cases} \quad (2.24)$$

for $\mathbf{x}_s = \mathbf{x} + \log(|\mathcal{M}|)$, then

$$\mathbb{P}(\hat{m} \notin \mathcal{M}^\circ(\mathbf{x})) \leq e^{-\mathbf{x}}. \quad (2.25)$$

Proof. By definition of $\mathcal{M}^\circ(\mathbf{x})$ and Proposition 2.2.1

$$\mathbb{P}(\hat{m} \notin \mathcal{M}^\circ(\mathbf{x})) \leq \sum_{m \notin \mathcal{M}^\circ(\mathbf{x})} \mathbb{P}(\hat{m} = m) \leq \sum_{m \notin \mathcal{M}^\circ(\mathbf{x})} \mathbb{P}(\tilde{\mathcal{R}}_m \leq \tilde{\mathcal{R}}_{m^*}) \leq \sum_{m \in \mathcal{M}} e^{-\mathbf{x}_s} \leq e^{-\mathbf{x}}.$$

One can conclude that if the m lies beyond the *insensitivity zone* $\mathcal{M}^\circ(\mathbf{x})$ around m^* , on which the difference $\mathcal{R}_m - \mathcal{R}_{m^*}$ is not sufficiently large, then the event $\hat{m} = m$ is very unlikely. The result (2.25) can be stated in the form that there exists a random set $\Omega(\mathbf{x})$ such that $\mathbb{P}(\Omega(\mathbf{x})) \geq 1 - e^{-\mathbf{x}}$, and on this set, it holds

$$\tilde{\mathcal{R}}_m > \tilde{\mathcal{R}}_{m^*}, \quad m \notin \mathcal{M}^\circ(\mathbf{x})$$

and hence $\hat{m} \in \mathcal{M}^\circ(\mathbf{x})$ on $\Omega(\mathbf{x})$.

2.2.4 A bound on the excess

Introduce another random set $\Omega_0(\mathbf{x})$ such that

$$\begin{aligned} \|\sigma^{-1} \Pi_{m,m^*} \varepsilon\|^2 &\leq \mathfrak{J}^+(m - m^*, \mathbf{x}_s), & m > m^*, \\ \|\sigma^{-1} \Pi_{m^*,m} \varepsilon\|^2 &\geq \mathfrak{J}^-(m^* - m, \mathbf{x}_s), & m < m^*. \end{aligned} \quad (2.26)$$

Lemma 2.2.3 implies that $\mathbb{P}(\Omega_0(\mathbf{x})) \geq 1 - e^{-\mathbf{x}}$.

The next question is what happens if $\hat{m} \in \mathcal{M}^\circ(\mathbf{x})$ and how big this set is. We will try to bound the loss difference $\varrho(\tilde{\mathbf{f}}_m, \mathbf{f}^*) - \varrho(\tilde{\mathbf{f}}_{m^*}, \mathbf{f}^*)$. It follows from (2.2) that for $m > m^*$

$$\sigma^{-2} \{ \varrho(\tilde{\mathbf{f}}_m, \mathbf{f}^*) - \varrho(\tilde{\mathbf{f}}_{m^*}, \mathbf{f}^*) \} = -\|\sigma^{-1} \mathbf{b}_{m,m^*}\|^2 + \|\sigma^{-1} \Pi_{m,m^*} \varepsilon\|^2.$$

This implies for $m \in \mathcal{M}_+(\mathbf{x})$ by (2.24) and (2.17) on $\Omega_0(\mathbf{x})$

$$\begin{aligned} &\sigma^{-2} \{ \varrho(\tilde{\mathbf{f}}_m, \mathbf{f}^*) - \varrho(\tilde{\mathbf{f}}_{m^*}, \mathbf{f}^*) \} \\ &\leq -\Delta^+(m - m^*, \mathbf{x}_s) + \mathfrak{J}^+(m - m^*, \mathbf{x}_s) \\ &\leq \mathfrak{J}^+(m - m^*, \mathbf{x}_s) - \mathfrak{J}^-(m - m^*, \mathbf{x}_s) + 2z_1(\mathbf{x}_s) \sqrt{m - m^*} \\ &= z^\pm(m - m^*, \mathbf{x}_s) + 2z_1(\mathbf{x}_s) \sqrt{m - m^*}; \end{aligned} \quad (2.27)$$

here $z^\pm(k, \mathbf{x}) = \mathfrak{z}^+(k, \mathbf{x}) - \mathfrak{z}^-(k, \mathbf{x})$.

Now we consider the similar difference for the parameter m from the insensitivity zone $\mathcal{M}_-(\mathbf{x})$ with $m < m^*$. It holds

$$\sigma^{-2}\{\varrho(\tilde{\mathbf{f}}_m, \mathbf{f}^*) - \varrho(\tilde{\mathbf{f}}_{m^*}, \mathbf{f}^*)\} = \|\sigma^{-1}\mathbf{b}_{m^*,m}\|^2 - \|\sigma^{-1}II_{m^*,m}\boldsymbol{\varepsilon}\|^2.$$

This implies for $m \in \mathcal{M}^\circ(\mathbf{x})$ by (2.24) and (2.18) on $\Omega_0(\mathbf{x})$

$$\begin{aligned} & \sigma^{-2}\{\varrho(\tilde{\mathbf{f}}_m, \mathbf{f}^*) - \varrho(\tilde{\mathbf{f}}_{m^*}, \mathbf{f}^*)\} \\ & \leq \|\sigma^{-1}\mathbf{b}_{m^*,m}\|^2 - \|\sigma^{-1}II_{m^*,m}\boldsymbol{\varepsilon}\|^2 \\ & \leq \Delta^-(m^* - m, \mathbf{x}_s) - \mathfrak{z}^-(m^* - m, \mathbf{x}_s) \\ & \leq \mathfrak{z}^+(m^* - m, \mathbf{x}_s) - \mathfrak{z}^-(m^* - m, \mathbf{x}_s) + 2z_1(\mathbf{x}_s)\sqrt{m^* - m} \\ & = z^\pm(m^* - m, \mathbf{x}_s) + 2z_1(\mathbf{x}_s)\sqrt{m^* - m}. \end{aligned} \tag{2.28}$$

Now we can summarize. Define the radius $R = R(\mathcal{M}^\circ(\mathbf{x}))$ of the set $\mathcal{M}^\circ(\mathbf{x})$:

$$R \stackrel{\text{def}}{=} \max_{m \in \mathcal{M}^\circ(\mathbf{x})} |m - m^*|.$$

Theorem 2.2.1. *Let \hat{m} be defined by (2.6) and m^* be the oracle choice from (2.4). Suppose that $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 I_n)$. Let $\mathbf{x}_s = \mathbf{x} + \log(|\mathcal{M}|)$. For the insensitivity zone $\mathcal{M}^\circ(\mathbf{x})$ from (2.24), it holds $\hat{m} \in \mathcal{M}^\circ(\mathbf{x})$ on a random set $\Omega(\mathbf{x})$ with $\mathbb{P}(\Omega(\mathbf{x})) \geq 1 - e^{-x}$. Moreover, on a random set $\Omega_0(\mathbf{x})$ with $\mathbb{P}(\Omega_0(\mathbf{x})) \geq 1 - 2e^{-x}$*

$$\sigma^{-2}\{\varrho(\tilde{\mathbf{f}}_{\hat{m}}, \mathbf{f}^*) - \varrho(\tilde{\mathbf{f}}_{m^*}, \mathbf{f}^*)\} \leq z^\pm(R, \mathbf{x}_s) + 2z_1(\mathbf{x}_s)\sqrt{R}.$$

Proof. The result follows from (2.27) and (2.28) by monotonicity of the function $z^\pm(k, \mathbf{x})$ in k and \mathbf{x} .

In view of $z^\pm(R, \mathbf{x}_s) \asymp \sqrt{R \mathbf{x}_s}$, we conclude that the data-driven choice of the parameter m leads to additional loss of order $\sigma^2 \sqrt{R \mathbf{x}_s}$. One can say that the model selection based on unbiased risk estimation works well if the size of the zone of insensitivity $R = R(\mathcal{M}^\circ(\mathbf{x}))$ is not too large compared with the loss and risk of the oracle $\tilde{\mathbf{f}}_{m^*}$.

Note that the loss $\varrho(\tilde{\mathbf{f}}_{m^*}, \mathbf{f}^*)$ fulfills

$$\varrho(\tilde{\mathbf{f}}_{m^*}, \mathbf{f}^*) = \|\mathbf{b}_{m^*}\|^2 + \|II_{m^*}\boldsymbol{\varepsilon}\|^2 \geq \|II_{m^*}\boldsymbol{\varepsilon}\|^2.$$

By Lemma 2.2.3, it is of order m^* .

Exercise 2.2.5. Let $\Omega_0(\mathbf{x})$ be a random set on which (2.26) holds. Show that $\mathbb{P}(\Omega_0(\mathbf{x})) \geq 1 - e^{-x}$ and for every m , the loss $\varrho(\tilde{\mathbf{f}}_m, \mathbf{f}^*)$ satisfies on $\Omega_0(\mathbf{x})$

$$\sigma^{-2} \varrho(\tilde{\mathbf{f}}_m, \mathbf{f}^*) \geq \delta^-(m, \mathbf{x}).$$

For the case when the value $R = \max_{m \in \mathcal{M}^\circ(\mathbf{x})} |m - m^*|$ is small relative to m^* , the loss of the estimate $\hat{\mathbf{f}} = \tilde{\mathbf{f}}_{\hat{m}}$ corresponding to the data-driven selector \hat{m} is not significantly larger than the loss of the oracle estimate $\tilde{\mathbf{f}}_{m^*}$. Unfortunately, the set $\mathcal{M}^\circ(\mathbf{x})$ can be very large if the risk \mathcal{R}_m is a flat function of m . The extreme case is given by the so called “noise reproducing” model. This model is described by the equations $|\theta_j^*|^2 \equiv \sigma^2$; see (2.3). One can easily check that

$$\begin{aligned} \|\mathbf{b}_{m,m^*}\|^2 &\equiv \sigma^2(m - m^*), & m > m^*, \\ \|\mathbf{b}_{m^*,m}\|^2 &\equiv \sigma^2(m^* - m), & m < m^*. \end{aligned}$$

This implies that the risk function \mathcal{R}_m is constant in m , and therefore, the set $\mathcal{M}^\circ(\mathbf{x})$ coincides with the whole set \mathcal{M} .

Exercise 2.2.6. Build an example in which the radius R is twice as large as the oracle risk \mathcal{R}_{m^*} .

2.3 The approach based on multiple testing. “Smallest accepted” rule

This section discusses the alternative approach to model selection based on the idea of multiple testing. Let m^* be a good choice in the sense of “bias-variance trade-off”. Now we aim to develop a procedure that would treat m^* as a good choice with a high probability. We interpret the model selection procedure as pairwise comparison: the “oracle” model wins in terms of the risk against all other models:

$$\mathcal{R}_{m^*} \leq \mathcal{R}_m, \quad m \neq m^*.$$

The selector \hat{m} suggests to apply the model which wins in term of the unbiased risk estimate $\tilde{\mathcal{R}}_m$:

$$\tilde{\mathcal{R}}_{\hat{m}} \leq \tilde{\mathcal{R}}_m, \quad m \neq \hat{m}.$$

Now we reconsider this approach in terms of hypothesis testing.

2.3.1 A LR test

Our null hypothesis will be that a model-candidate m° is “good” in the sense that it delivers a kind of bias-variance trade-off. This means that there is no reason for considering a larger model: the bias improvement will not be compensated by increase of model complexity. Now we interpret this check as a test of the model-candidate m° against any larger model $m > m^\circ$. The null hypothesis H_{m° means that $\theta_j^* \equiv 0$ for all $j > m^\circ$. The alternative is that there are significant coefficients θ_j^* for $m^\circ < j \leq m$.

Define

$$\Theta_m \stackrel{\text{def}}{=} \{\boldsymbol{\theta} = (\theta_1, \dots, \theta_m, 0, \dots, 0)^\top \in \mathbb{R}^p\}.$$

For the log-likelihood function $L(\boldsymbol{\theta}) = -(2\sigma^2)^{-1}\|\mathbf{Y} - \Psi^\top \boldsymbol{\theta}\|^2$, the likelihood ratio test statistic reads as

$$\begin{aligned} T_{m,m^\circ} &\stackrel{\text{def}}{=} \sup_{\boldsymbol{\theta} \in \Theta_m} L(\boldsymbol{\theta}) - \sup_{\boldsymbol{\theta} \in \Theta_{m^\circ}} L(\boldsymbol{\theta}) = \frac{1}{2\sigma^2} (\|\mathbf{Y} - \Psi^\top \tilde{\boldsymbol{\theta}}_{m^\circ}\|^2 - \|\mathbf{Y} - \Psi^\top \tilde{\boldsymbol{\theta}}_m\|^2) \\ &= \frac{1}{2\sigma^2} (\|\mathbf{Y} - \Pi_{m^\circ} \mathbf{Y}\|^2 - \|\mathbf{Y} - \Pi_m \mathbf{Y}\|^2) = \frac{1}{2\sigma^2} \|\Pi_{m,m^\circ} \mathbf{Y}\|^2. \end{aligned}$$

We aim to calibrate this test in a way that the null hypothesis of no significant bias \mathbf{b}_{m,m° between m° and m is not rejected if the bias is indeed insignificant. Significance can be measured by the energy $\sigma^{-2}\|\mathbf{b}_{m,m^\circ}\|^2$ of this bias. Namely, we say that the bias \mathbf{b}_{m,m° is not significant if

$$\sigma^{-2}\|\mathbf{b}_{m,m^\circ}\|^2 \leq \beta(m - m^\circ) \tag{2.29}$$

for a fixed value β . Remind that the definition of the oracle m^* yields the inequality $\|\mathbf{b}_{m,m^*}\|^2 \leq \sigma^2(m - m^*)$ corresponding to $\beta = 1$. We apply Lemma 2.2.4 to choose a critical value for the test statistic T_{m,m° . Define

$$\mathbf{z}_\beta(k, \mathbf{x}) \stackrel{\text{def}}{=} \mathfrak{z}^+(k, \beta k; \mathbf{x}),$$

where $\mathfrak{z}^+(k, \Delta; \mathbf{x})$ is the quantile of a non-central chi-squared from (2.12). Lemma 2.2.4 also implies

$$\mathbf{z}_\beta(k, \mathbf{x}) \leq \beta k + \mathfrak{z}^+(k, \mathbf{x}) + 2\sqrt{\beta k} z_1(\mathbf{x}). \tag{2.30}$$

These definitions imply the following statement.

Proposition 2.3.1. *Let $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 I_n)$. Then $\|\mathbf{b}_{m,m^\circ}\|^2 \leq \beta(m - m^\circ)$ implies*

$$\mathbb{P}(2T_{m,m^\circ} > \mathbf{z}_\beta(m - m^\circ, \mathbf{x})) \leq e^{-\mathbf{x}}.$$

Exercise 2.3.1. Prove Proposition 2.3.1.

2.3.2 Multiplicity correction

The hypothesis H_{m° is not rejected if each test T_{m,m° for $m > m^\circ$ does not reject the null. This requires a multiple testing procedure and a correction for multiplicity. A simple way is the Bonferroni correction: one increases each \mathbf{x} by the same value $q_{m^\circ} = \log(|\mathcal{M}(m^\circ)|) = \log(p - m^\circ)$. Here $\mathcal{M}(m^\circ) = \{m \in \mathcal{M} : m > m^\circ\}$. Then

$$\begin{aligned} & \mathbb{P}\left(\bigcup_{m \in \mathcal{M}(m^\circ)} \{2T_{m,m^\circ} > \mathbf{z}_\beta(m - m^\circ, \mathbf{x} + q_{m^\circ})\}\right) \\ & \leq \sum_{m \in \mathcal{M}(m^\circ)} e^{-\mathbf{x} - q_{m^\circ}} = |\mathcal{M}(m^\circ)| \exp\{-\mathbf{x} - \log(|\mathcal{M}(m^\circ)|)\} = e^{-\mathbf{x}}. \end{aligned}$$

However, the Bonferroni correction is known to be rather conservative especially if the test statistics T_{m,m° are correlated for different m . This is exactly the case under consideration. Another more careful way to choose the correction q_{m° is based on calibration for one special model.

First we consider the special case $\beta = 0$. Then the null hypothesis means that $\theta_j^* \equiv 0$ for all $j > m^\circ$. In this situation the relative bias $\mathbf{b}_{m,m^\circ} = \Pi_{m,m^\circ} \mathbf{f}^*$ is equal to zero, and does not depend on the first m° coefficients θ_j^* for $j \leq m^\circ$. This allows to define q_{m° by the condition

$$\mathbb{P}_0\left(\bigcup_{m \in \mathcal{M}(m^\circ)} \{2T_{m,m^\circ} > \mathfrak{z}^+(m - m^\circ, \mathbf{x} + q_{m^\circ})\}\right) = e^{-\mathbf{x}}, \quad (2.31)$$

where \mathbb{P}_0 is the measure corresponding to zero signal $\theta_j^* \equiv 0$, and $\mathfrak{z}^+(k, \mathbf{x})$ is the upper quantile of the χ_k^2 from Lemma 2.2.3. The meaning of (2.32) is that each particular test based on the test statistic $2T_{m,m^\circ}$ is performed at a higher level $e^{-\mathbf{x} - q_{m^\circ}} = A^{-1}e^{-\mathbf{x}}$ with $A = e^{q_{m^\circ}}$. The value q_{m° is the smallest number providing the familywise error probability $e^{-\mathbf{x}}$. The Bonferroni correction uses $A = \#\{\text{set of hypotheses}\}$ but this choice is conservative especially if the test statistics are strongly dependent.

In the case of β positive, consider another special “frontier” model $\mathbf{f}^* = \Psi^\top \boldsymbol{\theta}^*$ with the vector of parameters $\boldsymbol{\theta}^*$ with the equalities $\sigma^{-2} \|\mathbf{b}_{m,m^\circ}\|^2 = \beta(m - m^\circ)$ in place of inequalities in (2.29).

Exercise 2.3.2. For the case of an orthonormal design Ψ , find a vector $\boldsymbol{\theta}_\beta^*$ for which $\sigma^{-2} \|\mathbf{b}_{m,m^\circ}\|^2 \equiv \beta(m - m^\circ)$.

Now we calibrate the critical values \mathbf{z}_β for this special extreme model. Namely, we fix the smallest value $q_{m^\circ} = q_{m^\circ}(\beta)$ for which holds

$$\mathbb{P}_{\boldsymbol{\theta}_\beta^*}\left(\bigcup_{m \in \mathcal{M}_+(m^\circ)} \{2T_{m,m^\circ} > \mathbf{z}_\beta(m - m^\circ, \mathbf{x} + q_{m^\circ})\}\right) = e^{-\mathbf{x}}. \quad (2.32)$$

Suppose that such correction $q_{m^\circ} = q_{m^\circ}(\beta)$ is defined for each $m^\circ \in \mathcal{M}$.

Exercise 2.3.3. Consider the frontier model of Exercise 2.3.2. Let $q_{m^\circ}(\beta)$ be the corresponding value for the condition (2.32). Check that for each m° ,

$$q_{m^\circ}(\beta) \geq q_{m^\circ+1}(\beta)$$

yielding

$$\mathbf{z}_\beta(k, \mathbf{x} + q_{m^\circ}) \geq \mathbf{z}_\beta(k, \mathbf{x} + q_{m^\circ-1}).$$

Denote

$$\mathbf{x}_{m^\circ} = \mathbf{x} + q_{m^\circ}.$$

The suggested procedure selects the smallest m° which is accepted by the multiple test H_{m° against H_m for all $m > m^\circ$. This acceptance rule for the null hypothesis H_{m° reads as follows:

$$2T_{m,m^\circ} \leq \mathbf{z}_\beta(m - m^\circ, \mathbf{x}_{m^\circ}), \quad \forall m > m^\circ.$$

Now the “smallest accepted” procedure can be stated in the following form:

$$\begin{aligned} \hat{m} &\stackrel{\text{def}}{=} \min\{m^\circ \in \mathcal{M}: m^\circ \text{ is accepted}\} \\ &= \min\left\{m^\circ \in \mathcal{M}: \max_{m \in \mathcal{M}(m^\circ)} \{2T_{m,m^\circ} - \mathbf{z}_\beta(m - m^\circ, \mathbf{x}_{m^\circ})\} \leq 0\right\}. \end{aligned} \quad (2.33)$$

2.3.3 Definition of the oracle and propagation property

Further we aim at establishing the oracle inequality for this method. It is desirable to show that the data driven selector \hat{m} behaves as good as the oracle one. As the procedure does not rely on the quadratic risk we slightly extend the oracle definition. Let m^* be the oracle value which is now defined as the smallest value m^* satisfying the conditions $\sigma^{-2} \|\mathbf{b}_{m,m^*}\|^2 \leq \beta(m - m^*)$ for $m > m^*$:

$$m^* \stackrel{\text{def}}{=} \operatorname{argmin}\{m: \sigma^{-2} \|\mathbf{b}_{m,m^*}\|^2 \leq \beta(m - m^*), m > m^*\}. \quad (2.34)$$

This definition follows the structure of the null hypothesis considered in the procedure. For $\beta = 1$ it is consistent with the oracle definition based on the risk minimization. Now we aim at establishing the oracle property: the data-driven selector \hat{m} behaves essentially as good as the oracle choice m^* . The study is done in two steps. The first step is to check that the model m^* will be accepted with a high probability. This would mean that the selected model \hat{m} satisfies $\hat{m} \leq m^*$. The second step is in applying the test $T_{m^*,\hat{m}}$ for this situation.

Theorem 2.3.1. *Let m^* be defined by (2.34). Let also the selector \hat{m} be calibrated by (2.32). Then*

$$\mathbb{P}(m^* \text{ is not accepted}) \leq e^{-x}. \tag{2.35}$$

Proof. We use the following fact.

Lemma 2.3.1. *Let $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$. Let m° and \mathbf{x} be fixed.*

- *The values $\theta_1^*, \dots, \theta_{m^\circ}^*$ do not enter in the distribution of T_{m, m° .*
- *Consider the set $F_\beta(m^\circ)$ of all vectors \mathbf{f}^* with $\sigma^{-2} \|\mathbf{b}_{m, m^\circ}\|^2 \leq \beta(m - m^\circ)$. Within this set, the probability*

$$\mathbb{P}\left(\max_{m > m^\circ} \{2T_{m, m^\circ} - \mathbf{z}_\beta(m - m^\circ, \mathbf{x})\} > 0\right)$$

is maximized for the case $\sigma^{-2} \|\mathbf{b}_{m, m^\circ}\|^2 \equiv \beta(m - m^\circ)$.

To be done: Proof of Lemma 2.3.1

Now we check (2.35). The definition of m^* ensures that $\mathbf{f}^* \in F_\beta(m^*)$, that is, (2.34) is fulfilled. The last statement of Lemma 2.3.1 ensures that if the value q_{m^*} is fixed for the extreme model with $\sigma^{-2} \|\mathbf{b}_{m, m^*}\|^2 \equiv \beta(m - m^*)$, then

$$\mathbb{P}(m^* \text{ is not accepted}) = \mathbb{P}\left(\max_{m > m^*} \{2T_{m, m^*} - \mathbf{z}_\beta(m - m^*, \mathbf{x}_{m^*})\} > 0\right) \leq e^{-x}$$

because the similar inequality (with $m^\circ = m^*$) is fulfilled for the extreme model.

The “propagation” property (2.35) is very important and it is usually not fulfilled for classical procedures like unbiased risk estimation (SURE). The advantage of the proposed approach is that this property is in fact intrinsic for the method and is postulated by the calibration step.

2.3.4 A bound on the loss

It remains to clarify the situation if the selected model is smaller than m^* . This probability is not small but we can control the difference of the losses similarly to the SURE procedure. First we check that the new procedure is also able to identify a significant bias $\mathbf{b}_{m^*, m}$ measured by $\sigma^{-2} \|\mathbf{b}_{m^*, m}\|^2$.

Proposition 2.3.2. *Define $\mathbf{x}_s \stackrel{\text{def}}{=} \mathbf{x} + \log(m^*)$ and let*

$$\mathcal{M}^\circ(\mathbf{x}) = \left\{ m \leq m^* : \mathfrak{z}^-(m^* - m, \Delta_m; \mathbf{x}_s) \leq \mathbf{z}_\beta(m^* - m, \mathbf{x}_s) \right\}$$

with $\Delta_m \stackrel{\text{def}}{=} \sigma^{-2} \|\mathbf{b}_{m^, m}\|^2$. Then*

$$\mathbb{P}(\widehat{m} \notin \mathcal{M}^\circ(\mathbf{x})) \leq 2e^{-\mathbf{x}}.$$

Moreover, $m \in \mathcal{M}^\circ(\mathbf{x})$ is impossible if the bias $\Delta_m \stackrel{\text{def}}{=} \sigma^{-2} \|\mathbf{b}_{m^*,m}\|^2$ fulfills with $k = m^* - m$ and with $z^\pm(k, \mathbf{x}_s) = \mathfrak{z}^+(k, \mathbf{x}_s) - \mathfrak{z}^-(k, \mathbf{x}_s)$

$$\{\Delta_m^{1/2} - z_1(\mathbf{x}_s)\}^2 \geq z^\pm(k, \mathbf{x}_s) + \{\sqrt{\beta k} + z_1(\mathbf{x}_s)\}^2 \quad (2.36)$$

or a stronger condition

$$\Delta_m^{1/2} \geq \sqrt{z^\pm(k, \mathbf{x}_s)} + \sqrt{\beta k} + 2z_1(\mathbf{x}_s). \quad (2.37)$$

Proof. We already proved that $\mathbb{P}(\widehat{m} > m^*) \leq e^{-\mathbf{x}}$. It remains to study the case $\widehat{m} < m^*$. Lemma 2.2.4 yields for $m < m^*$

$$2T_{m^*,m} = \sigma^{-2} \|\Pi_{m^*,m} \mathbf{Y}\|^2 \geq \mathfrak{z}^-(m^* - m, \Delta_m; \mathbf{x})$$

with the probability at least $1 - e^{-\mathbf{x}}$. This implies a uniform bound

$$2T_{m^*,m} = \sigma^{-2} \|\Pi_{m^*,m} \mathbf{Y}\|^2 \geq \mathfrak{z}^-(m^* - m, \Delta_m; \mathbf{x}_s), \quad m < m^*$$

on a random set $\Omega_0(\mathbf{x})$ of probability at least $1 - e^{-\mathbf{x}}$. For $m \notin \mathcal{M}^\circ(\mathbf{x})$, this implies

$$2T_{m^*,m} > \mathbf{z}_\beta(m^* - m, \mathbf{x}_m),$$

which makes the event “ m is accepted” impossible on $\Omega_0(\mathbf{x})$ for $m \notin \mathcal{M}^\circ(\mathbf{x})$.

The bound $\mathfrak{z}^-(m^* - m, \Delta_m; \mathbf{x}_s) > \mathbf{z}_\beta(m^* - m, \mathbf{x}_m)$ is implicit in Δ_m . To make it explicit, we use the lower bound (2.13) of Lemma 2.2.4 and (2.30). In view of $\mathbf{x}_s \geq \mathbf{x}_m$ for all $m < m^*$, with $k = m^* - m$, the following inequality is sufficient for checking that $m \notin \mathcal{M}^\circ(\mathbf{x})$:

$$\Delta_m + \mathfrak{z}^-(k, \mathbf{x}_s) - 2\Delta_m^{1/2} z_1(\mathbf{x}_s) \geq \beta k + \mathfrak{z}^+(k, \mathbf{x}_s) + 2\sqrt{\beta k} z_1(\mathbf{x}_s), \quad (2.38)$$

which yields (2.36). The latter can be slightly simplified by using $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$: the bound (2.38) holds if

$$\Delta_m^{1/2} - z_1(\mathbf{x}_s) \geq \sqrt{z^\pm(m^* - m, \mathbf{x}_s)} + \sqrt{\beta(m^* - m)} + z_1(\mathbf{x}_s)$$

which coincides with (2.37).

We conclude that the “smallest accepted” rule leads to the choice of \widehat{m} in the insensitivity zone $\mathcal{M}^\circ(\mathbf{x})$ with a high probability. It remains to bound the loss difference $\varrho(\widetilde{\mathbf{f}}_m, \mathbf{f}^*) - \varrho(\widetilde{\mathbf{f}}_{m^*}, \mathbf{f}^*)$ for $m \in \mathcal{M}^\circ(\mathbf{x})$. We use that for $m < m^*$

$$\sigma^{-2}\{\varrho(\tilde{\mathbf{f}}_m, \mathbf{f}^*) - \varrho(\tilde{\mathbf{f}}_{m^*}, \mathbf{f}^*)\} = \|\sigma^{-1}\mathbf{b}_{m^*,m}\|^2 - \|\sigma^{-1}\mathbf{II}_{m^*,m}\boldsymbol{\varepsilon}\|^2.$$

By Lemma 2.2.3, the stochastic form $\|\sigma^{-1}\mathbf{II}_{m^*,m}\boldsymbol{\varepsilon}\|^2$ can be bounded from below by $\mathfrak{z}^-(m^* - m, \mathbf{x})$ with probability $1 - e^{-\mathbf{x}}$. This implies a uniform probability bound $\|\sigma^{-1}\mathbf{II}_{m^*,m}\boldsymbol{\varepsilon}\|^2 \geq \mathfrak{z}^-(m^* - m, \mathbf{x}_s)$ with $\mathbf{x}_s = \mathbf{x} + \log(m^*)$ for $m < m^*$. For $m \in \mathcal{M}^\circ(\mathbf{x})$, this implies

$$\sigma^{-2}\{\varrho(\tilde{\mathbf{f}}_m, \mathbf{f}^*) - \varrho(\tilde{\mathbf{f}}_{m^*}, \mathbf{f}^*)\} \leq \Delta_m - \mathfrak{z}^-(m^* - m, \mathbf{x}_s)$$

on a random set of probability at least $1 - 3e^{-\mathbf{x}}$, where Δ_m follows the bound (2.36) or (2.37).

We state the following result:

Theorem 2.3.2. *Let $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$. For the selector \hat{m} from (2.33) and the oracle m^* from (2.34), it holds*

$$\mathbb{P}(\hat{m} < m^*) \leq e^{-\mathbf{x}},$$

and for the insensitivity zone $\mathcal{M}^\circ(\mathbf{x})$, it follows on a random set $\Omega(\mathbf{x})$ with $\mathbb{P}(\Omega(\mathbf{x})) \geq 1 - 3e^{-\mathbf{x}}$

$$\begin{aligned} \sigma^{-2}\{\varrho(\tilde{\mathbf{f}}_m, \mathbf{f}^*) - \varrho(\tilde{\mathbf{f}}_{m^*}, \mathbf{f}^*)\} &\leq \max_{m \in \mathcal{M}^\circ(\mathbf{x})} \{\Delta_m - \mathfrak{z}^-(m^* - m, \mathbf{x}_s)\} \\ &\leq \max_{m \in \mathcal{M}^\circ(\mathbf{x})} \left\{ \left(\sqrt{z^\pm(m^* - m, \mathbf{x}_s)} + \sqrt{\beta(m^* - m)} + 2z_1(\mathbf{x}_s) \right)^2 - \mathfrak{z}^-(m^* - m, \mathbf{x}_s) \right\}. \end{aligned}$$

2.3.5 Role of β

The SURE method corresponds to $\beta = 1$. This choice is natural as long the squared risk is concerned. However, the procedure is meaningful for every $\beta > 0$. It even applies for $\beta = 0$. This is a very special case leading to a testing problem instead of estimation.

In general, consider the situation when we apply the SURE procedure with the noise level σ which is misspecified and the true noise level is $\sigma_0^2 \leq \sigma^2$. Then the true risk \mathcal{R}_m is

$$\mathcal{R}_m = \sigma_0^2 m + \|\mathbf{b}_m\|^2 = \beta^{-1} \sigma^2 m + \|\mathbf{b}_m\|^2$$

for $\beta = \sigma^2 / \sigma_0^2$. The use of the wrong noise level is equivalent to applying the procedure with such β .

The choice of a small β results in a strong condition on the bias, so the oracle choice will be shifted towards larger (more complex) models. With $\beta > 1$, one allows for more space for the bias, the procedure becomes more robust and tend to oversmooth the model, that is, to select \hat{m} which is smaller than the risk minimizer m^* .

Linear smoothers

Here we discuss the important situation when the number of predictors ψ_j and hence the number of parameters p in the linear model $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$ is not small relative to the sample size. Then the least square or the maximum likelihood approach meets serious problems. The first one relates to the numerical issues. The definition of the LSE $\tilde{\boldsymbol{\theta}}$ involves the inversion of the $p \times p$ matrix $\Psi\Psi^\top$ and such an inversion becomes a delicate task for p large. The other problem concerns the inference for the estimated parameter $\boldsymbol{\theta}^*$. The risk bound and the width of the confidence set are proportional to the parameter dimension p and thus, with large p , the inference statements become almost uninformative. In particular, if p is of order the sample size n , even consistency is not achievable. One faces a really critical situation. We already know that the MLE is the efficient estimate in the class of all unbiased estimates. At the same time it is highly inefficient in overparametrized models. The only way out of this situation is to sacrifice the unbiasedness property in favor of reducing the model complexity: some procedures can be more efficient than MLE even if they are biased. This section discusses one way of resolving these problems by regularization or shrinkage. To be more specific, for the rest of the section we consider the following setup. The observed vector \mathbf{Y} follows the model

$$\mathbf{Y} = \mathbf{f}^* + \boldsymbol{\varepsilon} \tag{3.1}$$

with a homogeneous error vector $\boldsymbol{\varepsilon}$: $\mathbb{E}\boldsymbol{\varepsilon} = 0$, $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}_n$. Noise misspecification is not considered in this section.

Furthermore, we assume a basis or a collection of basis vectors $\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_p$ is given with p large. This allows for approximating the response vector $\mathbf{f} = \mathbb{E}\mathbf{Y}$ in the form $\mathbf{f} = \Psi^\top \boldsymbol{\theta}^*$, or, equivalently,

$$\mathbf{f} = \theta_1^* \boldsymbol{\psi}_1 + \dots + \theta_p^* \boldsymbol{\psi}_p.$$

In many cases we will assume that the basis is already orthogonalized: $\Psi\Psi^\top = I_p$. The model (3.1) can be rewritten as

$$\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}, \quad \text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 I_n.$$

The MLE or oLSE of the parameter vector $\boldsymbol{\theta}^*$ for this model reads as

$$\tilde{\boldsymbol{\theta}} = (\Psi\Psi^\top)^{-1}\Psi\mathbf{Y}, \quad \tilde{\mathbf{f}} = \Psi^\top \tilde{\boldsymbol{\theta}} = \Psi^\top (\Psi\Psi^\top)^{-1}\Psi\mathbf{Y}.$$

If the matrix $\Psi\Psi^\top$ is degenerate or badly posed, computing the MLE $\tilde{\boldsymbol{\theta}}$ is a hard task. Below we discuss how this problem can be treated.

3.1 Regularization and ridge regression

Let R be a positive symmetric $p \times p$ matrix. Then the sum $\Psi\Psi^\top + R$ is positive symmetric as well and can be inverted whatever the matrix Ψ is. This suggests to replace $(\Psi\Psi^\top)^{-1}$ by $(\Psi\Psi^\top + R)^{-1}$ leading to the regularized least squares estimate $\tilde{\boldsymbol{\theta}}_R$ of the parameter vector $\boldsymbol{\theta}$ and the corresponding response estimate $\tilde{\mathbf{f}}_R$:

$$\tilde{\boldsymbol{\theta}}_R \stackrel{\text{def}}{=} (\Psi\Psi^\top + R)^{-1}\Psi\mathbf{Y}, \quad \tilde{\mathbf{f}}_R \stackrel{\text{def}}{=} \Psi^\top (\Psi\Psi^\top + R)^{-1}\Psi\mathbf{Y}. \quad (3.2)$$

Such a method is also called *ridge regression*. An example of choosing R is the multiple of the unit matrix: $R = \alpha I_p$ where $\alpha > 0$ and I_p stands for the unit matrix. This method is also called *Tikhonov regularization* and it results in the parameter estimate $\tilde{\boldsymbol{\theta}}_\alpha$ and the response estimate $\tilde{\mathbf{f}}_\alpha$:

$$\tilde{\boldsymbol{\theta}}_\alpha \stackrel{\text{def}}{=} (\Psi\Psi^\top + \alpha I_p)^{-1}\Psi\mathbf{Y}, \quad \tilde{\mathbf{f}}_\alpha \stackrel{\text{def}}{=} \Psi^\top (\Psi\Psi^\top + \alpha I_p)^{-1}\Psi\mathbf{Y}. \quad (3.3)$$

A proper choice of the matrix R for the ridge regression method (3.2) or the parameter α for the Tikhonov regularization (3.3) is an important issue. Below we discuss several approaches which lead to the estimate (3.2) with a specific choice of the matrix R . The properties of the estimates $\tilde{\boldsymbol{\theta}}_R$ and $\tilde{\mathbf{f}}_R$ will be studied in context of penalized likelihood estimation in the next section.

3.2 Penalized likelihood. Bias and variance

The estimate (3.2) can be obtained in a natural way within the (quasi) ML approach using the penalized least squares. The classical unpenalized method is based on minimizing the sum of residuals squared:

$$\tilde{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} L(\boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{arginf}} \|\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta}\|^2$$

with $L(\boldsymbol{\theta}) = \sigma^{-2} \|\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta}\|^2 / 2$. (Here we omit the terms which do not depend on $\boldsymbol{\theta}$.) Now we introduce an additional penalty on the objective function which penalizes for the complexity of the candidate vector $\boldsymbol{\theta}$ which is expressed by the value $\|G\boldsymbol{\theta}\|^2 / 2$ for a given symmetric matrix G . This choice of complexity measure implicitly assumes that the vector $\boldsymbol{\theta} \equiv 0$ has the smallest complexity equal to zero and this complexity increases with the norm of $G\boldsymbol{\theta}$. Define the *penalized log-likelihood*

$$\begin{aligned} L_G(\boldsymbol{\theta}) &\stackrel{\text{def}}{=} L(\boldsymbol{\theta}) - \|G\boldsymbol{\theta}\|^2 / 2 \\ &= -(2\sigma^2)^{-1} \|\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta}\|^2 - \|G\boldsymbol{\theta}\|^2 / 2 - (n/2) \log(2\pi\sigma^2). \end{aligned} \quad (3.4)$$

The penalized MLE reads as

$$\tilde{\boldsymbol{\theta}}_G = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} L_G(\boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \left\{ (2\sigma^2)^{-1} \|\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta}\|^2 + \|G\boldsymbol{\theta}\|^2 / 2 \right\}.$$

A straightforward calculus leads to the expression (3.2) for $\tilde{\boldsymbol{\theta}}_G$ with $R = \sigma^2 G^2$:

$$\tilde{\boldsymbol{\theta}}_G \stackrel{\text{def}}{=} (\boldsymbol{\Psi}\boldsymbol{\Psi}^\top + \sigma^2 G^2)^{-1} \boldsymbol{\Psi}\mathbf{Y}. \quad (3.5)$$

We see that $\tilde{\boldsymbol{\theta}}_G$ is again a linear estimate: $\tilde{\boldsymbol{\theta}}_G = \mathcal{S}_G \mathbf{Y}$ with $\mathcal{S}_G = (\boldsymbol{\Psi}\boldsymbol{\Psi}^\top + \sigma^2 G^2)^{-1} \boldsymbol{\Psi}$. The results of Section 1.3 explains that $\tilde{\boldsymbol{\theta}}_G$ in fact estimates the value $\boldsymbol{\theta}_G$ defined by

$$\begin{aligned} \boldsymbol{\theta}_G &= \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \mathbb{E} L_G(\boldsymbol{\theta}) \\ &= \underset{\boldsymbol{\theta}}{\operatorname{arginf}} \mathbb{E} \left\{ \|\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta}\|^2 + \sigma^2 \|G\boldsymbol{\theta}\|^2 \right\} \\ &= (\boldsymbol{\Psi}\boldsymbol{\Psi}^\top + \sigma^2 G^2)^{-1} \boldsymbol{\Psi} \mathbf{f}^* = \mathcal{S}_G \mathbf{f}^*. \end{aligned} \quad (3.6)$$

In particular, if $\mathbf{f}^* = \boldsymbol{\Psi}^\top \boldsymbol{\theta}^*$, then

$$\boldsymbol{\theta}_G = (\boldsymbol{\Psi}\boldsymbol{\Psi}^\top + \sigma^2 G^2)^{-1} \boldsymbol{\Psi}\boldsymbol{\Psi}^\top \boldsymbol{\theta}^* \quad (3.7)$$

and $\boldsymbol{\theta}_G \neq \boldsymbol{\theta}^*$ unless $G = 0$. In other words, the penalized MLE $\tilde{\boldsymbol{\theta}}_G$ is biased.

Exercise 3.2.1. Check that $\mathbb{E} \tilde{\boldsymbol{\theta}}_\alpha = \boldsymbol{\theta}_\alpha$ for $\boldsymbol{\theta}_\alpha = (\boldsymbol{\Psi}\boldsymbol{\Psi}^\top + \alpha I_p)^{-1} \boldsymbol{\Psi}\boldsymbol{\Psi}^\top \boldsymbol{\theta}^*$, the bias $\|\boldsymbol{\theta}_\alpha - \boldsymbol{\theta}^*\|$ grows with the regularization parameter α .

The penalized MLE $\tilde{\boldsymbol{\theta}}_G$ leads to the response estimate $\tilde{\mathbf{f}}_G = \boldsymbol{\Psi}^\top \tilde{\boldsymbol{\theta}}_G$.

Exercise 3.2.2. Check that the penalized ML approach leads to the response estimate

$$\tilde{\mathbf{f}}_G = \boldsymbol{\Psi}^\top \tilde{\boldsymbol{\theta}}_G = \boldsymbol{\Psi}^\top (\boldsymbol{\Psi}\boldsymbol{\Psi}^\top + \sigma^2 G^2)^{-1} \boldsymbol{\Psi}\mathbf{Y} = \Pi_G \mathbf{Y}$$

with $\Pi_G = \Psi^\top (\Psi \Psi^\top + \sigma^2 G^2)^{-1} \Psi$. Show that Π_G is a sub-projector in the sense that $\|\Pi_G \mathbf{u}\| \leq \|\mathbf{u}\|$ for any $\mathbf{u} \in \mathbb{R}^n$.

Exercise 3.2.3. Let Ψ be orthonormal: $\Psi \Psi^\top = I_p$. Then the penalized MLE $\tilde{\boldsymbol{\theta}}_G$ can be represented as

$$\tilde{\boldsymbol{\theta}}_G = (I_p + \sigma^2 G^2)^{-1} \mathbf{Z},$$

where $\mathbf{Z} = \Psi \mathbf{Y}$ is the vector of empirical Fourier coefficients. Specify the result for the case of a diagonal matrix $G = \text{diag}(g_1, \dots, g_p)$ and describe the corresponding response estimate $\tilde{\mathbf{f}}_G$.

The previous results indicate that introducing the penalization leads to some bias of estimation. One can ask about a benefit of using a penalized procedure. The next result shows that penalization decreases the variance of estimation and thus, makes the procedure more stable.

Theorem 3.2.1. Let $\tilde{\boldsymbol{\theta}}_G$ be a penalized MLE from (3.5). Then $\mathbb{E} \tilde{\boldsymbol{\theta}}_G = \boldsymbol{\theta}_G$, see (3.7), and under noise homogeneity $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 I_n$, it holds

$$\begin{aligned} \text{Var}(\tilde{\boldsymbol{\theta}}_G) &= (\sigma^{-2} \Psi \Psi^\top + G^2)^{-1} \sigma^{-2} \Psi \Psi^\top (\sigma^{-2} \Psi \Psi^\top + G^2)^{-1} \\ &= D_G^{-2} D^2 D_G^{-2} \end{aligned}$$

with $D_G^2 = \sigma^{-2} \Psi \Psi^\top + G^2$. In particular, $\text{Var}(\tilde{\boldsymbol{\theta}}_G) \leq D_G^{-2}$. If $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 I_n)$, then $\tilde{\boldsymbol{\theta}}_G$ is also normal: $\tilde{\boldsymbol{\theta}}_G \sim \mathcal{N}(\boldsymbol{\theta}_G, D_G^{-2} D^2 D_G^{-2})$.

Moreover, the bias $\|\boldsymbol{\theta}_G - \boldsymbol{\theta}^*\|$ monotonously increases in G^2 while the variance monotonously decreases with the penalization G .

Proof. The first two moments of $\tilde{\boldsymbol{\theta}}_G$ are computed from $\tilde{\boldsymbol{\theta}}_G = \mathcal{S}_G \mathbf{Y}$. Monotonicity of the bias and variance of $\tilde{\boldsymbol{\theta}}_G$ is proved below in Exercise 3.2.6.

Exercise 3.2.4. Let Ψ be orthonormal: $\Psi \Psi^\top = I_p$. Describe $\text{Var}(\tilde{\boldsymbol{\theta}}_G)$. Show that the variance decreases with the penalization G in the sense that $G_1 \geq G$ implies $\text{Var}(\tilde{\boldsymbol{\theta}}_{G_1}) \leq \text{Var}(\tilde{\boldsymbol{\theta}}_G)$.

Exercise 3.2.5. Let $\Psi \Psi^\top = I_p$ and let $G = \text{diag}(g_1, \dots, g_p)$ be a diagonal matrix. Compute the squared bias $\|\boldsymbol{\theta}_G - \boldsymbol{\theta}^*\|^2$ and show that it monotonously increases in each g_j for $j = 1, \dots, p$.

Exercise 3.2.6. Let G be a symmetric matrix and $\tilde{\boldsymbol{\theta}}_G$ the corresponding penalized MLE. Show that the variance $\text{Var}(\tilde{\boldsymbol{\theta}}_G)$ decreases while the bias $\|\boldsymbol{\theta}_G - \boldsymbol{\theta}^*\|$ increases in G^2 .

Hint: with $D^2 = \sigma^{-2}\Psi\Psi^\top$, show that for any vector $\mathbf{w} \in \mathbb{R}^p$ and $\mathbf{u} = D^{-1}\mathbf{w}$, it holds

$$\mathbf{w}^\top \text{Var}(\tilde{\boldsymbol{\theta}}_G)\mathbf{w} = \mathbf{u}^\top (I_p + D^{-1}G^2D^{-1})^{-2}\mathbf{u}$$

and this value decreases with G^2 because $I_p + D^{-1}G^2D^{-1}$ increases. Show in a similar way that

$$\|\boldsymbol{\theta}_G - \boldsymbol{\theta}^*\|^2 = \|(D^2 + G^2)^{-1}G^2\boldsymbol{\theta}^*\|^2 = \boldsymbol{\theta}^{*\top} \Gamma^{-1} \boldsymbol{\theta}^*$$

with $\Gamma = (I_p + G^{-2}D^2)(I_p + D^2G^{-2})$. Show that the matrix Γ monotonously increases and thus Γ^{-1} monotonously decreases as a function of the symmetric matrix $B = G^{-2}$.

Putting together the results about the bias and the variance of $\tilde{\boldsymbol{\theta}}_G$ yields the statement about the quadratic risk.

Theorem 3.2.2. *Assume the model $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$ with $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 I_n$. Then the estimate $\tilde{\boldsymbol{\theta}}_G$ fulfills*

$$\mathbb{E}\|\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}^*\|^2 = \|\boldsymbol{\theta}_G - \boldsymbol{\theta}^*\|^2 + \text{tr}(D_G^{-2}D^2D_G^{-2}).$$

This result is called the *bias-variance decomposition*. The choice of a proper regularization is usually based on this decomposition: one selects a regularization from a given class to provide the minimal possible risk. This approach is referred to as *bias-variance trade-off*.

3.3 Inference for the penalized MLE

Here we discuss some properties of the penalized MLE $\tilde{\boldsymbol{\theta}}_G$. In particular, we focus on the construction of confidence and concentration sets based on the penalized log-likelihood. We know that the regularized estimate $\tilde{\boldsymbol{\theta}}_G$ is the empirical counterpart of the value $\boldsymbol{\theta}_G$ which solves the regularized deterministic problem (3.6). We also know that the key results are expressed via the value of the supremum $\sup_{\boldsymbol{\theta}} L_G(\boldsymbol{\theta}) - L_G(\boldsymbol{\theta}_G)$. The next result extends Theorem 1.4.3 to the penalized likelihood.

Theorem 3.3.1. *Let $L_G(\boldsymbol{\theta})$ be the penalized log-likelihood from (3.4). Then*

$$2L_G(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G) = (\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G)^\top D_G^2 (\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G) \quad (3.8)$$

$$= \sigma^{-2} \boldsymbol{\varepsilon}^\top \Pi_G \boldsymbol{\varepsilon} \quad (3.9)$$

with $\Pi_G = \Psi^\top (\Psi\Psi^\top + \sigma^2 G^2)^{-1} \Psi$.

In general the matrix Π_G is not a projector and hence, $\sigma^{-2}\boldsymbol{\varepsilon}^\top \Pi_G \boldsymbol{\varepsilon}$ is not χ^2 -distributed, the chi-squared result does not apply.

Exercise 3.3.1. Prove (3.8).

Hint: apply the Taylor expansion to $L_G(\boldsymbol{\theta})$ at $\tilde{\boldsymbol{\theta}}_G$. Use that $\nabla L_G(\boldsymbol{\theta}_G) = 0$ and $-\nabla^2 L_G(\boldsymbol{\theta}) \equiv \sigma^{-2}\Psi\Psi^\top + G^2$.

Exercise 3.3.2. Prove (3.9).

Hint: show that $\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G = \mathcal{S}_G \boldsymbol{\varepsilon}$ with $\mathcal{S}_G = (\Psi\Psi^\top + \sigma^2 G^2)^{-1} \Psi$.

The straightforward corollaries of Theorem 3.3.1 are the concentration and confidence probabilities. Define the confidence set $\mathcal{E}_G(\mathfrak{z})$ for $\boldsymbol{\theta}_G$ as

$$\mathcal{E}_G(\mathfrak{z}) \stackrel{\text{def}}{=} \{\boldsymbol{\theta} : L_G(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}) \leq \mathfrak{z}\}.$$

The definition implies the following result for the coverage probability:

$$\mathbb{P}(\mathcal{E}_G(\mathfrak{z}) \not\ni \boldsymbol{\theta}_G) \leq \mathbb{P}(L_G(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G) > \mathfrak{z}).$$

Now the representation (3.9) for $L_G(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G)$ reduces the problem to a deviation bound for a quadratic form. We apply the general result of Section ??.

Theorem 3.3.2. *Let $L_G(\boldsymbol{\theta})$ be the penalized log-likelihood from (3.4) and let $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$. Then it holds with $\mathfrak{p}_G = \text{tr}(\Pi_G)$ and $v_G^2 = 2 \text{tr}(\Pi_G^2)$ that*

$$\mathbb{P}(2L_G(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G) > \mathfrak{p}_G + (2v_G \mathfrak{x}^{1/2}) \vee (6\mathfrak{x})) \leq \exp(-\mathfrak{x}).$$

Similarly one can state the concentration result. With $D_G^2 = \sigma^{-2}\Psi\Psi^\top + G^2$

$$2L_G(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G) = \|D_G(\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G)\|^2$$

and the result of Theorem 3.3.2 can be restated as the concentration bound:

$$\mathbb{P}(\|D_G(\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G)\|^2 > \mathfrak{p}_G + (2v_G \mathfrak{x}^{1/2}) \vee (6\mathfrak{x})) \leq \exp(-\mathfrak{x}).$$

In other words, $\tilde{\boldsymbol{\theta}}_G$ concentrates on the set $\mathcal{A}(\mathfrak{z}, \boldsymbol{\theta}_G) = \{\boldsymbol{\theta} : \|\boldsymbol{\theta} - \boldsymbol{\theta}_G\|^2 \leq 2\mathfrak{z}\}$ for $2\mathfrak{z} > \mathfrak{p}_G$.

3.4 Projection and shrinkage estimates

Consider a linear model $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$ in which the matrix Ψ is orthonormal in the sense $\Psi\Psi^\top = \mathbf{I}_p$. Then the multiplication with Ψ maps this model in the sequence space model $\mathbf{Z} = \boldsymbol{\theta}^* + \boldsymbol{\xi}$, where $\mathbf{Z} = \Psi\mathbf{Y} = (z_1, \dots, z_p)^\top$ is the vector of empirical Fourier

coefficients $z_j = \boldsymbol{\psi}_j^\top \mathbf{Y}$. The noise $\boldsymbol{\xi} = \boldsymbol{\Psi}\boldsymbol{\varepsilon}$ borrows the feature of the original noise $\boldsymbol{\varepsilon}$: if $\boldsymbol{\varepsilon}$ is zero mean and homogeneous, the same applies to $\boldsymbol{\xi}$. The number of coefficients p can be large or even infinite. To get a sensible estimate, one has to apply some regularization method. The simplest one is called *projection*: one just considers the first m empirical coefficients z_1, \dots, z_m and drop the others. The corresponding parameter estimate $\tilde{\boldsymbol{\theta}}_m$ reads as

$$\tilde{\theta}_{m,j} = \begin{cases} z_j & \text{if } j \leq m, \\ 0 & \text{otherwise.} \end{cases}$$

The response vector $\mathbf{f}^* = \mathbb{E}\mathbf{Y}$ is estimated by $\boldsymbol{\Psi}^\top \tilde{\boldsymbol{\theta}}_m$ leading to the representation

$$\tilde{\mathbf{f}}_m = z_1 \boldsymbol{\psi}_1 + \dots + z_m \boldsymbol{\psi}_m$$

with $z_j = \boldsymbol{\psi}_j^\top \mathbf{Y}$. A disadvantage of the projection method is that it either keeps each empirical coefficient z_m or completely discards it. An extension of the projection method is called *shrinkage*: one multiplies every empirical coefficient z_j with a factor $\alpha_j \in (0, 1)$. This leads to the *shrinkage* estimate $\tilde{\boldsymbol{\theta}}_\alpha$ with

$$\tilde{\theta}_{\alpha,j} = \alpha_j z_j.$$

Here $\boldsymbol{\alpha}$ stands for the vector of coefficients α_j for $j = 1, \dots, p$. A projection method is a special case of this shrinkage with α_j equal to one or zero. Another popular choice of the coefficients α_j is given by

$$\alpha_j = (1 - j/m)^\beta \mathbf{1}(j \leq m) \tag{3.10}$$

for some $\beta > 0$ and $m \leq p$. This choice ensures that the coefficients α_j smoothly approach zero as j approach the value m , and α_j vanish for $j > m$. In this case, the vector $\boldsymbol{\alpha}$ is completely specified by two parameters m and β . The projection method corresponds to $\beta = 0$. The design orthogonality $\boldsymbol{\Psi}\boldsymbol{\Psi}^\top = \mathbf{I}_p$ yields again that the estimation risk $\mathbb{E}\|\tilde{\boldsymbol{\theta}}_\alpha - \boldsymbol{\theta}^*\|^2$ coincides with the prediction risk $\mathbb{E}\|\tilde{\mathbf{f}}_\alpha - \mathbf{f}^*\|^2$.

Exercise 3.4.1. Let $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}_p$. The risk $\mathcal{R}(\tilde{\mathbf{f}}_\alpha)$ of the shrinkage estimate $\tilde{\mathbf{f}}_\alpha$ fulfills

$$\mathcal{R}(\tilde{\mathbf{f}}_\alpha) \stackrel{\text{def}}{=} \mathbb{E}\|\tilde{\mathbf{f}}_\alpha - \mathbf{f}^*\|^2 = \sum_{j=1}^p \theta_j^{*2} (1 - \alpha_j)^2 + \sum_{j=1}^p \alpha_j^2 \sigma^2.$$

Specify the cases of $\boldsymbol{\alpha} = \boldsymbol{\alpha}(m, \beta)$ from (3.10). Evaluate the variance term $\sum_j \alpha_j^2 \sigma^2$. Hint: approximate the sum over j by the integral $\int (1 - x/m)_+^{2\beta} dx$.

The oracle choice is again defined by risk minimization:

$$\boldsymbol{\alpha}^* \stackrel{\text{def}}{=} \underset{\boldsymbol{\alpha}}{\operatorname{argmin}} \mathcal{R}(\tilde{\boldsymbol{f}}_{\boldsymbol{\alpha}}),$$

where minimization is taken over the class of all considered coefficient vectors $\boldsymbol{\alpha}$.

One way of obtaining a shrinkage estimate in the sequence space model $\boldsymbol{Z} = \boldsymbol{\theta}^* + \boldsymbol{\xi}$ is by using a roughness penalization. Let G be a symmetric matrix. Consider the regularized estimate $\tilde{\boldsymbol{\theta}}_G$ from (3.2). The next result claims that if G is a diagonal matrix, then $\tilde{\boldsymbol{\theta}}_G$ is a shrinkage estimate. Moreover, a general penalized MLE can be represented as shrinkage by an orthogonal basis transformation.

Theorem 3.4.1. *Let G be a diagonal matrix, $G = \operatorname{diag}(g_1, \dots, g_p)$. The penalized MLE $\tilde{\boldsymbol{\theta}}_G$ in the sequence space model $\boldsymbol{Z} = \boldsymbol{\theta}^* + \boldsymbol{\xi}$ with $\boldsymbol{\xi} \sim \mathcal{N}(0, \sigma^2 \boldsymbol{I}_p)$ coincides with the shrinkage estimate $\tilde{\boldsymbol{\theta}}_{\boldsymbol{\alpha}}$ for $\alpha_j = (1 + \sigma^2 g_j^2)^{-1} \leq 1$. Moreover, a penalized MLE $\tilde{\boldsymbol{\theta}}_G$ for a general matrix G can be reduced to a shrinkage estimate by a basis transformation in the sequence space model.*

Proof. The first statement for a diagonal matrix G follows from the representation $\tilde{\boldsymbol{\theta}}_G = (\boldsymbol{I}_p + \sigma^2 G^2)^{-1} \boldsymbol{Z}$. Next, let U be an orthogonal transform leading to the diagonal representation $G^2 = U^\top D^2 U$ with $D^2 = \operatorname{diag}(g_1, \dots, g_p)$. Then

$$U \tilde{\boldsymbol{\theta}}_G = (\boldsymbol{I}_p + \sigma^2 D^2)^{-1} U \boldsymbol{Z}$$

that is, $U \tilde{\boldsymbol{\theta}}_G$ is a shrinkage estimate in the transformed model $U \boldsymbol{Z} = U \boldsymbol{\theta}^* + U \boldsymbol{\xi}$.

In other words, roughness penalization results in some kind of shrinkage. Interestingly, the inverse statement holds as well.

Exercise 3.4.2. Let $\tilde{\boldsymbol{\theta}}_{\boldsymbol{\alpha}}$ is a shrinkage estimate for a vector $\boldsymbol{\alpha} = (\alpha_j)$. Then there is a diagonal penalty matrix G such that $\tilde{\boldsymbol{\theta}}_{\boldsymbol{\alpha}} = \tilde{\boldsymbol{\theta}}_G$.

Hint: define the j th diagonal entry g_j by the equation $\alpha_j = (1 + \sigma^2 g_j^2)^{-1}$.

3.5 Smoothness constraints and roughness penalty approach

Another way of reducing the complexity of the estimation procedure is based on smoothness constraints. The notion of smoothness originates from regression estimation. A non-linear regression function f is expanded using a Fourier or some other functional basis and $\boldsymbol{\theta}^*$ is the corresponding vector of coefficients. Smoothness properties of the regression function imply certain rate of decay of the corresponding Fourier coefficients: the larger frequency is, the fewer amount of information about the regression function is

contained in the related coefficient. This leads to the natural idea to replace the original optimization problem over the whole parameter space with the constrained optimization over a subset of “smooth” parameter vectors. Here we consider one popular example of Sobolev smoothness constraints which effectively means that the s th derivative of the function \mathbf{f}^* has a bounded L_2 -norm. A general Sobolev ball can be defined using a diagonal matrix G :

$$\mathcal{B}_G(R) \stackrel{\text{def}}{=} \|\mathbf{G}\boldsymbol{\theta}\| \leq R.$$

Now we consider a constrained ML problem:

$$\tilde{\boldsymbol{\theta}}_{G,R} = \underset{\boldsymbol{\theta} \in \mathcal{B}_G(R)}{\operatorname{argmax}} L(\boldsymbol{\theta}) = \underset{\boldsymbol{\theta} \in \Theta: \|\mathbf{G}\boldsymbol{\theta}\| \leq R}{\operatorname{argmin}} \|\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta}\|^2. \quad (3.11)$$

The Lagrange multiplier method leads to an unconstrained problem

$$\tilde{\boldsymbol{\theta}}_{G,\lambda} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \{\|\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta}\|^2 + \lambda \|\mathbf{G}\boldsymbol{\theta}\|^2\}.$$

A proper choice of λ ensures that the solution $\tilde{\boldsymbol{\theta}}_{G,\lambda}$ belongs to $\mathcal{B}_G(R)$ and solves also the problem (3.11). So, the approach based on a Sobolev smoothness assumption, leads back to regularization and shrinkage.

3.6 Shrinkage in a linear inverse problem

This section extends the previous approaches to the situation with indirect observations. More precisely, we focus on the model

$$\mathbf{Y} = \mathbf{A}\mathbf{f}^* + \boldsymbol{\varepsilon},$$

where \mathbf{A} is a given linear operator (matrix) and \mathbf{f}^* is the target of analysis. With the obvious change of notation this problem can be put back in the general linear setup $\mathbf{Y} = \boldsymbol{\Psi}^\top \boldsymbol{\theta} + \boldsymbol{\varepsilon}$. The special focus is due to the facts that the target can be high dimensional or even functional and that the product $\mathbf{A}^\top \mathbf{A}$ is usually badly posed and its inversion is a hard task. Below we consider separately the cases when the spectral representation for this problem is available and the general case.

3.7 Spectral cut-off and spectral penalization. Diagonal estimates

Suppose that the eigenvectors of the matrix $\mathbf{A}^\top \mathbf{A}$ are available. This allows for reducing the model to the spectral representation by an orthogonal change of the coordinate

system: $\mathbf{Z} = \Lambda \mathbf{u} + \Lambda^{1/2} \boldsymbol{\xi}$ with a diagonal matrix $\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_p\}$ and a homogeneous noise $\text{Var}(\boldsymbol{\xi}) = \sigma^2 \mathbf{I}_p$; see Section 1.1.4. Below we assume without loss of generality that the eigenvalues λ_j are ordered and decrease with j . This spectral representation means that one observes empirical Fourier coefficients z_m described by the equation $z_j = \lambda_j u_j + \lambda_j^{1/2} \xi_j$ for $j = 1, \dots, p$. The LSE or qMLE estimate of the spectral parameter \mathbf{u} is given by

$$\tilde{\mathbf{u}} = \Lambda^{-1} \mathbf{Z} = (\lambda_1^{-1} z_1, \dots, \lambda_p^{-1} z_p)^\top.$$

Exercise 3.7.1. Consider the spectral representation $\mathbf{Z} = \Lambda \mathbf{u} + \Lambda^{1/2} \boldsymbol{\xi}$. The LSE $\tilde{\mathbf{u}}$ reads as $\tilde{\mathbf{u}} = \Lambda^{-1} \mathbf{Z}$.

If the dimension p of the model is high or, specifically, if the spectral values λ_j rapidly go to zero, it might be useful to only track few coefficients u_1, \dots, u_m and to set all the remaining ones to zero. The corresponding estimate $\tilde{\mathbf{u}}_m = (\tilde{u}_{m,1}, \dots, \tilde{u}_{m,p})^\top$ reads as

$$\tilde{u}_{m,j} \stackrel{\text{def}}{=} \begin{cases} \lambda_j^{-1} z_j & \text{if } j \leq m, \\ 0 & \text{otherwise.} \end{cases}$$

It is usually referred to as a *spectral cut-off* estimate.

Exercise 3.7.2. Consider the linear model $\mathbf{Y} = \mathbf{A} \mathbf{f}^* + \boldsymbol{\varepsilon}$. Let U be an orthogonal transform in \mathbb{R}^p providing $U \mathbf{A}^\top \mathbf{A} U^\top = \Lambda$ with a diagonal matrix Λ leading to the spectral representation for $\mathbf{Z} = U \mathbf{A} \mathbf{Y}$. Write the corresponding spectral cut-off estimate $\tilde{\mathbf{f}}_m$ for the original vector \mathbf{f}^* . Show that computing this estimate only requires to know the first m eigenvalues and eigenvectors of the matrix $\mathbf{A}^\top \mathbf{A}$.

Similarly to the direct case, a spectral cut-off can be extended to *spectral shrinkage*: one multiplies every empirical coefficient z_j with a factor $\alpha_j \in (0, 1)$. This leads to the *spectral shrinkage* estimate $\tilde{\mathbf{u}}_\alpha$ with $\tilde{u}_{\alpha,j} = \alpha_j \lambda_j^{-1} z_j$. Here α stands for the vector of coefficients α_j for $j = 1, \dots, p$. A spectral cut-off method is a special case of this shrinkage with α_j equal to one or zero.

Exercise 3.7.3. Specify the spectral shrinkage $\tilde{\mathbf{u}}_\alpha$ with a given vector α for the situation of Exercise 3.7.2.

The spectral cut-off method can be described as follows. Let $\boldsymbol{\psi}_1, \boldsymbol{\psi}_2, \dots$ be the intrinsic orthonormal basis of the problem composed of the standardized eigenvectors of

$A^\top A$ and leading to the spectral representation $\mathbf{Z} = \Lambda \mathbf{u} + \Lambda^{1/2} \boldsymbol{\xi}$ with the target vector \mathbf{u} . In terms of the original target \mathbf{f}^* , one is looking for a solution or an estimate in the form $\mathbf{f} = \sum_j u_j \boldsymbol{\psi}_j$. The design orthogonality allows to estimate every coefficient u_j independently of the others using the empirical Fourier coefficient $\boldsymbol{\psi}_j^\top \mathbf{Y}$. Namely, $\tilde{u}_j = \lambda_j^{-1} \boldsymbol{\psi}_j^\top \mathbf{Y} = \lambda_j^{-1} z_j$. The LSE procedure tries to recover \mathbf{f} as the full sum $\tilde{\mathbf{f}} = \sum_j \tilde{u}_j \boldsymbol{\psi}_j$. The projection method suggests to cut this sum at the index m : $\tilde{\mathbf{f}}_m = \sum_{j \leq m} \tilde{u}_j \boldsymbol{\psi}_j$, while the shrinkage procedure is based on downweighting the empirical coefficients \tilde{u}_j : $\tilde{\mathbf{f}}_\alpha = \sum_j \alpha_j \tilde{u}_j \boldsymbol{\psi}_j$.

Next we study the risk of the shrinkage method. Orthonormality of the basis $\boldsymbol{\psi}_j$ allows to represent the loss as $\|\tilde{\mathbf{u}}_\alpha - \mathbf{u}^*\|^2 = \|\tilde{\mathbf{f}}_\alpha - \mathbf{f}^*\|^2$. Under the noise homogeneity one obtains the following result.

Theorem 3.7.1. *Let $\mathbf{Z} = \Lambda \mathbf{u}^* + \Lambda^{1/2} \boldsymbol{\xi}$ with $\text{Var}(\boldsymbol{\xi}) = \sigma^2 \mathbf{I}_p$. It holds for the shrinkage estimate $\tilde{\mathbf{u}}_\alpha$*

$$\mathcal{R}(\tilde{\mathbf{u}}_\alpha) \stackrel{\text{def}}{=} \mathbb{E} \|\tilde{\mathbf{u}}_\alpha - \mathbf{u}^*\|^2 = \sum_{j=1}^p |\alpha_j - 1|^2 u_j^{*2} + \sum_{j=1}^p \alpha_j^2 \sigma^2 \lambda_j^{-1}.$$

Proof. The empirical Fourier coefficients z_j are uncorrelated and $\mathbb{E} z_j = \lambda_j u_j^*$, $\text{Var} z_j = \sigma^2 \lambda_j$. This implies

$$\mathbb{E} \|\tilde{\mathbf{u}}_\alpha - \mathbf{u}^*\|^2 = \sum_{j=1}^p \mathbb{E} |\alpha_j \lambda_j^{-1} z_j - u_j^*|^2 = \sum_{j=1}^p \{|\alpha_j - 1|^2 u_j^{*2} + \alpha_j^2 \sigma^2 \lambda_j^{-1}\}$$

as required.

Risk minimization leads to the oracle choice of the vector $\boldsymbol{\alpha}$ or

$$\boldsymbol{\alpha}^* = \underset{\boldsymbol{\alpha}}{\text{argmin}} \mathcal{R}(\tilde{\mathbf{u}}_\alpha)$$

where the minimum is taken over the set of all admissible vectors $\boldsymbol{\alpha}$.

Similar analysis can be done for the spectral cut-off method.

Exercise 3.7.4. The risk of the spectral cut-off estimate $\tilde{\mathbf{u}}_m$ fulfills

$$\mathcal{R}(\tilde{\mathbf{u}}_m) = \sum_{j=1}^m \lambda_j^{-1} \sigma^2 + \sum_{j=m+1}^p u_j^{*2}.$$

Specify the choice of the oracle cut-off index m^* .

Ordered model selection for linear smoothers

Consider the linear model $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$. Suppose that a family of linear smoothers $\tilde{\boldsymbol{\theta}}_m = \mathcal{S}_m \mathbf{Y}$ is given, where \mathcal{S}_m is for each $m \in \mathcal{M}$ a given $p \times n$ matrix. The task is to develop a data based model selector \hat{m} which performs nearly as good as the optimal choice which depends on the model and is not available.

4.1 Model and problem

This section defines our setup.

4.1.1 Loss and risk

First we have to define the loss and risk. We focus on the quadratic loss and risk with a weighting $q \times p$ matrix W :

$$\begin{aligned} \varrho_m &\stackrel{\text{def}}{=} \|W(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*)\|^2, \\ \mathcal{R}_m &\stackrel{\text{def}}{=} \mathbb{E} \|W(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*)\|^2. \end{aligned} \tag{4.1}$$

Of course, the loss and the risk depend on the choice of W . We drop this dependence but it is important to keep in mind the role of W in the definition of ϱ_m . Typical examples of W are as follows.

Estimation of the whole vector $\boldsymbol{\theta}^*$:

Let W be the identity matrix $W = \mathbf{I}_p$ with $q = p$. This means that the estimation loss is measured by the usual squared Euclidean norm $\|\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*\|^2$.

Prediction:

Let W be the square root of the total Fisher information matrix $\mathbb{F} = \sigma^{-2} \Psi \Psi^\top$, that is, $W^2 = \mathbb{F}$. Then minimization of the loss ϱ_m is equivalent to maximizing the log-likelihood $-(2\sigma^2)^{-1} \|\Psi^\top(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2$. Such type of loss is usually referred to as *prediction*

loss because it measures the fit and the prediction ability of the true model by the model with the parameter θ .

Semiparametric estimation:

Let the target of estimation is not the whole vector θ^* but its subvector θ_0^* of dimension q . The estimate $\Pi\tilde{\theta}_m$ is called the *profile maximum likelihood estimate*. The matrix W can be defined as the projector Π_0 on the θ_0^* subspace. The corresponding loss is equal to the squared Euclidean norm in this subspace:

$$\varrho_m = \|\Pi_0(\tilde{\theta}_m - \theta^*)\|^2.$$

Alternatively, one can select W^2 as the efficient Fisher information matrix defined by the relation

$$W^2 \stackrel{\text{def}}{=} \check{\mathbb{F}} = (\Pi_0 \mathbb{F}^{-1} \Pi_0^\top)^{-1}.$$

Linear functional estimation:

The choice of the weighting matrix W can be adjusted to the problem of estimating some functionals of the whole parameter θ^* . For instance, in the regression problem $EY_i = f(X_i)$ with the Fourier expansion the target function f can be represented as

$$f(x) = \sum_{j \geq 0} \theta_j^* \psi_j(x) = \sum_{j \geq 0} \{\theta_{2j}^* \cos(2\pi jx) + \theta_{2j+1}^* \sin(2\pi jx)\}.$$

The value of this function at zero coincides with the functional

$$f(0) = \sum_j \theta_{2j}^*.$$

The first derivative of this function leads to the functional

$$f'(0) = 2\pi \sum_{j \geq 0} j \theta_{2j+1}^*.$$

Exercise 4.1.1. Describe the linear functions for the values

- $f(0.5)$;
- $f''(1)$;
- $\int_0^1 f(x) dx$;
- $\int_0^1 \sin(2\pi x) f'(x) dx$.

The most important feature of the estimate $\tilde{\theta}_m$ is *linearity*. It allows us to study its properties similarly to above. The next result describes the risk decomposition of $\tilde{\theta}_m$.

Theorem 4.1.1. *Let $\mathbb{E}\varepsilon = 0$, $\text{Var}(\varepsilon) = \sigma^2 I_n$. Then*

$$\begin{aligned}\mathbb{E}\tilde{\boldsymbol{\theta}}_m &= \boldsymbol{\theta}_m^* = \mathcal{S}_m \mathbf{f}^*, \\ \text{Var}(\tilde{\boldsymbol{\theta}}_m) &= V_m^2 = \sigma^2 \mathcal{S}_m \mathcal{S}_m^\top, \\ \mathcal{R}_m &= \|W(\boldsymbol{\theta}_m^* - \boldsymbol{\theta}^*)\|^2 + \sigma^2 \text{tr}(W \mathcal{S}_m \mathcal{S}_m^\top W^\top) \\ &= \|W(\mathcal{S}_m - \mathcal{S})\mathbf{f}^*\|^2 + \sigma^2 \text{tr}(W \mathcal{S}_m \mathcal{S}_m^\top W^\top).\end{aligned}\tag{4.2}$$

The expression for the risk \mathcal{R}_m in (4.2) can be called the “bias-variance” decomposition.

4.1.2 Oracle choice. Bias-variance trade-off

Let the weighting matrix W be fixed and the loss and risk of $\tilde{\boldsymbol{\theta}}_m$ be defined by (4.1). The optimal (oracle) choice of the tuning parameter m can be defined as the risk minimizer:

$$m^* \stackrel{\text{def}}{=} \underset{m \in \mathcal{M}}{\text{argmin}} \mathcal{R}_m.$$

The *model selection* problem can be described as the choice of m by data which *mimics the oracle*, that is, we aim at constructing a selector \hat{m} leading to the adaptive estimate $\hat{\boldsymbol{\theta}} = \tilde{\boldsymbol{\theta}}_{\hat{m}}$ with the properties similar to the oracle estimate $\tilde{\boldsymbol{\theta}}_{m^*}$.

Below we discuss the *ordered case*. The parameter m is treated as complexity of the method $\tilde{\boldsymbol{\theta}}_m$ and this parameter may grows until the maximal value M . In some cases the set \mathcal{M} of possible m choices can be countable and/or continuous and even unbounded. For simplicity of presentation, we assume that \mathcal{M} is finite. We also assume without loss of generality that m is an integer non-negative number. Complexity can be naturally expressed via the variance of the stochastic term of the estimate $\tilde{\boldsymbol{\theta}}_m$: the larger m , the larger is the variance $\text{Var}(W\tilde{\boldsymbol{\theta}}_m)$. Due to the result of Theorem 4.1.1, this condition can be written as

$$W \mathcal{S}_m \mathcal{S}_m^\top W^\top \leq W \mathcal{S}_{m'} \mathcal{S}_{m'}^\top W^\top, \quad m' > m.\tag{4.3}$$

In the case of projection estimation with the identity matrix W , this variance is linear in m , $\text{Var}(\tilde{\boldsymbol{\theta}}_m) = \sigma^2 m$. In general the dependence of the variance term on m may be more complicated but the monotonicity constraint (4.3) has to be preserved.

Further, it is implicitly assumed that the bias term $\|W(\boldsymbol{\theta}^* - \boldsymbol{\theta}_m^*)\|^2$ becomes small when m increases. The smallest value $m = 0$ corresponds to the simplest (zero) model with probably a large bias, while m large ensures a good approximation quality $\boldsymbol{\theta}_m^* \approx \boldsymbol{\theta}^*$ and a small bias at cost of a big complexity measured by the variance term. However,

in general, in the contrary to the case of projection estimation, one cannot require that the bias term $\|W(\boldsymbol{\theta}^* - \boldsymbol{\theta}_m^*)\|$ monotonously decreases with m . One example is given by estimation-at-a-point problem.

Exercise 4.1.2. Build an example of a signal $\boldsymbol{\theta}^*$ in the sequence space model $Y_i = \theta_j^* + \varepsilon_j$ such that the bias $\|W(\boldsymbol{\theta}_m - \boldsymbol{\theta}^*)\|$ is not monotonous in m .

Hint: Consider, e.g., $\tilde{\boldsymbol{\theta}}_m$ being the projector on the first m coordinates and $W\boldsymbol{\theta} = \sum_j \theta_j$. Take $\boldsymbol{\theta}^*$ by alternating blocks of 1's and -1's with equal length.

4.2 Unbiased risk estimation

The empirical risk is obtained by substituting the data \mathbf{Y} in place of the true function \mathbf{f}^* . Let $\tilde{\boldsymbol{\theta}} = \mathcal{S}\mathbf{Y}$ be the estimate from the largest considered model which we assume to be unbiased: $\mathbb{E}\tilde{\boldsymbol{\theta}} = \boldsymbol{\theta}^* = \mathcal{S}\mathbf{f}^*$. The analog of the empirical risk is

$$\|W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}})\|^2 = \|W(\mathcal{S}_m - \mathcal{S})\mathbf{Y}\|^2.$$

Lemma 4.2.1. *Let $\mathbb{E}\boldsymbol{\varepsilon} = 0$, $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2\mathbf{I}_n$. Then*

$$\mathbb{E}\|W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}})\|^2 = \|W(\boldsymbol{\theta}_m^* - \boldsymbol{\theta}^*)\|^2 + \sigma^2 \text{tr}\{W(\mathcal{S} - \mathcal{S}_m)(\mathcal{S} - \mathcal{S}_m)^\top W^\top\}.$$

Proof. The use of $\mathbf{Y} = \mathbf{f}^* + \boldsymbol{\varepsilon}$ and $\mathbb{E}\boldsymbol{\varepsilon} = 0$ implies

$$\begin{aligned} \mathbb{E}\|W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}})\|^2 &= \mathbb{E}\|W(\mathcal{S}_m - \mathcal{S})(\mathbf{f}^* + \boldsymbol{\varepsilon})\|^2 \\ &= \|W(\mathcal{S}_m - \mathcal{S})\mathbf{f}^*\|^2 + \mathbb{E}\|W(\mathcal{S}_m - \mathcal{S})\boldsymbol{\varepsilon}\|^2. \end{aligned}$$

For the quadratic stochastic term, we use the trace representation and $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2\mathbf{I}_n$:

$$\begin{aligned} \mathbb{E}\|W(\mathcal{S}_m - \mathcal{S})\boldsymbol{\varepsilon}\|^2 &= \mathbb{E}[\text{tr}\{W(\mathcal{S}_m - \mathcal{S})\boldsymbol{\varepsilon}\}^\top \{W(\mathcal{S}_m - \mathcal{S})\boldsymbol{\varepsilon}\}] \\ &= \mathbb{E}[\text{tr}\{W(\mathcal{S}_m - \mathcal{S})\boldsymbol{\varepsilon}\} \{W(\mathcal{S}_m - \mathcal{S})\boldsymbol{\varepsilon}\}^\top] \\ &= \text{tr}\{W(\mathcal{S}_m - \mathcal{S})\mathbb{E}(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top)(\mathcal{S}_m - \mathcal{S})^\top W^\top\} \\ &= \sigma^2 \text{tr}\{W(\mathcal{S} - \mathcal{S}_m)(\mathcal{S} - \mathcal{S}_m)^\top W^\top\}. \end{aligned}$$

This proves the required decomposition.

Exercise 4.2.1. Derive a similar expansion for an inhomogeneous noise $\boldsymbol{\varepsilon}$ with $\text{Var}(\boldsymbol{\varepsilon}) = \boldsymbol{\Sigma}$.

Now we compare the risk \mathcal{R}_m of $\tilde{\boldsymbol{\theta}}_m$ with the expected empirical risk $\mathbb{E}\|W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}})\|^2$. Similarly to the projection case, the bias term $\|W(\mathcal{S}_m - \mathcal{S})\mathbf{f}^*\|^2$ of the risk \mathcal{R}_m and

of the expectation of the empirical risk are the same. The variance terms are different but known. So, one can make a correction and build an unbiased risk estimate $\tilde{\mathcal{R}}_m$. Indeed, by construction

$$\mathbb{E}\|W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}})\|^2 + \sigma^2 \operatorname{tr}\{W[\mathcal{S}_m \mathcal{S}_m^\top - (\mathcal{S} - \mathcal{S}_m)(\mathcal{S} - \mathcal{S}_m)^\top]W^\top\} = \mathcal{R}_m.$$

This suggests the following definition of the unbiased risk estimate $\tilde{\mathcal{R}}_m$:

$$\begin{aligned} \tilde{\mathcal{R}}_m &\stackrel{\text{def}}{=} \|W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}})\|^2 + 2\sigma^2 \operatorname{tr}(W \mathcal{S} \mathcal{S}_m^\top W^\top) \\ &= \|W(\mathcal{S}_m - \mathcal{S})\mathbf{Y}\|^2 + 2\sigma^2 \operatorname{tr}(W \mathcal{S} \mathcal{S}_m^\top W^\top). \end{aligned} \quad (4.4)$$

The term $\operatorname{tr}(W \mathcal{S} \mathcal{S}^\top W^\top)$ is omitted because it does not depend on m .

Exercise 4.2.2. Let $\operatorname{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}_n$ and $\mathcal{S}_m = \mathbf{I}_m$ is a projector on the subspace spanned by $\theta_1, \dots, \theta_m$.

1. Check that

$$W \mathcal{S} \mathcal{S}_m^\top W^\top = W \mathcal{S}_m \mathcal{S}_m^\top W^\top = W \mathcal{S}_m W^\top.$$

2. Compute $\tilde{\mathcal{R}}_m$ if W projects $\boldsymbol{\theta}$ onto the first q components.

The Stein unbiased risk estimate (SURE) \hat{m} is defined by minimization of $\tilde{\mathcal{R}}_m$ from (4.4)

$$\hat{m} = \underset{m}{\operatorname{argmin}} \tilde{\mathcal{R}}_m. \quad (4.5)$$

4.2.1 Zone of insensitivity

Here we discuss whether the SURE method (4.5) does the job of model selection. The study is again done by the analysis of pairwise comparison. Let the oracle choice m^* is defined by risk minimization. By construction and the definition of the oracle, it implies

$$\mathbb{E}(\tilde{\mathcal{R}}_m - \tilde{\mathcal{R}}_{m^*}) = \mathcal{R}_m - \mathcal{R}_{m^*} \geq 0.$$

So, one can expect that it works “in mean”. However, similarly to the projection case, we need a more detailed and rigorous analysis of the probability of selecting $\hat{m} = m \neq m^*$ which is related to the probability of the inequality $\tilde{\mathcal{R}}_m - \tilde{\mathcal{R}}_{m^*} < 0$. Intuitively, if the systematic part (expectation) $\mathcal{R}_m - \mathcal{R}_{m^*}$ is significantly large, the probability of selecting $\hat{m} = m$ is small. In the region where the risk difference $\mathcal{R}_m - \mathcal{R}_{m^*}$ is small, the selector \hat{m} most likely makes a random choice.

For a detailed study, we need to slightly extend the monotonicity condition. Namely, it is required that the variance of $W\tilde{\boldsymbol{\theta}}_m$ grows with m while the variance of $W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}})$ decreases with m . For the projection method these two conditions coincide.

We begin with the following lemma.

Lemma 4.2.2. *For each $m^\circ \neq m$, it holds*

$$\tilde{\mathcal{R}}_m - \tilde{\mathcal{R}}_{m^\circ} = \begin{cases} -\mathbf{Y}^\top \mathcal{D}_{m,m^\circ}^2 \mathbf{Y} + \mathcal{Q}_{m,m^\circ}, & m > m^\circ, \\ \mathbf{Y}^\top \mathcal{D}_{m^\circ,m}^2 \mathbf{Y} - \mathcal{Q}_{m^\circ,m}, & m < m^\circ, \end{cases} \quad (4.6)$$

where for $m > m^\circ$

$$\begin{aligned} \mathcal{D}_{m,m^\circ}^2 &\stackrel{\text{def}}{=} (\mathcal{S} - \mathcal{S}_{m^\circ})^\top W^\top W (\mathcal{S} - \mathcal{S}_{m^\circ}) - (\mathcal{S} - \mathcal{S}_m)^\top W^\top W (\mathcal{S} - \mathcal{S}_m), \\ \mathcal{Q}_{m,m^\circ} &\stackrel{\text{def}}{=} \sigma^2 \text{tr}(\mathcal{D}_{m,m^\circ}^2) + \sigma^2 \text{tr}(\mathcal{D}_{m^\circ,m}^2), \\ \mathcal{D}_{m,m^\circ}^2 &= \mathcal{S}_m^\top W^\top W \mathcal{S}_m - \mathcal{S}_{m^\circ}^\top W^\top W \mathcal{S}_{m^\circ}. \end{aligned} \quad (4.7)$$

Proof. The definition (4.4) yields

$$\begin{aligned} \tilde{\mathcal{R}}_m - \tilde{\mathcal{R}}_{m^\circ} &= \mathbf{Y}^\top (\mathcal{S}_m - \mathcal{S})^\top W^\top W (\mathcal{S}_m - \mathcal{S}) \mathbf{Y} + 2\sigma^2 \text{tr}(W \mathcal{S} \mathcal{S}_m^\top W^\top) \\ &\quad - \mathbf{Y}^\top (\mathcal{S}_{m^\circ} - \mathcal{S})^\top W^\top W (\mathcal{S}_{m^\circ} - \mathcal{S}) \mathbf{Y} - 2\sigma^2 \text{tr}\{W \mathcal{S} \mathcal{S}_{m^\circ}^\top W^\top\} \end{aligned}$$

and it remains to check that

$$2 \text{tr}\{W \mathcal{S} (\mathcal{S}_m - \mathcal{S}_{m^\circ})^\top W^\top\} = \text{tr}(\mathcal{D}_{m,m^\circ}^2) + \text{tr}(\mathcal{D}_{m^\circ,m}^2).$$

This identity follows directly from the definition.

Exercise 4.2.3. Let $\Psi \mathcal{S}_m$ be a projector in \mathbb{R}^n on a linear subspace \mathcal{L}_m of dimension m for a growing system $\mathcal{L}_1 \subset \mathcal{L}_2 \subset \dots \subset \mathcal{L}_M$:

1. Check the monotonicity condition $\mathcal{D}_{m,m^\circ}^2 \geq 0$ for all $m > m^\circ$;
2. Write explicitly the expansion (4.6);
3. Check for the case with $W = \Psi^\top$ that

$$\mathcal{D}_{m,m^\circ}^2 = D_{m,m^\circ}^2 = \mathcal{S}_m W W^\top \mathcal{S}_m^\top - \mathcal{S}_{m^\circ} W W^\top \mathcal{S}_{m^\circ}^\top;$$

4. Specify $\mathcal{D}_{m,m^\circ}^2$ and \mathcal{Q}_{m,m° for the situation when the loss function W only depends on the first coordinate: $W\boldsymbol{\theta} = \theta_1$.

Now we apply Lemma 4.2.2 to study the probability of the event $\{\tilde{\mathcal{R}}_m < \tilde{\mathcal{R}}_{m^*}\}$. The decomposition (4.6) allows to represent this event as

$$(\mathbf{f}^* + \boldsymbol{\varepsilon})^\top \mathcal{D}_{m,m^*}^2 (\mathbf{f}^* + \boldsymbol{\varepsilon}) \geq \sigma^2 \operatorname{tr}(\mathcal{D}_{m,m^*}^2) + \sigma^2 \operatorname{tr}(D_{m,m^*}^2), \quad m > m^*, \quad (4.8)$$

$$(\mathbf{f}^* + \boldsymbol{\varepsilon})^\top \mathcal{D}_{m^*,m}^2 (\mathbf{f}^* + \boldsymbol{\varepsilon}) \leq \sigma^2 \operatorname{tr}(\mathcal{D}_{m^*,m}^2) + \sigma^2 \operatorname{tr}(D_{m^*,m}^2), \quad m < m^*. \quad (4.9)$$

This probability can be studied in the same way with separate consideration for $m > m^*$ and $m < m^*$. In the region $m > m^*$, in view of the definition (4.7), the condition (4.8) can be rewritten as

$$\begin{aligned} & \boldsymbol{\varepsilon}^\top \mathcal{D}_{m,m^*}^2 \boldsymbol{\varepsilon} - \mathbb{E}(\boldsymbol{\varepsilon}^\top \mathcal{D}_{m,m^*}^2 \boldsymbol{\varepsilon}) + 2\boldsymbol{\varepsilon}^\top \mathcal{D}_{m,m^*}^2 \mathbf{f}^* \\ & \geq \sigma^2 \operatorname{tr}(D_{m,m^*}^2) - \mathbf{f}^{*\top} \mathcal{D}_{m,m^*}^2 \mathbf{f}^* = \mathcal{R}_m - \mathcal{R}_{m^*}. \end{aligned} \quad (4.10)$$

For typical ordering schemes, one can expect that the bias component $\mathbf{f}^{*\top} \mathcal{D}_{m,m^*}^2 \mathbf{f}^*$ in the risk difference $\mathcal{R}_m - \mathcal{R}_{m^*}$ remains bounded, while the variance term $\sigma^2 \operatorname{tr}(D_{m,m^*}^2)$ grows and tends to dominate the bias for $m \gg m^*$. The deviation bound of Theorem 3.3.2 can be used to check under which conditions the event in (4.8) or (4.10) can occur only with a very small probability. In particular, with a dominating probability the centered quadratic form in (4.10) can be bounded as follows

$$\boldsymbol{\varepsilon}^\top \mathcal{D}_{m,m^*}^2 \boldsymbol{\varepsilon} - \mathbb{E}(\boldsymbol{\varepsilon}^\top \mathcal{D}_{m,m^*}^2 \boldsymbol{\varepsilon}) \leq \mathbf{C} \sigma^2 \sqrt{\mathbf{x} \operatorname{tr}(\mathcal{D}_{m,m^*}^4)}$$

for a fixed constant \mathbf{C} . Similarly, the linear term $2\boldsymbol{\varepsilon}^\top \mathcal{D}_{m,m^*}^2 \mathbf{f}^*$ after standardization can be bounded by a quantile of the standard normal law: with a dominating probability

$$\boldsymbol{\varepsilon}^\top \mathcal{D}_{m,m^*}^2 \mathbf{f}^* \leq \sigma \|\mathcal{D}_{m,m^*}^2 \mathbf{f}^*\| z_1(\mathbf{x}).$$

One can conclude that the SURE selector does the job for values of $m > m^*$ for which

$$\operatorname{tr}(D_{m,m^*}^2) - \sigma^{-2} \mathbf{f}^{*\top} \mathcal{D}_{m,m^*}^2 \mathbf{f}^* \geq \mathbf{C} \sqrt{\mathbf{x} \operatorname{tr}(\mathcal{D}_{m,m^*}^4)} + \sigma^{-1} \|\mathcal{D}_{m,m^*}^2 \mathbf{f}^*\| z_1(\mathbf{x})$$

for some fixed constant \mathbf{C} . The region of insensitivity $\mathcal{M}^\circ(\mathbf{x})$ includes those $m > m^*$, for which this condition does not meet.

Similarly one can consider the case of $m < m^*$. Here the monotonicity condition ensures that the variance term $\operatorname{tr}(D_{m,m^*}^2)$ is bounded from above for all $m < m^*$. However, the definition of the oracle m^* yields a significant bias for $m < m^*$. The condition (4.9) can be represented as

$$\begin{aligned} & -\boldsymbol{\varepsilon}^\top \mathcal{D}_{m^*,m}^2 \boldsymbol{\varepsilon} + \mathbb{E}(\boldsymbol{\varepsilon}^\top \mathcal{D}_{m^*,m}^2 \boldsymbol{\varepsilon}) - 2\boldsymbol{\varepsilon}^\top \mathcal{D}_{m^*,m}^2 \mathbf{f}^* \\ & \geq -\sigma^2 \operatorname{tr}(D_{m^*,m}^2) + \mathbf{f}^{*\top} \mathcal{D}_{m^*,m}^2 \mathbf{f}^* = \mathcal{R}_m - \mathcal{R}_{m^*}. \end{aligned}$$

Here the zone of insensitivity corresponds to values of m for which the bias $\mathbf{f}^{*\top} \mathcal{D}_{m^*,m}^2 \mathbf{f}^*$ is not too large relative to the standard deviation of the quadratic form $\boldsymbol{\varepsilon}^\top \mathcal{D}_{m^*,m}^2 \boldsymbol{\varepsilon}$. The bound of Theorem 3.3.2 leads to a sufficient condition

$$\sigma^{-2} \mathbf{f}^{*\top} \mathcal{D}_{m^*,m}^2 \mathbf{f}^* - \text{tr}(D_{m^*,m}^2) \geq \mathfrak{C} \sqrt{\text{tr}(\mathcal{D}_{m^*,m}^4)} + \sigma^{-1} \|\mathcal{D}_{m^*,m}^2 \mathbf{f}^*\| z_1(\mathbf{x}).$$

Under this condition the probability of selecting such a value m is negligible.

In general, one can conclude that the selector based on the unbiased risk estimation, works well beyond the region of insensitivity which means that the risk difference $\mathcal{R}_m - \mathcal{R}_{m^*}$ is prominent relative to the standard deviation of the quadratic form $\boldsymbol{\varepsilon}^\top \mathcal{D}_{m^*,m}^2 \boldsymbol{\varepsilon}$.

4.2.2 An oracle bound

For stating an oracle bound on the loss and risk of $\widehat{\boldsymbol{\theta}} = \widetilde{\boldsymbol{\theta}}_{\widehat{m}}$, we also have to bound the excess, that is, the difference between the loss of the adaptive and the oracle procedure in the region of insensitivity. We only present a bound for each particular m .

Lemma 4.2.3. *It holds for $m > m^*$*

$$\varrho_m - \varrho_{m^*} = \boldsymbol{\varepsilon}^\top D_{m,m^*}^2 \boldsymbol{\varepsilon} - \mathbf{f}^{*\top} \mathcal{D}_{m,m^*}^2 \mathbf{f}^* + 2\mathbf{b}_{m,m^*}^\top \boldsymbol{\varepsilon} \quad (4.11)$$

with

$$\mathbf{b}_{m,m^*} \stackrel{\text{def}}{=} \mathbf{f}^{*\top} (\mathcal{S}_m - \mathcal{S})^\top W^\top W \mathcal{S}_m - \mathbf{f}^{*\top} (\mathcal{S}_{m^*} - \mathcal{S})^\top W^\top W \mathcal{S}_{m^*}.$$

Moreover,

$$\varrho_m - \varrho_{m^*} - (\mathcal{R}_m - \mathcal{R}_{m^*}) = \boldsymbol{\varepsilon}^\top D_{m,m^*}^2 \boldsymbol{\varepsilon} - \mathbb{E}(\boldsymbol{\varepsilon}^\top D_{m,m^*}^2 \boldsymbol{\varepsilon}) + 2\mathbf{b}_{m,m^*}^\top \boldsymbol{\varepsilon}. \quad (4.12)$$

Proof. The use of $\widetilde{\boldsymbol{\theta}}_m = \mathcal{S}_m \mathbf{Y} = \mathcal{S}_m(\mathbf{f}^* + \boldsymbol{\varepsilon})$ and $\boldsymbol{\theta}^* = \mathcal{S} \mathbf{f}^*$ for $m > m^*$ by (4.2)

$$\begin{aligned} \varrho_m - \varrho_{m^*} &= \|W(\widetilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*)\|^2 - \|W(\widetilde{\boldsymbol{\theta}}_{m^*} - \boldsymbol{\theta}^*)\|^2 \\ &= \|W \mathcal{S}_m \boldsymbol{\varepsilon}\|^2 - \|W \mathcal{S}_{m^*} \boldsymbol{\varepsilon}\|^2 + \|W(\mathcal{S}_m - \mathcal{S}) \mathbf{f}^*\|^2 - \|W(\mathcal{S}_{m^*} - \mathcal{S}) \mathbf{f}^*\|^2 \\ &\quad + 2\mathbf{f}^{*\top} (\mathcal{S}_m - \mathcal{S})^\top W^\top W \mathcal{S}_m \boldsymbol{\varepsilon} - 2\mathbf{f}^{*\top} (\mathcal{S}_{m^*} - \mathcal{S})^\top W^\top W \mathcal{S}_{m^*} \boldsymbol{\varepsilon} \\ &= \boldsymbol{\varepsilon}^\top D_{m,m^*}^2 \boldsymbol{\varepsilon} - \mathbf{f}^{*\top} \mathcal{D}_{m,m^*}^2 \mathbf{f}^* + 2\mathbf{b}_{m,m^*}^\top \boldsymbol{\varepsilon} \end{aligned} \quad (4.13)$$

and (4.11) follows. The assertion (4.12) holds by formula (4.2) for the risk \mathcal{R}_m .

Further, the deviation of the centered quadratic form $\boldsymbol{\varepsilon}^\top D_{m,m^*}^2 \boldsymbol{\varepsilon} - \mathbb{E}(\boldsymbol{\varepsilon}^\top D_{m,m^*}^2 \boldsymbol{\varepsilon})$ can be bounded with a probability $1 - e^{-x}$ as

$$|\boldsymbol{\varepsilon}^\top D_{m,m^*}^2 \boldsymbol{\varepsilon} - \mathbb{E}(\boldsymbol{\varepsilon}^\top D_{m,m^*}^2 \boldsymbol{\varepsilon})| \leq \mathfrak{C} \sigma^2 \sqrt{x \text{tr}(D_{m,m^*}^4)}$$

while the linear term in (4.13) can be bounded in the same sense by

$$|\mathbf{b}_{m,m^*}^\top \boldsymbol{\varepsilon}| \leq \sigma \|\mathbf{b}_{m,m^*}\| z_1(\mathbf{x}).$$

An increase of \mathbf{x} by $\log(\mathbb{N}(\mathbf{x})) = \log(\#\mathcal{M}^\circ(\mathbf{x}))$ yields a uniform bound over the region of insensitivity $\mathcal{M}^\circ(\mathbf{x})$. One can conclude that the risk of the adaptive estimate $\hat{\boldsymbol{\theta}} = \tilde{\boldsymbol{\theta}}_{\hat{m}}$ is close to the risk of the oracle estimate $\tilde{\boldsymbol{\theta}}_{m^*}$ if the range of $\sqrt{\text{tr}(D_{m,m^*}^4)}$ over the set $\mathcal{M}^\circ(\mathbf{x})$ is small relative to the oracle risk \mathcal{R}_{m^*} .

To be done: An oracle risk bound for the SURE method

4.3 Smallest accepted (SmA) method

This section discusses another approach based on multiple testing. Let H_{m° vs. H_m mean a hypothesis of no significant bias between the models m° and m . Instead of the likelihood ratio test, we consider the test statistic \mathbb{T}_{m,m° based on the weighted difference $W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ})$:

$$\mathbb{T}_{m,m^\circ} = \|W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ})\| = \|W(\mathcal{S}_m - \mathcal{S}_{m^\circ})\mathbf{Y}\|.$$

4.3.1 Decomposition of the test statistic. Bias and variance

Below we use the decomposition

$$\mathbb{T}_{m,m^\circ} = \|W(\mathcal{S}_m - \mathcal{S}_{m^\circ})(\mathbf{f}^* + \boldsymbol{\varepsilon})\| = \|\mathbf{b}_{m,m^\circ} + \boldsymbol{\xi}_{m,m^\circ}\|, \quad (4.14)$$

where

$$\begin{aligned} \mathbf{b}_{m,m^\circ} &\stackrel{\text{def}}{=} W(\mathcal{S}_m - \mathcal{S}_{m^\circ})\mathbf{f}^*, \\ \boldsymbol{\xi}_{m,m^\circ} &\stackrel{\text{def}}{=} W(\mathcal{S}_m - \mathcal{S}_{m^\circ})\boldsymbol{\varepsilon}. \end{aligned}$$

It obviously holds $\mathbb{E}\boldsymbol{\xi}_{m,m^\circ} = \mathbf{0}$. Introduce $q \times q$ -matrix \mathbb{V}_{m,m° as the variance of $W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ})$:

$$\mathbb{V}_{m,m^\circ} \stackrel{\text{def}}{=} \text{Var}\{W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ})\} = \text{Var}\{W(\mathcal{S}_m - \mathcal{S}_{m^\circ})\mathbf{Y}\}.$$

If the noise $\boldsymbol{\varepsilon}$ in the decomposition $\mathbf{Y} = \mathbf{f}^* + \boldsymbol{\varepsilon}$ is homogeneous with $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}_n$, it holds

$$\begin{aligned} \mathbb{V}_{m,m^\circ} &= W(\mathcal{S}_m - \mathcal{S}_{m^\circ}) \text{Var}(\boldsymbol{\varepsilon})(\mathcal{S}_m - \mathcal{S}_{m^\circ})^\top W^\top \\ &= \sigma^2 W(\mathcal{S}_m - \mathcal{S}_{m^\circ})(\mathcal{S}_m - \mathcal{S}_{m^\circ})^\top W^\top. \end{aligned}$$

Note that in the case of homogeneous noise the matrix \mathbb{V}_{m,m° is easily computable. The matrices W and \mathcal{S}_m are given by the method and only the noise variance has to be estimated from the data.

Further,

$$\mathbb{E} \mathbb{T}_{m,m^\circ}^2 = \|\mathbf{b}_{m,m^\circ}\|^2 + \mathbb{E}\|\boldsymbol{\xi}_{m,m^\circ}\|^2 = \|\mathbf{b}_{m,m^\circ}\|^2 + \text{tr}(\mathbb{V}_{m,m^\circ}). \quad (4.15)$$

The bias term $\mathbf{b}_{m,m^\circ} \stackrel{\text{def}}{=} W(\mathcal{S}_m - \mathcal{S}_{m^\circ})\mathbf{f}^*$ is significant if its squared norm is competitive with the variance term $\mathbf{p}_{m,m^\circ} = \text{tr}(\mathbb{V}_{m,m^\circ})$.

Below we proceed with the Gaussian errors $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$. Then $\boldsymbol{\xi}_{m,m^\circ} \sim \mathcal{N}(0, \mathbb{V}_{m,m^\circ})$ and its distribution is completely specified by the covariance matrix \mathbb{V}_{m,m° and it holds for any \mathbf{x}

$$\mathbb{P}\left(\|\boldsymbol{\xi}_{m,m^\circ}\| > z(\mathbb{V}_{m,m^\circ}, \mathbf{x})\right) = e^{-\mathbf{x}} \quad (4.16)$$

for the tail function $z(\mathbb{V}, \mathbf{x})$ of $\|\boldsymbol{\xi}\|$ with $\boldsymbol{\xi} \sim \mathcal{N}(0, \mathbb{V})$.

4.3.2 Multiplicity correction. A Bonferroni bound and Monte-Carlo method

Below we need a uniform in $m > m^\circ$ version of the probability bound (4.16). Let

$$\mathcal{M}^+(m^\circ) \stackrel{\text{def}}{=} \{m \in \mathcal{M} : m > m^\circ\}.$$

By $q_{m^\circ} = q_{m^\circ}(\mathbf{x})$ denote the corresponding multiplicity correction:

$$\mathbb{P}\left(\bigcup_{m \in \mathcal{M}^+(m^\circ)} \{\|\boldsymbol{\xi}_{m,m^\circ}\| \geq z(\mathbb{V}_{m,m^\circ}, \mathbf{x} + q_{m^\circ}(\mathbf{x}))\}\right) = e^{-\mathbf{x}}. \quad (4.17)$$

We write $\mathbf{x}_{m^\circ} \stackrel{\text{def}}{=} \mathbf{x} + q_{m^\circ}(\mathbf{x})$. A simple way of computing the multiplicity correction q_{m° is based on the Bonferroni bound: $q_{m^\circ} = \log(\#\mathcal{M}^+(m^\circ))$. Indeed

$$\begin{aligned} & \mathbb{P}\left(\bigcup_{m \in \mathcal{M}^+(m^\circ)} \{\|\boldsymbol{\xi}_{m,m^\circ}\| \geq z(\mathbb{V}_{m,m^\circ}, \mathbf{x} + \log(\#\mathcal{M}^+(m^\circ)))\}\right) \\ & \leq \sum_{m \in \mathcal{M}^+(m^\circ)} \mathbb{P}\left(\|\boldsymbol{\xi}_{m,m^\circ}\| \geq z(\mathbb{V}_{m,m^\circ}, \mathbf{x} + \log(\#\mathcal{M}^+(m^\circ)))\right) \\ & \leq \sum_{m \in \mathcal{M}^+(m^\circ)} e^{-\mathbf{x} - \log(\#\mathcal{M}^+(m^\circ))} = e^{-\mathbf{x}}. \end{aligned}$$

However, it is well known that the Bonferroni bound is very conservative and leads to very large correction q_{m° , especially if the random vectors $\boldsymbol{\xi}_{m,m^\circ}$ are strongly correlated. This is exactly the case under consideration. Note that the joint distribution of the $\boldsymbol{\xi}_{m,m^\circ}$'s is precisely known. An analytic bound for the correction $q_{m^\circ}(\mathbf{x})$ is hardly computable but one can use Monte-Carlo experiments: draw K samples $\boldsymbol{\varepsilon}^{(k)} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$ for $k = 1, \dots, K$, compute for each the stochastic vectors $\boldsymbol{\xi}_{m,m^\circ}^{(k)}$ for $m \in \mathcal{M}^+(m^\circ)$, and evaluate $q_{m^\circ}(\mathbf{x})$ as the smallest value for which the condition (4.16) is fulfilled under the Monte-Carlo empirical measure \mathbb{P}° in place of \mathbb{P} .

Exercise 4.3.1. Show that for each \mathbf{x} there is unique value $q_{m^\circ}(\mathbf{x})$ ensuring (4.17). Check whether q_{m° depends on \mathbf{x} .

4.3.3 Oracle choice

We say that m° is a good choice if there is no significant bias \mathbf{b}_{m,m° for any $m > m^\circ$. This condition can be quantified as

$$\|\mathbf{b}_{m,m^\circ}\|^2 \leq \beta^2 \mathbf{p}_{m,m^\circ}, \quad m > m^\circ \quad (4.18)$$

for a given β and $\mathbf{p}_{m,m^\circ} = \text{tr}(\mathbb{V}_{m,m^\circ}) = \mathbb{E}\|\boldsymbol{\xi}_{m,m^\circ}\|^2$. This condition can be viewed as the “bias-variance trade-off”. The parameter β controls the bias component in the risk due to decomposition (4.15). Now define the *oracle* m^* as the minimal m° with the property (4.18):

$$m^* \stackrel{\text{def}}{=} \min\left\{m^\circ: \max_{m>m^\circ} \{\|\mathbf{b}_{m,m^\circ}\|^2 - \beta^2 \mathbf{p}_{m,m^\circ}\} \leq 0\right\}. \quad (4.19)$$

Exercise 4.3.2. Check that the definition (4.19) uniquely defines the value m^* .

The proposed procedure based on the family of test statistics \mathbb{T}_{m,m° has to be validated by the condition that any good model in the sense (4.18) will be rejected with a very small nominal probability. In particular, this would apply to the oracle choice m^* .

4.3.4 Data-driven choice and the oracle inequality

Define the selector \hat{m} by the “smallest accepted” (SmA) rule. Namely, with $\mathbf{x}_{m^\circ} = \mathbf{x} + q_{m^\circ}(\mathbf{x})$; see (4.16), the acceptance rule reads as follows:

$$\begin{aligned} \{m^\circ \text{ is accepted}\} &\Leftrightarrow \left\{ \max_{m>m^\circ} \{\mathbb{T}_{m,m^\circ} - \mathbf{z}_{m,m^\circ}\} \leq 0 \right\}, \\ \mathbf{z}_{m,m^\circ} &\stackrel{\text{def}}{=} z(\mathbb{V}_{m,m^\circ}, \mathbf{x}_{m^\circ}) + \beta \sqrt{\text{tr}(\mathbb{V}_{m,m^\circ})}. \end{aligned}$$

The SmA rule is

$$\begin{aligned} \hat{m} &\stackrel{\text{def}}{=} \text{“smallest accepted”} \\ &= \min\left\{m^\circ: \max_{m>m^\circ} \{\mathbb{T}_{m,m^\circ} - \mathbf{z}_{m,m^\circ}\} \leq 0\right\}. \end{aligned} \quad (4.20)$$

Exercise 4.3.3. Check that (4.20) uniquely defines the value \hat{m} .

The bound (4.16) and the definition of the oracle m^* yield the desired propagation property:

$$\mathbb{P}(m^* \text{ is rejected}) \leq e^{-\mathbf{x}}. \quad (4.21)$$

In some sense, this property is built in the construction of the procedure and is fulfilled automatically. It ensures that with a large probability, the procedure does not reject the point m^* , this model is “good” and has to be accepted. Therefore, the selector \hat{m} typically takes its value in the region $m \leq m^*$. It remains to check the performance of the method in this region.

First we specify the zone of insensitivity for $m < m^*$. In this region we expect a bias component competitive with the variance of the oracle. If the bias $\Delta_{m^*,m} = \|\mathbf{b}_{m^*,m}\|$ is significantly large then such m will be accepted with only very small probability. We can use the bound

$$\begin{aligned} \mathbb{P}(\|\mathbf{b}_{m^*,m} + \boldsymbol{\xi}_{m^*,m}\| < \mathbf{z}_{m^*,m}) \\ \leq \mathbb{P}(\|\boldsymbol{\xi}_{m^*,m}\| > \|\mathbf{b}_{m^*,m}\| - \mathbf{z}_{m^*,m}). \end{aligned} \quad (4.22)$$

Acceptance of the model $m < m^*$ assumes that the check of $\mathbb{T}_{m^*,m}$ does not fail. If

$$\|\mathbf{b}_{m^*,m}\| \geq \mathbf{z}_{m^*,m} + z(\mathbb{V}_{m^*,m}, \mathbf{x}),$$

then (4.22) implies the upper bound for the probability of accepting this model m :

$$\begin{aligned} \mathbb{P}(m \text{ is accepted}) &\leq \mathbb{P}(\|\mathbf{b}_{m^*,m} + \boldsymbol{\xi}_{m^*,m}\| < \mathbf{z}_{m^*,m}) \\ &\leq \mathbb{P}(\|\boldsymbol{\xi}_{m^*,m}\| > z(\mathbb{V}_{m^*,m}, \mathbf{x})) \leq e^{-\mathbf{x}}. \end{aligned} \quad (4.23)$$

The definition of the zone of insensitivity requires a small multiplicity correction of the \mathbf{x} -level. Define $\bar{\mathbf{x}}_{m^*} \stackrel{\text{def}}{=} \mathbf{x} + \log(m^*)$ and

$$\mathcal{M}^\circ(\mathbf{x}) \stackrel{\text{def}}{=} \{m < m^* : \|\mathbf{b}_{m^*,m}\| \leq \mathbf{z}_{m^*,m} + z(\mathbb{V}_{m^*,m}, \bar{\mathbf{x}}_{m^*})\}. \quad (4.24)$$

Theorem 4.3.1. *Let the set $\mathcal{M}^\circ(\mathbf{x})$ be defined by (4.24). Then*

$$\mathbb{P}(\hat{m} < m^*, \hat{m} \notin \mathcal{M}^\circ(\mathbf{x})) \leq e^{-\mathbf{x}}.$$

Proof. It suffices to note that by (4.23)

$$\begin{aligned} \mathbb{P}(\hat{m} < m^*, \hat{m} \notin \mathcal{M}^\circ(\mathbf{x})) &\leq \sum_{m < m^*, m \notin \mathcal{M}^\circ(\mathbf{x})} \mathbb{P}(\|\mathbf{b}_{m^*,m} + \boldsymbol{\xi}_{m^*,m}\| < \mathbf{z}_{m^*,m}) \\ &\leq \sum_{m < m^*, m \notin \mathcal{M}^\circ(\mathbf{x})} e^{-\bar{\mathbf{x}}_{m^*}} \leq e^{-\mathbf{x}}. \end{aligned}$$

This implies the result because the cardinality of the set $\{m < m^*, m \notin \mathcal{M}^\circ(\mathbf{x})\}$ does not exceed m^* .

Note that the proof of this result is based on the crude Bonferroni upper bound and the definition of $\mathcal{M}^\circ(\mathbf{x})$ can be refined by choosing $\bar{\mathbf{x}}_{m^*}$ more carefully. However, this value only enters in the theoretical bound and is not used in the procedure, a fine tuning for this value is not required.

One can conclude that the SmA procedure typically selects some m from the set $\mathcal{M}^\circ(\mathbf{x})$. The situation is similar to the case of the SURE procedure with one essential difference: the just defined set $\mathcal{M}^\circ(\mathbf{x})$ is located from the left of the oracle, which enables us to guarantee some oracle quality of estimation.

If $\hat{m} = m < m^*$, the acceptance of \hat{m} implies by definition

$$\mathbb{T}_{m^*,m} \leq \mathbf{z}_{m^*,m}$$

or, equivalently

$$\|W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^*})\| \leq \mathbf{z}_{m^*,m}. \quad (4.25)$$

One can see that the variability of the adaptive estimate $\hat{\boldsymbol{\theta}} = \tilde{\boldsymbol{\theta}}_{\hat{m}}$ is bounded from above by a fixed deterministic quantity. This enables us to state the final oracle inequality.

Theorem 4.3.2. *Let \mathbf{z}_{m,m° ensure for each m° the condition (4.16). Then the propagation property (4.21) is fulfilled. The following oracle bound holds for the SmA estimate $\hat{\boldsymbol{\theta}} = \tilde{\boldsymbol{\theta}}_{\hat{m}}$ on a random set $\Omega(\mathbf{x})$ of dominating probability $1 - 2e^{-\mathbf{x}}$:*

$$\|W(\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}_{m^*})\| \leq \bar{\mathbf{z}}(m^*) \stackrel{\text{def}}{=} \max_{m \in \mathcal{M}^\circ(\mathbf{x})} \mathbf{z}_{m^*,m}. \quad (4.26)$$

This implies the probabilistic oracle bound: with probability at least $1 - 2e^{-\mathbf{x}}$

$$\|W(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \leq \|W(\tilde{\boldsymbol{\theta}}_{m^*} - \boldsymbol{\theta}^*)\| + \bar{\mathbf{z}}(m^*). \quad (4.27)$$

Proof. First we note that

$$\begin{aligned} \mathbb{P}(\hat{m} > m^*) &\leq e^{-\mathbf{x}}, \\ \mathbb{P}(\hat{m} \leq m^*, m \notin \mathcal{M}^\circ(\mathbf{x})) &\leq e^{-\mathbf{x}}. \end{aligned}$$

This implies (4.26) by (4.25). (4.27) follows now by the triangle inequality.

The result (4.27) is called the oracle bound because it compares the loss of the data-driven selector \hat{m} and of the optimal choice m^* .

The value $\bar{\mathbf{z}}(m^*)$ in the right hand-side of (4.27) can be viewed as ‘‘payment for adaptation’’. An interesting feature of the presented result is that not only the oracle quality but also the payment of adaptation depend upon the unknown value $\boldsymbol{\theta}^*$ and the

corresponding oracle choice m^* . In the worst case of the model with the flat risk \mathcal{R}_m , the set $\mathcal{M}^\circ(\mathbf{x})$ can coincide with the whole range $m \leq m^*$. Even in this case the bounds (4.25) and (4.27) are meaningful. However, the payment for adaptation $\bar{z}(m^*)$ in this case can be larger than the oracle risk.

4.3.5 Analysis of the payment for adaptation $\bar{z}(m^*)$

To be done: An upper bound on $\bar{z}(m^*)$

4.3.6 Choice of β and \mathbf{x}

The value β enters in the procedure, so it has to be selected by some rule. This value shows up in the definition of the oracle choice and it controls the bias-variance relation. A standard risk minimization suggests to take $\beta = 1$. However, we prefer to keep the option of selecting β open in dependence on the type of applications. An increase of β results in the corresponding increase of the critical values \mathbf{z}_{m,m° , the acceptance rule becomes more conservative. Therefore, the selected model \hat{m} is smaller for β large and is bigger for β small. One can say, big values of β result in oversmoothing. In general, if the underlying model $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$ is specified correctly, we recommend to take β small. Numerical results confirm this suggestion and they are best for $\beta \equiv 0$. However, if we aim at building a procedure which is robust against model misspecification then larger values of β can be recommended. This especially concerns the noise misspecification: if the real noise vector $\boldsymbol{\varepsilon}$ is not homogeneous independent Gaussian and if its covariance structure is not precisely known, one can compensate these issues by choosing a positive β , say $\beta = 1$.

Another tuning parameter entering in the procedure is the quantile level \mathbf{x} . As one can guess from the beginning, this value is not important for the method. A default choice $\mathbf{x} = 3$ does a good job in all conducted numerical results.

4.3.7 Power loss function

The probabilistic oracle bound of Theorem 4.3.2 provides some statement about typical behavior of the adaptive SmA estimate $\hat{\boldsymbol{\theta}} = \tilde{\boldsymbol{\theta}}_{\hat{m}}$. Unfortunately, this bound does not yield a risk bound for quadratic or polynomial losses: even if big loss occurs with a small probability, the related risk can be large. It is interesting that the SmA procedure can be easily tuned to secure an oracle risk bound.

For simplicity of notation, we only consider the quadratic risk

$$\mathcal{R}(\hat{\boldsymbol{\theta}}) \stackrel{\text{def}}{=} \mathbb{E} \|\mathbf{W}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2.$$

We aim at comparing the risk of the SmA procedure with the risk \mathcal{R}_{m^*} of the oracle estimate $\tilde{\boldsymbol{\theta}}_{m^*}$. Remind the representation

$$W(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*) = \boldsymbol{\xi}_m + \mathbf{b}_m$$

with $\boldsymbol{\xi}_m = W\mathcal{S}_m\boldsymbol{\varepsilon}$ and $\mathbf{b}_m = W(\mathcal{S}_m - \mathcal{S})\mathbf{f}^*$ yielding the risk decomposition

$$\mathcal{R}_m \stackrel{\text{def}}{=} \mathbb{E}\|W(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*)\|^2 = \mathbb{E}\|\boldsymbol{\xi}_m\|^2 + \|\mathbf{b}_m\|^2 = \mathbf{p}_m + \|\mathbf{b}_m\|^2$$

with $\mathbf{p}_m = \text{tr}(\mathbb{V}_m)$ and $\mathbb{V}_m = \text{Var}(\boldsymbol{\xi}_m)$. A similar decomposition holds for the quadratic risk of the difference $\tilde{\boldsymbol{\theta}}_{m'} - \tilde{\boldsymbol{\theta}}_m$ for any pair $m' > m$:

$$\mathbb{E}\|W(\tilde{\boldsymbol{\theta}}_{m'} - \tilde{\boldsymbol{\theta}}_m)\|^2 = \mathbb{E}\|\boldsymbol{\xi}_{m',m}\|^2 + \|\mathbf{b}_{m',m}\|^2 = \mathbf{p}_{m',m} + \|\mathbf{b}_{m',m}\|^2.$$

Usually the oracle choice is defined by minimizing the risk \mathcal{R}_m . For our analysis, we have to slightly modify the definition in the spirit of SmA oracle (4.19). To ensure a risk bound, it is required that not only the model m^* is “good” but also all the larger models $m > m^*$ are “good” as well.

$$m^* \stackrel{\text{def}}{=} \min\left\{m^\circ : \max_{m' > m \geq m^\circ} \{\|\mathbf{b}_{m',m}\|^2 - \beta^2 \mathbf{p}_{m',m}\} \leq 0\right\}. \quad (4.28)$$

Below we also suppose that the bias component $\|\mathbf{b}_m\|^2$ fulfills

$$\|\mathbf{b}_m\| \leq \|\mathbf{b}_{m^*}\|, \quad m > m^*. \quad (4.29)$$

Otherwise, one can define $\|\mathbf{b}_{m^*}\| \stackrel{\text{def}}{=} \max_{m \geq m^*} \|\mathbf{b}_m\|$.

Below we consider the SmA procedure with the critical values $\mathbf{z}_{m',m}$ and explain how this values can be adjusted to ensure a risk bound. For stating the result, we need a bit more detailed analysis of the SmA procedure in the propagation zone $m > m^*$. If such m is selected then $\hat{\boldsymbol{\theta}}$ coincides with $\tilde{\boldsymbol{\theta}}_m$ whose loss and risk can be much larger than ones for the oracle estimate $\tilde{\boldsymbol{\theta}}_{m^*}$ because of its larger complexity. There is no essential bias component in this zone. Therefore, the SmA procedure can be tuned in the situation when there is no signal and hence no bias at all:

$$\mathbb{T}_{m',m} = \|W(\tilde{\boldsymbol{\theta}}_{m'} - \tilde{\boldsymbol{\theta}}_m)\| = \|\boldsymbol{\xi}_{m',m}\|.$$

The analysis is based on a simple but important observation that if $\hat{m} = m > m^*$, then the good model $m^\circ = m - 1$ is rejected. The latter means that at least one check based on $\mathbb{T}_{m',m-1}$ fails. The same can be expressed as follows: the maximum of the r.v.'s $\mathbb{T}_{m',m-1} \mathbb{I}(\mathbb{T}_{m',m-1} > z_{m',m-1})$ is positive. Let $z(\mathbb{V}_{m',m-1}, \mathbf{x})$ be the tail function of the $\|\boldsymbol{\xi}_{m',m-1}\|$ for $\boldsymbol{\xi}_{m',m-1} \sim \mathcal{N}(0, \mathbb{V}_{m',m-1})$. For each m , we consider the expectation of

$\mathbf{p}_m^{-1} \|\boldsymbol{\xi}_m\|^2$ on the random set $\Omega_m(\mathbf{x})$ on which at least one of test statistics $\mathbb{T}_{m',m-1} = \|\boldsymbol{\xi}_{m',m-1}\|$ exceeds the critical value $z(\mathbb{V}_{m',m-1}, \mathbf{x})$:

$$\mathcal{R}_m^+(\mathbf{x}) \stackrel{\text{def}}{=} \mathbb{E} \left[(\mathbf{p}_m^{-1} \|\boldsymbol{\xi}_m\|^2 \vee 1) \mathbb{I} \left(\max_{m' \geq m} \{ \|\boldsymbol{\xi}_{m',m-1}\| - z(\mathbb{V}_{m',m-1}, \mathbf{x}) \} > 0 \right) \right].$$

Similarly one can consider any other power loss function by replacing $(\mathbf{p}_m^{-1/2} \|\boldsymbol{\xi}_m\|)^2$ with $(\mathbf{p}_m^{-1/2} \|\boldsymbol{\xi}_m\|)^q$. In particular, $q = 0$ yields the probability loss considered before.

Now we define the value \mathbf{x}_{m-1} in a way to control the related deviation risk $\mathcal{R}_m^+(\mathbf{x})$. Let α_m be a given decreasing sequence. Its choice will be discussed below. We fix each value \mathbf{x}_{m-1} such that

$$\mathcal{R}_m^+(\mathbf{x}_{m-1}) = \alpha_m. \quad (4.30)$$

It implies

$$\begin{aligned} \mathbb{E} \left[\|\boldsymbol{\xi}_m\|^2 \mathbb{I} \left(\max_{m' \geq m} \{ \|\boldsymbol{\xi}_{m',m-1}\| - z(\mathbb{V}_{m',m-1}, \mathbf{x}_{m-1}) \} > 0 \right) \right] &\leq \alpha_m \mathbf{p}_m, \\ \mathbb{P} \left(\max_{m' \geq m} \{ \|\boldsymbol{\xi}_{m',m-1}\| - z(\mathbb{V}_{m',m-1}, \mathbf{x}_{m-1}) \} > 0 \right) &\leq \alpha_m. \end{aligned} \quad (4.31)$$

Now define the critical values \mathbf{z}_{m,m° of the SmA procedure as

$$\mathbf{z}_{m,m^\circ} = z(\mathbb{V}_{m,m^\circ}, \mathbf{x}_m) + \beta \mathbf{p}_{m,m^\circ}^{1/2}, \quad (4.32)$$

that is,

$$\hat{m} = \min \left\{ m^\circ : \max_{m > m^\circ} \{ \mathbb{T}_{m,m^\circ} - \mathbf{z}_{m,m^\circ} \} \leq 0 \right\}. \quad (4.33)$$

It is worth mentioning that the procedure is the same, and even the critical values \mathbf{z}_{m,m° are given by the same formula, as in the case of probabilistic loss. The only difference is in the propagation condition (4.30) which is a bit stronger than a similar condition for indicator loss. This implies that the values \mathbf{x}_m and \mathbf{z}_{m,m° are a bit larger.

Theorem 4.3.3. *Let the SmA procedure (4.33) be applied with the critical values \mathbf{z}_{m,m° from (4.32) and the values \mathbf{x}_m defined by (4.30) with the coefficients α_m satisfying*

$$\sum_{m > m^*} \alpha_m \mathbf{p}_m \leq \bar{\alpha}_{m^*} \mathbf{p}_{m^*} \quad (4.34)$$

for some $\bar{\alpha}_{m^*}$. Then

$$\mathbb{E} \|W(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2 \leq 2\bar{\alpha}_{m^*} \mathcal{R}_{m^*} + (\mathcal{R}_{m^*}^{1/2} + \bar{\mathbf{z}}_{m^*})^2$$

with

$$\bar{\mathbf{z}}_{m^*} \stackrel{\text{def}}{=} \max_{m < m^*} \mathbf{z}_{m^*,m}.$$

Proof. Let us fix $m > m^*$ and $m' \geq m$. The definition (4.28) of the oracle m^* and the formula (4.32) for the critical value $\mathbf{z}_{m',m-1}$ implies for the test statistic $\mathbb{T}_{m',m-1} = \|\boldsymbol{\xi}_{m',m-1} + \mathbf{b}_{m',m-1}\|$

$$\{\mathbb{T}_{m',m-1} > \mathbf{z}_{m',m-1}\} \subseteq \{\|\boldsymbol{\xi}_{m',m-1}\| > z(\mathbb{V}_{m',m-1}, \mathbf{x}_{m-1})\}.$$

Now we can bound the risk of $\hat{\boldsymbol{\theta}}$ on the set $\hat{m} > m^*$. We use that for $\hat{m} = m > m^*$ in view of (4.29)

$$\begin{aligned} \|W(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2 &= \|W(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*)\|^2 = \|\boldsymbol{\xi}_m + \mathbf{b}_m\|^2 \\ &\leq 2\|\boldsymbol{\xi}_m\|^2 + 2\|\mathbf{b}_m\|^2 \leq 2\|\boldsymbol{\xi}_m\|^2 + 2\|\mathbf{b}_{m^*}\|^2 \end{aligned}$$

and it holds by (4.31) and monotonicity $\mathbf{p}_m > \mathbf{p}_{m^*}$

$$\begin{aligned} &\mathbb{E}\{\|W(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2 \mathbb{I}(\hat{m} > m^*)\} \\ &\leq 2 \sum_{m > m^*} \mathbb{E}\{(\|\boldsymbol{\xi}_m\|^2 + \|\mathbf{b}_{m^*}\|^2) \mathbb{I}(\hat{m} = m)\} \\ &\leq 2 \sum_{m > m^*} \mathbb{E}\{(\|\boldsymbol{\xi}_m\|^2 + \|\mathbf{b}_{m^*}\|^2) \mathbb{I}(m-1 \text{ is rejected})\} \\ &= 2 \sum_{m > m^*} \mathbb{E}\left[(\|\boldsymbol{\xi}_m\|^2 + \|\mathbf{b}_{m^*}\|^2) \mathbb{I}\left(\max_{m' \geq m} \{\|\boldsymbol{\xi}_{m',m-1}\| - z(\mathbb{V}_{m',m-1}, \mathbf{x}_m)\} > 0\right)\right] \\ &\leq 2 \sum_{m > m^*} \alpha_m (\mathbf{p}_m + \|\mathbf{b}_{m^*}\|^2) \leq 2\bar{\alpha}_{m^*} (\mathbf{p}_{m^*} + \|\mathbf{b}_{m^*}\|^2) = 2\bar{\alpha}_{m^*} \mathcal{R}_{m^*}. \end{aligned}$$

Here we have used that (4.34) and $\mathbf{p}_m > \mathbf{p}_{m^*}$ imply $\sum_{m > m^*} \alpha_m \leq \bar{\alpha}_{m^*}$. In the zone $\hat{m} = m < m^*$, we can still use the stability property:

$$\|W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^*})\| \leq \mathbf{z}_{m^*,m};$$

cf. (4.25). We conclude

$$\begin{aligned} \mathbb{E}\|W(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2 &\leq 2\bar{\alpha}_{m^*} \mathcal{R}_{m^*} + \mathbb{E}\{\|W(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2 \mathbb{I}(\hat{m} < m^*)\} \\ &\leq 2\bar{\alpha}_{m^*} \mathcal{R}_{m^*} + \mathbb{E}(\|W(\tilde{\boldsymbol{\theta}}_{m^*} - \boldsymbol{\theta}^*)\| + \bar{\mathbf{z}}_{m^*})^2 \\ &\leq 2\bar{\alpha}_{m^*} \mathcal{R}_{m^*} + (\mathcal{R}_{m^*}^{1/2} + \bar{\mathbf{z}}_{m^*})^2 \end{aligned}$$

as required.

Similarly to the probabilistic loss function, the result can be refined by considering the zone of insensitivity in the region $m < m^*$.

To be done: Define the zone of insensitivity. An oracle bound refined

The constants α_m have to fulfill (4.34). If \mathbf{p}_m satisfy

$$\sum_{m>m^*} (\mathbf{p}_{m^*}/\mathbf{p}_m)^a \leq \mathbf{C}$$

for some $a > 0$ and a fixed constant \mathbf{C} , then one can take

$$\alpha_m = \mathbf{p}_m^{-1-a}$$

yielding

$$\sum_{m>m^*} \alpha_m \mathbf{p}_m \leq \sum_{m>m^*} \mathbf{p}_m^{-a} \leq \mathbf{C} \mathbf{p}_{m^*}^{-a}.$$

This implies for \mathbf{x}_m a bound

To be done: An upper bound on \mathbf{x}_m

In a similar way one can bound the payment for adaptation.

To be done: An upper bound on $z_{m^*}^b$

To be done: Linear vs geometric growth of \mathbf{p}_m

4.3.8 Penalized MLE

Now we discuss how the procedure and the results can be specified for the case of penalized maximum likelihood estimation. We consider the regression model $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$ with a homogeneous errors $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}_p$. For simplicity of description, we assume “in the model” case when the true regression $\mathbf{f}^* = \mathbb{E}\mathbf{Y}$ indeed follows the linear parametric model $\mathbf{f}^* = \Psi^\top \boldsymbol{\theta}^*$ with the target parameter $\boldsymbol{\theta}^*$. Let us given a family of penalized procedures described by the roughness matrices G_m^2 :

$$\begin{aligned} \tilde{\boldsymbol{\theta}}_m &= \underset{\boldsymbol{\theta}}{\text{argmax}} \{L(\boldsymbol{\theta}) - \|G_m \boldsymbol{\theta}\|^2/2\} \\ &= \underset{\boldsymbol{\theta}}{\text{argmin}} \{\|\mathbf{Y} - \Psi^\top \boldsymbol{\theta}\|^2 + \sigma^2 \|G_m \boldsymbol{\theta}\|^2\}. \end{aligned}$$

The solution is given by

$$\tilde{\boldsymbol{\theta}}_m = \mathcal{S}_m \mathbf{Y} = (\Psi \Psi^\top + \sigma^2 G_m^2)^{-1} \Psi \mathbf{Y}.$$

The penalizing matrices G_m^2 are naturally ordered, the running example is given by $G_m^2 = \lambda_m G^2$ for a fixed matrix G^2 and a decreasing sequence λ_m . Note that small values of the coefficient λ correspond to a minor penalization and thus, to a large model.

For any two values $m^\circ < m$, the difference $\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ}$ can be represented as

$$\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ} = (\mathcal{S}_m - \mathcal{S}_{m^\circ}) \mathbf{Y} = \mathcal{S}_{m,m^\circ} \mathbf{Y}$$

with

$$\mathcal{S}_{m,m^\circ} = \{(\Psi\Psi^\top + \sigma^2 G_m^2)^{-1} - (\Psi\Psi^\top + \sigma^2 G_{m^\circ}^2)^{-1}\}\Psi.$$

We consider three options for the weighting matrix W : the prediction loss with $W = \Psi^\top$, the estimation loss with $W = I_p$, and functional estimation $W: \mathbb{R}^p \rightarrow \mathbb{R}$.

Prediction loss

In the case of prediction loss, it holds

$$\Psi^\top(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*) = \Pi_m \mathbf{Y} - \mathbf{f}^* = \boldsymbol{\xi}_m + \mathbf{b}_m,$$

where

$$\Pi_m \stackrel{\text{def}}{=} \Psi^\top(\Psi\Psi^\top + \sigma^2 G_m^2)^{-1}\Psi$$

is a subprojector in \mathbb{R}^n with $\Pi_m = \Pi_m^\top$ and $\|\Pi_m\| \leq 1$. The bias $\mathbf{b}_m = \Pi_m \mathbf{f}^* - \mathbf{f}^*$ and the stochastic component $\boldsymbol{\xi}_m$ are vectors in \mathbb{R}^n . The quadratic risk is given by

$$\mathbb{E}\|\Psi^\top(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*)\|^2 = \sigma^2 \text{tr}(\Pi_m^2) + \|\mathbf{b}_m\|^2.$$

Further, for each $m^\circ < m$

$$\Psi^\top(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ}) = (\Pi_m - \Pi_{m^\circ})\mathbf{Y} = \Pi_{m,m^\circ}\mathbf{Y}.$$

The assumption $G_m^2 \geq G_{m^\circ}^2$ implies the monotonicity $\Pi_{m,m^\circ} \geq 0$.

Below we consider the scaled test statistic \mathbb{T}_{m,m° which reads as

$$\mathbb{T}_{m,m^\circ} = \sigma^{-1}\|\Pi_{m,m^\circ}\mathbf{Y}\|.$$

The corresponding covariance matrix is

$$\mathbb{V}_{m,m^\circ} = \text{Var}(\sigma^{-1}\Pi_{m,m^\circ}\mathbf{Y}) = \Pi_{m,m^\circ}^2.$$

The test statistic and the acceptance rule can be represented as the set of inequalities

$$\sigma^{-1}\|\Pi_{m,m^\circ}\mathbf{Y}\| \leq z(\mathbb{V}_{m,m^\circ}, \mathbf{x}_{m^\circ}) + \beta\sqrt{\text{tr}(\mathbb{V}_{m,m^\circ})}, \quad m > m^\circ.$$

If the effective dimension $\mathbf{p}_{m,m^\circ} = \text{tr}(\mathbb{V}_{m,m^\circ})$ is larger than a certain constant then

$$\begin{aligned} z(\mathbb{V}_{m,m^\circ}, \mathbf{x}_{m^\circ}) &\approx \mathbf{p}_{m,m^\circ}^{1/2} + \mathbf{C}\mathbf{x}_{m^\circ}^{1/2}, \\ \mathbf{z}_\beta(\mathbb{V}_{m,m^\circ}, \mathbf{x}_{m^\circ}) &\approx (1 + \beta)\mathbf{p}_{m,m^\circ}^{1/2} + \mathbf{C}\mathbf{x}_{m^\circ}^{1/2} \end{aligned}$$

Estimation loss

Now we briefly discuss the situation with $W = I_p$. Then

$$\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^* = \mathcal{S}_m \mathbf{Y} - \boldsymbol{\theta}^* = \mathbf{b}_m + \boldsymbol{\xi}_m$$

where $\boldsymbol{\xi}_m = \mathcal{S}_m \boldsymbol{\varepsilon}$ and $\mathbf{b}_m = \mathcal{S}_m \mathbf{f}^* - \boldsymbol{\theta}^*$ are vectors in \mathbb{R}^p . The quadratic risk is given by

$$\mathbb{E} \|\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*\|^2 = \sigma^2 \text{tr}(\mathcal{S}_m \mathcal{S}_m^\top) + \|\mathbf{b}_m\|^2.$$

The assumption $G_{m^\circ}^2 \geq G_m^2$ implies the monotonicity $\mathcal{S}_m \mathcal{S}_m^\top \geq \mathcal{S}_{m^\circ} \mathcal{S}_{m^\circ}^\top$. Further, for each $m^\circ < m$

$$\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ} = (\mathcal{S}_m - \mathcal{S}_{m^\circ}) \mathbf{Y} = \mathcal{S}_{m,m^\circ} \mathbf{Y}.$$

The scaled test statistic \mathbb{T}_{m,m° reads as

$$\mathbb{T}_{m,m^\circ} = \sigma^{-1} \|\mathcal{S}_{m,m^\circ} \mathbf{Y}\|.$$

The corresponding covariance matrix is

$$\mathbb{V}_{m,m^\circ} = \text{Var}(\sigma^{-1} \mathcal{S}_{m,m^\circ} \mathbf{Y}) = \mathcal{S}_{m,m^\circ} \mathcal{S}_{m,m^\circ}^\top.$$

and the acceptance rule can be represented as the set of bounds

$$\sigma^{-1} \|\mathcal{S}_{m,m^\circ} \mathbf{Y}\| \leq z(\mathbb{V}_{m,m^\circ}, \mathbf{x}_{m^\circ}) + \beta \sqrt{\text{tr}(\mathbb{V}_{m,m^\circ})}, \quad m > m^\circ.$$

If the effective dimension $\mathfrak{p}_{m,m^\circ} = \text{tr}(\mathbb{V}_{m,m^\circ})$ is larger than a certain constant then

$$\begin{aligned} z(\mathbb{V}_{m,m^\circ}, \mathbf{x}_{m^\circ}) &\approx \mathfrak{p}_{m,m^\circ}^{1/2} + \mathbf{C} \mathbf{x}_{m^\circ}^{1/2}, \\ \mathbf{z}_\beta(\mathbb{V}_{m,m^\circ}, \mathbf{x}_{m^\circ}) &\approx (1 + \beta) \mathfrak{p}_{m,m^\circ}^{1/2} + \mathbf{C} \mathbf{x}_{m^\circ}^{1/2}. \end{aligned}$$

To be done: make this more precise

Estimation of a linear functional

Now we discuss how the general SmA procedure can be applied to the problem of linear functional estimation in linear regression. We focus on the profile approach when one constructs the estimate $\tilde{\boldsymbol{\theta}}_m = \mathcal{S}_m \mathbf{Y}$ and the functional is described via the weighting matrix W of rank one which can be treated as a vector in \mathbb{R}^p . The target of estimation is the value $\phi^* = W \boldsymbol{\theta}^*$. It holds for each estimate $\tilde{\phi}_m = W \tilde{\boldsymbol{\theta}}_m$

$$\tilde{\phi}_m - \phi^* = W(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*) = W\mathcal{S}_m\boldsymbol{\varepsilon} + W(\mathcal{S}_m\mathbf{f}^* - \boldsymbol{\theta}^*) = \xi_m + b_m,$$

where $\xi_m = W\mathcal{S}_m\boldsymbol{\varepsilon}$ is a zero mean random variable and $b_m = W(\mathcal{S}_m\mathbf{f}^* - \boldsymbol{\theta}^*)$ is the deterministic bias. The squared risk of $\tilde{\phi}_m$ is given by the usual bias-variance decomposition:

$$\mathcal{R}_m = \mathbb{E}(\tilde{\phi}_m - \phi^*)^2 = \mathbb{E}(\xi_m + b_m)^2 = b_m^2 + \text{Var}(\xi_m) = b_m^2 + s_m^2$$

with

$$\mathbb{V}_m = s_m^2 = \sigma^2 W\mathcal{S}_m\mathcal{S}_m^\top W^\top.$$

Monotonicity condition $\mathcal{S}_m\mathcal{S}_m^\top \geq \mathcal{S}_{m^\circ}\mathcal{S}_{m^\circ}^\top$ yields monotonicity $\mathbb{V}_m \geq \mathbb{V}_{m^\circ}$ for the functional estimate. For each pair $m^\circ < m$

$$\tilde{\phi}_m - \tilde{\phi}_{m^\circ} = W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ}) = W(\mathcal{S}_m - \mathcal{S}_{m^\circ})\mathbf{Y} = W\mathcal{S}_{m,m^\circ}\mathbf{Y}.$$

The variance of this difference reads as

$$\mathbb{V}_{m,m^\circ} = s_{m,m^\circ}^2 = \text{Var}(\xi_{m,m^\circ}) = \sigma^2 W\mathcal{S}_{m,m^\circ}\mathcal{S}_{m,m^\circ}^\top W^\top.$$

The scaled test statistic \mathbb{T}_{m,m° is given by

$$\mathbb{T}_{m,m^\circ} = s_{m,m^\circ}^{-1} |W\mathcal{S}_{m,m^\circ}\mathbf{Y}|.$$

This is the absolute value of a standard Gaussian random variable and its tail function is $z_1(\mathbf{x})$. The acceptance rule can be represented as the set of bounds

$$\mathbb{T}_{m,m^\circ} \leq z_1(\mathbf{x}_{m^\circ}) + \beta, \quad m > m^\circ.$$

Here $z_1(\mathbf{x})$ is the $1 - e^{-\mathbf{x}}$ quantile of the absolute value of a standard normal r.v. ξ

$$\mathbb{P}(|\xi| > z_1(\mathbf{x})) = e^{-\mathbf{x}}$$

and the multiplicity correction $\mathbf{x}_{m^\circ} = \mathbf{x} + q_{m^\circ}(\mathbf{x})$ is defined by the condition

$$\mathbb{P}\left(\bigcup_{m>m^\circ} \{s_{m,m^\circ}^{-1} |\xi_{m,m^\circ}| \geq z_1(\mathbf{x}_{m^\circ})\}\right) \leq e^{-\mathbf{x}}. \quad (4.35)$$

We define

$$\hat{m} = \text{smallest accepted}, \quad \hat{\phi} = \tilde{\phi}_{\hat{m}} = W\mathcal{S}_{\hat{m}}\mathbf{Y}.$$

The oracle choice m^* corresponds to the smallest index m° for which

$$|b_{m,m^\circ}| \leq \beta s_{m,m^\circ}, \quad m > m^\circ.$$

The construction ensures that $\mathbb{P}(\widehat{m} > m^*) \leq e^{-x}$ and the probability oracle bound reads as

$$|\widehat{\phi} - \widetilde{\phi}_{m^*}| \leq \bar{z}(m^*) = \max_{m \in \mathcal{M}^\circ(\mathbf{x})} \{z_1(\mathbf{x}_m) s_{m^*,m}\}.$$

This is a meaningful bound even in the worst case of a flat risk when $\mathcal{M}^\circ(\mathbf{x})$ coincides with the set of all $m \leq m^*$. Indeed, monotonicity condition allows to bound $z_1(\mathbf{x}_m) \leq z_1(\mathbf{x}_0)$ and $s_{m^*,m} \leq s_{m^*,0} \leq s_{m^*}$ with $s_{m^*} = \mathbb{V}_{m^*}^{1/2}$ where zero index corresponds to the smallest model with a trivial estimate. This yields

$$|\widehat{\phi} - \widetilde{\phi}_{m^*}| \leq z_1(\mathbf{x}_0) s_{m^*}.$$

The value $z_1(\mathbf{x}_0)$ here can be easily upper bounded by usual Bonferroni arguments. Indeed, the choice $\mathbf{x}_0 \leq \mathbf{x} + \log(\#\mathcal{M})$ yields the uniform bound (4.35) in a straightforward way. This implies $z_1(\mathbf{x}_0) \approx \sqrt{\mathbf{x}_0} \leq \sqrt{\mathbf{x} + \log(\#\mathcal{M})}$. However, as we already mentioned, the Bonferroni correction is very conservative. In many situation it can be improved to the bound $\mathbf{x}_0 \leq \mathbf{x} + \mathbf{C} \log \log(\#\mathcal{M})$. This upper bound is still pessimistic because it is computed for the full set \mathcal{M} . It cannot be improved in a general minimax setup: one can always build a model with a flat risk in which $\mathcal{M}^\circ(\mathbf{x})$ is nearly \mathcal{M} . However, for each particular configuration $\boldsymbol{\theta}^*$, the oracle result can be made more precise and accurate by considering the zone of insensitivity and using that $s_{m^*,m}$ is smaller than s_{m^*} for m close to m^* .

4.4 Lepski's method

This section briefly discusses one more method of model selection known as Lepski's method and its acceptance rule is a bit stronger than for the SmA method. Namely, a candidate model m° is accepted if all larger models are accepted and all checks for \mathbb{T}_{m,m° are fulfilled. This rule can be written as

$$\widehat{m} \stackrel{\text{def}}{=} \min\{m^\circ : \mathbb{T}_{m',m} \leq \mathbf{z}_{m',m,m^\circ} \text{ for all } m' > m \geq m^\circ\} \quad (4.36)$$

with properly selected critical values $\mathbf{z}_{m',m,m^\circ}$. An important advantage of this procedure is its algorithmic description: one starts with the largest model and performs at each step a check of the next smaller model m° . If one of the checks for \mathbb{T}_{m,m° fails, that is, if $\mathbb{T}_{m,m^\circ} > \mathbf{z}_{m,m^\circ}$ then the procedure terminates and selects the latest previously accepted model. In the contrary, the SmA procedure screens all models. Even if one large model is rejected, it still may happen that a smaller one is accepted.

The critical values z_{m',m,m° can again be fixed by a procedure based on multiplicity correction. Let $z(\mathbb{V}, \mathbf{x})$ be the tail function of the norm $\|\boldsymbol{\xi}\|$ with $\boldsymbol{\xi} \sim \mathcal{N}(0, \mathbb{V})$. Similarly to the SmA procedure, define for each \mathbf{x} and m the correction $q_m(\mathbf{x})$ which ensures the probability bound

$$\mathbb{P}\left(\bigcup_{m' \in \mathcal{M}^+(m)} \{\|\boldsymbol{\xi}_{m',m}\| \geq z(\mathbb{V}_{m',m}, \mathbf{x} + q_m(\mathbf{x}))\}\right) = e^{-\mathbf{x}}. \quad (4.37)$$

Further, with $q_m(\mathbf{x})$ fixed for each m , introduce for each m° an additional correction $q_{m^\circ}^+ = q_{m^\circ}^+(\mathbf{x})$ which ensures the uniformity in all $m \geq m^\circ$:

$$\mathbb{P}\left(\bigcup_{m \geq m^\circ} \bigcup_{m' \in \mathcal{M}^+(m)} \{\|\boldsymbol{\xi}_{m',m}\| \geq z(\mathbb{V}_{m',m}, \mathbf{x} + q_m(\mathbf{x}) + q_{m^\circ}^+(\mathbf{x}))\}\right) = e^{-\mathbf{x}}. \quad (4.38)$$

Exercise 4.4.1. Check that the conditions (4.37) and (4.38) uniquely define the values $q_m(\mathbf{x})$ and $q_{m^\circ}^+(\mathbf{x})$ for all m and \mathbf{x} .

Define

$$\mathbf{x}_{m,m^\circ} \stackrel{\text{def}}{=} \mathbf{x} + q_m(\mathbf{x}) + q_{m^\circ}^+(\mathbf{x})$$

and apply (4.36) with

$$\mathbf{z}_{m',m,m^\circ} \stackrel{\text{def}}{=} z(\mathbb{V}_{m',m}, \mathbf{x}_{m,m^\circ}) + \beta \sqrt{\text{tr}(\mathbb{V}_{m',m})}.$$

Further, let the oracle choice m^* be defined in a similar way by checking the bias $\mathbf{b}_{m',m}$ for all $m' > m \geq m^\circ$:

$$m^* \stackrel{\text{def}}{=} \min\{m^\circ \in \mathcal{M}: \|\mathbf{b}_{m',m}\|^2 \leq \beta^2 \text{tr}(\mathbb{V}_{m',m}) \quad \forall m' > m \geq m^\circ\}.$$

Now the construction ensures that

$$\mathbb{P}(m^* \text{ is accepted}) = \mathbb{P}(\widehat{m} \leq m^*) \geq 1 - e^{-\mathbf{x}}.$$

The procedure (4.36) also guarantees that for $\widehat{m} = m < m^*$, it holds

$$\mathbb{T}_{m^*,m} = \|W(\widetilde{\boldsymbol{\theta}}_m - \widetilde{\boldsymbol{\theta}}_{m^*})\| \leq z(\mathbb{V}_{m^*,m}, \mathbf{x}_{m,m}) + \beta \sqrt{\text{tr}(\mathbb{V}_{m^*,m})}.$$

This bound can be used for establishing the oracle inequality similarly to the SmA procedure.

4.5 Intersection of confidence sets (ICI) method

The *intersection of confidence intervals* (ICI) procedure is another method based on multiple comparison. Here one constructs for each m a confidence set \mathcal{E}_m for the target $\boldsymbol{\theta}^*$ based on the estimate $\tilde{\boldsymbol{\theta}}_m$ and then selects the smallest m° for which the intersection of \mathcal{E}_m for all $m \geq m^\circ$ is not empty.

The construction is again based on the decomposition

$$W(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*) = \boldsymbol{\xi}_m + \mathbf{b}_m$$

with the stochastic component $\boldsymbol{\xi}_m$ and the bias \mathbf{b}_m . If we ignore the bias component \mathbf{b}_m , the use of $\boldsymbol{\xi}_m \sim \mathcal{N}(0, \mathbb{V}_m)$ leads to the confidence set

$$\mathcal{E}_m(z_m) = \{\boldsymbol{\theta} : \|\mathbb{V}_m^{-1/2}W(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta})\| \leq z_m\}.$$

The width z_m will be defined a bit later. Now the procedure reads as

$$\hat{m} \stackrel{\text{def}}{=} \min \left\{ m^\circ : \bigcap_{m \in \mathcal{M}^+(m^\circ)} \mathcal{E}_m(z_m) \neq \emptyset \right\}.$$

A slightly weaker condition is that any two confidence sets $\mathcal{E}_m(z_m)$ and $\mathcal{E}_{m'}(z_{m'})$ with $m' > m \geq m^\circ$ overlap.

Exercise 4.5.1. Consider the univariate case $q = \text{rank}(W) = 1$. Show that the pairwise check $\mathcal{E}_{m'}(z_{m'}) \cap \mathcal{E}_m(z_m) \neq \emptyset$ can be rewritten as

$$|W(\tilde{\boldsymbol{\theta}}_{m'} - \tilde{\boldsymbol{\theta}}_m)| \leq z_{m'}s_{m'} + z_ms_m, \quad (4.39)$$

where $s_m^2 = \mathbb{V}_m$.

Exercise 4.5.2. For the univariate case $q = 1$, check that the following conditions are equivalent:

- $\bigcap_{m \in \mathcal{M}^+(m^\circ)} \mathcal{E}_m(z_m) \neq \emptyset$;
- $\mathcal{E}_{m'}(z_{m'}) \cap \mathcal{E}_m(z_m) \neq \emptyset$ for all $m' > m \geq m^\circ$.

Exercise 4.5.3. Check for the case $q = 2$ whether the conditions of Exercise 4.5.2 on bivariate confidence sets $\mathcal{E}_m(z_m)$ are equivalent.

The algorithmic check of the condition on nonempty overlap is quite hard if the target dimension q is larger than one. Below we only consider the case $q = 1$.

The condition (4.39) of an nonempty overlap is similar to the check in the SmA procedure based on $\mathbb{T}_{m',m}$. The SmA check seems to be more accurate because it accounts

for the joint distribution of the two estimates $\tilde{\boldsymbol{\theta}}_{m'}$ and $\tilde{\boldsymbol{\theta}}_m$: their difference is scaled by its standard deviation. The ICI check normalizes the difference by the weighted sum of standard deviations.

Now we discuss how the ICI procedure can be tuned by fixing the values \mathbf{z}_m . The approach stays the same: one selects the “smallest” values which ensure the propagation property. The simplest way is based on a uniform multiplicity correction. The scaled stochastic components $s_m^{-1}\xi_m$ are standard normal and for each the tail function is given by the Laplace function:

$$\mathbb{P}(|s_m^{-1}\xi_m| > z_1(\mathbf{x})) \leq e^{-\mathbf{x}}.$$

Now one can make an adjustment $q^* = q^*(\mathbf{x})$ to the level \mathbf{x} to ensure a uniform bound:

$$\mathbb{P}\left(\bigcup_{m \in \mathcal{M}} \{|s_m^{-1}\xi_m| > z_1(\mathbf{x} + q^*)\}\right) \leq e^{-\mathbf{x}}. \quad (4.40)$$

Further we make an additional correction for the bias and define with a fixed $\beta \geq 0$ the uniform confidence widths

$$\mathbf{z}_m \equiv \bar{z}_\beta(\mathbf{x}) \stackrel{\text{def}}{=} \beta + z_1(\mathbf{x} + q^*). \quad (4.41)$$

The ICI procedure reads now

$$\hat{m} = \min\{m^\circ : |W(\tilde{\boldsymbol{\theta}}_{m'} - \tilde{\boldsymbol{\theta}}_m)| \leq (s_{m'} + s_m)\bar{z}_\beta(\mathbf{x}), \forall m' > m \geq m^\circ\}. \quad (4.42)$$

The oracle choice m^* can be defined as

$$m^* \stackrel{\text{def}}{=} \min\{m : |b_m| \leq \beta s_m, \forall m \geq m^\circ\}. \quad (4.43)$$

Again and again, the construction is designed to ensure the propagation property

$$\mathbb{P}(m^* \text{ is accepted}) = \mathbb{P}(\hat{m} \leq m^*) \geq 1 - e^{-\mathbf{x}}.$$

This fact is an easy corollary of the definitions (4.42) and (4.43) and of the probability bound (4.40). Also the procedure secures the “stability after propagation” property: if $\hat{m} = m < m^*$, then by definition of \hat{m}

$$|W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^*})| \leq (s_m + s_{m^*})\bar{z}_\beta(\mathbf{x}).$$

Monotonicity condition ensures that $s_m \leq s_{m^*}$ for $m < m^*$, and we end up with the oracle inequality for the adaptive estimate $\hat{\boldsymbol{\theta}} = \tilde{\boldsymbol{\theta}}_{\hat{m}}$.

Theorem 4.5.1. *Consider the linear Gaussian model $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$ with $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$. Let the weighting matrix W be of rank one. Given a set of estimates $\tilde{\boldsymbol{\theta}}_m$, $m \in \mathcal{M}$, with $\text{Var}(W\tilde{\boldsymbol{\theta}}_m) = s_m^2$, define the selector \hat{m} by (4.42) for $\bar{z}_\beta(\mathbf{x})$ from (4.41) and (4.40). The adaptive estimate $\hat{\boldsymbol{\theta}} = \tilde{\boldsymbol{\theta}}_{\hat{m}}$ and the oracle estimate $\tilde{\boldsymbol{\theta}}_{m^*}$ with m^* from (4.43) satisfy*

$$\mathbb{P}\left(|W(\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}_{m^*})| > 2s_{m^*}\bar{z}_\beta(\mathbf{x})\right) \leq 2e^{-x}.$$

To be done: Adaptive vs oracle risk

Unordered case. Anisotropic sets and subset selection

The SmA method of the previous section is quite general and can be extended to many statistical models and problem. However, it essentially requires the ordered structure of the set of considered models/methods. This section discusses how the SmA procedure can be extended to some other setups without ordered structure. To distinguish ordered and unordered cases, we denote by $\mathcal{A} = \{\varkappa\}$ the set of all considered models. The basic idea is to assume a kind of partial ordering which enables to define an acceptance rule:

\varkappa° is accepted if it is not rejected against any larger model.

This rule allows to fix a set of accepted models. Further we need some global measure of complexity which can be used for final selection:

$\hat{\varkappa}$ is the *simplest accepted* model.

Below we illustrate how this method works in two important examples: anisotropic classes and *subset selection* problems.

5.1 Subset selection procedure

Consider a linear model $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$. The dimension p of the vector $\boldsymbol{\theta}$ can be very large and we implicitly assumes a kind of *sparse* structure:

most of $\boldsymbol{\theta}^*$ -entries are nearly zero and can be dropped, there is a relatively small subvector of $\boldsymbol{\theta}^*$ containing the important features.

The aim is to find this subset and to estimate the whole vector $\boldsymbol{\theta}^*$. As in the ordered case, one can separate between prediction $W = \Psi^\top$ and estimation $W = I_p$ loss. Of course, these two problems coincide in the sequence space model with $n = p$ and $\Psi = I_p$.

5.1.1 SmA procedure and multilevel synchronization

Let \varkappa mean a subset of the whole index set $\{1, 2, \dots, p\}$. We use the obvious notation $\varkappa \vee \varkappa'$ for the union, $\varkappa \wedge \varkappa'$ for the overlap of two subsets \varkappa and \varkappa' , $\varkappa' - \varkappa$ for the complement of \varkappa within \varkappa' . Further we consider the usual partial ordering: $\varkappa' > \varkappa$ means that $\varkappa \subseteq \varkappa'$. The subset \varkappa is good if there is no significant bias in its complement \varkappa^c . The approach is to design a procedure which rejects any such good model with a small probability. The proposed SmA rule will be again to select the smallest (in complexity) accepted model.

The acceptance rule is based on pairwise comparison with a family of tests $\mathbb{T}_{\varkappa, \varkappa^\circ}$ for $\varkappa > \varkappa^\circ$. The model-candidate \varkappa° is accepted if no of $\mathbb{T}_{\varkappa, \varkappa^\circ}$ rejects the hypothesis of “no bias”. Given the loss matrix W , the test statistic $\mathbb{T}_{\varkappa, \varkappa^\circ}$ reads as in the ordered case:

$$\mathbb{T}_{\varkappa, \varkappa^\circ} = \sigma^{-1} \|W(\tilde{\boldsymbol{\theta}}_\varkappa - \tilde{\boldsymbol{\theta}}_{\varkappa^\circ})\|. \quad (5.1)$$

The acceptance rule can be written as

$$\varkappa^\circ \text{ is accepted iff } \mathbb{T}_{\varkappa, \varkappa^\circ} \leq \mathbf{z}_{\varkappa, \varkappa^\circ} \quad \forall \varkappa > \varkappa^\circ. \quad (5.2)$$

Now we discuss how the critical values $\mathbf{z}_{\varkappa, \varkappa^\circ}$ can be fixed by *synchronization (multiplicity correction)* of the individual *tail functions*. We use the decomposition of the test statistic $\mathbb{T}_{\varkappa, \varkappa^\circ}$ from (5.1)

$$\mathbb{T}_{\varkappa, \varkappa^\circ} = \|\boldsymbol{\xi}_{\varkappa, \varkappa^\circ} + \mathbf{b}_{\varkappa, \varkappa^\circ}\|.$$

Define

$$\mathbf{p}_{\varkappa, \varkappa^\circ} = \text{tr}\{\text{Var}(\boldsymbol{\xi}_{\varkappa, \varkappa^\circ})\}.$$

Suppose we are given for each pair $\varkappa > \varkappa^\circ$ a tail function $z_{\varkappa, \varkappa^\circ}(\mathbf{x})$ of the noise component $\|\boldsymbol{\xi}_{\varkappa, \varkappa^\circ}\|$ providing

$$\mathbb{P}(\|\boldsymbol{\xi}_{\varkappa, \varkappa^\circ}\| > z_{\varkappa, \varkappa^\circ}(\mathbf{x})) \leq e^{-\mathbf{x}}.$$

This tail function can be used for testing the hypothesis of no significant bias component in the test statistic $\mathbb{T}_{\varkappa, \varkappa^\circ}$. The model-candidate \varkappa° is accepted by the SmA method if all such tests for $\varkappa > \varkappa^\circ$ do. To keep the overall test level, we have to synchronize all performed $\mathbb{T}_{\varkappa, \varkappa^\circ}$ -based tests by correcting for multiple check. The simplest way of multiplicity correction is by a uniform increase of the level \mathbf{x} to control the overall rejecting probability: define $q_{\varkappa^\circ}(\mathbf{x})$ by the condition

$$\mathbb{P}\left(\bigcup_{\mathcal{M}(\mathcal{M}^\circ)} \left\{ \|\boldsymbol{\xi}_{\mathcal{M},\mathcal{M}^\circ}\| > z_{\mathcal{M},\mathcal{M}^\circ}(\mathbf{x} + q_{\mathcal{M}^\circ}(\mathbf{x})) \right\}\right) \leq e^{-\mathbf{x}}.$$

Denote $\mathbf{x}_{\mathcal{M}^\circ} = \mathbf{x} + q_{\mathcal{M}^\circ}(\mathbf{x})$ and apply the acceptance rule (5.2) with $\mathbf{z}_{\mathcal{M},\mathcal{M}^\circ}$ equal to such defined $z_{\mathcal{M},\mathcal{M}^\circ}(\mathbf{x}_{\mathcal{M}^\circ})$ after a small bias correction:

$$\mathbf{z}_{\mathcal{M},\mathcal{M}^\circ} \stackrel{\text{def}}{=} z_{\mathcal{M},\mathcal{M}^\circ}(\mathbf{x}_{\mathcal{M}^\circ}) + \beta \mathbf{p}_{\mathcal{M},\mathcal{M}^\circ}^{1/2}.$$

One can use a more sophisticated *multilevel synchronization* procedure which accounts for the complexity of the alternative model \mathcal{M} . Let $|\mathcal{M}^\circ|$ mean the cardinality (complexity) of \mathcal{M}° . For each $\tau > 0$, denote by $\mathcal{M}_\tau(\mathcal{M}^\circ)$ the set of model $\mathcal{M} > \mathcal{M}^\circ$ with $|\mathcal{M}| \leq \tau + |\mathcal{M}^\circ|$. For a given growing sequence $\tau_1 < \tau_2 < \dots < \tau_{\mathcal{K}}$, we write $\mathcal{M}_k(\mathcal{M}^\circ)$ instead of $\mathcal{M}_{\tau_k}(\mathcal{M}^\circ)$. Obviously $\mathcal{M}_k(\mathcal{M}^\circ)$ grow to the set of models $\mathcal{M}(\mathcal{M}^\circ)$ containing \mathcal{M}° . Now one can define the correction step by step: first we fix the correction $q_{1,\mathcal{M}^\circ} = q_{1,\mathcal{M}^\circ}(\mathbf{x})$ for all $\mathcal{M} \in \mathcal{M}_1(\mathcal{M}^\circ)$

$$\mathbb{P}\left(\bigcup_{\mathcal{M} \in \mathcal{M}_1(\mathcal{M}^\circ)} \left\{ \|\boldsymbol{\xi}_{\mathcal{M},\mathcal{M}^\circ}\| > z_{\mathcal{M},\mathcal{M}^\circ}(\mathbf{x} + q_{1,\mathcal{M}^\circ}) \right\}\right) \leq e^{-\mathbf{x}}.$$

With such defined q_{1,\mathcal{M}° , define $q_{2,\mathcal{M}^\circ} = q_{2,\mathcal{M}^\circ}(\mathbf{x})$ such that

$$\mathbb{P}\left(\bigcup_{\mathcal{M} \in \mathcal{M}_2(\mathcal{M}^\circ)} \left\{ \|\boldsymbol{\xi}_{\mathcal{M},\mathcal{M}^\circ}\| > z_{\mathcal{M},\mathcal{M}^\circ}(\mathbf{x} + q_{1,\mathcal{M}^\circ} + q_{2,\mathcal{M}^\circ}) \right\}\right) \leq e^{-\mathbf{x}}.$$

We continue this way and define by induction: if $q_{1,\mathcal{M}^\circ}, \dots, q_{k-1,\mathcal{M}^\circ}$ are fixed then the correction for the set $\mathcal{M}_k(\mathcal{M}^\circ)$ is selected as the sum $\mathbf{x} + q_{1,\mathcal{M}^\circ} + \dots + q_{k,\mathcal{M}^\circ}$

$$\mathbb{P}\left(\bigcup_{\mathcal{M} \in \mathcal{M}_k(\mathcal{M}^\circ)} \left\{ \|\boldsymbol{\xi}_{\mathcal{M},\mathcal{M}^\circ}\| > z_{\mathcal{M},\mathcal{M}^\circ}(\mathbf{x} + q_{1,\mathcal{M}^\circ} + q_{2,\mathcal{M}^\circ} + \dots + q_{k,\mathcal{M}^\circ}) \right\}\right) \leq e^{-\mathbf{x}}. \quad (5.3)$$

For each $\mathcal{M} > \mathcal{M}^\circ$, there exists a unique $k = k(\mathcal{M})$ corresponding to the smallest set $\mathcal{M}_k(\mathcal{M}^\circ)$ containing \mathcal{M} . Finally, we define

$$\mathbf{z}_{\mathcal{M},\mathcal{M}^\circ} = z_{\mathcal{M},\mathcal{M}^\circ}(\mathbf{x} + q_{1,\mathcal{M}^\circ} + q_{2,\mathcal{M}^\circ} + \dots + q_{k,\mathcal{M}^\circ}) + \beta \mathbf{p}_{\mathcal{M},\mathcal{M}^\circ}^{1/2}, \quad k = k(\mathcal{M}). \quad (5.4)$$

Our selection rule chooses the smallest accepted model. It can be written as

$$\hat{\mathcal{M}} = \underset{\mathcal{M}^\circ \in \mathcal{M}}{\operatorname{argmin}} \left\{ |\mathcal{M}^\circ| : \mathbb{T}_{\mathcal{M},\mathcal{M}^\circ} \leq \mathbf{z}_{\mathcal{M},\mathcal{M}^\circ}, \quad \mathcal{M} \in \mathcal{M}(\mathcal{M}^\circ) \right\}. \quad (5.5)$$

If there are many such \mathcal{M}° , one can select arbitrarily among them. Now we define a good choice \mathcal{M}° as previously by “no significant bias” condition:

$$\|\mathbf{b}_{\mathcal{M},\mathcal{M}^\circ}\| \leq \beta \mathbf{p}_{\mathcal{M},\mathcal{M}^\circ}^{1/2} \quad \mathcal{M} \in \mathcal{M}(\mathcal{M}^\circ). \quad (5.6)$$

The construction ensures that a good model will be accepted with a high probability.

Theorem 5.1.1. *Let \varkappa° be a good model in the sense (5.6). Then it holds for the SmA procedure with the critical values $\mathbf{z}_{\varkappa, \varkappa^\circ}$ from (5.3) and (5.4)*

$$\mathbb{P}(\varkappa^\circ \text{ is rejected}) \leq e^{-x}.$$

Now we define the oracle choice \varkappa^* as the smallest (in complexity $|\varkappa^*|$) model under the constraint (5.6):

$$\varkappa^* \stackrel{\text{def}}{=} \operatorname{argmin}_{\varkappa^\circ \in \mathcal{M}} \{|\varkappa^\circ| : \|\mathbf{b}_{\varkappa, \varkappa^\circ}\| \leq \beta \mathbf{p}_{\varkappa, \varkappa^\circ}^{1/2}, \quad \varkappa \in \mathcal{M}(\varkappa^\circ)\}. \quad (5.7)$$

Again, this relation does not uniquely define the \varkappa^* value, if there are many \varkappa^* with this property, any of them can be taken. The oracle bound compares the risk of the oracle estimate $\mathcal{R}_{\varkappa^*}$ with the risk of the adaptive estimate $\hat{\boldsymbol{\theta}} = \tilde{\boldsymbol{\theta}}_{\hat{\varkappa}}$ for the SmA rule $\hat{\varkappa}$.

Theorem 5.1.2. *It holds on a random set of probability at least $1 - e^{-x}$*

$$\sigma^{-1} \|W(\tilde{\boldsymbol{\theta}}_{\varkappa^*} - \tilde{\boldsymbol{\theta}}_{\hat{\varkappa}})\| \leq \bar{\mathbf{z}}_{\varkappa^*}$$

with $k^* = |\varkappa^*|$ and $\bar{\mathbf{z}}_{\varkappa^*}$ defined by

$$\bar{\mathbf{z}}_{\varkappa^*} \stackrel{\text{def}}{=} \max_{|\varkappa| \leq k^*} (\mathbf{z}_{\varkappa \vee \varkappa^*, \varkappa^*} + \mathbf{z}_{\varkappa \vee \varkappa^*, \varkappa}).$$

Proof. The construction and the propagation property ensures that \varkappa^* is accepted with a high probability $1 - e^{-x}$. Below we focus on this case. Then the adaptive choice $\hat{\varkappa}$ from (5.5) has to fulfill

$$|\hat{\varkappa}| \leq |\varkappa^*|.$$

Due to partial ordering, the models \varkappa^* and $\hat{\varkappa}$ are not directly comparable. We use the model $\varkappa = \varkappa^* \vee \hat{\varkappa}$ which contains both and is the smallest one with this property. As \varkappa^* and $\hat{\varkappa}$ are both accepted, it holds

$$\sigma^{-1} \|W(\tilde{\boldsymbol{\theta}}_{\varkappa} - \tilde{\boldsymbol{\theta}}_{\hat{\varkappa}})\| \leq \mathbf{z}_{\varkappa, \hat{\varkappa}}, \quad \sigma^{-1} \|W(\tilde{\boldsymbol{\theta}}_{\varkappa} - \tilde{\boldsymbol{\theta}}_{\varkappa^*})\| \leq \mathbf{z}_{\varkappa, \varkappa^*}.$$

We conclude that

$$\sigma^{-1} \|W(\tilde{\boldsymbol{\theta}}_{\varkappa^*} - \tilde{\boldsymbol{\theta}}_{\hat{\varkappa}})\| \leq \mathbf{z}_{\varkappa, \hat{\varkappa}} + \mathbf{z}_{\varkappa, \varkappa^*}.$$

and the assertion follows.

5.1.2 Prediction loss

The case of prediction loss ($W = \Psi^\top$) in combination with projection estimates $\tilde{\boldsymbol{\theta}}_\varkappa$ allows to reduce the study to the sequence space model with $\Psi = I_p$. This dramatically simplifies the situation. Below we denote by Π_\varkappa the projector in \mathbb{R}^n onto the subspace corresponding to \varkappa . For a couple $\varkappa^\circ < \varkappa$, we also consider the projector $\Pi_{\varkappa, \varkappa^\circ} = \Pi_\varkappa - \Pi_{\varkappa^\circ} = \Pi_{\varkappa - \varkappa^\circ}$.

We also use that $\mathfrak{p}_\varkappa = \varkappa$ and

$$\mathfrak{p}_{\varkappa, \varkappa^\circ} = \mathbb{E} \|\boldsymbol{\xi}_{\varkappa, \varkappa^\circ}\|^2 = |\varkappa - \varkappa^\circ|.$$

Theorem 5.1.3. *For the regression model $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \varepsilon$ with homogeneous errors and for $W = \Psi^\top$, the tail functions $z_{\varkappa, \varkappa^\circ}(\mathbf{x})$ only depends on $|\varkappa - \varkappa^\circ|$. Moreover, for \mathbf{x} fixed, the multiplicity corrections $q_{k, \varkappa^\circ} = q_{k, \varkappa^\circ}(\mathbf{x})$ only depends on $p - |\varkappa^\circ|$ and on τ_k .*

Proof. We use that $\Psi^\top \tilde{\boldsymbol{\theta}}_\varkappa = \Pi_\varkappa \mathbf{Y}$ and

$$\mathbb{T}_{\varkappa, \varkappa^\circ}^2 = \sigma^{-2} \|\Psi^\top (\tilde{\boldsymbol{\theta}}_\varkappa - \tilde{\boldsymbol{\theta}}_{\varkappa^\circ})\|^2 = \sigma^{-2} \|\Pi_{\varkappa, \varkappa^\circ} \mathbf{Y}\|^2 \sim \chi_{|\varkappa - \varkappa^\circ|}^2.$$

To be done: complete the proof

The definition (5.7) of the oracle choice \varkappa^* can be restated as

$$\varkappa^* \stackrel{\text{def}}{=} \underset{\varkappa^\circ \in \mathcal{M}}{\operatorname{argmin}} \{ |\varkappa^\circ| : \|\mathbf{b}_{\varkappa, \varkappa^\circ}\| \leq \beta |\varkappa - \varkappa^\circ|^{1/2}, \quad \varkappa \in \mathcal{M}(\varkappa^\circ) \}. \quad (5.8)$$

Theorem 5.1.4. *The result of Theorem 5.1.2 holds with $\bar{z}_{\varkappa^*} = \bar{z}_{k^*}$ only depending on the cardinality $k^* = |\varkappa^*|$. Any \varkappa^* with this cardinality can be used in definition (5.8).*

To be done: An upper bound on q_{k, \varkappa°

To be done: An upper bound on \bar{z}_{k^*}

To be done: Algorithmic implementation

5.1.3 Estimation loss

The analysis for the problem of estimation loss $W = I_p$ is similar, but some nice features of the prediction loss do not apply here. We use

$$\begin{aligned} \tilde{\boldsymbol{\theta}}_\varkappa &= \mathcal{S}_\varkappa \mathbf{Y}, \\ \mathcal{S}_\varkappa &= (\Psi_\varkappa \Psi_\varkappa^\top)^{-1} \Psi_\varkappa. \end{aligned}$$

Therefore,

$$\mathbb{T}_{\varkappa, \varkappa^\circ} = \sigma^{-1} \|(\mathcal{S}_\varkappa - \mathcal{S}_{\varkappa^\circ})\mathbf{Y}\| = \|\boldsymbol{\xi}_{\varkappa, \varkappa^\circ} + \mathbf{b}_{\varkappa, \varkappa^\circ}\|.$$

If the bias component vanishes, the distribution of this test statistic is completely described by the matrices Ψ_\varkappa and Ψ_{\varkappa° but in a more complicated way than for the prediction loss. Also define

$$\mathbf{p}_{\varkappa, \varkappa^\circ} = \mathbb{E}\|\boldsymbol{\xi}_{\varkappa, \varkappa^\circ}\|^2 = \text{tr}\{(\mathcal{S}_\varkappa - \mathcal{S}_{\varkappa^\circ})(\mathcal{S}_\varkappa - \mathcal{S}_{\varkappa^\circ})^\top\}.$$

The general results of Theorems 5.1.1 and 5.1.2 apply here for such defined values $\mathbf{p}_{\varkappa, \varkappa^\circ}$. Unfortunately, the nice simplification of the formulation as in the prediction case is only possible if the design is orthonormal. Then the estimation and prediction problems coincide.

5.1.4 Linear functional estimation

Now we briefly discuss the case when W is a matrix of rank one which corresponds to estimation of a linear functional. An advantage of this situation is that each difference $W(\tilde{\boldsymbol{\theta}}_{\varkappa'} - \tilde{\boldsymbol{\theta}}_\varkappa)$ is univariate normal. This helps a lot in evaluating the distribution of each test statistic $\mathbb{T}_{\varkappa', \varkappa}$.

As in the ordered case, each estimate $\tilde{\phi}_\varkappa = W\tilde{\boldsymbol{\theta}}_\varkappa$ fulfills

$$\tilde{\phi}_\varkappa - \phi^* = W(\tilde{\boldsymbol{\theta}}_\varkappa - \boldsymbol{\theta}^*) = W\mathcal{S}_\varkappa\boldsymbol{\varepsilon} + W(\mathcal{S}_\varkappa\mathbf{f}^* - \boldsymbol{\theta}^*) = \xi_\varkappa + b_\varkappa,$$

where $\xi_\varkappa = W\mathcal{S}_\varkappa\boldsymbol{\varepsilon}$ is a zero mean random variable and $b_\varkappa = W(\mathcal{S}_\varkappa\mathbf{f}^* - \boldsymbol{\theta}^*)$ is the deterministic bias. The squared risk of $\tilde{\phi}_\varkappa$ is given by the usual bias-variance decomposition:

$$\mathcal{R}_\varkappa = \mathbb{E}(\tilde{\phi}_\varkappa - \phi^*)^2 = \mathbb{E}(\xi_\varkappa + b_\varkappa)^2 = b_\varkappa^2 + \text{Var}(\xi_\varkappa) = b_\varkappa^2 + s_\varkappa^2$$

with

$$s_\varkappa^2 = \sigma^2 W\mathcal{S}_\varkappa\mathcal{S}_\varkappa^\top W^\top.$$

Monotonicity condition $\mathcal{S}_\varkappa\mathcal{S}_\varkappa^\top \geq \mathcal{S}_{\varkappa^\circ}\mathcal{S}_{\varkappa^\circ}^\top$ for $\varkappa > \varkappa^\circ$ yields monotonicity $s_\varkappa \geq s_{\varkappa^\circ}$ for the functional estimate. For each pair $\varkappa^\circ < \varkappa$

$$\tilde{\phi}_\varkappa - \tilde{\phi}_{\varkappa^\circ} = W(\tilde{\boldsymbol{\theta}}_\varkappa - \tilde{\boldsymbol{\theta}}_{\varkappa^\circ}) = W(\mathcal{S}_\varkappa - \mathcal{S}_{\varkappa^\circ})\mathbf{Y} = W\mathcal{S}_{\varkappa, \varkappa^\circ}\mathbf{Y}.$$

The variance of this difference reads as

$$s_{\varkappa, \varkappa^\circ}^2 = \text{Var}(\xi_{\varkappa, \varkappa^\circ}) = \sigma^2 W\mathcal{S}_{\varkappa, \varkappa^\circ}\mathcal{S}_{\varkappa, \varkappa^\circ}^\top W^\top.$$

The scaled test statistic $\mathbb{T}_{\varkappa, \varkappa^\circ}$ is given for $\varkappa > \varkappa^\circ$ by

$$\mathbb{T}_{\mathcal{X}, \mathcal{X}^\circ} = s_{\mathcal{X}, \mathcal{X}^\circ}^{-1} |W\mathcal{S}_{\mathcal{X}, \mathcal{X}^\circ} \mathbf{Y}|.$$

One can use that the stochastic component $\xi_{\mathcal{X}, \mathcal{X}^\circ}$ of $s_{\mathcal{X}, \mathcal{X}^\circ}^{-1} W\mathcal{S}_{\mathcal{X}, \mathcal{X}^\circ} \mathbf{Y}$ is standard normal. Thus, the multiplicity corrections are computed from the same bound (5.3) with $z_1(\cdot)$ in place of $z_{\mathcal{X}, \mathcal{X}^\circ}(\cdot)$. Moreover, $\mathbf{p}_{\mathcal{X}, \mathcal{X}^\circ} \equiv 1$, and the oracle definition reads as

$$\mathcal{X}^* \stackrel{\text{def}}{=} \underset{\mathcal{X}^\circ \in \mathcal{M}}{\operatorname{argmin}} \{ |\mathcal{X}^\circ| : \|\mathbf{b}_{\mathcal{X}, \mathcal{X}^\circ}\| \leq \beta, \quad \mathcal{X} \in \mathcal{M}(\mathcal{X}^\circ) \}.$$

The general results of Theorems 5.1.1 and 5.1.2 apply here without any change. However, the involved values can be made more precise.

To be done: An upper bound on q_{k, \mathcal{X}°

To be done: An upper bound on \bar{z}_{k^*}

To be done: Algorithmic implementation

5.1.5 Subset selection problem

The presented oracle bound of Theorem 5.1.2 claims that the risk of the adaptive estimate $\widehat{\boldsymbol{\theta}}$ is linked to the risk of the oracle $\widetilde{\boldsymbol{\theta}}_{\mathcal{X}^*}$. However, it tells nothing about the selected set $\widehat{\mathcal{X}}$. Now we discuss this issue of choosing the active set of important features represented by non-zero entries of $\boldsymbol{\theta}^*$.

Note that this problem has to be put in a right way: one can suppose that $\boldsymbol{\theta}^*$ is sparse and only significant entries are non-zero. Alternatively, one tries to find an approximating model with another vector $\boldsymbol{\theta}^*$ having a sparse representation and delivering nearly the same approximation and prediction quality. We follow our oracle result and define \mathcal{X}^* by (5.7). This will be our target. We aim to find some sufficient conditions ensuring that $\widehat{\mathcal{X}} \approx \mathcal{X}^*$. We already know that the set \mathcal{X}^* will be accepted with a high probability. This particularly implies that $|\widehat{\mathcal{X}}| \leq |\mathcal{X}^*|$. So, the question under study is the probability of the situation when $\widehat{\mathcal{X}}$ selects some other features instead of those in \mathcal{X}^* .

For simplicity of notation, we consider the sequence space model with $\Psi = I_p$ and also assume $\sigma^2 = 1$. Our first result describes which candidate set \mathcal{X} will be rejected with a high probability. The definition of the oracle \mathcal{X}^* implies that norm of the remaining part $\|\Pi_{\mathcal{X}, \mathcal{X}^*} \boldsymbol{\theta}^*\|$ does not exceed $\beta |\mathcal{X} - \mathcal{X}^*|^{1/2}$ for any $\mathcal{X} \supset \mathcal{X}^*$. Let \mathcal{X} be another subset. We measure its departure from \mathcal{X}^* by the norm of $\Pi_{\mathcal{X}^* \setminus \mathcal{X}} \boldsymbol{\theta}^*$, which is the projection of the true signal onto components which are in \mathcal{X}^* but not in \mathcal{X} . Our result says that if this departure is large, such a candidate will be killed with a high probability. This result can be viewed as an extension of the zone-of-insensitivity result from the ordered case. Below we use the short notation for the tail functions of $\|\boldsymbol{\xi}_{\mathcal{X}, \mathcal{X}^\circ}\|$ with the proper multiplicity correction; cf (5.4):

$$z_{\mathcal{X}, \mathcal{X}^\circ} = z_{\mathcal{X}, \mathcal{X}^\circ}(\mathbf{x} + q_{1, \mathcal{X}^\circ} + q_{2, \mathcal{X}^\circ} + \dots + q_{k, \mathcal{X}^\circ}), \quad \mathcal{X} \in \mathcal{M}_k(\mathcal{X}^\circ).$$

Theorem 5.1.5. *Let \mathcal{X} be such that*

$$\|\mathbf{b}_{\mathcal{X} \vee \mathcal{X}^*}\| = \|\Pi_{\mathcal{X}^* \setminus \mathcal{X}} \boldsymbol{\theta}^*\| > 2z_{\mathcal{X} \vee \mathcal{X}^*} + \beta |\mathcal{X}^* \setminus \mathcal{X}|^{1/2}, \quad (5.9)$$

and let $\mathcal{M}^\circ(\mathcal{X}^*)$ be the collection of all such \mathcal{X} . Then

$$P(\text{any of } \mathcal{X} \in \mathcal{M}^\circ(\mathcal{X}^*) \text{ is accepted}) \leq e^{-\mathbf{x}}.$$

Proof. We just apply the acceptance rule for \mathcal{X} requiring

$$\|\tilde{\boldsymbol{\theta}}_{\mathcal{X} \vee \mathcal{X}^*} - \tilde{\boldsymbol{\theta}}_{\mathcal{X}}\| \leq \mathbf{z}_{\mathcal{X} \vee \mathcal{X}^*} = z_{\mathcal{X} \vee \mathcal{X}^*} + \beta |\mathcal{X}^* \setminus \mathcal{X}|^{1/2}. \quad (5.10)$$

The usual decomposition of $\tilde{\boldsymbol{\theta}}_{\mathcal{X}}$ implies on a dominating set $\Omega(\mathbf{x})$ by (5.3)

$$\begin{aligned} \|\tilde{\boldsymbol{\theta}}_{\mathcal{X} \vee \mathcal{X}^*} - \tilde{\boldsymbol{\theta}}_{\mathcal{X}}\| &\geq \|\mathbf{b}_{\mathcal{X} \vee \mathcal{X}^*}\| - \|\boldsymbol{\xi}_{\mathcal{X} \vee \mathcal{X}^*}\| \\ &\geq \|\mathbf{b}_{\mathcal{X} \vee \mathcal{X}^*}\| - z_{\mathcal{X} \vee \mathcal{X}^*}. \end{aligned} \quad (5.11)$$

It remains to check that the inequalities (5.10) and (5.11) are incompatible under (5.9).

This result can be directly applied to check for a subset \mathcal{X}_0^* of \mathcal{X}^* whether it will completely missed by $\hat{\mathcal{X}}$.

Corollary 5.1.1. *Let \mathcal{X}_0^* fulfill*

$$\|\Pi_{\mathcal{X}_0^*} \boldsymbol{\theta}^*\| \geq \bar{\mathbf{z}}_{\mathcal{X}^*}$$

with

$$\bar{\mathbf{z}}_{\mathcal{X}^*} \stackrel{\text{def}}{=} \max_{|\mathcal{X}| \leq |\mathcal{X}^*|} \{2z_{\mathcal{X} \vee \mathcal{X}^*} + \beta |\mathcal{X}^* \setminus \mathcal{X}|^{1/2}\}.$$

Then

$$P(\mathcal{X}_0^* \cap \hat{\mathcal{X}} = \emptyset) \leq e^{-\mathbf{x}}.$$

In particular, any coefficient θ_j^* of $\boldsymbol{\theta}^*$ with $|\theta_j^*| > \bar{\mathbf{z}}_{\mathcal{X}^*}$ will be included in $\hat{\mathcal{X}}$ with a probability at least $1 - e^{-\mathbf{x}}$.

5.2 Anisotropic models

This section studies the so called anisotropic models when one has a number tuning parameters to be selected, and each of them is ordered. In other words, \varkappa is a vector with two or more components, and we consider the set of the estimates $\tilde{\boldsymbol{\theta}}_{\varkappa}$. When only one component of \varkappa is varying and the other are fixed, the monotonicity assumption is assumed to be fulfilled. However, this only yields a componentwise partial ordering of the set \mathcal{M} of all considered models.

To simplify the presentation, we consider below the two dimensional case and a product structure. An extension to the general case is straightforward.

Let $\varkappa = (\varkappa_1, \varkappa_2)$ with $\varkappa_j \in \mathcal{M}_j$ for $j = 1, 2$ and $\mathcal{M} = \mathcal{M}_1 \times \mathcal{M}_2$. We write $\varkappa = (\varkappa_1, \varkappa_2) \geq \varkappa^\circ = (\varkappa_1^\circ, \varkappa_2^\circ)$ if $\varkappa_1 \geq \varkappa_1^\circ$ and $\varkappa_2 \geq \varkappa_2^\circ$. Assume we are given a collection of linear smoothers $\tilde{\boldsymbol{\theta}}_{\varkappa}$ with a partial ordering: if $\varkappa > \varkappa^\circ$ then

$$\text{Var}(W\tilde{\boldsymbol{\theta}}_{\varkappa}) > \text{Var}(W\tilde{\boldsymbol{\theta}}_{\varkappa^\circ}).$$

This particularly implies

$$\mathbf{p}_{\varkappa} \stackrel{\text{def}}{=} \text{tr}\{\text{Var}(W\tilde{\boldsymbol{\theta}}_{\varkappa})\} > \mathbf{p}_{\varkappa^\circ} \stackrel{\text{def}}{=} \text{tr}\{\text{Var}(W\tilde{\boldsymbol{\theta}}_{\varkappa^\circ})\}.$$

We aim at applying the SmA method to this special situation. The setup and notation of the previous section are kept. We focus on a linear regression model $\mathbf{Y} = \boldsymbol{\Psi}^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$ with homogeneous Gaussian error $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$ and consider a collection of pairwise test statistics

$$\mathbb{T}_{\varkappa, \varkappa^\circ} = \sigma^{-1} \|W(\tilde{\boldsymbol{\theta}}_{\varkappa} - \tilde{\boldsymbol{\theta}}_{\varkappa^\circ})\| = \|\boldsymbol{\xi}_{\varkappa, \varkappa^\circ} + \mathbf{b}_{\varkappa, \varkappa^\circ}\|.$$

As previously, define

$$\mathbf{p}_{\varkappa, \varkappa^\circ} = \text{tr}\{\text{Var}(\boldsymbol{\xi}_{\varkappa, \varkappa^\circ})\}.$$

The acceptance rule can be written as

$$\varkappa^\circ \text{ is accepted iff } \mathbb{T}_{\varkappa, \varkappa^\circ} \leq \mathbf{z}_{\varkappa, \varkappa^\circ} \quad \forall \varkappa \in \mathcal{M}(\varkappa^\circ), \quad (5.12)$$

where $\mathcal{M}(\varkappa^\circ) = \{\varkappa \in \mathcal{M}: \varkappa > \varkappa^\circ\}$. In words, \varkappa° is accepted if it is competitive with any larger model \varkappa . Now we discuss how the critical values $\mathbf{z}_{\varkappa, \varkappa^\circ}$ can be fixed by the multiplicity correction of the individual tail functions. Here we only discuss a uniform correction, however, it can be easily done in a multilevel form. Define $q_{\varkappa^\circ}(\mathbf{x})$ by the condition

$$\mathbb{P}\left(\bigcup_{\mathcal{M}(\mathcal{M}^\circ)} \left\{ \|\boldsymbol{\xi}_{\mathcal{M}, \mathcal{M}^\circ}\| > z_{\mathcal{M}, \mathcal{M}^\circ}(\mathbf{x} + q_{\mathcal{M}^\circ}(\mathbf{x})) \right\}\right) \leq e^{-x}. \quad (5.13)$$

Denote $\mathbf{x}_{\mathcal{M}^\circ} = \mathbf{x} + q_{\mathcal{M}^\circ}(\mathbf{x})$ and apply the acceptance rule (5.2) with $\mathbf{z}_{\mathcal{M}, \mathcal{M}^\circ}$ equal to such defined $z_{\mathcal{M}, \mathcal{M}^\circ}(\mathbf{x}_{\mathcal{M}^\circ})$ after a small bias correction:

$$\mathbf{z}_{\mathcal{M}, \mathcal{M}^\circ} \stackrel{\text{def}}{=} z_{\mathcal{M}, \mathcal{M}^\circ}(\mathbf{x}_{\mathcal{M}^\circ}) + \beta \mathbf{p}_{\mathcal{M}, \mathcal{M}^\circ}^{1/2}. \quad (5.14)$$

A good choice \mathcal{M}° can be defined as previously by “no significant bias” condition:

$$\|\mathbf{b}_{\mathcal{M}, \mathcal{M}^\circ}\| \leq \beta \mathbf{p}_{\mathcal{M}, \mathcal{M}^\circ}^{1/2} \quad \mathcal{M} \in \mathcal{M}(\mathcal{M}^\circ). \quad (5.15)$$

The construction ensures that a good model will be accepted with a high probability.

Theorem 5.2.1. *Let \mathcal{M}° be a good model in the sense (5.15). Then it holds for the acceptance rule (5.12) with the critical values $\mathbf{z}_{\mathcal{M}, \mathcal{M}^\circ}$ from (5.13) and (5.14)*

$$\mathbb{P}(\mathcal{M}^\circ \text{ is rejected}) \leq e^{-x}.$$

So, the construction allows to figure out a set of good models, each of them will be kept by the procedure with a high probability. It remains to introduce a natural ordering on the set of such good models. This can be done by the value $\mathbf{p}_{\mathcal{M}^\circ}$ which is proportional to $|\mathcal{M}^\circ|$ for the sequence space model. Define the oracle choice \mathcal{M}^* as the smallest (in complexity $\mathbf{p}_{\mathcal{M}^\circ}$) model under the constraint (5.15):

$$\mathcal{M}^* \stackrel{\text{def}}{=} \underset{\mathcal{M}^\circ \in \mathcal{M}}{\operatorname{argmin}} \{ \mathbf{p}_{\mathcal{M}^\circ} : \|\mathbf{b}_{\mathcal{M}, \mathcal{M}^\circ}\| \leq \beta \mathbf{p}_{\mathcal{M}, \mathcal{M}^\circ}^{1/2}, \quad \mathcal{M} \in \mathcal{M}(\mathcal{M}^\circ) \}.$$

Our selection rule chooses the smallest accepted model. It can be written as

$$\widehat{\mathcal{M}} = \underset{\mathcal{M}^\circ \in \mathcal{M}}{\operatorname{argmin}} \{ \mathbf{p}_{\mathcal{M}^\circ} : \mathbb{T}_{\mathcal{M}, \mathcal{M}^\circ} \leq \mathbf{z}_{\mathcal{M}, \mathcal{M}^\circ}, \quad \mathcal{M} \in \mathcal{M}(\mathcal{M}^\circ) \}. \quad (5.16)$$

The oracle bound compares the risk of the oracle estimate $\mathcal{R}_{\mathcal{M}^*}$ with the risk of the adaptive estimate $\widehat{\boldsymbol{\theta}} = \widetilde{\boldsymbol{\theta}}_{\widehat{\mathcal{M}}}$ for the SmA rule $\widehat{\mathcal{M}}$. Below for two given models \mathcal{M} and \mathcal{M}° , we denote by $\mathcal{M} \vee \mathcal{M}^\circ$ the smallest model which is larger than each:

$$\mathcal{M} \vee \mathcal{M}^\circ \stackrel{\text{def}}{=} (\mathcal{M}_1 \vee \mathcal{M}_1^\circ, \mathcal{M}_2 \vee \mathcal{M}_2^\circ).$$

Theorem 5.2.2. *It holds on a random set of probability at least $1 - e^{-x}$*

$$\sigma^{-1} \|W(\widetilde{\boldsymbol{\theta}}_{\mathcal{M}^*} - \widetilde{\boldsymbol{\theta}}_{\widehat{\mathcal{M}}})\| \leq \bar{\mathbf{z}}_{\mathcal{M}^*}$$

where $\bar{\mathbf{z}}_{\mathcal{M}^*}$ is defined by

$$\bar{\mathbf{z}}_{\mathcal{M}^*} \stackrel{\text{def}}{=} \max_{\mathbf{p}_{\mathcal{M}} \leq \mathbf{p}_{\mathcal{M}^*}} (\mathbf{z}_{\mathcal{M} \vee \mathcal{M}^*, \mathcal{M}^*} + \mathbf{z}_{\mathcal{M} \vee \mathcal{M}^*, \mathcal{M}}).$$

Proof. The construction and the propagation property ensures that \varkappa^* is accepted with a high probability $1 - e^{-x}$. Below we focus on this case. Then the adaptive choice $\widehat{\varkappa}$ from (5.16) has to fulfill

$$p_{\widehat{\varkappa}} \leq p_{\varkappa^*}.$$

Consider $\check{\varkappa} = \varkappa^* \vee \widehat{\varkappa}$ which contains both $\widehat{\varkappa}$ and \varkappa^* . As \varkappa^* and $\widehat{\varkappa}$ are both accepted, it holds

$$\sigma^{-1} \|W(\tilde{\theta}_{\check{\varkappa}} - \tilde{\theta}_{\widehat{\varkappa}})\| \leq z_{\check{\varkappa}, \widehat{\varkappa}}, \quad \sigma^{-1} \|W(\tilde{\theta}_{\check{\varkappa}} - \tilde{\theta}_{\varkappa^*})\| \leq z_{\check{\varkappa}, \varkappa^*}.$$

Therefore,

$$\sigma^{-1} \|W(\tilde{\theta}_{\varkappa^*} - \tilde{\theta}_{\widehat{\varkappa}})\| \leq z_{\check{\varkappa}, \widehat{\varkappa}} + z_{\check{\varkappa}, \varkappa^*}.$$

and the assertion follows.

The value \bar{z}_{\varkappa^*} is the ‘‘payment for adaptation’’, and it can be quite large relative to the oracle standard deviation $p_{\varkappa^*}^{1/2}$. The worst case is given by the anisotropic situation with the oracle of the form like $\varkappa^* = (\varkappa_1^*, \varkappa_{2,\max})$. In this case of a competitive model $\varkappa = (\varkappa_{1,\max}^*, \varkappa_2)$, the maximum of \varkappa^* and \varkappa is the largest possible model

$$\varkappa^* \vee \varkappa = (\varkappa_{1,\max}, \varkappa_{2,\max})$$

and the corresponding critical value $z_{\varkappa \vee \varkappa^*, \varkappa^*}$ is very large.

To be done: In the isotropic case with $\varkappa_1^* \asymp \varkappa_2^*$, this problem disappears.

To be done: One can refine the result by considering the zone of insensitivity $\mathcal{M}^\circ(\varkappa^*)$.

Fisher and Wilks expansion

This chapter presents two prominent results of classical parametric statistics, namely the Fisher and Wilks Theorems, in a non-classical framework. The main features to be addressed here are a finite sample setup with large parameter dimension and a possible model misspecification.

First we specify our set-up. Let \mathbf{Y} denote the observed data and \mathbb{P} mean their distribution. A general parametric assumption (PA) means that \mathbb{P} belongs to p -dimensional family $(\mathbb{P}_\theta, \theta \in \Theta \subseteq \mathbb{R}^p)$ dominated by a measure μ_0 . This family yields the log-likelihood function $L(\theta) = L(\mathbf{Y}, \theta) \stackrel{\text{def}}{=} \log \frac{d\mathbb{P}_\theta}{d\mu_0}(\mathbf{Y})$. The PA can be misspecified, so, in general, $L(\theta)$ is a *quasi log-likelihood*. The classical likelihood principle suggests to estimate θ by maximizing the function $L(\theta)$:

$$\tilde{\theta} \stackrel{\text{def}}{=} \operatorname{argmax}_{\theta \in \Theta} L(\theta). \quad (6.1)$$

If $\mathbb{P} \notin (\mathbb{P}_\theta)$, then the (quasi) MLE estimate $\tilde{\theta}$ from (6.1) is still meaningful and it appears to be an estimate of the value θ^* defined by maximizing the expected value of $L(\theta)$:

$$\theta^* \stackrel{\text{def}}{=} \operatorname{argmax}_{\theta \in \Theta} \mathbb{E}L(\theta).$$

Here θ^* is the true value in the parametric situation and can be viewed as the parameter of the best parametric fit in the general case. The study is non-asymptotic, that is, we proceed with only one sample \mathbf{Y} . One can easily extend it to an asymptotic setup in which the data, its distribution, the parameter space and the parametric family depend on the asymptotic parameter like the sample size. One example is given below in Section 6.4 for the case of an i.i.d. sample.

The Fisher expansion of the qMLE $\tilde{\theta}$ is given as follows:

$$D(\tilde{\theta} - \theta^*) \approx \xi \stackrel{\text{def}}{=} D^{-1} \nabla L(\theta^*),$$

where $\nabla L(\boldsymbol{\theta}) = \frac{dL}{d\boldsymbol{\theta}}(\boldsymbol{\theta})$ and $D^2 \stackrel{\text{def}}{=} -\nabla^2 \mathbb{E}L(\boldsymbol{\theta}^*)$ is the analog of the total Fisher information matrix. In classical situations, the standardized score $\boldsymbol{\xi}$ is asymptotically standard normal yielding asymptotic root-n normality and efficiency of the MLE $\tilde{\boldsymbol{\theta}}$. Theorem 6.3.2 carefully describes how the error of this expansion depends on the parameter dimension p and the regularity of the model. The Wilks expansion means

$$L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*) \approx \|\boldsymbol{\xi}\|^2/2.$$

Again, if the vector $\boldsymbol{\xi}$ is asymptotically standard normal, the expansion yields the classical χ_p^2 asymptotic distribution for the excess $L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*)$.

The whole study is nonasymptotic and all “small” terms are carefully described. This helps to understand how the parameter dimension is involved and particularly to address the question of a *critical dimension*; see Section 6.4 which specifies the result to the i.i.d. case with n observations and links the obtained results to the classical literature.

6.1 Main results

This section presents our main results which include the Fisher and Wilks expansions for a non-classical and non-asymptotic framework. First we present the frequentist results: concentration and large deviation properties of the maximum likelihood estimator $\tilde{\boldsymbol{\theta}}$, the Fisher expansion for the difference $\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}$ and the Wilks expansion for the excess $L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta})$. The results are stated in a concise way, all the terms are given explicitly. Surprisingly, the leading terms in all bounds are sharp, in particular, the classical results on asymptotic efficiency can be easily derived from the obtained expansions.

Introduce the notation $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^*)$ for the (quasi) log-likelihood ratio. The main step in the approach is the following *uniform local bracketing result*:

$$\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \Delta \leq L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \leq \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \Delta, \quad \boldsymbol{\theta} \in \Theta_0. \quad (6.2)$$

Here $\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is a quadratic in $\boldsymbol{\theta} - \boldsymbol{\theta}^*$ expression, Δ is a small error only depending on Θ_0 which is a local vicinity of the central point $\boldsymbol{\theta}^*$. This result can be viewed as an extension of the famous Le Cam *local asymptotic normality* (LAN) condition. The LAN condition postulates an approximation of the log-likelihood $L(\boldsymbol{\theta})$ by a nearly Gaussian process; see e.g. Ibragimov and Khas'minskij (1981) or Kleijn and van der Vaart (2012) for an extension of this condition (stochastic LAN). The bracketing bound (6.2) requires only some general conditions listed in Section 6.2. A model misspecification case is included. Similarly to the LAN theory, the bracketing result has a number of remarkable corollaries like the Wilks and Fisher Theorems; see Theorems 6.3.2 and 6.3.3.

For making a precise statement, we have to specify the ingredients of the bracketing device. The most important one is a symmetric positive $p \times p$ -matrix D^2 . In typical situations, it can be defined as the negative Hessian of the expected log-likelihood: $D^2 = -\nabla^2 \mathbb{E}L(\boldsymbol{\theta}^*)$. Also one has to specify a radius \mathbf{r}_0 entering in the definition of the local vicinity $\Theta_0(\mathbf{r}_0)$ of the central point $\boldsymbol{\theta}^*$: $\Theta_0(\mathbf{r}_0) = \{\boldsymbol{\theta} : \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}_0\}$. The bracketing result (6.2) can be stated for $\Theta_0 = \Theta_0(\mathbf{r}_0)$ with

$$\begin{aligned} \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) &\stackrel{\text{def}}{=} (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla L(\boldsymbol{\theta}^*) - \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2 \\ &= \boldsymbol{\xi}^\top D(\boldsymbol{\theta} - \boldsymbol{\theta}^*) - \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2 \end{aligned}$$

and

$$\boldsymbol{\xi} \stackrel{\text{def}}{=} D^{-1} \nabla L(\boldsymbol{\theta}^*). \quad (6.3)$$

The construction is essentially changed relative to Spokoiny (2012) (in fact, it is simplified) by using only one matrix D^2 while Spokoiny (2012) used three matrices D_ϵ^2 , $D_{\underline{\epsilon}}^2$, and V^2 . The bracketing bound (6.2) becomes useful if the error Δ is relatively small and can be neglected.

6.2 Conditions

This section collects the conditions which are systematically used in the text. The conditions are quite general and seem to be non-restrictive; see the discussion at the end of the section. We mainly require some regularity and smoothness of the log-likelihood process $L(\boldsymbol{\theta})$. With $D^2 = -\nabla^2 \mathbb{E}L(\boldsymbol{\theta}^*)$, define the local elliptic sets $\Theta_0(\mathbf{r})$ as

$$\Theta_0(\mathbf{r}) \stackrel{\text{def}}{=} \{\boldsymbol{\theta} \in \Theta : \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}\}.$$

We distinguish between local and global conditions. Local ones are stated on $\Theta_0(\mathbf{r}_0)$, while the global one corresponds to $\mathbf{r} \geq \mathbf{r}_0$, where the value \mathbf{r}_0 will be specified later.

The first condition requires that the expected log-likelihood $\mathbb{E}L(\boldsymbol{\theta})$ is twice continuously differentiable.

(\mathcal{L}_0) For each $\mathbf{r} \leq \mathbf{r}_0$, there is a constant $\delta(\mathbf{r}) \leq 1/2$ such that it holds for any $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$ and $D^2(\boldsymbol{\theta}) \stackrel{\text{def}}{=} -\nabla^2 \mathbb{E}L(\boldsymbol{\theta})$:

$$\|D^{-1}D^2(\boldsymbol{\theta})D^{-1} - I_p\|_\infty \leq \delta(\mathbf{r}).$$

Under (\mathcal{L}_0), it follows from the second order Taylor expansion at $\boldsymbol{\theta}^*$:

$$|-2\mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2| \leq \delta(\mathbf{r})\|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2, \quad \boldsymbol{\theta} \in \Theta_0(\mathbf{r}).$$

For $\mathbf{r} > \mathbf{r}_0$, we need a global identification property which ensures that the deterministic component $\mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ of the log-likelihood is competitive with the variation of the stochastic component.

(\mathcal{L}) *There exists $\mathbf{b}(\mathbf{r}) > 0$ such that $\mathbf{r}\mathbf{b}(\mathbf{r})$ is non-decreasing and*

$$\frac{-2\mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)}{\|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2} \geq \mathbf{b}(\mathbf{r}), \quad \forall \mathbf{r} \geq \mathbf{r}_0, \boldsymbol{\theta} \in \Theta_0(\mathbf{r}).$$

Now we consider the stochastic component of the process $L(\boldsymbol{\theta})$:

$$\zeta(\boldsymbol{\theta}) \stackrel{\text{def}}{=} L(\boldsymbol{\theta}) - \mathbb{E}L(\boldsymbol{\theta}).$$

We assume that it is twice differentiable and denote by $\nabla\zeta(\boldsymbol{\theta})$ its gradient and by $\nabla^2\zeta(\boldsymbol{\theta})$ its Hessian matrix.

(\mathbf{ED}_0) *There exist a positive symmetric matrix V^2 , and constants $\mathbf{g} > 0$, $\nu_0 \geq 1$ such that $\text{Var}\{\nabla\zeta(\boldsymbol{\theta}^*)\} \leq V^2$ and*

$$\sup_{\boldsymbol{\gamma} \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{\gamma}^\top \nabla\zeta(\boldsymbol{\theta}^*)}{\|V\boldsymbol{\gamma}\|} \right\} \leq \nu_0^2 \lambda^2 / 2, \quad |\lambda| \leq \mathbf{g}.$$

(\mathbf{ED}_2) *There exist a value $\omega > 0$ and for each $\mathbf{r} > 0$, a constant $\mathbf{g}(\mathbf{r}) > 0$ such that it holds for any $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$:*

$$\sup_{\boldsymbol{\gamma}_1, \boldsymbol{\gamma}_2 \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega} \frac{\boldsymbol{\gamma}_1^\top \nabla^2\zeta(\boldsymbol{\theta}) \boldsymbol{\gamma}_2}{\|D\boldsymbol{\gamma}_1\| \cdot \|D\boldsymbol{\gamma}_2\|} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq \mathbf{g}(\mathbf{r}).$$

Below we only need that the constant $\mathbf{g}(\mathbf{r})$ is larger than $\mathbf{C}p$ for a fixed constant \mathbf{C} and all \mathbf{r} .

The *identifiability condition* relates the matrices V^2 and D^2 .

(\mathcal{I}) *There is a constant $\mathbf{a} > 0$ such that*

$$\mathbf{a}^2 D^2 \geq V^2.$$

Remark 6.2.1. The conditions involve some constants. We distinguish between important constants and technical ones. The impact of the important constants is shown in our results, the list includes $\delta(\mathbf{r})$, ω , and \mathbf{a} . The constant \mathbf{a} can be viewed as the largest eigenvalue of $B = D^{-1}V^2D^{-1}$ and it enters in the definition of the upper quantile function $q(B, \mathbf{x})$ for $\|\boldsymbol{\xi}\|$; see Proposition 6.5.3 below. The other constants like ν_0 or $\mathbf{g}(\mathbf{r})$ are technical. The constant ν_0 is introduced for convenience only, it can be omitted by rescaling the matrix V . In the asymptotic setup it can usually be selected very close to one.

Remark 6.2.2. We briefly comment how restrictive the imposed conditions are. Spokoiny (2012), Section 5.1, considered in details the i.i.d. case and presented some mild sufficient conditions on the parametric family which imply the above general conditions. Condition (ED_0) requires some exponential moments of the observations (errors). Usually one only assumes some finite moments of the errors; cf. Ibragimov and Khas'minskij (1981), Chapter 2. Our condition is a bit more restrictive but it allows to obtain some finite sample bounds. Condition (\mathcal{L}_0) only requires some regularity of the considered parametric family and is not restrictive. Conditions (ED_2) with $\mathbf{g}(\mathbf{r}) \equiv \mathbf{g} > 0$ and (\mathcal{L}) with $\mathbf{b}(\mathbf{r}) \equiv \mathbf{b} > 0$ are easy to verify if the parameter set Θ is compact and the sample size n exceeds $\mathcal{C}p$ for a fixed constant \mathcal{C} . It suffices to check a usual identifiability condition that the value $\mathcal{I}EL(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ does not vanish for $\boldsymbol{\theta} \neq \boldsymbol{\theta}^*$.

The regression and generalized regression models are included as well; cf. Ghosal (1999, 2000) or Kim (2006). Spokoiny (2012), Section 5.2, argued that (ED_2) is automatically fulfilled for a generalized linear model, while (ED_0) requires that regression errors have to fulfill some exponential moments conditions. If this condition is too restrictive and a more stable (robust) estimation procedure is desirable, one can apply the LAD-type contrast leading to median regression. Spokoiny (2012), Section 5.3, showed for the case of linear median regression that all the required conditions are fulfilled automatically if the sample size n exceeds $\mathcal{C}p$ for a fixed constant \mathcal{C} . Spokoiny et al. (2013) applied this approach for local polynomial quantile regression. Zaitsev et al. (2013) applied the approach to the problem of regression with Gaussian process where the unknown parameters enter in the likelihood in a rather complicated way.

6.3 Properties of the MLE $\tilde{\theta}$

This section collects the main results about the MLE $\tilde{\theta}$. We begin by a large deviation bound which ensures a small probability of the event $\tilde{\theta} \notin \Theta_0(\mathbf{r}_0)$. Then we present the Fisher and Wilks expansions. The formulation involves two growing functions of the argument \mathbf{x} : $q(B, \mathbf{x})$ and $q_{\mathbb{H}}(\mathbf{x})$. The functions are given analytically and only depend on the parameters of the model. The function $q(B, \mathbf{x})$ describes the quantiles of the norm of the vector $\|\boldsymbol{\xi}\|$ from (6.3). The definition is given in (6.25). Further, the function $q_{\mathbb{H}}(\mathbf{x})$ is related to the entropy of the parameter space and it is given by (6.13). In typical situations one can use the upper bound $q^2(\mathbf{x}) \leq \mathcal{C}(p + \mathbf{x})$ for both functions. The first result explains the choice of \mathbf{r}_0 ensuring with a high probability that $\tilde{\theta} \in \Theta_0(\mathbf{r}_0)$.

Theorem 6.3.1. *Suppose (ED_0) and (ED_2) , (\mathcal{L}_0) , (\mathcal{L}) , and (\mathcal{I}) . Let also the function $\mathbf{b}(\mathbf{r})$ in (\mathcal{L}) satisfy*

$$\mathbf{b}(\mathbf{r}) \mathbf{r} \geq 2q(B, \mathbf{x}) + 2\varrho(\mathbf{r}, \mathbf{x}), \quad \mathbf{r} > \mathbf{r}_0, \quad (6.4)$$

where

$$\varrho(\mathbf{r}, \mathbf{x}) \stackrel{\text{def}}{=} 6\nu_0 q_{\mathbb{H}}(\mathbf{x} + \log(2\mathbf{r}/\mathbf{r}_0)) \omega. \quad (6.5)$$

Then

$$P(\tilde{\boldsymbol{\theta}} \notin \Theta_0(\mathbf{r}_0)) \leq 3e^{-x}.$$

Remark 6.3.1. The radius \mathbf{r}_0 has to fulfill (6.4). This condition is easy to check in typical situations. One can use that $\mathbf{b}(\mathbf{r}_0) \geq 1 - \delta(\mathbf{r}_0) \approx 1$, that the constant ω is small, and $\mathbf{r}\mathbf{b}(\mathbf{r})$ grows with \mathbf{r} . A simple rule $\mathbf{r}_0 \geq (2 + \delta)q(B, \mathbf{x})$ for some $\delta > 0$ works in most of cases.

Now we state the result about the Fisher expansion for the qMLE $\tilde{\boldsymbol{\theta}}$.

Theorem 6.3.2. *Suppose the conditions of Theorem 6.3.1. On a random set $\Omega(\mathbf{x})$ of a dominating probability at least $1 - 4e^{-x}$, it holds*

$$\|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \boldsymbol{\xi}\| \leq \diamond(\mathbf{r}_0, \mathbf{x}), \quad (6.6)$$

where for the function $q_{\mathbb{H}}(\mathbf{x})$ given by (6.13), the value $\diamond(\mathbf{r}, \mathbf{x})$ is given by

$$\diamond(\mathbf{r}, \mathbf{x}) \stackrel{\text{def}}{=} \{\delta(\mathbf{r}) + 6\nu_0 q_{\mathbb{H}}(\mathbf{x}) \omega\} \mathbf{r}. \quad (6.7)$$

Our version of the Wilks result can be stated in the following form.

Theorem 6.3.3. *Suppose the conditions of Theorem 6.3.1. On a random set $\Omega(\mathbf{x})$ of a dominating probability at least $1 - 5e^{-x}$, it holds with $\diamond(\mathbf{r}_0, \mathbf{x})$ from (6.7)*

$$\begin{aligned} |2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2| &\leq 2\mathbf{r}_0 \diamond(\mathbf{r}_0, \mathbf{x}) + \diamond^2(\mathbf{r}_0, \mathbf{x}), \\ |2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \|\boldsymbol{\xi}\|^2| &\leq 2\mathbf{r}_0 \diamond(\mathbf{r}_0, \mathbf{x}) + \diamond^2(\mathbf{r}_0, \mathbf{x}). \end{aligned} \quad (6.8)$$

Furthermore,

$$\begin{aligned} \left| \sqrt{2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)} - \|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \right| &\leq 2\diamond(\mathbf{r}_0, \mathbf{x}), \\ \left| \sqrt{2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)} - \|\boldsymbol{\xi}\| \right| &\leq 3\diamond(\mathbf{r}_0, \mathbf{x}), \end{aligned} \quad (6.9)$$

and for any $\boldsymbol{\theta}^\circ \in \Theta_0(\mathbf{r}_0)$, it holds on $\Omega(\mathbf{x})$

$$\begin{aligned} \left| \{2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^\circ)\}^{1/2} - \|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^\circ)\| \right| &\leq 4\diamond(\mathbf{r}_0, \mathbf{x}), \\ \left| \{2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^\circ)\}^{1/2} - \|\boldsymbol{\xi} + D(\boldsymbol{\theta}^* - \boldsymbol{\theta}^\circ)\| \right| &\leq 5\diamond(\mathbf{r}_0, \mathbf{x}). \end{aligned}$$

Remark 6.3.2. The classical Fisher and Wilks results describe asymptotic behavior of the MLE $\tilde{\boldsymbol{\theta}}$ and of the excess $L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)$. The whole derivations are based on expansions similar to (6.6) and (6.8) and on the limiting behavior of the vector $\boldsymbol{\xi}$ and its squared norm. Under standard assumptions in the regression or i.i.d. setup the vector $\boldsymbol{\xi}$ is standard normal and $\|\boldsymbol{\xi}\|^2$ is asymptotically χ^2 with p degrees of freedom. The asymptotic distribution of the MLE $\tilde{\boldsymbol{\theta}}$ or of the excess $L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)$ can be used for building the confidence sets or for test critical values. However, the use of asymptotic arguments is limited and faces serious problems in practical applications.

This especially concerns the likelihood ratio statistic $L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)$. It is well recognized that the accuracy of χ^2 -approximation of the tails of $L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)$ is very poor and a reasonable quality requires a huge sample size. If the parameter dimension grows with n this problem becomes even more crucial. The qualitative tail behavior of $\|\boldsymbol{\xi}\|^2$ is described in Proposition 6.5.3 but the upper bounds given there appear to be too conservative for practical use.

Remark 6.3.3. Another issue is a possible model misspecification. The expansions (6.8) and (6.9) apply, even if $L(\boldsymbol{\theta})$ is a quasi-log-likelihood function. However, the covariance matrix $V^2 = \text{Var}\{\nabla L(\boldsymbol{\theta}^*)\}$ of the score does not necessarily coincide with the information matrix D^2 . Then the covariance matrix of the vector $\boldsymbol{\xi}$ follows the famous “sandwich” formula $\text{Var}(\boldsymbol{\xi}) = D^{-1}V^2D^{-1}$, and the distribution of the squared norm $\|\boldsymbol{\xi}\|^2$ depends on the unknown covariance matrix V^2 .

The results presented above focus on the expansions of the MLE $\tilde{\boldsymbol{\theta}}$ and on the excess $L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)$. Numerical results (not presented here) indicate that the accuracy of the expansions (6.6) and (6.8) is very reasonable even for moderate and small samples and it is stable w.r.t. possible model misspecifications. It seems that such expansions are of independent interest and can be used for many further statistical tasks.

6.4 Critical dimension. Examples

This section discusses the issue of the critical dimension of the parameter space Θ . This particularly allows to consider the case of a growing parameter dimension. It appears that the error of expansions is different for different types of problems under consideration. The large deviation result of Theorem 6.3.1 for the MLE $\tilde{\boldsymbol{\theta}}$ requires (6.4). In particular, the value \mathbf{r}_0^2 should be at least $\mathbf{C}(p + \mathbf{x})$. This result yields consistency of the MLE if the neighborhood $\Theta_0(\mathbf{r}_0)$ is small, or, equivalently, $(p + \mathbf{x})^{1/2}\|D^{-1}\|_\infty$ is small. The Fisher expansion (6.6) of Theorem 6.3.2 requires that the error term $\diamond(\mathbf{r}_0, \mathbf{x})$ from (6.5) is small. Finally, the Wilks expansion holds if the error term $\Delta(\mathbf{r}_0, \mathbf{x})$ is small. Now we specify the results for some popular statistical models.

6.4.1 Linear and generalized linear models

In the case of a linear model $\mathbf{Y} = \Psi^\top \boldsymbol{\theta} + \boldsymbol{\varepsilon}$ with a given design $p \times n$ matrix Ψ under the assumption of Gaussian noise $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \Sigma)$, the standard calculus leads to the log-likelihood

$$L(\boldsymbol{\theta}) = -\frac{1}{2}(\mathbf{Y} - \Psi^\top \boldsymbol{\theta})^\top \Sigma^{-1}(\mathbf{Y} - \Psi^\top \boldsymbol{\theta}) + R,$$

where the remainder R does not depend on $\boldsymbol{\theta}$. Moreover, $L(\boldsymbol{\theta})$ is quadratic in $\boldsymbol{\theta}$ and its Hessian is constant: $\nabla^2 \mathbb{E}L(\boldsymbol{\theta}) = -\Psi \Sigma^{-1} \Psi^\top$. One can summarize as follows: with $\mathbb{E}\mathbf{Y} = \mathbf{f}$

$$\begin{aligned} D^2 &\stackrel{\text{def}}{=} -\nabla^2 \mathbb{E}L(\boldsymbol{\theta}) = \Psi \Sigma^{-1} \Psi^\top, & \tilde{\boldsymbol{\theta}} &= D^{-2} \Psi \Sigma^{-1} \mathbf{Y}, \\ \boldsymbol{\xi} &\stackrel{\text{def}}{=} D^{-1} \nabla L(\boldsymbol{\theta}^*) = D^{-1} \Psi \Sigma^{-1} (\mathbf{Y} - \mathbf{f}). & \boldsymbol{\theta}^* &= D^{-2} \Psi \Sigma^{-1} \mathbf{f}, \end{aligned}$$

Moreover,

$$D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \equiv \boldsymbol{\xi}, \quad L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*) \equiv \|\boldsymbol{\xi}\|^2/2.$$

All these results are straightforward, the last one is obtained by the Taylor expansion of the second order around $\tilde{\boldsymbol{\theta}}$ with the use of $\nabla L(\tilde{\boldsymbol{\theta}}) = 0$:

$$L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*) = -\frac{1}{2}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)^\top \nabla^2 L(\tilde{\boldsymbol{\theta}})(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) = -\frac{1}{2} \|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2 = -\frac{1}{2} \|\boldsymbol{\xi}\|^2.$$

The presented derivations mean that the Fisher and Wilks expansions are *identities*, they apply for *any sample size* without *any conditions*, and are only based on *quadraticity* of the likelihood function $L(\boldsymbol{\theta})$ in $\boldsymbol{\theta}$. The *true distribution* of \mathbf{Y} can be whatever and it is not involved at all. There is *no dimensional restrictions*. However, for *inference*, the parametric assumption is important. It only concerns the *distribution of $\boldsymbol{\xi}$* . Let $\text{Var}(\mathbf{Y}) = \Sigma_0 \neq \Sigma$. Then with $D^2 = \Psi \Sigma^{-1} \Psi^\top$

$$\text{Var}\{\nabla L(\boldsymbol{\theta}^*)\} = \text{Var}\{\Psi \Sigma^{-1} \mathbf{Y}\} = \Psi \Sigma^{-1} \Sigma_0 \Sigma^{-1} \Psi^\top \stackrel{\text{def}}{=} V^2 \neq D^2$$

which leads to the famous *sandwich formula*

$$\text{Var}(\boldsymbol{\xi}) = \text{Var}\{D^{-1} \nabla L(\boldsymbol{\theta}^*)\} = D^{-1} V^2 D^{-1} \neq I_p.$$

6.4.2 Generalized linear models (GLM)

Let $\mathbf{Y} = (Y_1, \dots, Y_n)^\top \sim \mathcal{I}P$ be a sample of independent r.v.s. The parametric GLM model is given by $Y_i \sim P_{\Psi_i^\top \boldsymbol{\theta}} \in (P_{\mathbf{v}})$, where Ψ_i are given factors in \mathbb{R}^p , $\boldsymbol{\theta} \in \mathbb{R}^p$

is the unknown parameter in \mathbb{R}^p , and $(P_{\mathbf{v}})$ is an exponential family with canonical parametrization yielding the log-density $\ell(y, \mathbf{v}) = y\mathbf{v} - d(\mathbf{v})$ for a convex function $d(\mathbf{v})$. The MLE $\tilde{\boldsymbol{\theta}}$ and the target $\boldsymbol{\theta}^*$ for this GLM read as

$$\begin{aligned}\tilde{\boldsymbol{\theta}} &= \operatorname{argmax}_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{i=1}^n \{Y_i \Psi_i^\top \boldsymbol{\theta} - d(\Psi_i^\top \boldsymbol{\theta})\}, \\ \boldsymbol{\theta}^* &= \operatorname{argmax}_{\boldsymbol{\theta}} \mathbb{E}L(\boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{i=1}^n \{f_i \Psi_i^\top \boldsymbol{\theta} - d(\Psi_i^\top \boldsymbol{\theta})\}\end{aligned}$$

with $f_i = \mathbb{E}Y_i$. An important feature of a GLM is that the stochastic component $\zeta(\boldsymbol{\theta})$ of $L(\boldsymbol{\theta})$ is *linear in $\boldsymbol{\theta}$* : with $\varepsilon_i = Y_i - \mathbb{E}Y_i$

$$\zeta(\boldsymbol{\theta}) = L(\boldsymbol{\theta}) - \mathbb{E}L(\boldsymbol{\theta}) = \left(\sum_{i=1}^n \varepsilon_i \Psi_i \right)^\top \boldsymbol{\theta}, \quad \nabla \zeta(\boldsymbol{\theta}) = \sum_{i=1}^n \varepsilon_i \Psi_i.$$

In the contrary to the linear case, the Fisher information matrix D^2 depends on the true data distribution via the target $\boldsymbol{\theta}^*$:

$$D^2 = \sum_i \Psi_i \Psi_i^\top d''(\Psi_i^\top \boldsymbol{\theta}^*).$$

The vector $\boldsymbol{\xi}$ is given by

$$\boldsymbol{\xi} = D^{-1} \nabla \zeta(\boldsymbol{\theta}^*) = D^{-1} \sum_{i=1}^n \varepsilon_i \Psi_i.$$

Here is a list of sufficient conditions which ensure our general conditions from Section 6.2: the functions $d''(\Psi_i^\top \boldsymbol{\theta})$ are uniformly continuous in $\boldsymbol{\theta}$ and $i \leq n$, for some matrices \mathbf{v}_i^2 and fixed constants $\mathbf{C}_0, \lambda_0 > 0$

$$\mathbb{E} \exp\{\lambda_0 \mathbf{v}_i^{-1} \varepsilon_i\} \leq \mathbf{C}_0, \quad i = 1, \dots, n,$$

the sample size n is larger than $\mathbf{C}p$ for a prescribed constant \mathbf{C} , and the matrix $V^2 = \sum_i \mathbf{v}_i^2$ fulfills

$$V^2 \leq \mathbf{a}^2 D^2.$$

The details can be found in [Spokoiny \(2012\)](#).

As already mentioned, the error of Fisher and Wilks approximations only include the term $\delta(\mathbf{r}_0)$ because the stochastic term is linear by definition. By using the higher order expansion of the function $d(\cdot)$, [Portnoy \(1988\)](#) showed that the error $\diamond(\mathbf{r}_0, \mathbf{x})$ can be improved under the true the parametric assumption, however, the case of a model misspecification is not include. Under mild regularity conditions on the design Ψ and on

smoothness of $d(\cdot)$, the consistency result applies with $\mathbf{r}_0^2 = \mathbf{C}(p + \mathbf{x})$, and the bound $\delta(\mathbf{r}_0) \asymp \mathbf{r}_0/\sqrt{n}$ yields the errors $\diamond(\mathbf{r}_0, \mathbf{x}) \asymp \mathbf{r}_0\delta(\mathbf{r}_0) \leq \mathbf{C}(p + \mathbf{x})/\sqrt{n}$ and $\Delta(\mathbf{r}_0, \mathbf{x}) \asymp \mathbf{r}_0^2\delta(\mathbf{r}_0) \leq \mathbf{C}(p + \mathbf{x})^{3/2}/\sqrt{n}$.

6.4.3 I.i.d. case

Now we consider the asymptotic setup in which $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ is an i.i.d. sample from a measure P . The parametric assumption means that $P \in (P_\theta, \theta \in \Theta)$ for a given family of marginal measures (P_θ) . We admit that the parametric family depends on n and the parameter dimension $p = p_n$ grows to infinity with the sample size n . We also allow that the parametric assumption can be misspecified.

Similarly to Spokoiny (2012), our general conditions can be transformed into the conditions on the family (P_θ) and the marginal measure P . It suffices to check that the first and second derivatives of the log-density function $\ell(y, \theta) = \log dP_\theta/d\mu_0(y)$ have exponential moments and the expectation $\mathbb{E}\ell(Y_1, \theta)$ is three times continuously differentiable in θ . See Section 5.1 in Spokoiny (2012) for more details.

Also select $\mathbf{x} = \mathbf{x}_n$ depending on n and growing slowly with n , for instance, $\mathbf{x}_n = \log n$. The matrix D^2 satisfies $D^2 = n\mathbb{F}_{\theta^*}$, where \mathbb{F}_{θ^*} is the Fisher information of (P_θ) at θ^* if the parametric assumption holds.

The bracketing bound and the large deviation result from Spokoiny (2012) and from Section 6.3 apply if the sample size n fulfills $n \geq \mathbf{C}p_n$ for a fixed constant \mathbf{C} . It appears that the Fisher and Wilks results require stronger conditions. Indeed, in the regular i.i.d. case it holds $\delta(\mathbf{r}_0) \asymp \mathbf{r}_0/\sqrt{n}$, $q_{\mathbb{H}}^2(\mathbf{x}_n) \asymp p_n + \mathbf{x}_n$, and $\omega \asymp 1/\sqrt{n}$. The radius \mathbf{r}_0 should fulfill $\mathbf{r}_0^2 \geq \mathbf{C}p_n$ to ensure the large deviation result. This yields

$$\diamond(\mathbf{r}_0, \mathbf{x}_n) \leq \{\delta(\mathbf{r}_0) + 3\nu_0 q_{\mathbb{H}}(\mathbf{x}_n)\omega\}\mathbf{r}_0 \leq \mathbf{C}p_n/\sqrt{n}.$$

Similarly

$$\Delta(\mathbf{r}_0, \mathbf{x}_n) \leq \{\delta(\mathbf{r}_0) + 6\nu_0 q_{\mathbb{H}}(\mathbf{x}_n)\omega\}\mathbf{r}_0^2 \leq \mathbf{C}\sqrt{p_n^3/n}.$$

One can conclude that the consistency result is valid under $p_n/n \rightarrow 0$, the Fisher expansion requires $p_n^2/n \rightarrow 0$, while the Wilks is applicable under $p_n^3/n \rightarrow 0$ as $n \rightarrow \infty$.

Theorem 6.4.1. *Suppose the conditions of Theorem 5.4 in Spokoiny (2012). If $p_n^2/n \rightarrow 0$, then the Fisher expansion (6.6) of Theorem 6.3.2 holds with the error term $\diamond(\mathbf{r}_0, \mathbf{x}_n) \rightarrow 0$. Let also $p_n^3/n \rightarrow 0$. Then the error term $\Delta(\mathbf{r}_0, \mathbf{x}_n)$ in the Wilks expansion (6.8) satisfies $\Delta(\mathbf{r}_0, \mathbf{x}_n) \rightarrow 0$.*

Existing statistical literature addresses the issue of a growing parameter dimension in different set-ups. The classical results by Portnoy (1984, 1985, 1986) provide some

constraints on parameter dimension for consistency and asymptotic normality of the M-estimator for regression models. Our results are consistent with the conclusion of that papers. [Mammen \(1993, 1996\)](#) discussed the validity of bootstrap procedures in linear models with many parameters. The obtained results are valid under $p_n^{3/2}/n \rightarrow 0$, however are limited to the testing problem for linear models.

The setup with growing parameter dimension is naturally used in sieve nonparametric estimation when a nonparametric model is approximated by a sequence of parametric ones. We mention papers by [Shen and Wong \(1994\)](#); [Shen \(1997\)](#), [Birgé and Massart \(1993\)](#), [Van de Geer \(1993\)](#); [van de Geer \(2002\)](#). Some minimal smoothness assumptions are normally imposed on the underlying nonparametric function which ensure that the parameter dimension of a sieve is smaller in order than the sample size.

6.5 Some auxiliary results and proofs

This section collects some auxiliary results about the behavior of the posterior measures which might be of independent interest.

6.5.1 Local linear approximation of the gradient of the log-likelihood

The principle step of the proof is a bound on the local linear approximation of the gradient $\nabla L(\boldsymbol{\theta})$. Below we study separately its stochastic and deterministic components coming from the decomposition $L(\boldsymbol{\theta}) = \mathbb{E}L(\boldsymbol{\theta}) + \zeta(\boldsymbol{\theta})$. With $D^2 = -\nabla^2 \mathbb{E}L(\boldsymbol{\theta}^*)$, this leads to the decomposition

$$\begin{aligned} \chi(\boldsymbol{\theta}, \boldsymbol{\theta}^*) &\stackrel{\text{def}}{=} D^{-1}\{\nabla L(\boldsymbol{\theta}) - \nabla L(\boldsymbol{\theta}^*) + D^2(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\} \\ &= D^{-1}\{\nabla \zeta(\boldsymbol{\theta}) - \nabla \zeta(\boldsymbol{\theta}^*) + \nabla \mathbb{E}L(\boldsymbol{\theta}) - \nabla \mathbb{E}L(\boldsymbol{\theta}^*) + D^2(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\}. \end{aligned} \quad (6.10)$$

First we check the deterministic part. For any $\boldsymbol{\theta}$ with $\|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}$

$$\begin{aligned} \mathbb{E}[\chi(\boldsymbol{\theta}, \boldsymbol{\theta}^*)] &\stackrel{\text{def}}{=} D^{-1}\{\nabla \mathbb{E}L(\boldsymbol{\theta}) - \nabla \mathbb{E}L(\boldsymbol{\theta}^*)\} + D(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \\ &= \{I_p - D^{-1}D^2(\boldsymbol{\theta}^\circ)D^{-1}\}D(\boldsymbol{\theta} - \boldsymbol{\theta}^*), \end{aligned}$$

where $\boldsymbol{\theta}^\circ$ is a point on the line connecting $\boldsymbol{\theta}^*$ and $\boldsymbol{\theta}$. This implies by (\mathcal{L}_0)

$$\mathbb{E}[\chi(\boldsymbol{\theta}, \boldsymbol{\theta}^*)] \leq \|I_p - D^{-1}D^2(\boldsymbol{\theta}^\circ)D^{-1}\|_\infty \mathbf{r} \leq \delta(\mathbf{r})\mathbf{r}. \quad (6.11)$$

Now we study the stochastic part. Consider the vector process

$$\mathcal{U}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \stackrel{\text{def}}{=} D^{-1}\{\nabla \zeta(\boldsymbol{\theta}) - \nabla \zeta(\boldsymbol{\theta}^*)\}.$$

It is convenient to change the variable by $\mathbf{v} = D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)$ and consider the vector process $\mathcal{Y}(\mathbf{v}) = \mathcal{U}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$. It obviously holds $\nabla \mathcal{Y}(\mathbf{v}) = D^{-1} \nabla^2 \zeta(\boldsymbol{\theta}) D^{-1}$. Moreover, for any $\boldsymbol{\gamma}_1, \boldsymbol{\gamma}_2 \in \mathbb{R}^p$ with $\|\boldsymbol{\gamma}_1\| = \|\boldsymbol{\gamma}_2\| = 1$, condition (ED_2) implies

$$\log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega} \boldsymbol{\gamma}_1^\top \nabla \mathcal{Y}(\mathbf{v}) \boldsymbol{\gamma}_2 \right\} = \log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega} \boldsymbol{\gamma}_1^\top D^{-1} \nabla^2 \zeta(\boldsymbol{\theta}) D^{-1} \boldsymbol{\gamma}_2 \right\} \leq \frac{\nu_0^2 \lambda^2}{2}.$$

Define $\Upsilon_o(\mathbf{r}) \stackrel{\text{def}}{=} \{\mathbf{v} : \|\mathbf{v}\| \leq \mathbf{r}\}$. Then

$$\sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \|\mathcal{U}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\| = \sup_{\mathbf{v} \in \Upsilon_o(\mathbf{r})} \|\mathcal{Y}(\mathbf{v})\|. \quad (6.12)$$

Theorem 6.5.3 yields

$$\sup_{\mathbf{v} \in \Upsilon_o(\mathbf{r})} \|\mathcal{Y}(\mathbf{v})\| \leq 6\nu_0 q_{\mathbb{H}}(\mathbf{x}) \omega \mathbf{r}$$

on a set of a dominating probability at least $1 - e^{-\mathbf{x}}$, where the function $q_{\mathbb{H}}(\mathbf{x})$ is given by the following rules:

$$q_{\mathbb{H}}(\mathbf{x}) = \begin{cases} \sqrt{\mathbb{H}_2 + 2\mathbf{x}}, & \text{if } \mathbb{H}_2 + 2\mathbf{x} \leq \mathbf{g}^2, \\ \mathbf{g}^{-1}\mathbf{x} + \frac{1}{2}(\mathbf{g}^{-1}\mathbb{H}_2 + \mathbf{g}), & \text{if } \mathbb{H}_2 + 2\mathbf{x} > \mathbf{g}^2. \end{cases} \quad (6.13)$$

Here $\mathbb{H}_2 = 4p$ and $\mathbb{H}_1 = 2p^{1/2}$; see Theorem 6.5.1 in the Appendix.

Putting together the bounds (6.11) and (6.12) imply the following result.

Proposition 6.5.1. *Suppose that the matrix $D^2(\boldsymbol{\theta}) \stackrel{\text{def}}{=} -\nabla^2 \mathbb{E}L(\boldsymbol{\theta})$ fulfills the condition (\mathcal{L}_0) and let also (ED_2) be fulfilled on $\Theta_0(\mathbf{r})$ for any fixed \mathbf{r} . Then*

$$\mathbb{P} \left\{ \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \left\| D^{-1} \{ \nabla L(\boldsymbol{\theta}) - \nabla L(\boldsymbol{\theta}^*) \} + D(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \right\| \geq \diamond(\mathbf{r}, \mathbf{x}) \right\} \leq e^{-\mathbf{x}},$$

where

$$\diamond(\mathbf{r}, \mathbf{x}) \stackrel{\text{def}}{=} \{ \delta(\mathbf{r}) + 6\nu_0 q_{\mathbb{H}}(\mathbf{x}) \omega \} \mathbf{r}. \quad (6.14)$$

The result of Proposition 6.5.1 can be extended to the differences $\mathcal{U}(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) = D^{-1} \{ \nabla \zeta(\boldsymbol{\theta}) - \nabla \zeta(\boldsymbol{\theta}^\circ) \}$: on a set of probability at least $1 - e^{-\mathbf{x}}$, it holds for any $\boldsymbol{\theta}, \boldsymbol{\theta}^\circ \in \Theta_0(\mathbf{r})$ and $\chi(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) = D^{-1} \{ \nabla L(\boldsymbol{\theta}) - \nabla L(\boldsymbol{\theta}^\circ) \} + D(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)$

$$\begin{aligned} \mathbb{E}[\chi(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ)] &\leq \delta(\mathbf{r}) \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\| \leq 2\mathbf{r} \delta(\mathbf{r}), \\ \|\mathcal{U}(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ)\| &\leq \|\mathcal{U}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\| + \|\mathcal{U}(\boldsymbol{\theta}^\circ, \boldsymbol{\theta}^*)\| \leq 2\mathbf{r} \varrho(\mathbf{x}), \\ \|\chi(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ)\| &\leq 2 \diamond(\mathbf{r}, \mathbf{x}). \end{aligned} \quad (6.15)$$

6.5.2 Local quadratic approximation of the log-likelihood

As the next step, we derive a uniform deviation bound on the error of a quadratic approximation $\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) = (\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top \nabla L(\boldsymbol{\theta}^\circ) - \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\|^2/2$ of $L(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ)$:

$$\begin{aligned} \alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) &\stackrel{\text{def}}{=} L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^\circ) - (\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top \nabla L(\boldsymbol{\theta}^\circ) + \frac{1}{2} \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\|^2 \\ &= L(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) - \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) \end{aligned}$$

in all $\boldsymbol{\theta}, \boldsymbol{\theta}^\circ \in \Theta_0$, where Θ_0 is some vicinity of a fixed point $\boldsymbol{\theta}^*$. With $\boldsymbol{\theta}^\circ$ fixed, the gradient $\nabla \alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) \stackrel{\text{def}}{=} \frac{d}{d\boldsymbol{\theta}} \alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ)$ fulfills

$$\nabla \alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) = \nabla L(\boldsymbol{\theta}) - \nabla L(\boldsymbol{\theta}^\circ) + D^2(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ) = D \chi(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ);$$

cf. (6.10). This implies

$$\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) = (\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top \nabla \alpha(\boldsymbol{\theta}', \boldsymbol{\theta}^\circ),$$

where $\boldsymbol{\theta}'$ is a point on the line connecting $\boldsymbol{\theta}$ and $\boldsymbol{\theta}^\circ$. Further,

$$|\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ)| = |(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top D D^{-1} \nabla \alpha(\boldsymbol{\theta}', \boldsymbol{\theta}^\circ)| \leq \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\| \sup_{\boldsymbol{\theta}' \in \Theta_0(\mathbf{r})} |\chi(\boldsymbol{\theta}', \boldsymbol{\theta}^\circ)|.$$

and one can apply (6.15). This yields the following result.

Proposition 6.5.2. *Suppose (\mathcal{L}_0) , (ED_0) , and (ED_2) . For each \mathbf{r} , it holds on a random set $\Omega(\mathbf{r}, \mathbf{x})$ of a dominating probability at least $1 - e^{-\mathbf{x}}$, it holds with any $\boldsymbol{\theta}, \boldsymbol{\theta}^\circ \in \Theta_0(\mathbf{r})$*

$$\begin{aligned} \frac{|\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^*)|}{\|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|} &\leq \diamond(\mathbf{r}, \mathbf{x}), & |\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^*)| &\leq \mathbf{r} \diamond(\mathbf{r}, \mathbf{x}), \\ \frac{|\alpha(\boldsymbol{\theta}^*, \boldsymbol{\theta})|}{\|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|} &\leq 2\diamond(\mathbf{r}, \mathbf{x}), & |\alpha(\boldsymbol{\theta}^*, \boldsymbol{\theta})| &\leq 2\mathbf{r} \diamond(\mathbf{r}, \mathbf{x}), \\ \frac{|\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ)|}{\|D(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\|} &\leq 2\diamond(\mathbf{r}, \mathbf{x}), & |\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ)| &\leq 4\mathbf{r} \diamond(\mathbf{r}, \mathbf{x}), \end{aligned}$$

where $\diamond(\mathbf{r}, \mathbf{x})$ is from (6.14).

6.5.3 Proof of Theorem 6.3.1

By definition $\sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0)} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq 0$. So, it suffices to check that $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) < 0$ for all $\boldsymbol{\theta} \in \Theta \setminus \Theta_0(\mathbf{r}_0)$. The proof is based on the following bound: for each \mathbf{r}

$$\mathbb{P} \left(\sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} |\zeta(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla \zeta(\boldsymbol{\theta}^*)| \geq 3\nu_0 q_{\mathbb{H}}(\mathbf{x}) \omega \mathbf{r} \right) \leq e^{-\mathbf{x}}.$$

This bound is a special case of the general result from Theorem 6.5.1 below. It implies by Theorem 6.5.2 with $\rho = 1/2$ on a set $\Omega(\mathbf{x})$ of probability at least $1 - e^{-\mathbf{x}}$ that for all $\mathbf{r} \geq \mathbf{r}_0$ and all $\boldsymbol{\theta}$ with $\|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}$

$$|\zeta(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla \zeta(\boldsymbol{\theta}^*)| \leq \varrho(\mathbf{r}, \mathbf{x}) \mathbf{r},$$

where

$$\varrho(\mathbf{r}, \mathbf{x}) = 6\nu_0 q_{\mathbb{H}}(\mathbf{x} + \log(2\mathbf{r}/\mathbf{r}_0)) \omega.$$

The use of $\nabla \mathbb{E}L(\boldsymbol{\theta}^*) = 0$ yields

$$\sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} |L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla L(\boldsymbol{\theta}^*)| \leq \varrho(\mathbf{r}, \mathbf{x}) \mathbf{r}.$$

By Proposition 6.5.3, the vector $\boldsymbol{\xi} = D^{-1} \nabla \zeta(\boldsymbol{\theta}^*)$ fulfills $\mathbb{P}(\|\boldsymbol{\xi}\| \geq q(B, \mathbf{x})) \leq 2e^{-\mathbf{x}}$. We ignore here the negligible term of order $e^{-\mathbf{x}c}$. The condition $\|\boldsymbol{\xi}\| \leq q(B, \mathbf{x})$ implies for $\mathbf{r} \geq \mathbf{r}_0$

$$\begin{aligned} & \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} |(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla L(\boldsymbol{\theta}^*)| \\ & \leq \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \times \|D^{-1} \nabla L(\boldsymbol{\theta}^*)\| = \mathbf{r} \|\boldsymbol{\xi}\| \leq q(B, \mathbf{x}) \mathbf{r}. \end{aligned}$$

Condition (\mathcal{L}) implies $-2\mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq \mathbf{r}^2 \mathbf{b}(\mathbf{r})$ for each $\boldsymbol{\theta}$ with $\|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| = \mathbf{r}$. We conclude that the condition

$$\mathbf{r} \mathbf{b}(\mathbf{r}) \geq 2q(B, \mathbf{x}) + 2\varrho(\mathbf{r}, \mathbf{x}), \quad \mathbf{r} > \mathbf{r}_0,$$

ensure $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) < 0$ for all $\boldsymbol{\theta} \notin \Theta_0(\mathbf{r}_0)$ with a dominating probability.

6.5.4 Proof of Theorem 6.3.2

Let \mathbf{r}_0 be selected to ensure that $\mathbb{P}\{\tilde{\boldsymbol{\theta}} \notin \Theta_0(\mathbf{r}_0)\} \leq e^{-\mathbf{x}}$. Furthermore, the definition of $\tilde{\boldsymbol{\theta}}$ yields $\nabla L(\tilde{\boldsymbol{\theta}}) = 0$ and

$$\chi(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) = -D^{-1} \nabla L(\boldsymbol{\theta}^*) + D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*).$$

Now Proposition 6.5.1 implies on a set of a dominating probability

$$\|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \boldsymbol{\xi}\| \leq \diamond(\mathbf{r}, \mathbf{x}) \tag{6.16}$$

and the assertion follows.

6.5.5 Proof of Theorem 6.3.3

We apply the result of Proposition 6.5.2 on a random set of dominating probability $1 - 2e^{-x}$ on which $\tilde{\boldsymbol{\theta}} \in \Theta_0(\mathbf{r}_0)$ and the inequalities from that proposition are fulfilled. For the special case with $\boldsymbol{\theta}^\circ = \tilde{\boldsymbol{\theta}}$ and $\boldsymbol{\theta} = \boldsymbol{\theta}^*$ we obtain in view of $\nabla L(\tilde{\boldsymbol{\theta}}) = 0$ that

$$\left| L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2/2 \right| \leq |\alpha(\boldsymbol{\theta}^*, \tilde{\boldsymbol{\theta}})| \leq 2\mathbf{r}_0 \diamond(\mathbf{r}_0, \mathbf{x}). \quad (6.17)$$

Furthermore, with $\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}$ and $\boldsymbol{\theta}^\circ = \boldsymbol{\theta}^*$

$$\left| L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \boldsymbol{\xi}^\top D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) + \|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2/2 \right| = |\alpha(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)| \leq \mathbf{r}_0 \diamond(\mathbf{r}_0, \mathbf{x})$$

which implies

$$\left| L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \|\boldsymbol{\xi}\|^2/2 + \|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \boldsymbol{\xi}\|^2/2 \right| \leq \mathbf{r}_0 \diamond(\mathbf{r}_0, \mathbf{x}).$$

Now it follows by (6.16) that

$$\left| L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \|\boldsymbol{\xi}\|^2/2 \right| \leq \mathbf{r}_0 \diamond(\mathbf{r}_0, \mathbf{x}) + \diamond^2(\mathbf{r}_0, \mathbf{x})/2.$$

For the squared root of the excess, (6.17) implies

$$\begin{aligned} \left| \{2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)\}^{1/2} - \|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \right| &\leq \frac{|2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2|}{\|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|} \\ &\leq \frac{2|\alpha(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)|}{\|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|} \leq \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0)} \frac{2|\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^*)|}{\|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|} \leq 2 \diamond(\mathbf{r}_0, \mathbf{x}). \end{aligned} \quad (6.18)$$

Similarly, for any $\boldsymbol{\theta}^\circ \in \Theta_0(\mathbf{r})$, it holds

$$\left| \{2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^\circ)\}^{1/2} - \|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^\circ)\| \right| \leq \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0)} \frac{2|\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ)|}{\|D(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\|} \leq 4 \diamond(\mathbf{r}_0, \mathbf{x}).$$

The Fisher expansion (6.16) allows to replace in (6.18) the norm of the standardized error $D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)$ with the norm of the normalized score $\boldsymbol{\xi}$. This completes the proof of Theorem 6.3.3.

6.5.6 An entropy bound for the maximum of a random process

We use one general result on the upper bound for the maximum of a centered random process in the form of Spokoiny (2012); see Corollary 7.2 in the supplement of that paper.

Here we discuss the special case when Υ is an open subset in \mathbb{R}^p , the stochastic process $\mathcal{U}(\mathbf{v})$ is absolutely continuous and its gradient $\nabla \mathcal{U}(\mathbf{v}) \stackrel{\text{def}}{=} d\mathcal{U}(\mathbf{v})/d\mathbf{v}$ has bounded exponential moments.

(**ED**) *There exist $\mathbf{g} > 0$, $\nu_0 \geq 1$, and a symmetric non-negative matrix H_0 such that for any $\lambda \leq \mathbf{g}$ and any unit vector $\boldsymbol{\gamma} \in \mathbb{R}^p$, it holds*

$$\log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{\gamma}^\top \nabla \mathcal{U}(\mathbf{v})}{\|H_0 \boldsymbol{\gamma}\|} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}.$$

We consider the local sets of the elliptic form $\Upsilon_\circ(\mathbf{r}) \stackrel{\text{def}}{=} \{\mathbf{v} : \|H_0(\mathbf{v} - \mathbf{v}^*)\| \leq \mathbf{r}\}$.

Theorem 6.5.1 (Spokoiny (2012)). *Let (ED) hold with some $\mathbf{g} > 0$, and a matrix H_0 . For any $\mathbf{x} \geq 0$ and any $\mathbf{r} > 0$*

$$\mathbb{P} \left\{ \sup_{\mathbf{v} \in \Upsilon_\circ(\mathbf{r})} |\mathcal{U}(\mathbf{v}, \mathbf{v}^*)| \geq 3\nu_0 \mathbf{r} q_{\mathbb{H}}(\mathbf{x}) \right\} \leq e^{-\mathbf{x}},$$

where $q_{\mathbb{H}}(\mathbf{x})$ is given by the following rule:

$$q_{\mathbb{H}}(\mathbf{x}) = \begin{cases} \sqrt{\mathbb{H} + 2\mathbf{x}} & \text{if } \mathbb{H} + 2\mathbf{x} \leq \mathbf{g}^2, \\ \mathbf{g}^{-1}\mathbf{x} + \frac{1}{2}(\mathbf{g}^{-1}\mathbb{H} + \mathbf{g}) & \text{if } \mathbb{H} + 2\mathbf{x} > \mathbf{g}^2, \end{cases} \quad (6.19)$$

with $\mathbb{H} = 4p$.

Due to the result of Theorem 6.5.1, the bound for the maximum of $\mathcal{U}(\mathbf{v}, \mathbf{v}^*)$ over $\mathbf{v} \in \mathcal{B}_{\mathbf{r}}(\mathbf{v}^*)$ grows linearly in \mathbf{r} . So, its applications to situations with $\mathbf{r} \gg \mathbb{Q}_1(\Upsilon^\circ)$ are limited. The next result shows that introducing a negative drift helps to state a uniform in \mathbf{r} local probability bound. Namely, the bound for the process $\mathcal{U}(\mathbf{v}, \mathbf{v}^*) - f(d(\mathbf{v}, \mathbf{v}^*))$ for some function $f(\mathbf{r})$ over a ball $\mathcal{B}_{\mathbf{r}}(\mathbf{v}^*)$ around the point \mathbf{v}^* does not depend on \mathbf{r} . Here the generic chaining arguments are accomplished with the slicing technique. The idea is for a given $\mathbf{r}^* > 1$ to split the ball $\mathcal{B}_{\mathbf{r}^*}(\mathbf{v}^*)$ into the slices $\mathcal{B}_{\mathbf{r}_k}(\mathbf{v}^*) \setminus \mathcal{B}_{\mathbf{r}_{k-1}}(\mathbf{v}^*)$ and to apply Theorem 6.5.1 to each slice separately.

Theorem 6.5.2. *Let \mathbf{r}^* be such that (ED) holds on $\mathcal{B}_{\mathbf{r}^*}(\mathbf{v}^*)$. Given $\mathbf{r}_0 < \mathbf{r}^*$, let a monotonous function $f(\mathbf{r}, \mathbf{r}_0)$ fulfill*

$$f(\mathbf{r}, \mathbf{r}_0) \geq 3\nu_0 \mathbf{r} q_{\mathbb{H}}(\mathbf{x} + \log(\mathbf{r}/\mathbf{r}_0)), \quad \mathbf{r}_0 \leq \mathbf{r} \leq \mathbf{r}^*, \quad (6.20)$$

where the function $q_{\mathbb{H}}(\cdot)$ is given by (6.19). Then it holds for any $\rho < 1$

$$\mathbb{P} \left(\sup_{\mathbf{r}_0 \leq \mathbf{r} \leq \mathbf{r}^*} \sup_{\mathbf{v} \in \Upsilon_\circ(\mathbf{r})} \{\mathcal{U}(\mathbf{v}, \mathbf{v}^*) - f(\rho^{-1}\mathbf{r}, \mathbf{r}_0)\} \geq 0 \right) \leq \frac{\rho}{1-\rho} e^{-\mathbf{x}}.$$

Remark 6.5.1. Formally the bound applies even with $\mathbf{r}^* = \infty$ provided that (ED) is fulfilled on the whole set Υ° .

Remark 6.5.2. If $\mathbf{g} = \infty$, then $q_{\mathbb{H}}(\mathbf{x}) = \sqrt{2\mathbf{x} + 4p}$ and the condition (6.20) on the drift simplifies to $(3\nu_0 \mathbf{r})^{-1} f(\mathbf{r}, \mathbf{r}_0) \geq \sqrt{2\mathbf{x} + 4p + 2 \log(\mathbf{r}/\mathbf{r}_0)}$.

Proof. By (6.20) and Theorem 6.5.1 for any $\mathbf{r} > \mathbf{r}_0$

$$\begin{aligned} & \mathbb{P}\left(\sup_{\mathbf{v} \in \mathcal{B}_{\mathbf{r}}(\mathbf{v}^*) \setminus \mathcal{B}_{\rho \mathbf{r}}(\mathbf{v}^*)} \{\mathcal{U}(\mathbf{v}, \mathbf{v}^*) - f(\mathbf{r}, \mathbf{r}_0)\} \geq 0\right) \\ & \leq \mathbb{P}\left(\frac{1}{3\nu_0 \mathbf{r}} \sup_{\mathbf{v} \in \mathcal{B}_{\mathbf{r}}(\mathbf{v}^*)} \mathcal{U}(\mathbf{v}, \mathbf{v}^*) \geq q_{\mathbb{H}}(\mathbf{x} + \log(\mathbf{r}/\mathbf{r}_0))\right) \leq \frac{\mathbf{r}_0}{\mathbf{r}} e^{-\mathbf{x}}. \end{aligned} \quad (6.21)$$

Now define $\mathbf{r}_k = \mathbf{r}_0 \rho^{-k}$ for $k = 0, 1, 2, \dots$. Define also $k^* \stackrel{\text{def}}{=} \log(\mathbf{r}^*/\mathbf{r}_0) + 1$. It follows from (6.21) that

$$\begin{aligned} & \mathbb{P}\left(\sup_{\mathbf{v} \in \mathcal{B}_{\mathbf{r}^*}(\mathbf{v}^*) \setminus \mathcal{B}_{\mathbf{r}_0}(\mathbf{v}^*)} \left\{ \mathcal{U}(\mathbf{v}, \mathbf{v}^*) - f(\rho^{-1}d(\mathbf{v}, \mathbf{v}^*), \mathbf{r}_0) \right\} \geq 0\right) \\ & \leq \sum_{k=1}^{k^*} \mathbb{P}\left(\frac{1}{\mathbf{r}_k} \sup_{\mathbf{v} \in \mathcal{B}_{\mathbf{r}_k}(\mathbf{v}^*) \setminus \mathcal{B}_{\mathbf{r}_{k-1}}(\mathbf{v}^*)} \left\{ \mathcal{U}(\mathbf{v}, \mathbf{v}^*) - f(\mathbf{r}_k, \mathbf{r}_0) \right\} \geq 0\right) \\ & \leq e^{-\mathbf{x}} \sum_{k=1}^{k^*} \rho^k \leq \frac{\rho}{1-\rho} e^{-\mathbf{x}} \end{aligned}$$

as required.

6.5.7 A bound for the norm of a vector random process

Let $\mathcal{Y}(\mathbf{v})$, $\mathbf{v} \in \Upsilon$, be a smooth centered random vector process with values in \mathbb{R}^q , where $\Upsilon \subseteq \mathbb{R}^p$. Let also $\mathcal{Y}(\mathbf{v}^*) = 0$ for a fixed point $\mathbf{v}^* \in \Upsilon$. Without loss of generality assume $\mathbf{v}^* = 0$. We aim to bound the maximum of the norm $\|\mathcal{Y}(\mathbf{v})\|$ over a vicinity Υ_{\circ} of \mathbf{v}^* . Suppose that $\mathcal{Y}(\mathbf{v})$ satisfies for each $\boldsymbol{\gamma} \in \mathbb{R}^p$ and $\boldsymbol{\alpha} \in \mathbb{R}^q$ with $\|\boldsymbol{\gamma}\| = \|\boldsymbol{\alpha}\| = 1$

$$\sup_{\mathbf{v} \in \Upsilon} \log \mathbb{E} \exp\left\{ \lambda \boldsymbol{\gamma}^{\top} \nabla \mathcal{Y}(\mathbf{v}) \boldsymbol{\alpha} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad \lambda^2 \leq 2\mathbf{g}^2. \quad (6.22)$$

Condition (6.22) implies for any $\mathbf{v} \in \Upsilon_{\circ}$ with $\|\mathbf{v}\| \leq \mathbf{r}$ and $\|\boldsymbol{\gamma}\| = 1$ in view of $\mathcal{Y}(\mathbf{v}^*) = 0$

$$\log \mathbb{E} \exp\left\{ \frac{\lambda}{\mathbf{r}} \boldsymbol{\gamma}^{\top} \mathcal{Y}(\mathbf{v}) \right\} \leq \frac{\nu_0^2 \lambda^2 \|\mathbf{v}\|^2}{2\mathbf{r}^2}, \quad \lambda^2 \leq 2\mathbf{g}^2; \quad (6.23)$$

In what follows, we use the representation

$$\|\mathcal{Y}(\mathbf{v})\| = \sup_{\|\mathbf{u}\| \leq \mathbf{r}} \frac{1}{\mathbf{r}} \mathbf{u}^{\top} \mathcal{Y}(\mathbf{v}).$$

This implies for $\Upsilon_{\circ}(\mathbf{r}) = \{\mathbf{v} \in \Upsilon : \|\mathbf{v} - \mathbf{v}^*\| \leq \mathbf{r}\}$

$$\sup_{\mathbf{v} \in \Upsilon_{\circ}(\mathbf{r})} \|\mathcal{Y}(\mathbf{v})\| = \sup_{\mathbf{v} \in \Upsilon_{\circ}(\mathbf{r})} \sup_{\|\mathbf{u}\| \leq \mathbf{r}} \frac{1}{\mathbf{r}} \mathbf{u}^{\top} \mathcal{Y}(\mathbf{v}).$$

Consider a bivariate process $\mathbf{u}^\top \mathfrak{y}(\mathbf{v})$ of $\mathbf{u} \in \mathbb{R}^q$ and $\mathbf{v} \in \Upsilon \subset \mathbb{R}^p$. By definition $\mathbb{E}\mathbf{u}^\top \mathfrak{y}(\mathbf{v}) = 0$. Further, $\nabla_{\mathbf{u}}[\mathbf{u}^\top \mathfrak{y}(\mathbf{v})] = \mathfrak{y}(\mathbf{v})$ while $\nabla_{\mathbf{v}}[\mathbf{u}^\top \mathfrak{y}(\mathbf{v})] = \mathbf{u}^\top \nabla \mathfrak{y}(\mathbf{v}) = \|\mathbf{u}\| \boldsymbol{\gamma}^\top \nabla \mathfrak{y}(\mathbf{v})$ for $\boldsymbol{\gamma} = \mathbf{u}/\|\mathbf{u}\|$. Suppose that $\mathbf{u} \in \mathbb{R}^q$ and $\mathbf{v} \in \Upsilon$ are such that $\|\mathbf{u}\|^2 + \|\mathbf{v}\|^2 \leq 2\mathbf{r}^2$. By the Hölder inequality, (6.23), and (6.22), it holds for $\|\boldsymbol{\gamma}\| = \|\boldsymbol{\alpha}\| = 1$ and $\mathbf{v} \in \Upsilon_{\circ}(\mathbf{r})$

$$\begin{aligned} & \log \mathbb{E} \exp \left\{ \frac{\lambda}{2\mathbf{r}} (\boldsymbol{\gamma}, \boldsymbol{\alpha})^\top \nabla [\mathbf{u}^\top \mathfrak{y}(\mathbf{v})] \right\} \\ & \leq \frac{1}{2} \log \mathbb{E} \exp \left\{ \frac{\lambda}{\mathbf{r}} \boldsymbol{\gamma}^\top \mathfrak{y}(\mathbf{v}) \right\} + \frac{1}{2} \log \mathbb{E} \exp \left\{ \frac{\lambda}{\mathbf{r}} \mathbf{u}^\top \nabla \mathfrak{y}(\mathbf{v}) \boldsymbol{\alpha} \right\} \\ & \leq \frac{1}{2} \log \mathbb{E} \exp \left\{ \frac{\lambda}{\mathbf{r}} \boldsymbol{\gamma}^\top \mathfrak{y}(\mathbf{v}) \right\} + \frac{1}{2} \log \mathbb{E} \exp \left\{ \frac{\lambda}{\mathbf{r}} \|\mathbf{u}\| \boldsymbol{\gamma}^\top \nabla \mathfrak{y}(\mathbf{v}) \boldsymbol{\alpha} \right\} \\ & \leq \frac{\nu_0^2 \lambda^2}{4\mathbf{r}^2} (\|\mathbf{v}\|^2 + \|\mathbf{u}\|^2) \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq \mathbf{g}. \end{aligned}$$

We summarize our findings in the following theorem.

Theorem 6.5.3. *Let a random p -vector process $\mathfrak{y}(\mathbf{v})$ for $\mathbf{v} \in \Upsilon \subseteq \mathbb{R}^p$ fulfill $\mathfrak{y}(\mathbf{v}^*) = 0$, $\mathbb{E}\mathfrak{y}(\mathbf{v}) \equiv 0$, and the condition (6.22) be satisfied. Then for each \mathbf{r} and any $\mathbf{x} \geq 1/2$, it holds*

$$\mathbb{P} \left\{ \sup_{\mathbf{v} \in \Upsilon_{\circ}(\mathbf{r})} \|\mathfrak{y}(\mathbf{v})\| > 6\nu_0 \mathbf{r} q_{\mathbb{H}}(\mathbf{x}) \right\} \leq e^{-\mathbf{x}},$$

where $q_{\mathbb{H}}(\mathbf{x})$ is given by (6.19).

Proof. The results follow from Theorem 6.5.1 applied to the process $\mathbf{u}^\top \mathfrak{y}(\mathbf{v})/2$ of the variable $(\mathbf{u}, \mathbf{v}) \in \mathbb{R}^{p+q}$.

6.5.8 A deviation bound for the quadratic form $\|\boldsymbol{\xi}\|^2$

This section presents a bound for a quadratic form $\|\boldsymbol{\xi}\|^2$ where $\boldsymbol{\xi} = D^{-1} \nabla \zeta(\boldsymbol{\theta}^*)$. The result only uses the condition (ED_0) which we restate in a slightly different form. For $B = D^{-1} V^2 D^{-1}$, define

$$\mathfrak{p}_B \stackrel{\text{def}}{=} \text{tr}(B), \quad \mathfrak{v}_B^2 \stackrel{\text{def}}{=} 2 \text{tr}(B^2), \quad \lambda_B \stackrel{\text{def}}{=} \lambda_{\max}(B).$$

Note that $\mathfrak{p}_B = \mathbb{E}\|\boldsymbol{\xi}\|^2$. Moreover, if $\boldsymbol{\xi}$ is a Gaussian vector then $\mathfrak{v}_B^2 = \text{Var}(\|\boldsymbol{\xi}\|^2)$. If $V^2 = D^2$, then $\lambda_B = 1$. The condition (ED_0) means that the vector $V^{-1} \nabla \zeta(\boldsymbol{\theta}^*)$ fulfills the following exponential moment condition:

$$\log \mathbb{E} \exp(\boldsymbol{\gamma}^\top V^{-1} \nabla \zeta(\boldsymbol{\theta}^*)) \leq \|\boldsymbol{\gamma}\|^2/2, \quad \boldsymbol{\gamma} \in \mathbb{R}^p, \|\boldsymbol{\gamma}\| \leq \mathbf{g}.$$

Here ν_0 is set to one. Spokoiny (2012) argued how the case of any $\nu_0 \geq 1$ can be reduced to $\nu_0 \approx 1$ by a slight change of scale and reducing the value \mathbf{g} which is typically large. For ease of presentation, suppose that $\mathbf{g}^2 \geq 2\mathbf{p}_B$. The other case only changes the constants in the inequalities. Note that $\|\boldsymbol{\xi}\|^2 = \boldsymbol{\eta}^\top B \boldsymbol{\eta}$. Define $\mu_c = 2/3$ and

$$\begin{aligned} \mathbf{g}_c &\stackrel{\text{def}}{=} \sqrt{\mathbf{g}^2 - \mu_c \mathbf{p}_B}, \\ 2\mathbf{x}_c &\stackrel{\text{def}}{=} (\mathbf{g}^2/\mu_c - \mathbf{p}_B)/\lambda_B + \log \det(I_p - \mu_c B/\lambda_B). \end{aligned} \quad (6.24)$$

Proposition 6.5.3 (Spokoiny (2012)). *Let (ED_0) hold with $\nu_0 = 1$ and $\mathbf{g}^2 \geq 2\mathbf{p}_B$. Then for each $\mathbf{x} > 0$*

$$\mathbb{P}(\|\boldsymbol{\xi}\| \geq q(B, \mathbf{x})) \leq 2e^{-\mathbf{x}} + 8.4e^{-\mathbf{x}_c},$$

where $q(B, \mathbf{x})$ is defined by

$$q^2(B, \mathbf{x}) \stackrel{\text{def}}{=} \begin{cases} \mathbf{p}_B + 2\mathbf{v}_B \mathbf{x}^{1/2}, & \mathbf{x} \leq \mathbf{v}_B/(18\lambda_B), \\ \mathbf{p}_B + 6\lambda_B \mathbf{x}, & \mathbf{v}_B/(18\lambda_B) < \mathbf{x} \leq \mathbf{x}_c, \\ |\mathbf{y}_c + 2\lambda_B(\mathbf{x} - \mathbf{x}_c)/\mathbf{g}_c|^2, & \mathbf{x} > \mathbf{x}_c. \end{cases} \quad (6.25)$$

with $\mathbf{y}_c^2 \leq \mathbf{p}_B + 6\lambda_B \mathbf{x}_c$.

Depending on the value \mathbf{x} , we observe three types of tail behavior of the quadratic form $\|\boldsymbol{\xi}\|^2$. The sub-Gaussian regime for $\mathbf{x} \leq \mathbf{v}_B/(18\lambda_B)$ and the Poissonian regime for $\mathbf{x} \leq \mathbf{x}_c$ are similar to the case of a Gaussian quadratic form. The value \mathbf{x}_c from (6.24) is of order \mathbf{g}^2 . In all our results we suppose that \mathbf{g}^2 and hence, \mathbf{x}_c is sufficiently large and the quadratic form $\|\boldsymbol{\xi}\|^2$ can be bounded with a dominating probability by $\mathbf{p}_B + 6\lambda_B \mathbf{x}$ for a proper \mathbf{x} . We refer to Spokoiny (2012) for the proof of this and related results, further discussion and references.

6.6 Some results for the normal law

This section collects some simple but useful facts about the properties of the multivariate standard normal distribution. Many similar results can be found in the literature, we present the proofs to keep the presentation self-contained.

6.6.1 Deviation bounds

This section collects some deviation bounds on the norm or quadratic form of a standard normal vector. Everywhere in this section $\boldsymbol{\gamma}$ means a standard normal vector in \mathbb{R}^p .

Lemma 6.6.1. *Let $\mu \in (0, 1)$. Then for any vector $\boldsymbol{\lambda} \in \mathbb{R}^p$ with $\|\boldsymbol{\lambda}\|^2 \leq p$ and any $\mathbf{r} > 0$*

$$\log \mathbb{E}\{\exp(\boldsymbol{\lambda}^\top \boldsymbol{\gamma}) \mathbb{I}(\|\boldsymbol{\gamma}\| > \mathbf{r})\} \leq -\frac{1-\mu}{2}\mathbf{r}^2 + \frac{1}{2\mu}\|\boldsymbol{\lambda}\|^2 + \frac{p}{2}\log(\mu^{-1}). \quad (6.26)$$

Moreover, if $\mathbf{r}^2 \geq 6p + 4\mathbf{x}$, then

$$\mathbb{E}\left\{\exp(\boldsymbol{\lambda}^\top \boldsymbol{\gamma}) \mathbb{I}(\|\boldsymbol{\gamma}\| \leq \mathbf{r})\right\} \geq e^{\|\boldsymbol{\lambda}\|^2/2}(1 - e^{-\mathbf{x}}). \quad (6.27)$$

Proof. We use that for $\mu < 1$

$$\mathbb{E}\{\exp(\boldsymbol{\lambda}^\top \boldsymbol{\gamma}) \mathbb{I}(\|\boldsymbol{\gamma}\| > \mathbf{r})\} \leq e^{-(1-\mu)\mathbf{r}^2/2} \mathbb{E} \exp\{\boldsymbol{\lambda}^\top \boldsymbol{\gamma} + (1-\mu)\|\boldsymbol{\gamma}\|^2/2\}.$$

It holds

$$\begin{aligned} \mathbb{E} \exp\{\boldsymbol{\lambda}^\top \boldsymbol{\gamma} + (1-\mu)\|\boldsymbol{\gamma}\|^2/2\} &= (2\pi)^{-p/2} \int \exp\{\boldsymbol{\lambda}^\top \boldsymbol{\gamma} - \mu\|\boldsymbol{\gamma}\|^2/2\} d\boldsymbol{\gamma} \\ &= \mu^{-p/2} \exp(\mu^{-1}\|\boldsymbol{\lambda}\|^2/2) \end{aligned}$$

and (6.26) follows.

Now we apply this result with $\mu = 1/2$. In view of $\mathbb{E} \exp(\boldsymbol{\lambda}^\top \boldsymbol{\gamma}) = e^{\|\boldsymbol{\lambda}\|^2/2}$, $\mathbf{r}^2 \geq 6p + 4\mathbf{x}$, and $2 + \log(2) < 3$, it follows for $\|\boldsymbol{\lambda}\|^2 \leq p$

$$\begin{aligned} e^{-\|\boldsymbol{\lambda}\|^2/2} \mathbb{E}\{\exp(\boldsymbol{\lambda}^\top \boldsymbol{\gamma}) \mathbb{I}(\|\boldsymbol{\gamma}\| \leq \mathbf{r})\} \\ \geq 1 - \exp(-\mathbf{r}^2/4 + p + (p/2)\log(2)) \geq 1 - \exp(-\mathbf{x}) \end{aligned}$$

which implies (6.27).

Lemma 6.6.2. *For any $\mathbf{u} \in \mathbb{R}^p$, any unit vector $\mathbf{a} \in \mathbb{R}^p$, and any $q > 0$, it holds*

$$\mathbb{P}(\|\boldsymbol{\gamma} - \mathbf{u}\| \geq q) \leq \exp\{-q^2/4 + p/2 + \|\mathbf{u}\|^2/2\}, \quad (6.28)$$

$$\mathbb{E}\{|\boldsymbol{\gamma}^\top \mathbf{a}|^2 \mathbb{I}(\|\boldsymbol{\gamma} - \mathbf{u}\| \geq q)\} \leq (2 + |\mathbf{u}^\top \mathbf{a}|^2) \exp\{-q^2/4 + p/2 + \|\mathbf{u}\|^2/2\}. \quad (6.29)$$

Proof. By the exponential Chebyshev inequality, for any $\lambda < 1$

$$\begin{aligned} \mathbb{P}(\|\boldsymbol{\gamma} - \mathbf{u}\| \geq q) &\leq \exp(-\lambda q^2/2) \mathbb{E} \exp(\lambda \|\boldsymbol{\gamma} - \mathbf{u}\|^2/2) \\ &= \exp\left\{-\frac{\lambda q^2}{2} - \frac{p}{2} \log(1-\lambda) + \frac{\lambda}{2(1-\lambda)} \|\mathbf{u}\|^2\right\}. \end{aligned}$$

In particular, with $\lambda = 1/2$, this implies (6.28). Further, for $\|\mathbf{a}\| = 1$

$$\begin{aligned} \mathbb{E}\{|\boldsymbol{\gamma}^\top \mathbf{a}|^2 \mathbb{I}(\|\boldsymbol{\gamma} - \mathbf{u}\| \geq q)\} &\leq \exp(-q^2/4) \mathbb{E}\{|\boldsymbol{\gamma}^\top \mathbf{a}|^2 \exp(\|\boldsymbol{\gamma} - \mathbf{u}\|^2/4)\} \\ &\leq (2 + |\mathbf{u}^\top \mathbf{a}|^2) \exp(-q^2/4 + p/2 + \|\mathbf{u}\|^2/2) \end{aligned}$$

and (6.29) follows.

The next result explains the concentration effect for the norm $\|\boldsymbol{\xi}\|^2$ of a Gaussian vector. We use a version from Spokoiny (2012).

Lemma 6.6.3. *For each \mathbf{x} ,*

$$\mathbb{P}(\|\boldsymbol{\gamma}\| \geq q(p, \mathbf{x})) \leq \exp(-\mathbf{x}), \quad \mathbb{P}(\|\boldsymbol{\gamma}\| \leq q_1(p, \mathbf{x})) \leq \exp(-\mathbf{x}), \quad (6.30)$$

where

$$q^2(p, \mathbf{x}) \stackrel{\text{def}}{=} p + \sqrt{6.6p\mathbf{x}} \vee (6.6\mathbf{x}), \quad q_1^2(p, \mathbf{x}) \stackrel{\text{def}}{=} p - 2\sqrt{p\mathbf{x}}.$$

Proof. The exponential Chebyshev inequality implies for any positive $\lambda \leq 1/2$ in view of $-\log(1 - \lambda) \leq \lambda + \lambda^2$

$$\begin{aligned} \mathbb{P}(\|\boldsymbol{\gamma}\|^2 > p + \sqrt{2p\mathbf{x}}) &\leq \exp\left\{-\frac{\lambda}{2}(p + \sqrt{2p\mathbf{x}})\right\} \mathbb{E} \exp\left(\frac{\lambda}{2}\|\boldsymbol{\gamma}\|^2\right) \\ &\leq \exp\left\{-\frac{\lambda}{2}(p + \sqrt{2p\mathbf{x}})\right\} \mathbb{E} \exp\left(\frac{\lambda}{2}\|\boldsymbol{\gamma}\|^2\right) \\ &= \exp\left\{-\frac{\lambda}{2}(p + \sqrt{2p\mathbf{x}}) - \frac{p}{2}\log(1 - \lambda)\right\} \\ &\leq \exp\{-\lambda\sqrt{p\mathbf{x}/2} + p\lambda^2/2\}. \end{aligned}$$

The choice $\lambda = \sqrt{\mathbf{x}/(2p)}$ yields the result (6.30). Similarly

$$\begin{aligned} \mathbb{P}(\|\boldsymbol{\gamma}\|^2 - p < -\sqrt{2p\mathbf{x}}) &\leq \exp\left\{-\frac{\lambda}{2}(\sqrt{2p\mathbf{x}} - p)\right\} \mathbb{E} \exp\left(-\frac{\lambda}{2}\|\boldsymbol{\gamma}\|^2\right) \\ &\leq \exp\left\{-\frac{\lambda}{2}(\sqrt{2p\mathbf{x}} - p)\right\} \mathbb{E} \exp\left(-\frac{\lambda}{2}\|\boldsymbol{\gamma}\|^2\right) \\ &= \exp\left\{-\frac{\lambda}{2}(\sqrt{2p\mathbf{x}} - p) - \frac{p}{2}\log(1 + \lambda)\right\} \\ &\leq \exp\{-\lambda\sqrt{p\mathbf{x}/2} + p\lambda^2/4\}. \end{aligned}$$

Here the choice $\lambda = \sqrt{2\mathbf{x}/p}$ does the job.

A similar bound can be obtained for a norm of the vector $B\boldsymbol{\xi}$ where B is some given deterministic matrix. For notational simplicity we assume that B is symmetric. Otherwise one should replace it with $(B^\top B)^{1/2}$.

Theorem 6.6.1. *Let $\boldsymbol{\xi}$ be standard normal in \mathbb{R}^p . Then for every $\mathbf{x} > 0$ and any symmetric matrix B , it holds with $\mathbf{p} = \text{tr}(B^2)$, $\mathbf{v}^2 = 2\text{tr}(B^4)$, and $a^* = \|B\|_\infty$*

$$\mathbb{P}(\|B\boldsymbol{\xi}\|^2 > \mathbf{p} + (2\mathbf{v}\mathbf{x}^{1/2}) \vee (6a^*\mathbf{x})) \leq \exp(-\mathbf{x}).$$

Proof. The matrix B^2 can be represented as $U^\top \text{diag}(a_1, \dots, a_p)U$ for an orthogonal matrix U . The vector $\tilde{\xi} = U\xi$ is also standard normal and $\|B\xi\|^2 = \tilde{\xi}^\top UB^2U^\top\tilde{\xi}$. This means that one can reduce the situation to the case of a diagonal matrix $B^2 = \text{diag}(a_1, \dots, a_p)$. We can also assume without loss of generality that $a_1 \geq a_2 \geq \dots \geq a_p$. The expressions for the quantities \mathbf{p} and \mathbf{v}^2 simplifies to

$$\begin{aligned}\mathbf{p} &= \text{tr}(B^2) = a_1 + \dots + a_p, \\ \mathbf{v}^2 &= 2 \text{tr}(B^4) = 2(a_1^2 + \dots + a_p^2).\end{aligned}$$

Moreover, rescaling the matrix B^2 by a_1 reduces the situation to the case with $a_1 = 1$.

Lemma 6.6.4. *It holds*

$$\mathbb{E}\|B\xi\|^2 = \text{tr}(B^2), \quad \text{Var}(\|B\xi\|^2) = 2 \text{tr}(B^4).$$

Moreover, for $\mu < 1$

$$\mathbb{E} \exp\{\mu\|B\xi\|^2/2\} = \det(1 - \mu B^2)^{-1/2} = \prod_{i=1}^p (1 - \mu a_i)^{-1/2}. \quad (6.31)$$

Proof. If B^2 is diagonal, then $\|B\xi\|^2 = \sum_i a_i \xi_i^2$ and the summands $a_i \xi_i^2$ are independent. It remains to note that $\mathbb{E}(a_i \xi_i^2) = a_i$, $\text{Var}(a_i \xi_i^2) = 2a_i^2$, and for $\mu a_i < 1$,

$$\mathbb{E} \exp\{\mu a_i \xi_i^2/2\} = (1 - \mu a_i)^{-1/2}$$

yielding (6.31).

Given u , fix $\mu < 1$. The exponential Markov inequality yields

$$\begin{aligned}\mathbb{P}(\|B\xi\|^2 > \mathbf{p} + u) &\leq \exp\left\{-\frac{\mu(\mathbf{p} + u)}{2}\right\} \mathbb{E} \exp\left(\frac{\mu\|B\xi\|^2}{2}\right) \\ &\leq \exp\left\{-\frac{\mu u}{2} - \frac{1}{2} \sum_{i=1}^p [\mu a_i + \log(1 - \mu a_i)]\right\}.\end{aligned}$$

We start with the case when $\mathbf{x}^{1/2} \leq \mathbf{v}/3$. Then $u = 2\mathbf{x}^{1/2}\mathbf{v}$ fulfills $u \leq 2\mathbf{v}^2/3$. Define $\mu = u/\mathbf{v}^2 \leq 2/3$ and use that $t + \log(1 - t) \geq -t^2$ for $t \leq 2/3$. This implies

$$\begin{aligned}\mathbb{P}(\|B\xi\|^2 > \mathbf{p} + u) &\leq \exp\left\{-\frac{\mu u}{2} + \frac{1}{2} \sum_{i=1}^p \mu^2 a_i^2\right\} = \exp(-u^2/(4\mathbf{v}^2)) = e^{-\mathbf{x}}.\end{aligned} \quad (6.32)$$

Next, let $\mathbf{x}^{1/2} > \mathbf{v}/3$. Set $\mu = 2/3$. It holds similarly to the above

$$\sum_{i=1}^p [\mu a_i + \log(1 - \mu a_i)] \geq - \sum_{i=1}^p \mu^2 a_i^2 \geq -2v^2/9 \geq -2x.$$

Now, for $u = 6x$ and $\mu u/2 = 2x$, (6.32) implies

$$\mathbb{P}(\|B\xi\|^2 > p + u) \leq \exp\{-(2x - x)\} = \exp(-x)$$

as required.

6.6.2 Gaussian comparison via KL-divergence and Pinsker's inequality

Suppose that two p -dimensional zero mean Gaussian vectors $\xi \sim \mathcal{N}(0, \Sigma)$ and $\xi^b \sim \mathcal{N}(0, \Sigma^b)$ are given. Let also T map \mathbb{R}^p to \mathbb{R}^M and $\mathbf{X} = T(\xi)$ and $\mathbf{Y} = T(\xi^b)$. We aim to bound the distance between distributions of \mathbf{X} and \mathbf{Y} under the conditions

$$\|\Sigma^{-1/2} \Sigma^b \Sigma^{-1/2} - I_p\| \leq \epsilon \leq 1/2, \quad \text{tr}(\Sigma^{-1/2} \Sigma^b \Sigma^{-1/2} - I_p)^2 \leq \square^2 \quad (6.33)$$

for some $\epsilon \leq 1/2$ and $\square \geq 0$. The next lemma bounds from above the Kullback-Leibler divergence between two normal distributions.

Lemma 6.6.5. *Let $\mathbb{P}_0 = \mathcal{N}(\mathbf{b}, \Sigma)$ and $\mathbb{P}_1 = \mathcal{N}(\mathbf{b}^b, \Sigma^b)$ some non-degenerated matrices Σ and Σ^b . If*

$$\|\Sigma^{-1/2} \Sigma^b \Sigma^{-1/2} - I_p\| \leq 1/2, \quad \text{tr}\left\{(\Sigma^{-1/2} \Sigma^b \Sigma^{-1/2} - I_p)^2\right\} \leq \square^2,$$

then

$$\mathcal{K}(\mathbb{P}_0, \mathbb{P}_1) = -\mathbb{E}_0 \log \frac{d\mathbb{P}_1}{d\mathbb{P}_0} \leq \frac{\square^2}{2} + \frac{1}{2}(\mathbf{b} - \mathbf{b}^b)^\top \Sigma^b (\mathbf{b} - \mathbf{b}^b).$$

For any measurable set $A \subset \mathbb{R}^p$, it holds

$$|\mathbb{P}_0(A) - \mathbb{P}_1(A)| \leq \sqrt{\mathcal{K}(\mathbb{P}_0, \mathbb{P}_1)/2}.$$

Proof. The change of variables $\mathbf{u} = \Sigma^{-1/2}(\mathbf{x} - \mathbf{b})$ reduces the general case to the situation when \mathbb{P}_0 is standard normal in \mathbb{R}^p while $\mathbb{P}_1 = \mathcal{N}(\boldsymbol{\beta}, B)$ with $\boldsymbol{\beta} = \Sigma^{1/2}(\mathbf{b}^b - \mathbf{b})$ and $B \stackrel{\text{def}}{=} \Sigma^{-1/2} \Sigma^b \Sigma^{-1/2}$

$$2 \log \frac{d\mathbb{P}_1}{d\mathbb{P}_0}(\boldsymbol{\gamma}) = \log \det(B) - (\boldsymbol{\gamma} - \boldsymbol{\beta})^\top B (\boldsymbol{\gamma} - \boldsymbol{\beta}) + \|\boldsymbol{\gamma}\|^2$$

with $\boldsymbol{\gamma}$ standard normal and

$$2\mathcal{K}(\mathbb{P}_0, \mathbb{P}_1) = -2\mathbb{E}_0 \log \frac{d\mathbb{P}_1}{d\mathbb{P}_0} = -\log \det(B) + \text{tr}(B - I_p) + \boldsymbol{\beta}^\top B \boldsymbol{\beta}. \quad (6.34)$$

Let a_j be the j th eigenvalue of $B - I_p$. The condition $\|B - I_p\| \leq 1/2$ yields $|a_j| \leq 1/2$ and

$$\begin{aligned} 2\mathcal{K}(\mathbb{P}_0, \mathbb{P}_1) &= \boldsymbol{\beta}^\top B \boldsymbol{\beta} + \sum_{j=1}^p \{a_j - \log(1 + a_j)\} \\ &\leq \boldsymbol{\beta}^\top B \boldsymbol{\beta} + \sum_{j=1}^p a_j^2 \\ &\leq \boldsymbol{\beta}^\top B \boldsymbol{\beta} + \text{tr}(B - I_p)^2 \leq \boldsymbol{\beta}^\top B \boldsymbol{\beta} + \square^2. \end{aligned}$$

This implies by Pinsker's inequality

$$\sup_A |\mathbb{P}_0(A) - \mathbb{P}_1(A)| \leq \sqrt{\frac{1}{2}\mathcal{K}(\mathbb{P}_0, \mathbb{P}_1)} \leq \frac{1}{2}\sqrt{\square^2 + \boldsymbol{\beta}^\top B \boldsymbol{\beta}} \quad (6.35)$$

as required.

Theorem 6.6.2. *Let two p -dimensional zero mean Gaussian vectors $\boldsymbol{\xi} \sim \mathcal{N}(0, \Sigma)$ and $\boldsymbol{\xi}^b \sim \mathcal{N}(0, \Sigma^b)$ be given, and (6.33) holds. Then for any mapping $T: \mathbb{R}^p \rightarrow \mathbb{R}^M$ and any set of values (q_η) , the random vectors $\mathbf{X} = T(\boldsymbol{\xi})$ and $\mathbf{Y} = T(\boldsymbol{\xi}^b)$ fulfill*

$$|\mathbb{P}(\max_\eta X_\eta - q_\eta > 0) - \mathbb{P}(\max_\eta Y_\eta - q_\eta > 0)| \leq \square/2.$$

Proof. We simply apply the result of the lemma to the set $A = \{\mathbf{x} \in \mathbb{R}^p : T(\mathbf{x}) \leq \mathbf{z}\}$.

As a corollary, for the special case with $\boldsymbol{\beta} \equiv 0$, we bound for any Borel set $A \subset \mathbb{R}^M$

$$|\mathbb{P}(T(\boldsymbol{\xi}) \in A) - \mathbb{P}(T(\boldsymbol{\xi}^b) \in A)| \leq \square/2.$$

Notice that the operator norm bound

$$\|\Sigma^{-1/2}\Sigma^b\Sigma^{-1/2} - I_p\| \leq \epsilon$$

implies for $B = \Sigma^{-1/2}\Sigma^b\Sigma^{-1/2}$

$$\text{tr}(B - I_p)^2 \leq p\epsilon^2, \quad \boldsymbol{\beta}^\top B \boldsymbol{\beta} \leq (1 + \epsilon)\|\boldsymbol{\beta}\|^2.$$

Interestingly, this method can be used for obtaining an anti-concentration bound in the case of a homogeneous mapping $T: \mathbb{R}^p \rightarrow \mathbb{R}^M$.

Theorem 6.6.3. *Let $\boldsymbol{\xi} \sim \mathcal{N}(0, \Sigma)$ be a Gaussian vector in \mathbb{R}^p . For any homogeneous mapping $T: \mathbb{R}^p \rightarrow \mathbb{R}^M$, and for any $q > 0$ and Δ satisfying $0 \leq \Delta/q \leq 1$, it holds*

$$\mathbb{P}(\max_\eta T_\eta(\boldsymbol{\xi}) \geq q) - \mathbb{P}(\max_\eta T_\eta(\boldsymbol{\xi}) \geq q + \Delta) \leq \Delta q^{-1}\sqrt{p/2}.$$

Moreover, if $\boldsymbol{\xi}^b \sim \mathcal{N}(0, \Sigma^b)$ is another Gaussian vector and (6.33) holds with $\epsilon \leq 1/2$ and some $\square \geq 0$, then

$$|\mathbb{P}(\max_{\boldsymbol{\eta}} T_{\boldsymbol{\eta}}(\boldsymbol{\xi}) \geq q) - \mathbb{P}(\max_{\boldsymbol{\eta}} T_{\boldsymbol{\eta}}(\boldsymbol{\xi}^b) \geq q + \Delta)| \leq \square/2 + \Delta q^{-1} \sqrt{p/2}. \quad (6.36)$$

Proof. Given q and Δ , define $\boldsymbol{\xi}^b = q/(q + \Delta) \boldsymbol{\xi}$. It holds by homogeneity of T

$$\mathbb{P}(\max_{\boldsymbol{\eta}} T_{\boldsymbol{\eta}}(\boldsymbol{\xi}) \geq q + \Delta) = \mathbb{P}(\max_{\boldsymbol{\eta}} T_{\boldsymbol{\eta}}(\boldsymbol{\xi}^b) \geq q).$$

It is obvious that $\text{Var}(\boldsymbol{\xi}^b) = (1 + \Delta/q)^{-2} \Sigma$. Now it holds for the KL-divergence between $\boldsymbol{\xi}$ and $\boldsymbol{\xi}^b$

$$\mathcal{K}(\mathbb{P}_{\boldsymbol{\xi}}, \mathbb{P}_{\boldsymbol{\xi}^b}) = \frac{p}{2} \{2\Delta/q + (\Delta/q)^2 - 2 \log(1 + \Delta/q)\} \leq p(\Delta/q)^2. \quad (6.37)$$

Here we used that $\log(1 + \rho) \leq \rho - \rho^2/2$ for $\rho \leq 1$. Now Pinsker's bound (6.35) implies

$$\begin{aligned} & |\mathbb{P}(\max_{\boldsymbol{\eta}} T_{\boldsymbol{\eta}}(\boldsymbol{\xi}) \geq q) - \mathbb{P}(\max_{\boldsymbol{\eta}} T_{\boldsymbol{\eta}}(\boldsymbol{\xi}^b) \geq q + \Delta)| \\ & \leq \mathbb{P}(\max_{\boldsymbol{\eta}} T_{\boldsymbol{\eta}}(\boldsymbol{\xi}) \geq q) - \mathbb{P}(\max_{\boldsymbol{\eta}} T_{\boldsymbol{\eta}}(\boldsymbol{\xi}) \geq q + \Delta) \\ & \quad + |\mathbb{P}(\max_{\boldsymbol{\eta}} T_{\boldsymbol{\eta}}(\boldsymbol{\xi}) \geq q + \Delta) - \mathbb{P}(\max_{\boldsymbol{\eta}} T_{\boldsymbol{\eta}}(\boldsymbol{\xi}^b) \geq q + \Delta)| \\ & \leq \square/2 + \Delta q^{-1} \sqrt{p/2} \end{aligned}$$

and (6.36) follows.

We also present a simple corollary of the above result which concerns the change in the expectation $\mathbb{E}f(\boldsymbol{\xi})$ for a bounded function f .

Lemma 6.6.6. *Let $\boldsymbol{\xi} \sim \mathcal{N}(0, \Sigma)$ and $\boldsymbol{\xi}^b \sim \mathcal{N}(0, \Sigma^b)$, where Σ, Σ^b satisfy (6.33). For any function f on \mathbb{R}^p with $|f(\mathbf{x})| \leq 1$, and any $\delta > 0$ it holds*

$$|\mathbb{E}f(\boldsymbol{\xi}) - \mathbb{E}f(\boldsymbol{\xi}^b)| \leq \square. \quad (6.38)$$

Also, for any $\delta \geq 0$

$$|\mathbb{E}f(\boldsymbol{\xi}) - \mathbb{E}f((1 + \delta)\boldsymbol{\xi})| \leq \delta \sqrt{2p}. \quad (6.39)$$

Proof. in view of $|f(\mathbf{x})| \leq 1$, it holds

$$|\mathbb{E}f(\boldsymbol{\xi}) - \mathbb{E}f(\boldsymbol{\xi}^b)| \leq \int |f(\mathbf{x})| \cdot |\phi_{\boldsymbol{\xi}}(\mathbf{x}) - \phi_{\boldsymbol{\xi}^b}(\mathbf{x})| d\mathbf{x} \leq \int |\phi_{\boldsymbol{\xi}}(\mathbf{x}) - \phi_{\boldsymbol{\xi}^b}(\mathbf{x})| d\mathbf{x}.$$

One more use of Pinsker's inequality yields

$$\int |\phi_{\boldsymbol{\xi}}(\mathbf{x}) - \phi_{\boldsymbol{\xi}^b}(\mathbf{x})| d\mathbf{x} = 2\|\mathbb{P}_{\boldsymbol{\xi}} - \mathbb{P}_{\boldsymbol{\xi}^b}\|_{TV} \leq \sqrt{2\mathcal{K}(\mathbb{P}_{\boldsymbol{\xi}}, \mathbb{P}_{\boldsymbol{\xi}^b})},$$

and the assertion (6.39) follows by $2\mathcal{K}(\mathbb{P}_{\boldsymbol{\xi}}, \mathbb{P}_{\boldsymbol{\xi}^b}) \leq \square^2$. It remains to note that for $\Sigma^b = (1 + \delta)^2 \Sigma$, it holds $\mathcal{K}(\mathbb{P}_{\boldsymbol{\xi}}, \mathbb{P}_{\boldsymbol{\xi}^b}) \leq \delta^2 p$; see (6.37).

Sieve Model Selection

This chapter discusses the problem of sieve model selection in the situation when no prior information about the underlying noise distribution is available. The SmA procedure from Section 2.3 requires the set of critical values \mathbf{z}_{m,m° to be fixed. Here we discuss how this can be done in a data driven way with a resampling procedure.

7.1 Sieve SmA procedure

We consider a general parametric setup $\mathbf{Y} \sim \mathbb{P} \in (\mathbb{P}_\theta)$, where the parameter θ is high or infinite dimensional. The sieve approximations assumes that there is a growing sequence of subspaces $\Theta_1 \subset \Theta_2 \subset \dots$, one fixes a proper value m and applied the MLE $\tilde{\theta}_m$ obtained by maximization of the log-likelihood $L(\theta)$ over Θ_m :

$$\tilde{\theta}_m \stackrel{\text{def}}{=} \operatorname{argmax}_{\theta \in \Theta_m} L(\theta). \quad (7.1)$$

The main issue in applying this approach is the choice of the model parameter m . This problem was discussed in details for linear models with a quadratic log-likelihood function in Chapters 2 and 4. Now we aim at extending the SmA approach to the general sieve maximum likelihood setup.

Define the sieve target

$$\theta_m^* = \operatorname{argmax}_{\theta \in \Theta_m} \mathbb{E}L(\theta).$$

The Fisher Theorem claims that the sieve MLE $\tilde{\theta}_m$ estimates θ_m^* with the parametric accuracy corresponding to the parameter dimension p_m of the subset Θ_m . The value θ_m^* differs in general from θ^* yielding some sieve bias which decreases with m . At the same time, the use of large m leads to a rather complicated parametric problem (7.1) with p_m parameters. Therefore, a proper model choice has to balance the parametric complexity within the sieve model Θ_m and the bias occurring by replacing the full model by its approximation.

Similarly to the linear case we suppose to be given by a loss weighting matrix W and measure the loss of estimation by $\|W(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*)\|$. The quadratic risk is defined by its expectation

$$\mathbb{E}\|W(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*)\|^2.$$

More generally, one can consider polynomial loss function $\|W(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*)\|^q$ for a given $q \geq 0$. For the quadratic risk, one can use the bias-variance decomposition

$$\mathbb{E}\|W(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*)\|^2 = \|W(\mathbb{E}\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*)\|^2 + \text{tr}\{\text{Var}(W\tilde{\boldsymbol{\theta}}_m)\}.$$

Below we use a slightly different decomposition based on the Fisher expansion. Define

$$\begin{aligned} D_m^2 &= -\nabla_m^2 \mathbb{E}L(\boldsymbol{\theta}^*), \\ \nabla &= \nabla L(\boldsymbol{\theta}^*). \end{aligned}$$

The Fisher expansion for the largest sieve $m = \mathbf{M}$ yields a similar statement for each $m < \mathbf{M}$: on a random set of probability at least $1 - e^{-x}$

$$\|D_m(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}_m^*) - D_m^{-1}\nabla\| \leq \diamond(\mathbf{x})$$

where $\diamond(\mathbf{x}) = \diamond(\mathbf{r}_{\mathbf{M}}, \mathbf{x})$ is the error of approximation in the \mathbf{M} sieve.

Now consider the estimation loss with a weighting matrix W satisfying

$$\|WD_m^{-1}\| \leq 1.$$

Then

$$\|W(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}_m^*) - WD_m^{-2}\nabla\| = \|WD_m^{-1}\{D_m(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}_m^*) - D_m^{-1}\nabla\}\| \leq \diamond(\mathbf{x})$$

For any two $m > m^\circ$,

$$\|W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ}) - W(\boldsymbol{\theta}_m^* - \boldsymbol{\theta}_{m^\circ}^*) - W(D_m^{-2} - D_{m^\circ}^{-2})\nabla\| \leq 2\diamond(\mathbf{x}).$$

This expansion can be rewritten as

$$\|W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ}) - \mathbf{b}_{m,m^\circ} - \boldsymbol{\xi}_{m,m^\circ}\| \leq 2\diamond(\mathbf{x}) \quad (7.2)$$

with

$$\begin{aligned} \mathbf{b}_{m,m^\circ} &\stackrel{\text{def}}{=} W(\boldsymbol{\theta}_m^* - \boldsymbol{\theta}_{m^\circ}^*), \\ \boldsymbol{\xi}_{m,m^\circ} &\stackrel{\text{def}}{=} W(D_m^{-2} - D_{m^\circ}^{-2})\nabla. \end{aligned} \quad (7.3)$$

In the linear case, the expansion $W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ}) - \mathbf{b}_{m,m^\circ} - \boldsymbol{\xi}_{m,m^\circ} = 0$ is exact. The formula (7.2) is an extension to a general nonlinear regular case.

Now we suppose that the error term is small enough and proceed as if this expansion is identity. The SmA procedure selects the “smallest accepted” with the acceptance rule

$$m^\circ \text{ is accepted if } \|W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ})\| \leq \mathbf{z}_{m,m^\circ} \quad \forall m \in \mathcal{M}(m^\circ).$$

The critical values \mathbf{z}_{m,m° have to be selected to ensure the propagation property: if there is no bias for $m > m^\circ$ then the procedure should not reject m° .

Below we discuss how these values can be selected by resampling methods.

7.2 Resampling methods for parameter tuning in generalized regression

This section explains the choice of the critical values \mathbf{z}_{m,m° for the generalized regression model by a multiplier bootstrap procedure.

7.2.1 Generalized regression

We consider a sample $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ with independent observations Y_i .

Our parametric assumption concerns the marginal distribution of each observation Y_i and the structure of the regression function f .

We suppose that the distribution P_i of each observation Y_i belongs to a given family $\mathcal{P} = (P_{\mathbf{v}})$ and the value of this parameter is an unknown function f of the regressor \mathbf{X}_i . We write this relation in the form

$$Y_i \sim P_{f(\mathbf{X}_i)}. \quad (7.4)$$

Further we suppose the family \mathcal{P} to be regular and dominated by a sigma-finite measure μ_0 . By $\ell(y, \mathbf{v})$ we denote the corresponding log-density function: $\ell(y, \mathbf{v}) \stackrel{\text{def}}{=} \log \frac{dP_{\mathbf{v}}}{d\mu_0}(y)$.

The regression function f will be modelled by a linear expansion

$$f(\mathbf{x}) = f(\mathbf{x}, \boldsymbol{\theta}) = \sum_{j=1}^p \theta_j^* \psi_j(\mathbf{x}) \quad (7.5)$$

for a given basis system $\{\psi_j\}$. For simplicity this basis will be considered finite: $j \leq p$ for a finite p . We write this expansion in the vector form $\mathbf{f}^* = \boldsymbol{\Psi}^\top \boldsymbol{\theta}^*$.

Our parametric assumptions (7.4) about the marginal distribution of each Y_i and (7.5) about the structure of the regression function f lead to the log-likelihood function

$$L(\boldsymbol{\theta}) = \sum_{i=1}^n \ell(Y_i, f(X_i, \boldsymbol{\theta})) = \sum_{i=1}^n \ell(Y_i, \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}) \quad (7.6)$$

and the MLE $\tilde{\boldsymbol{\theta}}$ maximizes $L(\boldsymbol{\theta})$ over the large set of all feasible $\boldsymbol{\theta}$ -values. The sieve approach leads to the family of estimates $\tilde{\boldsymbol{\theta}}_m$ each of them is defined by restricting the parameter set to a subset Θ_m in which only first m components of $\boldsymbol{\theta}$ are varying. This means that for $\boldsymbol{\theta} \in \Theta_m$, the expansion (7.5) reads as

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum_{j=1}^m \theta_j^* \psi_j(\mathbf{x}). \quad (7.7)$$

Exercise 7.2.1. Write the sieve MLE $\tilde{\boldsymbol{\theta}}_m$ for the generalized linear regression (7.7) in

- Poisson regression with canonical parameter $\boldsymbol{v} = 1/\lambda$, where λ is the Poisson intensity parameter;
- logit (Bernoulli canonical) model

Describe in each case the Fisher matrices D_m , the score ∇_m , and the stochastic term $\boldsymbol{\xi}_{m,m^\circ}$ in expansion (7.3).

7.2.2 Multiplier bootstrap

We consider now the SmA procedure based on the family of estimates $\tilde{\boldsymbol{\theta}}_m$ and discuss how the critical values \mathbf{z}_{m,m° can be selected in a data-driven way. In what follows we suppose that the data sample \mathbf{Y} is fixed as well as the corresponding feature vectors $\boldsymbol{\Psi}_i$. This means that the critical values will be computed given data and are in general data-dependent.

Introduce the central object of analysis - the weighted log-likelihood $L^b(\boldsymbol{\theta})$ defined via the family of random weights w_i^b :

$$L^b(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \sum_{i=1}^n \ell(Y_i, f(X_i, \boldsymbol{\theta})) w_i^b. \quad (7.8)$$

The weights w_i^b are assumed i.i.d. given the data \mathbf{Y} with $\mathbb{E}w_i^b = \text{Var} w_i^b = 1$. Two leading examples of such weights are Gaussian weights with $w_i^b \sim \mathcal{N}(1, 1)$ and the exponential weights $w_i^b \sim \text{Exp}(1)$. Note the the expression (7.8) looks very similar to (7.6) but there is an essential difference in the probabilistic nature of them. The original log-likelihood is a function of the random data \mathbf{Y} , its distribution is described via the unknown data distribution. In the expression (7.8) the data \mathbf{Y} are considered as fixed and non-random, the only random element there is a collection of weights w_i^b with known distribution. In particular, if w_i^b are i.i.d. normal then $L^b(\boldsymbol{\theta})$ is also normal as a linear combination of normal r.v.s. One can conclude that $L(\boldsymbol{\theta})$ and $L^b(\boldsymbol{\theta})$ are living on two different probability spaces and have very different properties. The link between these two worlds (of real and of bootstrap ones) is given by a very simple observation

$$\mathbb{E}^b L^b(\boldsymbol{\theta}) \equiv L(\boldsymbol{\theta}).$$

Here and everywhere below \mathbb{P}^b means the distribution of the weights w_i^b given the data \mathbf{Y} . By \mathbb{E}^b we denote the corresponding expectation. The main benefit of considering the measure \mathbb{P}^b is that it is completely known.

Now we interpret $L^b(\boldsymbol{\theta})$ as a log-likelihood process and define

$$\tilde{\boldsymbol{\theta}}^b \stackrel{\text{def}}{=} \underset{\boldsymbol{\theta}}{\operatorname{argmax}} L^b(\boldsymbol{\theta}). \quad (7.9)$$

Similarly one can define $\tilde{\boldsymbol{\theta}}_m^b$ by maximizing $L^b(\boldsymbol{\theta})$ over Θ_m :

$$\tilde{\boldsymbol{\theta}}_m^b \stackrel{\text{def}}{=} \underset{\boldsymbol{\theta} \in \Theta_m}{\operatorname{argmax}} L^b(\boldsymbol{\theta}). \quad (7.10)$$

The value $\tilde{\boldsymbol{\theta}}_m^b$ formally depends on both the data \mathbf{Y} and the weights w_i^b but we consider its conditional distribution given the data \mathbf{Y} with the hope that this distribution somehow reflects the original distribution of $\tilde{\boldsymbol{\theta}}_m$. In fact, for running the SmA procedure we only need to know the distribution (tail function) of stochastic parts $\boldsymbol{\xi}_{m,m^\circ}$ of considered test statistics $W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ})$. Define their analog in the bootstrap world:

$$\mathbb{T}_{m,m^\circ}^b \stackrel{\text{def}}{=} \|W(\tilde{\boldsymbol{\theta}}_m^b - \tilde{\boldsymbol{\theta}}_{m^\circ}^b)\|.$$

Now we use the great advantage of considering the bootstrap world: the distribution of the $\tilde{\boldsymbol{\theta}}_m^b$ is known, in particular, one can compute the expectation $\mathbb{E}^b(\tilde{\boldsymbol{\theta}}_m^b)$ or use the knowledge of the true value in the bootstrap model which coincides with the real world estimate $\tilde{\boldsymbol{\theta}}_m$. The corresponding stochastic component $\boldsymbol{\xi}_{m,m^\circ}^b$ is given by

$$\boldsymbol{\xi}_{m,m^\circ}^b \stackrel{\text{def}}{=} W\{\tilde{\boldsymbol{\theta}}_m^b - \tilde{\boldsymbol{\theta}}_{m^\circ}^b - \mathbb{E}^b(\tilde{\boldsymbol{\theta}}_m^b - \tilde{\boldsymbol{\theta}}_{m^\circ}^b)\}.$$

Now define the tail function $z_{m,m^\circ}^b(\mathbf{x})$ by the relation

$$\mathbb{P}^b(\|\boldsymbol{\xi}_{m,m^\circ}^b\| \geq z_{m,m^\circ}^b(\mathbf{x})) = e^{-\mathbf{x}} \quad (7.11)$$

and further proceed as in the case of known functions z_{m,m° 's. In particular, the multiplicity correction $q_{m^\circ} = q_{m^\circ}(\mathbf{x})$ is defined as the smallest value q satisfying the relation

$$\mathbb{P}^b\left(\max_{m > m^\circ} \left\{ \|\boldsymbol{\xi}_{m,m^\circ}^b\| - z_{m,m^\circ}^b(\mathbf{x} + q) \right\} \geq 0\right) \leq e^{-\mathbf{x}}. \quad (7.12)$$

Finally the SmA procedure is applied with

$$\mathbf{z}_{m,m^\circ} \stackrel{\text{def}}{=} z_{m,m^\circ}^b(\mathbf{x} + q_{m^\circ}) + \beta \sqrt{\mathbb{P}_{m,m^\circ}^b}. \quad (7.13)$$

7.2.3 Numerical issues

It was mentioned many times that the joint distribution of the test statistics \mathbb{T}_{m,m°^b under \mathbb{P}^b is known. However, its analytic study is a hard task even if the weights w_i^b are normal. Similarly to linear regression case, one can use a numerical Monte Carlo procedure for evaluating the tail functions of \mathbb{T}_{m,m°^b and the multiplicity corrections q_{m° . The approach can be described as follows:

- generate B samples of weights $\mathbf{w}^b = (w_i^b)$;
- for each sample, compute and store $\tilde{\boldsymbol{\theta}}_m^b = \tilde{\boldsymbol{\theta}}_m^b(\mathbf{w}^b)$ for all m ;
- compute the bootstrap empirical means

$$\mathbb{E}^b(\tilde{\boldsymbol{\theta}}_m^b) \stackrel{\text{def}}{=} \frac{1}{B} \sum_{\mathbf{w}^b} \tilde{\boldsymbol{\theta}}_m^b(\mathbf{w}^b);$$

- For each $m > m^\circ$ and all feasible \mathbf{x} , compute the tail functions $z_{m,m^\circ}(\mathbf{x})$ by the relations (7.11) when \mathbb{P}^b is replaced by its bootstrap empirical distribution:

$$\frac{1}{B} \sum_{\mathbf{w}^b} \mathbb{I}(\|\boldsymbol{\xi}_{m,m^\circ}^b(\mathbf{w}^b)\| \geq z_{m,m^\circ}^b(\mathbf{x})) \leq e^{-\mathbf{x}}$$

with

$$\boldsymbol{\xi}_{m,m^\circ}^b(\mathbf{w}^b) \stackrel{\text{def}}{=} W \left\{ \tilde{\boldsymbol{\theta}}_m^b(\mathbf{w}^b) - \tilde{\boldsymbol{\theta}}_{m^\circ}^b(\mathbf{w}^b) - \mathbb{E}^b(\tilde{\boldsymbol{\theta}}_m^b - \tilde{\boldsymbol{\theta}}_{m^\circ}^b) \right\}.$$

Alternatively,

$$\boldsymbol{\xi}_{m,m^\circ}^b(\mathbf{w}^b) \stackrel{\text{def}}{=} W \left\{ \tilde{\boldsymbol{\theta}}_m^b(\mathbf{w}^b) - \tilde{\boldsymbol{\theta}}_{m^\circ}^b(\mathbf{w}^b) - (\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ}) \right\}.$$

Also estimate

$$\mathbf{P}_{m,m^\circ}^b = \text{Var}^b(\boldsymbol{\xi}_{m,m^\circ}^b) \approx \text{tr} \left\{ \frac{1}{B} \sum_{\mathbf{w}^b} \boldsymbol{\xi}_{m,m^\circ}^b(\mathbf{w}^b) \boldsymbol{\xi}_{m,m^\circ}^b(\mathbf{w}^b)^\top \right\}.$$

- Compute for each m° the multiplicity corrections q_{m° as the smallest value with

$$\frac{1}{B} \sum_{\mathbf{w}^b} \mathbb{I} \left(\max_{m > m^\circ} \left\{ \|\boldsymbol{\xi}_{m,m^\circ}^b(\mathbf{w}^b)\| - z_{m,m^\circ}^b(\mathbf{x} + q_{m^\circ}) \right\} \geq 0 \right) \leq e^{-\mathbf{x}}.$$

Complexity of this procedure is mainly determined by the complexity of computing the family of estimates $\tilde{\boldsymbol{\theta}}_m^b(\mathbf{w}^b)$ for each bootstrap sample \mathbf{w}^b . This has to be done many many times to reduce the Monte-Carlo error. Each estimate $\tilde{\boldsymbol{\theta}}_m^b$ is given implicitly via the optimization problem (7.10). This can be a hard task especially if the parameter dimension p is large. Note however one important feature of the bootstrap world: we know the true value, which the estimate $\tilde{\boldsymbol{\theta}}$ computed from the original data. This true

value is actually there target of estimation from the resampled data and automatically it is a very good starting point of the estimation procedure; see further details of reducing the computational burden below in Section ??.

In the next section we discuss how the proposed procedure can be justified. First we consider the linear Gaussian case. Then we extend the results to the case of linear models with non-Gaussian errors. Finally we consider the general case of sieve model selection and show how the bootstrap validity can be derived from the Fisher expansions (7.2) in the asymptotic sense for a reasonably large sample size n .

7.3 Why does it work? Linear Gaussian case

This section studies the validity of the bootstrap procedure for a linear Gaussian model

$$\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, I_n).$$

In words, the model assumes linear systematic dependence $\mathbb{E}\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^*$ and standard Gaussian noise $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, I_n)$. In reality, both assumptions can be misspecified. We aim to check to which extent the proposed procedure is robust w.r.t. possible model misspecifications. First we make the analysis for just one estimate, then we discuss how the results apply to the SmA model selection procedure.

7.3.1 Small modeling bias condition

This section discusses in the setup of linear Gaussian regression whether the bootstrap procedure does the job if the model assumptions are misspecified. Our study admits that the model is misspecified: the linear dependence can be violated and we admit a inhomogeneous Gaussian noise, that is, $\boldsymbol{\varepsilon}$ be zero mean normal with a diagonal covariance matrix $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$. The main question is whether the bootstrap procedure mimics well the noise distribution and is robust with respect to a moderate deviation from the linear modeling assumption. The final result presents some sufficient conditions on the deviation $\mathbb{E}\mathbf{Y} - \Psi^\top \boldsymbol{\theta}^*$ and on the variance coefficients σ_i^2 which ensure the bootstrap validity.

Let $\tilde{\boldsymbol{\theta}} = \mathcal{S}\mathbf{Y} = (\Psi\Psi^\top)^{-1}\Psi\mathbf{Y}$ be the qMLE for our linear Gaussian parametric model. The corresponding target value is $\boldsymbol{\theta}^* = (\Psi\Psi^\top)^{-1}\Psi\mathbf{f}^*$. Below we suppose that $\boldsymbol{\theta}^* = 0$ which does not restrict generality, because this case can be obtained by reparametrization (a linear shift of the parameter $\boldsymbol{\theta}$). For the real world estimate $\tilde{\boldsymbol{\theta}} = D^{-2}\Psi\mathbf{Y}$ holds

$$\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^* = D^{-2}\Psi(\mathbf{Y} - \Psi^\top \boldsymbol{\theta}^*) = D^{-2}\Psi\boldsymbol{\varepsilon} = D^{-2}\nabla$$

with $\nabla = \Psi \boldsymbol{\varepsilon}$. Under the assumption of an inhomogeneous Gaussian noise with $\boldsymbol{\varepsilon}$ being zero mean normal with a diagonal covariance matrix $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$, the score $\nabla = \Psi \boldsymbol{\varepsilon}$ is zero mean normal as well and

$$V^2 \stackrel{\text{def}}{=} \text{Var}(\nabla) = \mathbb{E}(\nabla \nabla^\top) = \Psi \text{Var}(\boldsymbol{\varepsilon}) \Psi^\top = \Psi \Sigma \Psi^\top. \quad (7.14)$$

Now we check what happens in the bootstrap world. One can easily see that

$$L^b(\boldsymbol{\theta}) = -\frac{1}{2} \sum_i (Y_i - \Psi_i^\top \boldsymbol{\theta})^2 w_i^b + R,$$

where the remainder R does not depend on $\boldsymbol{\theta}$. The corresponding estimates $\tilde{\boldsymbol{\theta}}^b$ read as

$$\tilde{\boldsymbol{\theta}}^b = (\Psi \mathcal{W}^b \Psi^\top)^{-1} \Psi \mathcal{W}^b \mathbf{Y}.$$

Alternatively one can consider the estimate

$$\tilde{\boldsymbol{\theta}}^b = (\Psi \Psi^\top)^{-1} \Psi \mathcal{W}^b \mathbf{Y} = D^{-2} \Psi \mathcal{W}^b \mathbf{Y} \quad (7.15)$$

with $D^2 = \Psi \Psi^\top$. These two versions are very close to each other because the matrix $\Psi \Psi^\top$ and its bootstrap counterpart $\Psi \mathcal{W}^b \Psi^\top$ are close to each other. The version (7.15) is preferable for the bootstrap procedure because one can use the same matrix $D^{-2} \Psi$ for all bootstrap runs, otherwise one has to recompute it for each \mathcal{W}^b . The theoretical study is also simpler for the version (7.15). By construction $\mathbb{E}^b(w_i^b) = 1$ yielding $\mathbb{E}^b \mathcal{W}^b = I_n$ and

$$\tilde{\boldsymbol{\theta}}^b - \tilde{\boldsymbol{\theta}} = D^{-2} \Psi (\mathcal{W}^b - I_p) \mathbf{Y} = D^{-2} \nabla^b,$$

where for $\mathcal{E}^b = \mathcal{W}^b - I_p$

$$\nabla^b \stackrel{\text{def}}{=} \Psi \mathcal{E}^b \mathbf{Y} \quad (7.16)$$

and $\mathbb{E}^b(\nabla^b) = 0$.

Further, the model equation $\mathbf{Y} = \mathbf{f}^* + \boldsymbol{\varepsilon}$ yields the decomposition of the bootstrap score $\nabla^b = \Psi \mathcal{E}^b \mathbf{Y}$ from (7.16) as

$$\nabla^b = \Psi \mathcal{E}^b \boldsymbol{\varepsilon} + \Psi \mathcal{E}^b \mathbf{f}^*.$$

Below we show that the first term $\Psi \mathcal{E}^b \boldsymbol{\varepsilon}$ nicely mimics in distribution the original score ∇ , while the second term is small under the so called ‘‘small modeling bias’’ condition. Let the weights w_i^b be normal $\mathcal{N}(1, 1)$. Then each $e_i^b = w_i^b - 1$ is standard normal and $\Psi \mathcal{E}^b \boldsymbol{\varepsilon}$ is also zero mean normal under \mathbb{P}^b with

$$\text{Var}^b(\Psi \mathcal{E}^b \varepsilon) = \mathbb{E}^b \{ \Psi \mathcal{E}^b \varepsilon (\Psi \mathcal{E}^b \varepsilon)^\top \} = \Psi \text{Var}(\mathcal{E}^b \varepsilon) \Psi^\top = \Psi \text{diag}(\varepsilon \cdot \varepsilon) \Psi^\top. \quad (7.17)$$

Exercise 7.3.1. Let $\mathcal{E}^b \varepsilon$ be the vector with the entries $w_i^b \varepsilon_i$.

- Check that

$$\text{Var}^b(\mathcal{E}^b \varepsilon) = \text{diag}(\varepsilon \cdot \varepsilon),$$

where $\text{diag}(\varepsilon \cdot \varepsilon) = \text{diag}(\varepsilon_1^2, \dots, \varepsilon_n^2)$ is the diagonal matrix with the entries ε_i^2 .

- Check the formula (7.17)

For the study we need a kind of Lindeberg condition which requires for each vector Ψ_i to be small relative to the standardized sum $\Sigma^{1/2}$. More precisely, define

$$\delta_\Psi \stackrel{\text{def}}{=} \max_{i=1, \dots, n} \|V^{-1} \Psi_i\| \sigma_i. \quad (7.18)$$

Exercise 7.3.2. Consider the case of a regular design, when all the Ψ_i 's belongs to a compact subset \mathcal{X} of \mathbb{R}^p , the matrix

$$d_\Psi^2 \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \Psi_i \Psi_i^\top$$

is non-degenerate, and the value

$$a_\Psi \stackrel{\text{def}}{=} \max_{i=1, \dots, n} \|d_\Psi^{-1} \Psi_i\| \quad (7.19)$$

is finite. Moreover, let the ratio $\lambda_{\max}(\Sigma)/\lambda_{\min}(\Sigma)$ of the maximal and minimal values of σ_i^2 's be bounded by a_Σ^2 , that is, for all $i, j \leq n$

$$\sigma_i/\sigma_j \leq a_\Sigma. \quad (7.20)$$

Check that δ_Ψ from (7.18) fulfills

$$\delta_\Psi \leq a_\Psi a_\Sigma n^{-1/2}. \quad (7.21)$$

Two centered Gaussian distributions with different covariance matrices can be compared via the Pinsker inequality. The general result of Theorem 6.6.2 claims that these two distributions are close to each other if the quantity

$$\begin{aligned} \square^2 &\stackrel{\text{def}}{=} \text{tr} \left[\left\{ V^{-1} \text{Var}^b(\Psi \mathcal{E}^b \varepsilon) V^{-1} - I_p \right\}^2 \right] \\ &= \text{tr} \left[\left\{ V^{-1} \Psi \text{diag}(\varepsilon \cdot \varepsilon) \Psi^\top V^{-1} - I_p \right\}^2 \right] \end{aligned}$$

is small. For a symmetric $p \times p$ matrix A , the trace $\text{tr}(A^2)$ of A^2 coincides with its squared Frobenius norm and it can be bounded from above using the operator norm (maximal eigenvalue) $\|A\|$ of the matrix A :

$$\text{tr}(A^2) \leq p\|A\|^2.$$

We apply this bound for the matrix \mathbf{A} in the form

$$\mathbf{A} \stackrel{\text{def}}{=} V^{-1}\Psi \text{diag}(\boldsymbol{\varepsilon} \cdot \boldsymbol{\varepsilon})\Psi^\top V^{-1} - I_p.$$

Now we use that $\tilde{\boldsymbol{\varepsilon}} = \Sigma^{-1/2}\boldsymbol{\varepsilon}$ is a standard normal vector in \mathbb{R}^n . Define

$$\mathbf{U}^\top \stackrel{\text{def}}{=} V^{-1}\Psi \Sigma^{1/2}. \quad (7.22)$$

Exercise 7.3.3. Check whether the definition (7.14) implies that the columns $\mathbf{u}_j \in \mathbb{R}^n$ of the matrix \mathbf{U} are orthonormal:

$$\mathbf{u}_i^\top \mathbf{u}_j = \mathbb{I}(i = j).$$

Suppose that

$$\max_{i=1,\dots,n} \|V^{-1}\Psi_i\| \sigma_i \leq \delta_\Psi$$

Then for any unit vector $\boldsymbol{\gamma} \in \mathbb{R}^p$, the vector $\mathbf{u} = \mathbf{U}\boldsymbol{\gamma} \in \mathbb{R}^n$ fulfills $\|\mathbf{u}\| = 1$ and

$$\|\mathbf{u}\|_\infty \leq \delta_\Psi. \quad (7.23)$$

Hint: suppose that all $\sigma_i \equiv 1$. Then the Cauchy-Schwartz inequality, each component u_i of \mathbf{u} satisfies

$$u_i = \boldsymbol{\gamma}^\top V^{-1}\Psi_i \leq \|\boldsymbol{\gamma}\| \|V^{-1}\Psi_i\|.$$

With \mathbf{U} from (7.22), one can write

$$\mathbf{A} = \mathbf{U}^\top \text{diag}\{\tilde{\boldsymbol{\varepsilon}} \cdot \tilde{\boldsymbol{\varepsilon}} - 1\}\mathbf{U}. \quad (7.24)$$

The spectral norm of A can be easily bounded in the univariate case with $p = 1$. Then \mathbf{U} coincides with a unit vector $\mathbf{u} \in \mathbb{R}^n$ satisfying $\|\mathbf{u}\|_\infty \leq \delta_\Psi$ and the value

$$\mathbf{u}^\top \text{diag}\{\tilde{\boldsymbol{\varepsilon}} \cdot \tilde{\boldsymbol{\varepsilon}} - 1\}\mathbf{u} = \sum_{i=1}^n u_i^2 (\tilde{\varepsilon}_i^2 - 1)$$

can be bounded using the lower and upper bounds for centered Gaussian quadratic forms; cf. Theorem 6.6.1: with $\mathbf{v}_\mathbf{u}^2 = 2 \sum_i u_i^4$ and $a^* = \max u_i^2 = \|\mathbf{u}\|_\infty^2$

$$\mathbb{P}\left(\left|\sum_{i=1}^n u_i^2(\tilde{\varepsilon}_i^2 - 1)\right| \geq \sqrt{4v_{\mathbf{u}}^2 \mathbf{x}} \vee (6\|\mathbf{u}\|_{\infty}^2 \mathbf{x})\right) \leq e^{-\mathbf{x}}.$$

One can also bound

$$v_{\mathbf{u}}^2 = \sum_{i=1}^n u_i^4 \leq \|\mathbf{u}\|_{\infty}^2 \sum_{i=1}^n u_i^2 = \|\mathbf{u}\|_{\infty}^2 \leq \delta_{\Psi}^2.$$

Therefore, with a high probability

$$\mathbf{u}^{\top} \text{diag}\{\tilde{\varepsilon} \cdot \tilde{\varepsilon} - 1\} \mathbf{u} \leq (2\delta_{\Psi} \mathbf{x}^{1/2}) \vee (6\delta_{\Psi}^2 \mathbf{x}). \quad (7.25)$$

Exercise 7.3.4. Complete the proof of (7.25).

In the matrix case one has to increase the tail level \mathbf{x} by $\log p$. Putting together all obtained bounds yields for $3\mathbf{x}^{1/2}\delta_{\Psi} \leq 1$ with probability at least $1 - 2e^{-\mathbf{x}}$

$$\|A\| = \|\mathbf{U}^{\top} \text{diag}\{\tilde{\varepsilon} \cdot \tilde{\varepsilon} - 1\} \mathbf{U}\| \leq 2\delta_{\Psi}(\mathbf{x} + \log p)^{1/2}. \quad (7.26)$$

To be done: complete the case $p > 1$.

The bootstrap procedure in the unbiased case is justified if $\text{tr}(A^2)$ is small. In view of the bound $\text{tr}(A^2) \leq p\|A\|^2$ and (7.26), it follows from the condition “ $p\delta_{\Psi}^2(\mathbf{x} + \log p)$ is small”. Under regular noise and regular design, δ_{Ψ} is of order $n^{-1/2}$; see (7.21). Therefore, it suffices that $n^{-1}p \log(p)$ is small.

Now we consider the general case and try to evaluate the impact of the bias $\mathbf{f}^* - \Psi^{\top} \boldsymbol{\theta}^*$ which under $\boldsymbol{\theta}^* = 0$ coincides with \mathbf{f}^* . For $\nabla^b = \Psi \mathcal{E}^b \mathbf{Y} = \Psi \mathcal{E}^b(\boldsymbol{\varepsilon} + \mathbf{f}^*)$, it holds with \mathbf{U} from (7.22)

$$V^{-1} \text{Var}^b(\nabla^b) V^{-1} = \mathbf{U}^{\top} \text{diag}\{(\tilde{\varepsilon} + \mathbf{B}) \cdot (\tilde{\varepsilon} + \mathbf{B})\} \mathbf{U}, \quad (7.27)$$

where

$$\mathbf{B} = \Sigma^{-1/2} \mathbf{f}^* = \Sigma^{-1/2}(\mathbf{f}^* - \Pi \mathbf{f}^*).$$

Then

$$\begin{aligned} \mathbf{A} &= \mathbf{U}^{\top} \text{diag}\{(\tilde{\varepsilon} + \mathbf{B}) \cdot (\tilde{\varepsilon} + \mathbf{B}) - 1\} \mathbf{U} \\ &= \mathbf{U}^{\top} \text{diag}\{\tilde{\varepsilon} \cdot \tilde{\varepsilon} - 1\} \mathbf{U} + \mathbf{U}^{\top} \text{diag}\{\mathbf{B} \cdot \mathbf{B}\} \mathbf{U} + 2\mathbf{U}^{\top} \text{diag}\{\tilde{\varepsilon} \cdot \mathbf{B}\} \mathbf{U}. \end{aligned} \quad (7.28)$$

The quadratic bias term can be easily evaluated. Indeed, for any unit vector $\mathbf{u} = (u_1, \dots, u_n)^{\top} \in \mathbb{R}^n$ with $\|\mathbf{u}\|_{\infty} \leq \delta_{\Psi}$, it holds

$$\mathbf{u}^\top \text{diag}\{\mathbf{B} \cdot \mathbf{B}\} \mathbf{u} = \sum_{i=1}^n u_i^2 b_i^2 \leq \|\mathbf{u}\|_\infty^2 \|\mathbf{B}\|^2 \leq \delta_\Psi^2 \|\mathbf{B}\|^2. \quad (7.29)$$

This also implies by (7.23)

$$\begin{aligned} \|\mathbf{U} \text{diag}\{\mathbf{B} \cdot \mathbf{B}\} \mathbf{U}^\top\| &= \sup_{\boldsymbol{\gamma} \in \mathbb{R}^p: \|\boldsymbol{\gamma}\|=1} \boldsymbol{\gamma}^\top \mathbf{U} \text{diag}\{\mathbf{B} \cdot \mathbf{B}\} \mathbf{U}^\top \boldsymbol{\gamma} \\ &\leq \sup_{\mathbf{u} \in \mathbb{R}^n: \|\mathbf{u}\|=1, \|\mathbf{u}\|_\infty \leq \delta_\Psi} \mathbf{u}^\top \text{diag}\{\mathbf{B} \cdot \mathbf{B}\} \mathbf{u} \leq \delta_\Psi^2 \|\mathbf{B}\|^2. \end{aligned} \quad (7.30)$$

Exercise 7.3.5. Check (7.27) and the inequalities (7.29) and (7.30).

The cross term can be easily bounded in the univariate case with $p = 1$. Then \mathbf{U} coincides with a unit vector $\mathbf{u} \in \mathbb{R}^n$ satisfying $\|\mathbf{u}\|_\infty \leq \delta_\Psi$ and

$$\begin{aligned} \mathbf{u}^\top \text{diag}\{\mathbf{B} \cdot \tilde{\boldsymbol{\varepsilon}}\} \mathbf{u} &= \sum_{i=1}^n u_i^2 b_i \tilde{\varepsilon}_i \sim \mathcal{N}(0, v_u^2), \\ v_u^2 &\stackrel{\text{def}}{=} \sum_{i=1}^n u_i^4 b_i^2 \leq \|\mathbf{u}\|_\infty^4 \|\mathbf{B}\|^2 \leq \delta_\Psi^4 \|\mathbf{B}\|^2. \end{aligned}$$

This implies that on a dominating set of probability $1 - e^{-x}$

$$\mathbf{u}^\top \text{diag}\{\mathbf{B} \cdot \tilde{\boldsymbol{\varepsilon}}\} \mathbf{u} \leq \delta_\Psi^2 \|\mathbf{B}\| z_1(\mathbf{x}). \quad (7.31)$$

In the general case $p > 1$, the result continues to apply with \mathbf{x} increased by $\log(p)$.

Putting together (7.28), (7.26), (7.30), and (7.31) yields on a random set of probability at least $1 - 3e^{-x}$

$$\|\mathbf{A}\| \leq 2(\mathbf{x} + \log p)^{1/2} \delta_\Psi + \delta_\Psi^2 \|\mathbf{B}\|^2 + 2\delta_\Psi^2 \|\mathbf{B}\| z_1(\mathbf{x} + \log p). \quad (7.32)$$

Now we evaluate the distance \square between the distribution of $\nabla = \Psi \boldsymbol{\varepsilon}$ under \mathbb{P} and of $\nabla^b = \Psi \mathcal{E}^b \mathbf{Y}$ under \mathbb{P}^b using that δ_Ψ is a small number. Remind that under regular noise and regular design, it is of order $n^{-1/2}$; see (7.21). In view of $\text{tr}(\mathbf{A}^2) \leq p \|\mathbf{A}\|^2$ and of (7.32), we ensure that \square is small if $n^{-1}p(\mathbf{x} + \log p)$ is small and $p^{1/2}n^{-1} \|\mathbf{B}\|^2$ is small as well. Under a homogeneous noise with the variance σ^2 , it holds $\|\mathbf{B}\|^2 = \sigma^{-2} \|\mathbf{f}^*\|^2$. A similar bound holds under ‘‘regular noise’’ condition (7.20). If this condition is fulfilled, it suffices that

$$\frac{p^{1/2}}{n\sigma^2} \|\mathbf{f}^* - \Psi^\top \boldsymbol{\theta}^*\|^2 \text{ is small}$$

with $\bar{\sigma}^2 \stackrel{\text{def}}{=} n^{-1} \text{tr} \Sigma$.

We are prepared to state the following result.

Theorem 7.3.1. Let $\mathbf{Y} = \mathbf{f}^* + \boldsymbol{\varepsilon}$ be a Gaussian vector in \mathbb{R}^n with independent components, $\mathbf{Y} \sim \mathcal{N}(\mathbf{f}^*, \Sigma)$ for a diagonal matrix $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$. Let also Ψ be a $p \times n$ feature matrix such that the $p \times p$ -matrix $V^2 = \Psi \Sigma \Psi^\top$ is non-degenerated and the value

$$\delta_\Psi = \max_{i=1, \dots, n} \|V^{-1} \Psi_i\| \sigma_i \quad (7.33)$$

is small. Denote by $\Pi = \Psi^\top (\Psi \Psi^\top)^{-1} \Psi$ the projector on the feature space and define

$$\mathbf{B} = \Sigma^{-1/2} (\mathbf{f}^* - \Pi \mathbf{f}^*).$$

Further, let $\mathcal{W}^b = \text{diag}\{w_1^b, \dots, w_n^b\}$ be a diagonal matrix of bootstrap $\mathcal{N}(1, 1)$ multipliers and $\mathcal{E}^b = \mathcal{W}^b - \mathbb{E}^b \mathcal{W}^b$. Then on a set $\Omega(\mathbf{x})$ with $\mathbb{P}(\Omega(\mathbf{x})) \geq 1 - 3e^{-x}$, the distribution \mathbb{Q} of the score $\nabla = \Psi \boldsymbol{\varepsilon}$ and the conditional distribution \mathbb{Q}^b of the bootstrap score $\nabla^b = \Psi \mathcal{E}^b \mathbf{Y}$ given \mathbf{Y} are related by

$$\|\mathbb{Q} - \mathbb{Q}^b\|_{TV} \leq \sqrt{p/2} \Delta \quad (7.34)$$

with

$$\Delta = 2(\mathbf{x} + \log p)^{1/2} \delta_\Psi + \delta_\Psi^2 \|\mathbf{B}\|^2 + 2\delta_\Psi^2 \|\mathbf{B}\| z_1(\mathbf{x} + \log p). \quad (7.35)$$

The same bounds applies to the distance between the distribution of the qMLE $\tilde{\boldsymbol{\theta}}$ and of the bootstrap estimate $\tilde{\boldsymbol{\theta}}^b = (\Psi \Psi^\top)^{-1} \Psi \mathcal{W}^b \mathbf{Y}$.

Furthermore, if the design Ψ and the noise $\boldsymbol{\varepsilon}$ are regular in the sense (7.19) and (7.20) with some fixed constants a_Ψ and a_Σ , then

$$\delta_\Psi \leq a_\Psi a_\Sigma n^{-1/2}.$$

The statement (7.34) yields

$$\sup_A |\mathbb{Q}(A) - \mathbb{Q}^b(A)| \leq \sqrt{p/2} \Delta, \quad (7.36)$$

where sup is taken over all measurable subsets of \mathbb{R}^p . This result has a number of remarkable corollaries. First we specify the result to the considered estimation loss $\|W(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|$.

Corollary 7.3.1. Under the conditions of Theorem 7.3.1, it holds on a set $\Omega(\mathbf{x})$ of dominating probability

$$\sup_{z>0} \left| \mathbb{P}(\|W(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \leq z) - \mathbb{P}^b(\|W(\tilde{\boldsymbol{\theta}}^b - \tilde{\boldsymbol{\theta}}\| \leq z) \right| \leq \sqrt{p/2} \Delta. \quad (7.37)$$

Exercise 7.3.6. Derive (7.37) as a special case of (7.36).

Our next corollary concerns the “no-bias” case with $\mathbf{B} \equiv 0$.

Corollary 7.3.2. *Assume the conditions of Theorem 7.3.1 including the regularity conditions (7.19) and (7.20). If the linear parametric assumption is correct, that is, if $\mathbf{f}^* = \Psi^\top \boldsymbol{\theta}^*$, then the bootstrap procedure is justified with the accuracy*

$$\square \asymp \sqrt{n^{-1}p \log(p)}.$$

In particular, the procedure is asymptotically valid if the parameter dimension $p = p_n$ satisfies

$$n^{-1}p_n \log(p_n) \rightarrow 0, \quad n \rightarrow \infty.$$

One can conclude from (7.34) and (7.35) that the bootstrap procedure does the job if the value

$$p\delta_{\Psi}^2 \|\mathbf{B}\|^2 \text{ is small.}$$

In the regular case it can be rewritten as “ $p\|\mathbf{B}\|^2/n$ is small”. As a special case we present the following corollary.

Corollary 7.3.3. *Assume the conditions of Theorem 7.3.1 including the regularity conditions (7.19) and (7.20). If \mathbf{f}^* is Sobolev-smooth with the parameter s , that is, if*

$$\frac{1}{n} \|\mathbf{B}\|^2 \leq Cp^{-2s}$$

then the bootstrap procedure is justified for $s > 1/2$ if $p = p_n \rightarrow \infty$ as $n \rightarrow \infty$.

7.3.2 The “large bias” case

The result of Theorem 7.3.1 claims that if the modeling bias $\|\mathbf{B}\|^2 = \|\Sigma^{-1/2}(\mathbf{f}^* - \Pi \mathbf{f}^*)\|^2$ is small then the bootstrap mimics well the real distribution. The next question is whether we can apply the bootstrap when the bias is large. The answer is a bit surprising: yes, we can apply, the bootstrap can be validated but the corresponding bootstrap quantiles are larger than the quantiles for the original distributions. In other words, the bootstrap procedure in the “large bias” situation becomes conservative and yields higher coverage probability than the nominal one. The reason is that the bias in the original data transfers to the additional variance in the bootstrap world. More exactly, the vector $V^{-1}\nabla = V^{-1}\Psi\boldsymbol{\varepsilon}$ is standard normal in \mathbb{R}^p under the underlying measure \mathbb{P} . At the same time, the vector $V^{-1}\nabla^b = V^{-1}\Psi\mathcal{E}^b\mathbf{Y}$ is Gaussian zero mean under \mathbb{P}^b but its variance is larger than the identity if the bias \mathbf{B} is significant.

Define

$$\mathcal{V}^2 \stackrel{\text{def}}{=} V^2 + \Psi \text{diag}\{(\mathbf{f}^* - \Pi \mathbf{f}^*) \cdot (\mathbf{f}^* - \Pi \mathbf{f}^*)\} \Psi^\top. \quad (7.38)$$

Then we can approximate the variance of ∇^b by \mathcal{V}^2 , and the latter can be significantly larger than V^2 depending on the bias $\mathbf{f}^* - \Pi \mathbf{f}^*$.

Theorem 7.3.2. *Suppose the conditions of Theorem 7.3.1. Then on a set $\Omega(\mathbf{x})$ of dominating probability, it holds with \mathcal{V}^2 from (7.38)*

$$\|\text{Var}^b(\mathcal{V}^{-1} \nabla^b) - I_p\| \leq 2\delta_\Psi(\mathbf{x} + \log p)^{1/2}. \quad (7.39)$$

Moreover, if

$$\|\mathcal{V}^{-1} V^2 \mathcal{V}^{-1}\| \leq 1 - 2\delta_\Psi(\mathbf{x} + \log p)^{1/2}, \quad (7.40)$$

then on the same set $\Omega(\mathbf{x})$

$$\text{Var}^b(V^{-1} \nabla^b) \geq I_p.$$

Proof. The second statement follows directly from the first one and the definition of \mathcal{V}^2 .

As the variance of ∇^b is larger than the variance of ∇ , the related probability of any centrally symmetric convex set A in \mathbb{R}^p is larger under \mathbb{P} than under \mathbb{P}^b .

Corollary 7.3.4. *For any convex centrally symmetric set $A \subset \mathbb{R}^p$, it holds on a set $\Omega(\mathbf{x})$ of dominating probability*

$$\mathbb{Q}(A) - \mathbb{Q}^b(A) \geq \delta_\Psi \sqrt{2p(\mathbf{x} + \log p)}. \quad (7.41)$$

Moreover, under (7.40)

$$\mathbb{Q}(A) - \mathbb{Q}^b(A) \geq 0. \quad (7.42)$$

Exercise 7.3.7. Check that (7.39) implies (7.41).

In particular, one can bound from above the deviation probabilities.

Corollary 7.3.5. *For any $W: \mathbb{R}^p \rightarrow \mathbb{R}^q$, it holds on a set $\Omega(\mathbf{x})$ of dominating probability*

$$\sup_{z>0} \left\{ \mathbb{P}(\|W(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \leq z) - \mathbb{P}^b(\|W(\tilde{\boldsymbol{\theta}}^b - \tilde{\boldsymbol{\theta}}\| \leq z) \right\} \geq \delta_\Psi \sqrt{2p(\mathbf{x} + \log p)}. \quad (7.43)$$

Moreover, under (7.40)

$$\sup_{z>0} \left\{ \mathbb{P}(\|W(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \leq z) - \mathbb{P}^b(\|W(\tilde{\boldsymbol{\theta}}^b - \tilde{\boldsymbol{\theta}}\| \leq z) \right\} \geq 0. \quad (7.44)$$

Exercise 7.3.8. Derive (7.43) from (7.41) and (7.44) from (7.42).

The inequality (7.44) implies that the $1 - \alpha$ quantile of $\|W(\tilde{\boldsymbol{\theta}}^b - \tilde{\boldsymbol{\theta}})\|$ is systematically larger than the corresponding quantile of $\|W(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|$, and thus, the bootstrap procedure becomes conservative in the “large bias” situation.

7.3.3 Bootstrap and the SmA procedure

Now we come back to the model selection setup. We again assume a linear Gaussian model with an inhomogeneous noise and $\mathbf{Y} = \mathbf{f}^* + \boldsymbol{\varepsilon}$, $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \Sigma)$.

Let $(\tilde{\boldsymbol{\theta}}_m, m \in \mathcal{M})$ be the considered ordered family of estimates $\tilde{\boldsymbol{\theta}}_m = \mathcal{S}_m \mathbf{Y}$ with $\mathcal{S}_m = (\Psi_m \Psi_m^\top)^{-1} \Psi_m$ for $\Psi_m = \Pi_m \Psi$, $m \in \mathcal{M}$. For notational simplicity $\mathcal{M} = \{1, 2, \dots, p\}$. In particular, p is the largest considered parameter dimension, and here it is supposed to be finite. Below we identify Ψ and Ψ_p .

Let also $m^\circ \in \mathcal{M}$ be fixed. We measure the corresponding modeling bias by the vector

$$\mathbf{f}^* - \Psi_{m^\circ}^\top \boldsymbol{\theta}_{m^\circ}^* = \mathbf{f}^* - \Pi_{m^\circ} \mathbf{f}^*.$$

The aim is to compare the joint distribution of scaled differences $W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ})$ for all $m > m^\circ$ with the analogous distribution of their bootstrap counterparts. We use the decomposition (4.14)

$$W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ}) = \mathbf{b}_{m,m^\circ} + \boldsymbol{\xi}_{m,m^\circ}$$

with

$$\begin{aligned} \mathbf{b}_{m,m^\circ} &= W(\Pi_m \mathbf{f}^* - \Pi_{m^\circ} \mathbf{f}^*), \\ \boldsymbol{\xi}_{m,m^\circ} &= W(D_m^{-2} \Pi_m - D_{m^\circ}^{-2} \Pi_{m^\circ}) \nabla, \end{aligned} \tag{7.45}$$

and $\nabla = \Psi \boldsymbol{\varepsilon}$. The bias term \mathbf{b}_{m,m° can be well controlled if m° is “good”.

Now we switch to the bootstrap counterpart. The main problem is a possible systematic component in \mathbf{Y} . We aim to design a sensitive procedure which detects the bias of the smoothing method \mathcal{S}_m as precise as possible. The use of \mathbf{Y} in the bootstrap world leads for $m^\circ < m^*$ to quantiles of $\boldsymbol{\xi}_{m,m^\circ}^b$ which are much larger than the corresponding quantiles of $\boldsymbol{\xi}_{m,m^\circ}$ due to the large bias component. This would result in a procedure which heavily oversmooths the data. A possible way out of this problem is based on using an under smoothing pilot $\tilde{\boldsymbol{\theta}}$ which removes the bias but preserves the variance. Typically one can use the largest considered model $m = p$ and the corresponding estimate $\tilde{\boldsymbol{\theta}} = \tilde{\boldsymbol{\theta}}_p = \mathcal{S}_p \mathbf{Y}$. Now we replace the data by the residuals for the pilot fit:

$$\check{\boldsymbol{\varepsilon}} \stackrel{\text{def}}{=} \mathbf{Y} - \Psi^\top \tilde{\boldsymbol{\theta}}_p = \mathbf{Y} - \Pi \mathbf{Y}$$

with $\Pi = \Pi_p$. Below we use that for any $m \leq p$, it holds

$$\Psi_m \check{\boldsymbol{\varepsilon}} = \Psi_m (\mathbf{I} - \Pi) \mathbf{Y} = 0. \quad (7.46)$$

Now we apply these residuals in place of the original data yielding the family of bootstrap residuals

$$\boldsymbol{\zeta}_m^b = \mathcal{S}_m \mathcal{W}^b \check{\boldsymbol{\varepsilon}}$$

for a diagonal weighting matrix $\mathcal{W}^b = \text{diag}(w_1^b, \dots, w_n^b)$ of bootstrap weights. If $\mathcal{E}^b = \mathcal{W}^b - \mathbb{E}^b \mathcal{W}^b$, then (7.46) and the definition imply

$$\boldsymbol{\zeta}_m^b = (\Psi_m \Psi_m^\top)^{-1} \Psi_m \mathcal{E}^b \check{\boldsymbol{\varepsilon}} = D_m^{-2} \Pi_m \check{\nabla}^b,$$

where $D_m^2 = \Psi_m \Psi_m^\top$ and

$$\check{\nabla}^b = \Psi \mathcal{E}^b \check{\boldsymbol{\varepsilon}} = \Psi \mathcal{E}^b (\mathbf{I} - \Pi) \mathbf{Y}.$$

For the differences $\boldsymbol{\zeta}_m^b - \boldsymbol{\zeta}_{m^\circ}^b$, $m > m^\circ$, this yields:

$$\boldsymbol{\zeta}_m^b - \boldsymbol{\zeta}_{m^\circ}^b = (D_m^{-2} \Pi_m - D_{m^\circ}^{-2} \Pi_{m^\circ}) \check{\nabla}^b.$$

The bootstrap procedure involves the stochastic component obtained by multiplying with the weighting matrix W :

$$\boldsymbol{\xi}_{m,m^\circ}^b = W (\boldsymbol{\zeta}_m^b - \boldsymbol{\zeta}_{m^\circ}^b) = W (D_m^{-2} \Pi_m - D_{m^\circ}^{-2} \Pi_{m^\circ}) \check{\nabla}^b \quad (7.47)$$

for all $m > m^\circ$. One can see that each $\boldsymbol{\xi}_{m,m^\circ}^b$ in (7.45) is a deterministic functions of the vector ∇ while $\boldsymbol{\xi}_{m,m^\circ}^b$ is the same function of $\check{\nabla}^b$. Therefore, it suffices to compare the distribution of ∇ with the conditional distribution of $\check{\nabla}^b$ given \mathbf{Y} and show that the latter mimics well the former. This question was studied in the previous section in terms of the quantity δ_Ψ defined by (7.33) for the model p , and of the bias $\mathbf{B} = \Sigma^{-1/2} (\mathbf{f}^* - \Pi \mathbf{f}^*)$. The result of Theorem 7.3.1 does not directly applies here, because we use $\check{\nabla}^b$ in place of ∇^b . However, the main arguments of the approach can still be used: given \mathbf{Y} , the vector $\check{\nabla}^b$ is Gaussian zero mean with the covariance matrix

$$\Psi \text{diag}\{\check{\boldsymbol{\varepsilon}} \cdot \check{\boldsymbol{\varepsilon}}\} \Psi^\top.$$

If this random matrix is close in probability to the covariance matrix $\Psi \Sigma \Psi^\top$ of the score ∇ , then the approach of Section 7.3.1 is still validated. A linear transformation

$\Psi \rightarrow \Psi \Sigma^{1/2}$ reduces the study to the case with $\Sigma = \text{Var}(\varepsilon) = I_n$. It obviously holds in this situation

$$\begin{aligned} \mathbb{E}(\check{\varepsilon}\check{\varepsilon}^\top) &= (I - \Pi)\mathbb{E}(\mathbf{Y}\mathbf{Y}^\top)(I - \Pi) \\ &= (\mathbf{f}^* - \Pi\mathbf{f}^*)(\mathbf{f}^* - \Pi\mathbf{f}^*)^\top + (I - \Pi)\Sigma(I - \Pi) = \mathbf{B}\mathbf{B}^\top + I - \Pi \end{aligned}$$

with $\Pi = \Pi_p$. Now one can check that each diagonal element Π_{ii} of Π satisfies

$$|\Pi_{ii}| \leq \delta_\Psi^2. \quad (7.48)$$

Exercise 7.3.9. Check (7.48) with δ_Ψ from (7.33) for the case $\sigma_i \equiv 1$.

This implies the operator-norm inequality

$$\|I_n - \text{diag}\{\check{\varepsilon} \cdot \check{\varepsilon}\}\| \leq \delta_\Psi^2.$$

Now we can summarize that under the small modeling bias condition $\mathbf{B} \approx 0$, the covariance of $\check{\nabla}^b$ is close in mean to the score covariance $\Psi\Psi^\top$ up to the error of order δ_Ψ . Moreover, the bias term $\mathbf{B}\mathbf{B}^\top$ can be treated exactly as in Section 7.3.2. These informal considerations motivate the following result.

Theorem 7.3.3. *Let $\mathbf{Y} = \mathbf{f}^* + \varepsilon$ be a Gaussian vector in \mathbb{R}^n with independent components, $\mathbf{Y} \sim \mathcal{N}(\mathbf{f}^*, \Sigma)$ for $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$. Let also Ψ be a $p \times n$ feature matrix such that the $p \times p$ -matrix $V^2 = \Psi\Sigma\Psi^\top$ is non-degenerated and the value*

$$\delta_\Psi = \max_{i=1, \dots, n} \|V^{-1}\Psi_i\|\sigma_i$$

is small. Define

$$\mathbf{B} = \Sigma^{-1/2}(\mathbf{f}^* - \Pi\mathbf{f}^*).$$

Given $m^\circ < p$, let the values $z_{m, m^\circ}^b(\mathbf{x})$ and q_{m° be fixed for the bootstrap distribution due to (7.11) and (7.12). Then, with Δ from (7.35), it holds

$$\mathbb{P}\left(\max_{m > m^\circ} \left\{ \|\xi_{m, m^\circ}\| - z_{m, m^\circ}^b(\mathbf{x} + q_{m^\circ}) \right\} \geq 0 \right) \leq 4e^{-\mathbf{x}} + \Delta.$$

*Proof. **To be done:***

The result of this theorem justifies the SmA procedure with the bootstrap-based critical values (7.13) in the asymptotic sense.

7.4 Linear non-Gaussian case and GAR

This section briefly comment why the bootstrap procedure can be validated even if the true error distribution is not Gaussian. This means that we again consider the linear Gaussian likelihood and the corresponding qMLE $\tilde{\boldsymbol{\theta}}$ is given by $\tilde{\boldsymbol{\theta}} = D^{-2}\Psi\mathbf{Y}$, the errors $\boldsymbol{\varepsilon}$ are independent but no more Gaussian. The discussion of the previous section shows that the most challenging step of analysis is to check that two vectors $\nabla = \Psi\boldsymbol{\varepsilon}$ and $\nabla^b = \Psi\mathcal{E}^b\mathbf{Y}$ have a similar distribution under the corresponding measures. In the Gaussian case, both vectors are normal zero mean and it suffices to compare their covariance matrices. In the non-Gaussian case the situation is more involved. A nice feature of Gaussian bootstrap multipliers is that the distribution of $\nabla^b = \Psi\mathcal{E}^b\mathbf{Y}$ given \mathbf{Y} is again Gaussian, and this fact does not rely on the true data distribution. It is entirely due to the construction of the bootstrap multipliers: ∇^b is normal because it is a linear combination of standard normal weights $e_i^b = w_i^b - 1$. The real score $\nabla = \Psi\boldsymbol{\varepsilon}$ is again a linear combination of errors ε_i , however these errors can be non-normal. In fact, in typical applications, there is no reason to assume that the errors are exactly normal. However, $\Psi\boldsymbol{\varepsilon}$ can be viewed as a linear combination of the errors ε_i . In combination with the condition that the value δ_Ψ from (7.33) is small, the central limit theorem applies and the zero mean standardized vector $V^{-1}\nabla$ is nearly standard normal under some further regularity and moment conditions. This allows to extend the result on bootstrap validity to the non-Gaussian case in some asymptotic sense. In the univariate case with $p = 1$ one can use the famous Berry-Esseen theorem, which can be also extended to the multivariate case in various special setups.

7.5 Sieve Generalized Linear regression

In a special case of a generalized liner model, \mathcal{P} is an exponential family with canonical parametrization. Then $\ell(y, \mathbf{v}) = y\mathbf{v} - d(\mathbf{v})$ and

$$L(\boldsymbol{\theta}) = \sum_{i=1}^n \ell(Y_i, f(X_i)) = \sum_{i=1}^n \{Y_i \Psi_i^\top \boldsymbol{\theta} - d(\Psi_i^\top \boldsymbol{\theta})\} = \mathbf{Y}^\top \Psi^\top \boldsymbol{\theta} - A(\boldsymbol{\theta})$$

with

$$A(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \sum_{i=1}^n d(\Psi_i^\top \boldsymbol{\theta}).$$

Also define

$$D^2(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \nabla^2 A(\boldsymbol{\theta}) = \sum_i \Psi_i \Psi_i^\top d''(\Psi_i^\top \boldsymbol{\theta}).$$

The estimate $\tilde{\boldsymbol{\theta}}$ can be computed by the iterative Newton procedure: start with any $\tilde{\boldsymbol{\theta}}_0$, e.g. LSE, and then compute

$$\tilde{\boldsymbol{\theta}}_{k+1} = \tilde{\boldsymbol{\theta}}_k + D^{-2}(\tilde{\boldsymbol{\theta}}_k)\Psi(\mathbf{Y} - \Psi^\top \tilde{\boldsymbol{\theta}}_k), \quad k = 1, 2, \dots$$

Under standard conditions this procedure converges very fast after just few iterations.

The target $\boldsymbol{\theta}^*$ is described by the optimizing the value $\mathbb{E}L(\boldsymbol{\theta})$. In view of $\mathbb{E}Y_i = f(\mathbf{X}_i)$

$$\boldsymbol{\theta}^* \stackrel{\text{def}}{=} \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \mathbb{E}L(\boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_i \{f(\mathbf{X}_i)\Psi_i^\top \boldsymbol{\theta} - d(\Psi_i^\top \boldsymbol{\theta})\}$$

which leads to the normal equation

$$\sum_{i=1}^n \{f(\mathbf{X}_i) - d'(\Psi_i^\top \boldsymbol{\theta})\}\Psi_i = 0.$$

One also has

$$\mathbb{F} = D^2 = -\nabla^2 \mathbb{E}L(\boldsymbol{\theta}^*) = \sum_i \Psi_i \Psi_i^\top d''(\Psi_i^\top \boldsymbol{\theta}^*) = \Psi \mathbf{d}''(\Psi^\top \boldsymbol{\theta}^*) \Psi^\top,$$

where $\mathbf{d}''(\mathbf{f})$ is a $n \times n$ diagonal matrix with the diagonal entries $d''(f_i)$.

Further, for the stochastic component $\zeta(\boldsymbol{\theta}) = L(\boldsymbol{\theta}) - \mathbb{E}L(\boldsymbol{\theta})$, it holds with $\varepsilon_i = Y_i - \mathbb{E}Y_i$

$$\nabla \zeta(\boldsymbol{\theta}) = \sum_i \varepsilon_i \Psi_i = \Psi \boldsymbol{\varepsilon} \tag{7.49}$$

and its covariance matrix fulfills

$$V^2 = \operatorname{Var}\{\nabla \zeta(\boldsymbol{\theta})\} = \sum_i \operatorname{Var}(Y_i)\Psi_i \Psi_i^\top = \Psi \operatorname{Var}(\boldsymbol{\varepsilon})\Psi^\top. \tag{7.50}$$

If the assumption on marginal distributions is correct, that is, if $Y_i \sim P_{f(\mathbf{X}_i)}$, then

$$\operatorname{Var}(Y_i) = d''(f(\mathbf{X}_i)). \tag{7.51}$$

By comparing the equations (7.50) and (7.51), one can conclude that D^2 and V^2 coincides if $d''(\Psi_i^\top \boldsymbol{\theta}^*) = d''(f(\mathbf{X}_i))$, that is, if PA is correct. Otherwise these two matrices may differ from each other. The Fisher expansion reads as

$$\|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \boldsymbol{\xi}\| \leq \diamond(\mathbf{x}) \tag{7.52}$$

on a dominating set of probability at least $1 - e^{-\mathbf{x}}$. Here in view of (7.49)

$$\boldsymbol{\xi} \stackrel{\text{def}}{=} D^{-1} \nabla \zeta(\boldsymbol{\theta}^*) = D^{-1} \Psi \boldsymbol{\varepsilon}.$$

7.5.1 Sieve MLE

Now we consider the sieve MLE setup when a sequence of subsets Θ_m of growing dimension is given and for each m the MLE $\tilde{\boldsymbol{\theta}}_m$ is computed. The corresponding quadratic risk \mathcal{R}_m is

$$\mathcal{R}_m = \mathbb{E} \|W(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}^*)\|^2.$$

The Fisher expansion (7.52) restricted to the sieve Θ_m yields

$$\|D_m(\tilde{\boldsymbol{\theta}}_m - \boldsymbol{\theta}_m^*) - \boldsymbol{\xi}_m\| \leq \diamond(\mathbf{x}), \quad (7.53)$$

where

$$\begin{aligned} D_m^2 &\stackrel{\text{def}}{=} -\nabla_m^2 \mathbb{E}L(\boldsymbol{\theta}^*) = \Psi_m \mathbf{d}''(\Psi_m^\top \boldsymbol{\theta}^*) \Psi_m, \\ \boldsymbol{\xi}_m &\stackrel{\text{def}}{=} D_m^{-1} \nabla_m \zeta(\boldsymbol{\theta}^*) = D_m^{-1} \Psi_m \boldsymbol{\varepsilon}, \\ \boldsymbol{\theta}_m^* &\stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta_m} \mathbb{E}L(\boldsymbol{\theta}). \end{aligned}$$

Alternatively one can define

$$\boldsymbol{\theta}_m^* = \Pi_m \boldsymbol{\theta}^*$$

the expansion (7.53) continues to apply.

The test statistic \mathbb{T}_{m,m° can be written as

$$\mathbb{T}_{m,m^\circ} = \|W(\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ})\|$$

and the decomposition (7.53) implies

$$\left| \mathbb{T}_{m,m^\circ} - \|\mathbf{b}_{m,m^\circ} + \boldsymbol{\xi}_{m,m^\circ}\| \right| \leq 2\diamond(\mathbf{x})$$

with

$$\begin{aligned} \mathbf{b}_{m,m^\circ} &= W(\boldsymbol{\theta}_m^* - \boldsymbol{\theta}_{m^\circ}^*), \\ \boldsymbol{\xi}_{m,m^\circ} &= W(D_m^{-2} \nabla_m - D_{m^\circ}^{-2} \nabla_{m^\circ}) = W(D_m^{-2} \Pi_m - D_{m^\circ}^{-2} \Pi_{m^\circ}) \nabla \end{aligned} \quad (7.54)$$

for $\nabla = \Psi \boldsymbol{\varepsilon}$.

7.5.2 Bootstrap counterpart

Now we check what happens in the bootstrap world. It holds

$$L^b(\boldsymbol{\theta}) = \sum_{i=1}^n \ell_i(\boldsymbol{\theta}) w_i^b = \sum_{i=1}^n \{Y_i \boldsymbol{\Psi}_i^\top \boldsymbol{\theta} - d(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta})\} w_i^b = \boldsymbol{\theta}^\top \boldsymbol{\Psi} \mathcal{W}^b \mathbf{Y} - A^b(\boldsymbol{\theta})$$

with

$$A^b(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \sum_{i=1}^n d(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta}) w_i^b.$$

One can use that $\mathbb{E}^b A^b(\boldsymbol{\theta}) = A$ and even more, the matrix $A^b(\boldsymbol{\theta})$ is close to its expectation $A(\boldsymbol{\theta})$ for n large. This suggests to replace $A^b(\boldsymbol{\theta})$ by $A(\boldsymbol{\theta})$ in our analysis. The same applies to its Hessian $D^2(\boldsymbol{\theta})$.

This implies similarly to the linear case that the stochastic part of the bootstrap-world estimate $\tilde{\boldsymbol{\theta}}^b$ is $D^{-2} \check{\nabla}^b$ for $\check{\nabla}^b = \boldsymbol{\Psi} \mathcal{E}^b \check{\boldsymbol{\varepsilon}}$ for $\check{\boldsymbol{\varepsilon}} = \mathbf{Y} - \boldsymbol{\Psi}^\top \tilde{\boldsymbol{\theta}}$. One can see that the question of bootstrap validity is again reduced to comparing of two distributions: of ∇ w.r.t. the original measure \mathbb{P} and of $\check{\nabla}^b$ w.r.t. to the bootstrap measure \mathbb{P}^b given \mathbf{Y} . If the bootstrap weights w_i^b are Gaussian then the bootstrap score $\check{\nabla}^b = \boldsymbol{\Psi} \mathcal{E}^b \check{\boldsymbol{\varepsilon}}$ is Gaussian as well, because it is a linear combination of the $e_i^b = w_i^b - 1$'s which are i.i.d. standard normal. Unfortunately, the real world score ∇ is a linear combination of the errors ε_i which in general are not Gaussian. Therefore, we cannot assume that the score ∇ is Gaussian. However, the desirable justification of the bootstrap procedure can be obtained by GAR arguments for the score vector ∇ .

7.5.3 Bootstrap for the SmA procedure

Now we return to the SmA procedure for the sieve GLM estimation scheme. For each m , we consider the sieve estimate $\tilde{\boldsymbol{\theta}}_m$ and its bootstrap counterpart $\tilde{\boldsymbol{\theta}}_m^b$ and try to design the procedure by its desired performance in the bootstrap world where we know the true value $\tilde{\boldsymbol{\theta}}_m$ in each sieve. We use the Fisher expansion as the main tool. Similarly to the real world case, it allows to decompose the difference $\tilde{\boldsymbol{\theta}}_m^b - \tilde{\boldsymbol{\theta}}_m$ into the deterministic and stochastic part. For the SmA procedure, we have to consider the pairs $\tilde{\boldsymbol{\theta}}_m - \tilde{\boldsymbol{\theta}}_{m^\circ}$ for all $m > m^\circ$ and their bootstrap versions. The counterpart of the real world expansion (7.54) looks as

$$\boldsymbol{\xi}_{m,m^\circ}^b = W(D_m^{-2} \Pi_m - D_{m^\circ}^{-2} \Pi_{m^\circ}) \check{\nabla}^b$$

with $\check{\nabla}^b = \boldsymbol{\Psi} \mathcal{E}^b \check{\boldsymbol{\varepsilon}}$. Again, as in the linear case, for m° fixed, all the $\boldsymbol{\xi}_{m,m^\circ}^b$ are the deterministic linear functions of the score $\check{\nabla}^b$, and this is exactly as for the real world stochastic vectors $\boldsymbol{\xi}_{m,m^\circ}$. Therefore, it suffices to compare the distribution of ∇ and of $\check{\nabla}^b$. These two distributions do not depend on m° and they are close to each other in probability if the small modeling bias condition is fulfilled for the largest considered

model p . Therefore, the bootstrap-adjusted critical values \mathbf{z}_{m,m° can be used in the real world, and the central propagation result continues to hold: any good model will be accepted with a high probability.

SmA and parameter tuning in high dimensional regression

Extending the methods and results on model selection to situation with a large or even huge parameter dimension is one of the main challenge of modern statistics. An important requirement to any such method is an automatic parameter tuning. If a proposed procedure involves some tuning parameters without explaining their automatic choice, then one problem is just replaced by another. This chapter focuses a special problem of subset selection for linear models with unknown noise structure. We aim to bring together the SmA procedure of Section 5.1 and the resampling idea of Chapter 7. In this chapter we restrict ourselves to a linear model: the observation vector $\mathbf{Y} \in \mathbb{R}^n$ is described by the equation

$$\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$$

for a given dictionary Ψ . We allow that the dictionary is overcomplete, that is, the dimension p of the vector $\boldsymbol{\theta}^*$ can be much larger than the number of observations n . Some structural assumptions are necessary to make the problem of recovering $\boldsymbol{\theta}^*$ meaningful. As in Chapter 7 we assume that $\boldsymbol{\theta}^*$ is sparse or can be approximated by a sparse vector. For a subset \varkappa , by $\tilde{\boldsymbol{\theta}}_\varkappa$ we denote the corresponding LSE $\tilde{\boldsymbol{\theta}}_\varkappa$:

$$\tilde{\boldsymbol{\theta}}_\varkappa = (\Psi_\varkappa \Psi_\varkappa^\top)^{-1} \Psi_\varkappa \mathbf{Y} = D_\varkappa^{-2} \Pi_\varkappa \Psi \mathbf{Y},$$

where $\Psi_\varkappa = \Pi_\varkappa \Psi$ with Π_\varkappa being a projector on the \varkappa -subspace of the $\boldsymbol{\theta}$ -space and $\nabla = \Psi \boldsymbol{\varepsilon}$ is the score vector in \mathbb{R}^p . By D_\varkappa^{-2} we denote the pseudo inverse of the matrix $D_\varkappa^2 = \Psi_\varkappa \Psi_\varkappa^\top$. Given a weighing loss matrix W , the SmA procedure can be applied as soon as the family of tail functions $z_{\varkappa, \varkappa^\circ}(\mathbf{x})$ for the norm of the stochastic component

$$\boldsymbol{\xi}_{\varkappa, \varkappa^\circ} = W(D_\varkappa^{-2} - D_{\varkappa^\circ}^{-2}) \nabla = \mathcal{S}_{\varkappa, \varkappa^\circ} \nabla$$

is fixed. Prior information about the noise $\boldsymbol{\varepsilon}$ makes this possible. Here we discuss how this information can be recovered from the data using a resampling method.

8.1 SmA subset selection in high dimensional regression

This section discusses the parameter choice in the problem of subset selection for a high dimensional linear regression model with unknown noise structure. The SmA procedure of Section 5.1 will be extended by the bootstrap step for choosing the critical values $\mathbf{z}_{\varkappa, \varkappa^\circ}$. We follow the approach of previous chapters and consider for a given sample \mathbf{Y} a family of Gaussian multipliers $\mathbf{w}^b = (w_i^b)$. We also need a pilot estimate $\tilde{\boldsymbol{\theta}}$ which removes most of systematic part from the data \mathbf{Y} . Further we proceed with the residuals $\check{\boldsymbol{\varepsilon}} = \mathbf{Y} - \Psi^\top \tilde{\boldsymbol{\theta}}$ for the bootstrap step. For each pair $\varkappa > \varkappa^\circ$, the bootstrap counterpart $\boldsymbol{\xi}_{\varkappa, \varkappa^\circ}^b$ of $\boldsymbol{\xi}_{\varkappa, \varkappa^\circ}$ looks as

$$\boldsymbol{\xi}_{\varkappa, \varkappa^\circ}^b = W(D_\varkappa^{-2} \Pi_\varkappa - D_{\varkappa^\circ}^{-2} \Pi_{\varkappa^\circ}) \Psi \mathcal{E}^b \check{\boldsymbol{\varepsilon}},$$

where \mathcal{E}^b is the diagonal matrix of the centered bootstrap multipliers; cf. (7.47). The proposed approach uses the bootstrap tail function $z_{\varkappa, \varkappa^\circ}(\mathbf{x})$ as a proxy for the true one:

$$\mathbb{P}^b(\|\boldsymbol{\xi}_{\varkappa, \varkappa^\circ}^b\| > z_{\varkappa, \varkappa^\circ}(\mathbf{x})) \leq e^{-x}.$$

After each tail function $z_{\varkappa, \varkappa^\circ}(\mathbf{x})$ is built, one can apply either uniform or multilevel synchronization of such tail functions for fixing the set of critical values $\mathbf{z}_{\varkappa, \varkappa^\circ}$. The main argument in validating this bootstrap procedure in the ordered case was that all real-world stochastic terms $\boldsymbol{\xi}_{\varkappa, \varkappa^\circ}$ and bootstrap-world stochastic components $\boldsymbol{\xi}_{\varkappa, \varkappa^\circ}^b$ are deterministic linear functions of the corresponding scores: $\check{\nabla}^b = \Psi \mathcal{E}^b \check{\boldsymbol{\varepsilon}}$ in the bootstrap world and $\nabla = \Psi \boldsymbol{\varepsilon}$ in the underlying model. This is still the case.

If the errors $\boldsymbol{\varepsilon}$ are Gaussian then the score vector $\nabla = \Psi \boldsymbol{\varepsilon}$ is Gaussian as well. The bootstrap score $\check{\nabla}^b$ is Gaussian under the bootstrap measure \mathbb{P}^b by construction. So, validation of the bootstrap procedure can now be restated as comparison of two Gaussian distributions in a rather special sense. Below we admit a possibly inhomogeneous noise. The only necessary assumptions are independence of the errors ε_i and a kind of the Lindeberg condition on the rows of the matrix Ψ . If $\psi_j = (\psi_{i,j})$ denotes the j th row of Ψ , then each component $\nabla_j = \psi_j \boldsymbol{\varepsilon}$ of the score $\nabla = \Psi \boldsymbol{\varepsilon}$ reads

$$\nabla_j = \sum_{i=1}^n \psi_{i,j} \varepsilon_i$$

Under the assumption that the errors ε_i are Gaussian, the same holds for ∇_j :

$$\nabla_j \sim \mathcal{N}(0, v_j^2), \quad v_j^2 = \sum_{i=1}^n \psi_{i,j}^2 \sigma_i^2.$$

The corresponding bootstrap-score component can be represented as

$$\nabla_j^b = \sum_{i=1}^n \psi_{i,j} \check{\varepsilon}_i w_i^b,$$

where $\check{\varepsilon}_i$ are the components of $\check{\varepsilon}$. With Gaussian multipliers, it is normal as well conditioned on the data \mathbf{Y} :

$$\nabla_j^b \sim \mathcal{N}(0, v_j^{b^2}), \quad v_j^{b^2} = \sum_{i=1}^n \psi_{i,j}^2 \check{\varepsilon}_i^2.$$

Even if we ignore the systematic component in the residuals $\check{\varepsilon}_i$ and use ε_i in place of $\check{\varepsilon}_i$, we can only hope that two covariances v_j^2 and $v_j^{b^2}$ are close to each other with high probability under Lindeberg type conditions on ε_i . The same applies to cross-covariance of the different component of ∇ and similar components of $\check{\nabla}^b$. Therefore, the problem can be reduced to comparing of two high dimensional Gaussian measures with similar covariance structure. Unfortunately the dimension p of the vectors ∇ and $\check{\nabla}^b$ can be very large. The tools of previous chapters based on the Pinsker inequality hardly apply here, because the dimension p enters in the error bound. In regular situation the error term is of order $\sqrt{p/n}$, and this value is not small if p is larger than n . One needs other technique which applies in a very high dimensional space.

8.1.1 Norm comparison for a family of Gaussian vectors

The SmA procedure involves a multiple comparison of many test statistics each of them is based on the norm of a Gaussian vector. Suppose we are given a family of vectors $\{\boldsymbol{\xi}_{\varkappa}, \varkappa \in \mathcal{M}_m\}$ each of dimension m . For each $\boldsymbol{\xi}_{\varkappa}$ suppose that the corresponding critical value z_{\varkappa} is fixed in a way that

$$\mathbb{P}\left(\bigcup_{\varkappa \in \mathcal{M}_m} \{\|\boldsymbol{\xi}_{\varkappa}\| > z_{\varkappa}\}\right) \leq e^{-x}.$$

We now want to control a similar deviation probability for another family of vectors $\{\boldsymbol{\xi}_{\varkappa}^b\}$ whose covariance structure is close to that of $\{\boldsymbol{\xi}_{\varkappa}\}$. To justify the bootstrap procedure in the case of high parameter dimension, we need to reduce this problem to the problem of comparison of maxima for two large Gaussian vectors.

Let \mathcal{S}_m denote a unit sphere in \mathbb{R}^m . We use that for any vector $\boldsymbol{\xi}$ in \mathbb{R}^m , it holds

$$\|\boldsymbol{\xi}\| = \sup_{\boldsymbol{\gamma} \in \mathcal{S}_m} \boldsymbol{\gamma}^\top \boldsymbol{\xi}.$$

Further we have to replace the maximum over the whole sphere by the maximum over a finite subset. Given $\delta > 0$, consider a finite δ -net $\mathcal{S}_m(\delta)$ in \mathcal{S}_m . It is obvious that

$$(1 - \delta)\|\boldsymbol{\xi}\| \leq \max_{\boldsymbol{\gamma} \in \mathcal{S}_m(\delta)} \boldsymbol{\gamma}^\top \boldsymbol{\xi} \leq \|\boldsymbol{\xi}\|. \quad (8.1)$$

Putting together for all \varkappa implies

$$\begin{aligned} \mathbb{P}\left(\bigcup_{\varkappa \in \mathcal{M}_m} \bigcup_{\gamma \in \mathcal{S}_m} \left\{ \frac{1}{\mathbf{z}_\varkappa} \gamma^\top \boldsymbol{\xi}_\varkappa > 1 \right\}\right) &\leq \mathbb{P}\left(\bigcup_{\varkappa \in \mathcal{M}_m} \{\|\boldsymbol{\xi}_\varkappa\| > \mathbf{z}_\varkappa\}\right) \\ &\leq \mathbb{P}\left(\bigcup_{\varkappa \in \mathcal{M}_m} \bigcup_{\gamma \in \mathcal{S}_m(\delta)} \left\{ \frac{1}{\mathbf{z}_\varkappa} \gamma^\top \boldsymbol{\xi}_\varkappa > 1 - \delta \right\}\right) \end{aligned}$$

Now introduce the vector \mathbf{X} of dimension $\mathbb{M}_m = \mathbb{M}_m(\delta) = |\mathcal{S}_m(\delta)| \times |\mathcal{M}_m|$ with the entries $\mathbf{z}_\varkappa^{-1} \gamma^\top \boldsymbol{\xi}_\varkappa$ for $\gamma \in \mathcal{S}_m(\delta)$ and $\varkappa \in \mathcal{M}_m$.

Below we consider a special case when $\boldsymbol{\xi}_\varkappa = \Pi_\varkappa \boldsymbol{\xi}$ for some linear mapping $\Pi_\varkappa: \mathbb{R}^p \rightarrow \mathbb{R}^m$ for $\varkappa \in \mathcal{M}_m$. Suppose also that another Gaussian zero mean vector $\boldsymbol{\xi}^b$ is given and $\boldsymbol{\xi}_\varkappa^b = \Pi_\varkappa \boldsymbol{\xi}^b$. Build a vector \mathbf{X}^b out of the $\boldsymbol{\xi}_\varkappa^b$'s in the same way as \mathbf{X} was constructed out of the $\boldsymbol{\xi}_\varkappa$'s. As a next step we evaluate the distance between two covariance operators for \mathbf{X} and \mathbf{X}^b . Let $\Sigma = \text{Var}(\boldsymbol{\xi})$. Obviously, for each two pairs (\varkappa, γ) and (\varkappa_1, γ_1) ,

$$\begin{aligned} \mathbb{E}[\gamma^\top \boldsymbol{\xi}_\varkappa \gamma_1^\top \boldsymbol{\xi}_{\varkappa_1}] &= \mathbb{E}[\gamma^\top \boldsymbol{\xi}_\varkappa \boldsymbol{\xi}_{\varkappa_1}^\top \gamma_1] = \mathbb{E}[\gamma^\top \Pi_\varkappa \boldsymbol{\xi} \boldsymbol{\xi}^\top \Pi_{\varkappa_1}^\top \gamma_1] \\ &= \gamma^\top \Pi_\varkappa \Sigma \Pi_{\varkappa_1}^\top \gamma_1 \leq \|\Pi_\varkappa \Sigma \Pi_{\varkappa_1}^\top\|. \end{aligned} \quad (8.2)$$

A similar formula holds for the covariance operator of $\boldsymbol{\xi}^b$. Below we denote

$$\square_m \stackrel{\text{def}}{=} \max_{\varkappa, \varkappa_1 \in \mathcal{M}_m} \|\Pi_\varkappa (\Sigma - \Sigma^b) \Pi_{\varkappa_1}^\top\|. \quad (8.3)$$

To be done: Exp-Bernstein inequality implies $\square_m \leq \mathbf{C}n^{-1/2} \log(p)$

Then (8.2) and (8.3) imply for $\mathbf{z}_\varkappa \geq 1$

$$\|\Sigma_{\mathbf{X}} - \Sigma_{\mathbf{X}^b}\|_\infty \leq \square_m.$$

This and the result (8.15) of Theorem 8.2.2 imply

$$\begin{aligned} &\mathbb{P}\left(\max_{\varkappa \in \mathcal{M}_m} \max_{\gamma \in \mathcal{S}_m(\delta)} \frac{1}{\mathbf{z}_\varkappa} \gamma^\top \boldsymbol{\xi}_\varkappa \geq 1\right) \\ &\leq \mathbb{P}\left(\max_{\varkappa \in \mathcal{M}_m} \max_{\gamma \in \mathcal{S}_m(\delta)} \frac{1}{\mathbf{z}_\varkappa} \gamma^\top \boldsymbol{\xi}_\varkappa^b \geq 1 - 2\Delta\right) + 2\Delta^{-2} \{\log(\mathbb{M}_m) + 1\} \square_m. \end{aligned}$$

Together with the norm approximation (8.1) this yields

$$\begin{aligned} &\mathbb{P}\left(\max_{\varkappa \in \mathcal{M}_m} \frac{1}{\mathbf{z}_\varkappa} \|\boldsymbol{\xi}_\varkappa\| \geq 1\right) \\ &\leq \mathbb{P}\left(\max_{\varkappa \in \mathcal{M}_m} \frac{1}{\mathbf{z}_\varkappa} \|\boldsymbol{\xi}_\varkappa^b\| \geq (1 - \delta)(1 - 2\Delta)\right) + 2\Delta^{-2} \{\log(\mathbb{M}_m) + 1\} \square_m. \end{aligned} \quad (8.4)$$

It remains to account for δ . The cardinality $|\mathcal{S}_m(\delta)|$ can be roughly upper bounded by $(1 + 2\delta^{-1})^m$, see Lemma 5.2 in www-personal.umich.edu/~romanv/papers/non-asymptotic-rmt-plain.pdf, yielding

$$\mathbb{M}_m < |\mathcal{M}_m| (1 + 2\delta^{-1})^m. \quad (8.5)$$

The above calculus with $\delta = \Delta$ imply the following result.

Theorem 8.1.1. *Let $\boldsymbol{\xi}$ and $\boldsymbol{\xi}^b$ be two zero mean Gaussian vectors in \mathbb{R}^p , and let $\{\Pi_{\mathcal{z}}, \mathcal{z} \in \mathcal{M}_m\}$ be a collection of linear mappings $\mathbb{R}^M \rightarrow \mathbb{R}^m$. For the quantity \square_m from (8.3) and for any set of critical values $\mathbf{z}_{\mathcal{z}} \geq 1$, it holds with any $\Delta < 1/3$*

$$\begin{aligned} & \mathbb{P} \left(\max_{\mathcal{z} \in \mathcal{M}_m} \frac{1}{\mathbf{z}_{\mathcal{z}}} \|\boldsymbol{\xi}_{\mathcal{z}}\| \geq 1 \right) \\ & \leq \mathbb{P} \left(\max_{\mathcal{z} \in \mathcal{M}_m} \frac{1}{\mathbf{z}_{\mathcal{z}}} \|\boldsymbol{\xi}_{\mathcal{z}}^b\| \geq 1 - 3\Delta \right) + 2\Delta^{-2} \left\{ \log(|\mathcal{M}_m|) + m \log(1 + 2/\Delta) \right\} \square_m. \end{aligned}$$

Proof. Apply (8.4) and (8.5) and use that $(1 - \Delta)(1 - 2\Delta) \geq 1 - 3\Delta$.

If \mathcal{M}_m is the set of all subsets of the full index set $\{1, \dots, p\}$, then by the Stirling formula for $m!$

$$\log |\mathcal{M}_m| \leq \log \binom{p}{m} \leq \log(p^m/m!) \leq m \log(ep/m).$$

In particular, with $\Delta < 1/3$ we obtain

$$\begin{aligned} & \mathbb{P} \left(\max_{\mathcal{z} \in \mathcal{M}_m} \frac{1}{\mathbf{z}_{\mathcal{z}}} \|\boldsymbol{\xi}_{\mathcal{z}}\| \geq 1 \right) \\ & \leq \mathbb{P}^b \left(\max_{\mathcal{z} \in \mathcal{M}_m} \frac{1}{\mathbf{z}_{\mathcal{z}}} \|\boldsymbol{\xi}_{\mathcal{z}}^b\| \geq 1 - 3\Delta \right) + 2\Delta^{-2} m \log \left(\frac{2ep}{m\Delta} \right) \square_m. \end{aligned}$$

To be done: put altogether and find the bound on n , m , and p

8.2 Gaussian comparison in high dimension

If the dimension p of the vector $\boldsymbol{\theta}$ is large, the approach of Section 7.3 based on the Pinsker inequality faces a crucial problem: the error term is proportional to $p^{1/2}$ and can be very large. One needs a different technique which allows comparing two Gaussian measures in a high dimensional space in terms of the corresponding covariance operators.

8.2.1 Stein identity, Slepian bridge, and Gaussian comparison

Below for a $\mathbb{M} \times \mathbb{M}$ matrix A , we denote

$$\begin{aligned} \|A\| &= \sup_{\|\mathbf{u}\|=1} \|A\mathbf{u}\|, & \|A\|_{\infty} &= \max_{i,j} |a_{i,j}|, \\ \|A\|_1 &= \sum_{i,j} |a_{i,j}|, & \|A\|_{Fr}^2 &= \sum_{i,j} a_{i,j}^2. \end{aligned}$$

Lemma 8.2.1. Let \mathbf{X} and \mathbf{Y} be two zero mean Gaussian vectors in $\mathbb{R}^{\mathbb{M}}$ with $\Sigma_{\mathbf{X}} = \text{Var}(\mathbf{X})$ and $\Sigma_{\mathbf{Y}} = \text{Var}(\mathbf{Y})$. Let also $f(\mathbf{x})$ be a smooth function on $\mathbb{R}^{\mathbb{M}}$. Then

$$\epsilon \stackrel{\text{def}}{=} |\mathbb{E}f(\mathbf{X}) - \mathbb{E}f(\mathbf{Y})| \leq \frac{1}{2} \|\Sigma_{\mathbf{X}} - \Sigma_{\mathbf{Y}}\|_{\infty} \|\nabla^2 f\|_{1,\infty}, \quad (8.6)$$

where $\|\nabla^2 f\|_{1,\infty} \stackrel{\text{def}}{=} \sup_{\mathbf{x}} \|\nabla^2 f(\mathbf{x})\|_1$.

Proof. Without loss of generality assume that \mathbf{X} and \mathbf{Y} are given on the same probability space and independent. For each $t \in [0, 1]$, define

$$\begin{aligned} \mathbf{Z}(t) &\stackrel{\text{def}}{=} \sqrt{t} \mathbf{X} + \sqrt{1-t} \mathbf{Y}, \\ \Psi(t) &\stackrel{\text{def}}{=} \mathbb{E}f(\mathbf{Z}(t)) = \mathbb{E}f(\sqrt{t} \mathbf{X} + \sqrt{1-t} \mathbf{Y}). \end{aligned}$$

Obviously

$$\epsilon = |\Psi(1) - \Psi(0)| = \left| \int_0^1 \Psi'(t) dt \right|. \quad (8.7)$$

Further,

$$\Psi'(t) = \mathbb{E}[\nabla f(\mathbf{Z}(t))^\top \mathbf{Z}'(t)] = \frac{1}{2} \mathbb{E}[\{t^{-1/2} \mathbf{X} - (1-t)^{-1/2} \mathbf{Y}\}^\top \nabla f(\mathbf{Z}(t))].$$

To compute this expectation, we apply the *Stein identity*. Let \mathbf{W} be a zero mean Gaussian vector in $\mathbb{R}^{\mathbb{M}}$. Then for any C^1 function $s: \mathbb{R}^{\mathbb{M}} \rightarrow \mathbb{R}^{\mathbb{M}}$, it holds

$$\mathbb{E}[\mathbf{W} s(\mathbf{W})] = \text{Var}(\mathbf{W}) \mathbb{E}[\nabla s(\mathbf{W})]. \quad (8.8)$$

Exercise 8.2.1. Prove (8.9) for standard normal \mathbf{W} using integration by part:

$$\int_{\mathbb{R}^{\mathbb{M}}} s(\mathbf{w}) \mathbf{w} e^{-\|\mathbf{w}\|^2/2} d\mathbf{w} = \int_{\mathbb{R}^{\mathbb{M}}} \nabla s(\mathbf{w}) e^{-\|\mathbf{w}\|^2/2} d\mathbf{w}.$$

Reduce the case of a Gaussian zero mean $\mathbf{w} \sim \mathcal{N}(0, \Sigma)$ with a positive symmetric matrix Σ to the case $\Sigma = I_{\mathbb{M}}$.

This results can be directly extended to any C^1 vector function $\mathbf{s}: \mathbb{R}^{\mathbb{M}} \rightarrow \mathbb{R}^q$: it holds

$$\mathbb{E}[\mathbf{W} \mathbf{s}(\mathbf{W})^\top] = \text{Var}(\mathbf{W}) \mathbb{E}[\nabla \mathbf{s}(\mathbf{W})^\top]. \quad (8.9)$$

Here $\nabla \mathbf{s}(\mathbf{w})^\top$ means the $p \times q$ matrix with the entries $\frac{d}{d\theta_j} s_m(\mathbf{w})$ for $j = 1, \dots, p$ and $m = 1, \dots, q$.

Exercise 8.2.2. Derive (8.9) by applying (8.8) columnwise.

The identity (8.9) is used with $\mathbf{W} = (\mathbf{X}^\top, \mathbf{Y}^\top)^\top$ and $\mathbf{s}(\mathbf{w}) = \nabla f(\mathbf{z}(t))$ for $\mathbf{z}(t) = \sqrt{t}\mathbf{x} + \sqrt{1-t}\mathbf{y}$. Independence of \mathbf{X} and \mathbf{Y} implies

$$\text{Var}(\mathbf{W}) = \begin{pmatrix} \Sigma_{\mathbf{X}} & 0 \\ 0 & \Sigma_{\mathbf{Y}} \end{pmatrix}.$$

Also $\nabla \mathbf{s}(\mathbf{w}) = (t^{1/2} \nabla^2 f(\mathbf{z}(t)), (1-t)^{1/2} \nabla^2 f(\mathbf{z}(t)))^\top$ and by (8.9)

$$\begin{aligned} \mathbb{E}[\nabla f(\mathbf{Z}(t))\mathbf{X}^\top] &= t^{1/2} \Sigma_{\mathbf{X}} \mathbb{E}[\nabla^2 f(\mathbf{Z}(t))] \\ \mathbb{E}[\nabla f(\mathbf{Z}(t))\mathbf{Y}^\top] &= (1-t)^{1/2} \Sigma_{\mathbf{Y}} \mathbb{E}[\nabla^2 f(\mathbf{Z}(t))], \end{aligned}$$

This and (8.12) imply

$$\begin{aligned} |\Psi'(t)| &\leq \frac{1}{2} \left| \text{tr}\{(\Sigma_{\mathbf{X}} - \Sigma_{\mathbf{Y}}) \mathbb{E}[\nabla^2 f(\mathbf{Z}(t))]\} \right| \\ &\leq \frac{1}{2} \|\Sigma_{\mathbf{X}} - \Sigma_{\mathbf{Y}}\|_\infty \|\mathbb{E}[\nabla^2 f(\mathbf{Z}(t))]\|_1 \leq \frac{1}{2} \|\Sigma_{\mathbf{X}} - \Sigma_{\mathbf{Y}}\|_\infty \|\nabla^2 f\|_{1,\infty}. \end{aligned}$$

Now the assertion follows from (8.7).

Now we apply the obtained bound to $f(\mathbf{x}) \stackrel{\text{def}}{=} g(\Delta^{-1}h_\beta(\mathbf{x}))$, where $g(z)$ is a smooth univariate function with bounded first and second derivatives, and the *smooth maximum* function: for some $\beta > 0$

$$h_\beta(\mathbf{x}) = \beta^{-1} \log\left(\sum_j e^{\beta x_j}\right). \quad (8.10)$$

Lemma 8.2.2. *Let \mathbf{X} and \mathbf{Y} be two zero mean Gaussian vectors in $\mathbb{R}^{\mathbb{M}}$ with $\Sigma_{\mathbf{X}} = \text{Var}(\mathbf{X})$ and $\Sigma_{\mathbf{Y}} = \text{Var}(\mathbf{Y})$. For a univariate function $g(z)$ with bounded first and second derivatives, and $h_\beta(\mathbf{x})$ from (8.10)*

$$|\mathbb{E}g(\Delta^{-1}h_\beta(\mathbf{X})) - \mathbb{E}g(\Delta^{-1}h_\beta(\mathbf{Y}))| \leq \left(\frac{\beta\|g'\|_\infty}{\Delta} + \frac{\|g''\|_\infty}{2\Delta^2}\right) \|\Sigma_{\mathbf{X}} - \Sigma_{\mathbf{Y}}\|_\infty. \quad (8.11)$$

Proof. It holds for $f(\mathbf{x}) \stackrel{\text{def}}{=} g(\Delta^{-1}h_\beta(\mathbf{x}))$

$$\begin{aligned} \nabla f(\mathbf{x}) &= \Delta^{-1}g'(\Delta^{-1}h_\beta(\mathbf{x}))\nabla h_\beta(\mathbf{x}), \\ \nabla^2 f(\mathbf{x}) &= \Delta^{-1}g'(\Delta^{-1}h_\beta(\mathbf{x}))\nabla^2 h_\beta(\mathbf{x}) + \Delta^{-2}g''(\Delta^{-1}h_\beta(\mathbf{x}))\nabla h_\beta(\mathbf{x})\nabla h_\beta(\mathbf{x})^\top. \end{aligned}$$

Also for any \mathbf{x} by direct calculus

$$\begin{aligned} \|\nabla h_\beta(\mathbf{x})\|_1 &= 1, \\ \|\nabla^2 h_\beta(\mathbf{x})\|_1 &\leq 2\beta. \end{aligned} \quad (8.12)$$

This implies

$$\begin{aligned}\|\nabla f(\mathbf{x})\|_1 &\leq \Delta^{-1}\|g'\|_\infty \times \|\nabla h_\beta(\mathbf{x})\|_1 \leq \Delta^{-1}\|g'\|_\infty, \\ \|\nabla^2 f(\mathbf{x})\|_1 &\leq \Delta^{-1}\|g'\|_\infty \times \|\nabla^2 h_\beta(\mathbf{x})\|_1 + \Delta^{-2}\|g''\|_\infty \times \|\nabla h_\beta(\mathbf{x})\|_1^2 \\ &\leq 2\Delta^{-1}\beta\|g'\|_\infty + \Delta^{-2}\|g''\|_\infty.\end{aligned}$$

Now (8.11) follows from (8.6).

A particular choice of the function g is given by

$$g(z) \stackrel{\text{def}}{=} \begin{cases} 2u^2, & u \in [0, 1/2], \\ 1 - 2(1 - u)^2, & u \in [1/2, 1]. \\ 0 & \text{otherwise.} \end{cases} \quad (8.13)$$

Obviously $|g'(u)| \leq 2$, $|g''(u)| \leq 4$ for all u . Then $\|g'_\Delta\|_\infty \leq 2\Delta^{-1}$, $\|g''_\Delta\|_\infty \leq 4\Delta^{-2}$. We conclude with the following bound.

Theorem 8.2.1. *Let \mathbf{X} and \mathbf{Y} be two zero mean Gaussian vectors in \mathbb{R}^M with $\Sigma_{\mathbf{X}} = \text{Var}(\mathbf{X})$ and $\Sigma_{\mathbf{Y}} = \text{Var}(\mathbf{Y})$. Then with $g(\cdot)$ given by (8.13), it holds for any $\Delta > 0$ and $\beta > 0$*

$$|\mathbb{E}g(\Delta^{-1}h_\beta(\mathbf{X})) - \mathbb{E}g(\Delta^{-1}h_\beta(\mathbf{Y}))| \leq 2(\beta\Delta^{-1} + \Delta^{-2})\|\Sigma_{\mathbf{X}} - \Sigma_{\mathbf{Y}}\|_\infty. \quad (8.14)$$

8.2.2 Comparing of the maximum of Gaussians

Let $\mathbf{X} = (X_j)$ and $\mathbf{Y} = (Y_j)$ be two zero mean Gaussian vectors in \mathbb{R}^M with $\Sigma_{\mathbf{X}} = \text{Var}(\mathbf{X})$ and $\Sigma_{\mathbf{Y}} = \text{Var}(\mathbf{Y})$, and let $\square = \|\Sigma_{\mathbf{X}} - \Sigma_{\mathbf{Y}}\|_\infty$. Now we aim at comparing the distributions of $\max_j X_j$ and $\max_j Y_j$. We use that the smooth maximum h_β fulfills

$$\max_j x_j \leq h_\beta(\mathbf{x}) \leq \max_j x_j + \beta^{-1} \log(M).$$

As the indicator function $\mathbb{I}(z \geq 0)$ is not differentiable, we approximate it by a smooth function g_Δ . Namely, select a two times differentiable function g with $g(u) = 0$ for $u \leq 0$, $g(u) = 1$ for $u \geq 1$, and $g(u)$ monotonously grows from zero to one when u grows from zero to one. Define also $g_\Delta(u) = g(\Delta^{-1}u)$ for $\Delta > 0$. With $\Delta = \beta^{-1} \log(M)$

$$g_\Delta \circ h_\beta(\mathbf{x} - \mathbf{\Delta}) \leq \mathbb{I}(\max_j x_j > 0) \leq g_\Delta \circ h_\beta(\mathbf{x} + \mathbf{\Delta}).$$

Here $\mathbf{\Delta}$ is the vector with all entries equal to Δ . Indeed, $g_\Delta(z) \in [0, 1]$ for any z . If $x_j \geq 0$ for some j , then $h_\beta(\mathbf{x} + \mathbf{\Delta}) \geq \Delta$ and hence,

$$g_\Delta \circ h_\beta(\mathbf{x} + \mathbf{\Delta}) \geq g(\Delta/\Delta) = g(1) = 1.$$

Similarly, if $\max_j x_j \leq 0$, then due to $\Delta = \beta^{-1} \log(\mathbb{M})$

$$h_\beta(\mathbf{x} - \mathbf{\Delta}) \leq \max_j (x_j - \Delta) + \beta^{-1} \log(\mathbb{M}) \leq 0$$

and $g_\Delta \circ h_\beta(\mathbf{x} - \mathbf{\Delta}) = 0$. This and (8.14) yield the bound

$$\begin{aligned} \mathbb{P}(\max_j X_j > 0) &\leq \mathbb{E}[g_\Delta \circ h_\beta(\mathbf{X} + \mathbf{\Delta})] \\ &\leq \mathbb{E}[g_\Delta \circ h_\beta(\mathbf{Y} + \mathbf{\Delta})] + 2(\beta\Delta^{-1} + \Delta^{-2}) \square \\ &\leq \mathbb{P}(\max_j Y_j > -2\Delta) + 2\Delta^{-2} \{\log(\mathbb{M}) + 1\} \square. \end{aligned}$$

Similarly one can approximate any indicator $\mathbb{I}(z \geq z_0)$ by shifting the function g .

Theorem 8.2.2. *Let \mathbf{X} and \mathbf{Y} be two zero mean Gaussian vectors in $\mathbb{R}^{\mathbb{M}}$ with $\Sigma_{\mathbf{X}} = \text{Var}(\mathbf{X})$ and $\Sigma_{\mathbf{Y}} = \text{Var}(\mathbf{Y})$.*

$$\square \stackrel{\text{def}}{=} \|\Sigma_{\mathbf{X}} - \Sigma_{\mathbf{Y}}\|_\infty,$$

it holds for any Δ and z

$$\mathbb{P}(\max_j X_j > z) \leq \mathbb{P}(\max_j Y_j > z - 2\Delta) + 2\Delta^{-2} \{\log(\mathbb{M}) + 1\} \square. \quad (8.15)$$

8.2.3 Anti-concentration for Gaussian maxima

This section explains how one can compare the distribution of two maxima using the anti-concentration bound. The obtained results allow to bound the probability $\mathbb{P}(\max_j X_j > 0)$ by a similar probability $\mathbb{P}(\max_j Y_j > -2\Delta)$ from above and $\mathbb{P}(\max_j Y_j > 2\Delta)$ from below up to the error term $2\Delta^{-2} \{\log(\mathbb{M}) + 1\} \square$. The next question is whether we can replace 2Δ or -2Δ by zero without an essential change of probability. In other words, we have to bound the difference

$$\mathbb{P}(\max_j Y_j > -2\Delta) - \mathbb{P}(\max_j Y_j > 0).$$

The following theorem provides bounds on the Lévy concentration function of the maximum of a Gaussian random vector in $\mathbb{R}^{\mathbb{M}}$, where the terminology is borrowed from [Chernozhukov et al. \(2013\)](#). The Lévy concentration function of a real valued random variable ξ is defined for $\varepsilon > 0$ as

$$\mathcal{L}(\xi, \varepsilon) = \sup_{x \in \mathbb{R}} \mathbb{P}(|\xi - x| \leq \varepsilon).$$

Theorem 8.2.3 (Anti-concentration). *Let $(X_1, \dots, X_{p_n})^\top$ be a centered Gaussian random vector in $\mathbb{R}^{\mathbb{M}}$ with $\sigma_j^2 = \mathbb{E}[X_j^2] > 0$ for all $1 \leq j \leq \mathbb{M}$. Moreover, let $\underline{\sigma} = \min_{1 \leq j \leq \mathbb{M}} \sigma_j$, $\bar{\sigma} = \max_{1 \leq j \leq \mathbb{M}} \sigma_j$, and $a_{\mathbb{M}} = \mathbb{E}[\max_{1 \leq j \leq \mathbb{M}} (X_j/\sigma_j)]$.*

1. *If the variances are all equal, namely $\underline{\sigma} = \bar{\sigma} = \sigma$, then for every $\epsilon > 0$,*

$$\mathcal{L}\left(\max_{1 \leq j \leq \mathbb{M}} X_j, \epsilon\right) \leq 4\epsilon(a_{\mathbb{M}} + 1)/\sigma;$$

2. *If the variances are not equal, namely $\underline{\sigma} < \bar{\sigma}$, then for every $\epsilon > 0$,*

$$\mathcal{L}\left(\max_{1 \leq j \leq \mathbb{M}} X_j, \epsilon\right) \leq \mathbf{C}\epsilon\{a_{\mathbb{M}} + 1 \vee \log(\underline{\sigma}/\epsilon)\}$$

where $\mathbf{C} > 0$ depends only on $\underline{\sigma}$ and $\bar{\sigma}$.

To compare the distribution of two maxima, we use the anti-concentration bound: if $\text{Var}(Y_j) \equiv \sigma^2$

$$\mathbb{P}(\max_j Y_j > 0) - \mathbb{P}(\max_j Y_j > -2\Delta) \leq 8\Delta(a_{\mathbb{M}} + 1)/\sigma,$$

where $a_{\mathbb{M}} \stackrel{\text{def}}{=} \mathbb{E} \max_j |Y_j/\sigma| \leq (2 \log \mathbb{M})^{1/2}$. If the variances $\sigma_j^2 \stackrel{\text{def}}{=} \text{Var}(Y_j)$ are unequal then

$$\mathbb{P}(\max_j Y_j > 0) - \mathbb{P}(\max_j Y_j > -2\Delta) \leq \mathbf{C}\Delta\sqrt{\log(\mathbb{M}/\Delta)}.$$

We now apply all the inequalities with the following choice: with $\Delta = b^{-1} = \beta^{-1} \log(\mathbb{M})$ and $\mathbb{Q} = \mathbb{M}/\Delta$

$$b = \square^{-1/3} \{\log(\mathbb{Q})\}^{-1/6}.$$

It follows by (8.14)

$$\begin{aligned} & \left| \mathbb{P}(\max_j X_j > 0) - \mathbb{P}(\max_j Y_j > 0) \right| \\ & \leq (2\beta\Delta^{-1} + 2\Delta^{-2})\square + \mathbf{C}\Delta\sqrt{\log(\mathbb{Q})} \\ & \leq \mathbf{C}\square^{1/3} \log^{2/3}(\mathbb{Q}) + \mathbf{C}\square^{1/3} \log^{2/3}(\mathbb{Q}) \\ & \leq \mathbf{C}\square^{1/3} \log^{2/3}(\mathbb{Q}). \end{aligned}$$

The definition yields $\mathbb{Q} = \mathbb{M}/\Delta \leq \mathbb{M}/\square^{1/3}$. We conclude with the following result.

Theorem 8.2.4. *Let \mathbf{X} and \mathbf{Y} be two zero mean Gaussian vectors in $\mathbb{R}^{\mathbb{M}}$ with $\Sigma_{\mathbf{X}} = \text{Var}(\mathbf{X})$ and $\Sigma_{\mathbf{Y}} = \text{Var}(\mathbf{Y})$. With*

$$\square \stackrel{\text{def}}{=} \|\Sigma_{\mathbf{X}} - \Sigma_{\mathbf{Y}}\|_{\infty},$$

it holds for any Δ

$$\left| \mathbb{P}(\max_j X_j > 0) - \mathbb{P}(\max_j Y_j > 0) \right| \leq \mathbf{C} \Delta^{1/3} \log^{2/3}(\mathbb{M}/\Delta^{1/3}).$$

Penalized model selection

This chapter discusses a class of procedures which can be represented as penalized minimization of the empirical risk. We consider the linear model $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$ in which the parameter dimension p can be very large. The empirical risk is just the squared norm of the difference $\mathbf{Y} - \Psi^\top \boldsymbol{\theta}$. The penalized procedure tries to minimize this empirical risk penalized by the complexity of the vector $\boldsymbol{\theta}$ used for prediction:

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \{ \|\mathbf{Y} - \Psi^\top \boldsymbol{\theta}\|^2 + \operatorname{pen}(\boldsymbol{\theta}) \} \quad (9.1)$$

for a penalty function $\operatorname{pen}(\boldsymbol{\theta})$.

The roughness penalty has been already discussed in Chapter 3. The resulting estimate is again linear and the general approach of linear model selection continues to apply. Here we consider two special choices of $\operatorname{pen}(\boldsymbol{\theta})$ which are essentially non-linear. The complexity penalty $\operatorname{pen}(\boldsymbol{\theta}) = C\|\boldsymbol{\theta}\|_0$ just counts the number of non-zero $\boldsymbol{\theta}$ -coefficients. The sparse penalty $\operatorname{pen}(\boldsymbol{\theta}) = C\|\boldsymbol{\theta}\|_q^q$ use the q -norm of $\boldsymbol{\theta}$ for some $q < 2$. The most popular sparse penalty corresponds to $q = 1$.

9.1 Complexity penalization

This section considers the important special case of penalization by complexity. The famous Akaike criteria is a special case of such penalization. We offer another viewpoint based on the SmA idea. As previously in Section 2.1, for a subset \varkappa of the index set $\{1, \dots, p\}$, the estimate $\tilde{\boldsymbol{\theta}}_\varkappa$ is the corresponding projection MLE:

$$\tilde{\boldsymbol{\theta}}_\varkappa = (\Psi_\varkappa \Psi_\varkappa^\top)^{-1} \Psi_\varkappa \mathbf{Y}.$$

Theorem 9.1.1. *Let $\operatorname{pen}(\boldsymbol{\theta}) = C\|\boldsymbol{\theta}\|_0$. Then the solution $\hat{\boldsymbol{\theta}}$ of the problem (9.1) satisfies*

$$\hat{\boldsymbol{\theta}} = \tilde{\boldsymbol{\theta}}_{\hat{\varkappa}},$$

where

$$\hat{\varkappa} = \underset{\varkappa}{\operatorname{argmin}} \{ \|\tilde{\boldsymbol{\varepsilon}}_{\varkappa}\|^2 + \mathbf{C}|\varkappa| \}, \quad (9.2)$$

with

$$\tilde{\boldsymbol{\varepsilon}}_{\varkappa} = \mathbf{Y} - \boldsymbol{\Psi}_{\varkappa}^{\top} \tilde{\boldsymbol{\theta}}_{\varkappa} = (\mathbf{I}_n - \boldsymbol{\Pi}_{\varkappa}) \mathbf{Y}$$

and $|\varkappa|$ means the cardinality of \varkappa or, equivalently the number of coefficients in the support set \varkappa .

The unbiased risk estimation procedure corresponds to the special choice of constant $\mathbf{C} = 2\sigma^2$. This results in selecting a proper subset \varkappa which provides a reasonable fit under complexity constraint:

$$\hat{\varkappa} = \underset{\varkappa}{\operatorname{argmin}} \{ \|\tilde{\boldsymbol{\varepsilon}}_{\varkappa}\|^2 + 2\sigma^2|\varkappa| \}.$$

This procedure faces two essential problems when the dimension p becomes large. One of them is algorithmic: the procedure requires to compute the MLE $\tilde{\boldsymbol{\theta}}_{\varkappa}$ and the related empirical risk for any subset \varkappa ; such a problem is in general NP-hard and can be solved in very special cases, e.g. if the design matrix $\boldsymbol{\Psi}$ is orthogonal. Then the procedure can be reduced to thresholding of individual Fourier coefficients $\tilde{\theta}_j = \psi_j \mathbf{Y}$, where ψ_j denotes the j th row of $\boldsymbol{\Psi}$.

Theorem 9.1.2. *Let $n \geq p$ and the matrix $\boldsymbol{\Psi}$ be orthonormal, that is, $\boldsymbol{\Psi}\boldsymbol{\Psi}^T = \mathbf{I}_p$. Then the active set $\hat{\varkappa}$ from (9.2) is given by hard thresholding*

$$\hat{\varkappa} = \{j: |\tilde{\theta}_j| \geq \lambda\}$$

for a proper $\lambda > 0$. The corresponding hard thresholding estimate $\hat{\boldsymbol{\theta}} = (\hat{\theta}_j)$ reads as

$$\hat{\theta}_j = \begin{cases} \psi_j \mathbf{Y}, & |\psi_j \mathbf{Y}| > \lambda, \\ 0, & \text{otherwise.} \end{cases}$$

The other problem is statistical. First of all, the procedure assumes a homogeneous noise and requires that the noise variance σ^2 is given. Second, the penalization in the form $2\sigma^2|\varkappa|$ is too mild and does not ensure a proper model selection for p large. The reason is that there are very many models to select between, this number is exponential in p , and therefore, the price for this choice has to be much larger than in the ordered case.

Below we discuss several ways of solving the statistical problem. First we consider the penalized model selection procedure (9.1) or equivalently (9.2) with a data-driven constant \mathbf{C} . Then we extend it to a more sophisticated non-linear choice of the penalty function $\text{pen}(\boldsymbol{\theta})$. Finally we discuss the saddle-point bivariate model selection.

The approach is based on the propagation idea: if the model is “good” in the sense that it provides a reasonable data fit, it should be competitive against larger models: no reason to increase the complexity if you already have a proper prediction ability. A typical situation is as follows: we have a model-candidate \mathcal{K}° which is not too complex, that is, $|\mathcal{K}^\circ|$ is small relative to the sample size n and the total dimension p . One says in such cases that \mathcal{K}° is “sparse”. Further, this choice \mathcal{K}° is good if the coefficients θ_j^* for $j \notin \mathcal{K}^\circ$ are nearly zero. We aim at designing a data-driven procedure such that the criterium (9.2) keeps \mathcal{K}° alive in competition with all larger models. This precisely means that

$$\|\tilde{\varepsilon}_{\mathcal{K}}\|^2 + \mathbf{C}|\mathcal{K}| \geq \|\tilde{\varepsilon}_{\mathcal{K}^\circ}\|^2 + \mathbf{C}|\mathcal{K}^\circ|$$

This inequality has to be verified for all $\mathcal{K} > \mathcal{K}^\circ$ with a high probability. Rearranging yields in view of $\tilde{\varepsilon}_{\mathcal{K}} = (\mathbf{I}_n - \Pi_{\mathcal{K}})\mathbf{Y}$

$$\|\tilde{\varepsilon}_{\mathcal{K}^\circ}\|^2 - \|\tilde{\varepsilon}_{\mathcal{K}}\|^2 = \|\Pi_{\mathcal{K}, \mathcal{K}^\circ}\mathbf{Y}\|^2 \leq \mathbf{C}(|\mathcal{K}| - |\mathcal{K}^\circ|).$$

If we knew the noise distribution then we can fix the constant \mathbf{C} in the “pure noise” situation. Indeed, the underlying structural assumption means that there is no significant signal θ_j^* for $j \notin \mathcal{K}^\circ$, and hence, one can simply ignore such signal and consider $\theta_j^* \equiv 0$ in the complement of \mathcal{K}° :

$$\mathbb{P}\left(\bigcup_{\mathcal{K} > \mathcal{K}^\circ} \{\|\Pi_{\mathcal{K}, \mathcal{K}^\circ}\boldsymbol{\varepsilon}\|^2 > \mathbf{C}_0(|\mathcal{K}| - |\mathcal{K}^\circ|)\}\right) \leq e^{-x} \quad (9.3)$$

for a constant \mathbf{C}_0 . It is obvious that this condition becomes stronger if the set \mathcal{K}° is taken smaller. The hardest case corresponds to the empty set \mathcal{K}° , yielding the constraint

$$\mathbb{P}\left(\bigcup_{\mathcal{K}} \{\|\Pi_{\mathcal{K}}\boldsymbol{\varepsilon}\|^2 > \mathbf{C}_0|\mathcal{K}|\}\right) \leq e^{-x}. \quad (9.4)$$

If the dimension of \mathcal{K} only slightly higher than the dimension of \mathcal{K}° , the inequality $\|\Pi_{\mathcal{K}, \mathcal{K}^\circ}\boldsymbol{\varepsilon}\|^2 > \mathbf{C}_0(|\mathcal{K}| - |\mathcal{K}^\circ|)$ would require a very large constant \mathbf{C}_0 . At the same time, this constant rapidly stabilizes if $|\mathcal{K}| - |\mathcal{K}^\circ|$ exceeds some prescribed value. This suggests to extend the condition (9.6) (resp. (9.4)): given $\tau \geq 1$

$$\mathbb{P}\left(\bigcup_{\mathcal{K} \in \mathcal{M}_\tau(\mathcal{K}^\circ)} \{\|\Pi_{\mathcal{K}, \mathcal{K}^\circ}\boldsymbol{\varepsilon}\|^2 > \mathbf{C}_0(|\mathcal{K}| - |\mathcal{K}^\circ|)\}\right) \leq e^{-x}. \quad (9.5)$$

Here $\mathcal{M}_\tau(\mathcal{K}^\circ) = \{\mathcal{K} : \mathcal{K} > \mathcal{K}^\circ, |\mathcal{K}| \geq |\mathcal{K}^\circ| + \tau\}$.

Now suppose that such a constant $\mathbf{C}_0 = \mathbf{C}_0(\tau)$ is fixed. Then the procedure can be applied with

$$\mathbf{C} \stackrel{\text{def}}{=} \mathbf{C}_0 + \beta,$$

where β appears in the definition of a “good” choice: \varkappa° is good if

$$\|\Pi_{\varkappa, \varkappa^\circ} \mathbf{f}^*\|^2 \leq \beta(|\varkappa| - |\varkappa^\circ|).$$

To be done: An upper bound on \mathbf{C}_0

Apply the arguments from Section 8.1.1 and matrix Bernstein inequality from Section 1.6;

To be done: Choice of τ by concentration of $\|\Pi_{\varkappa, \varkappa^\circ} \boldsymbol{\varepsilon}\|^2$

Bound the difference between $\|\Pi_{\varkappa, \varkappa^\circ} \boldsymbol{\varepsilon}\|^2$ and its expectation; Use matrix Bernstein from Section 1.6.

To be done: Oracle inequality for $\hat{\varkappa}$ -choice

The construction ensures with a dominating probability that $\hat{\varkappa} \leq \varkappa^* + \tau$.

To be done: the quantity $\|\Pi_{\varkappa, \varkappa^\circ} \mathbf{Y}\|^2$ is in the Gaussian case non-central χ^2 . make it accurate

Now we discuss the situation when the noise distribution is unknown. Then the relation (9.4) cannot be used for fixing the constant \mathbf{C}_0 . Below we discuss the bootstrap based choice of this constant.

9.1.1 Bootstrap based tuning

As in Section 7.3.3 we fix some model-candidate \varkappa° and consider the family of estimates

$$\tilde{\boldsymbol{\theta}}_\varkappa = (\Psi_\varkappa \Psi_\varkappa^\top)^{-1} \Psi_\varkappa \mathbf{Y}$$

for $\varkappa > \varkappa^\circ$. We also suppose a pilot estimate $\tilde{\boldsymbol{\theta}}$ to be given which provides a reasonable data fit but is probably too volatile. Define the residuals $\check{\boldsymbol{\varepsilon}} = \mathbf{Y} - \Psi^\top \tilde{\boldsymbol{\theta}}$. The procedure follows the same path as in the ordered case. For each \varkappa we compute and store the corresponding MLE $\tilde{\boldsymbol{\theta}}_\varkappa$ and a collection of the bootstrap-based stockstic vectors $\boldsymbol{\zeta}_\varkappa^b$:

$$\boldsymbol{\zeta}_\varkappa^b = (\Psi_\varkappa \Psi_\varkappa^\top)^{-1} \Psi_\varkappa \mathcal{E}^b \check{\boldsymbol{\varepsilon}}.$$

Further we can fix the value \mathbf{C}_0 using the bootstrap differences

$$\boldsymbol{\xi}_{\varkappa, \varkappa^\circ}^b = W(\boldsymbol{\zeta}_\varkappa^b - \boldsymbol{\zeta}_{\varkappa^\circ}^b)$$

for the weighting loss matrix W . The bootstrap critical values can be computed from these differences by the bootstrap analog of the propagation condition (9.6)

$$\mathbb{P}^b \left(\bigcup_{\varkappa \in \mathcal{M}_\tau(\varkappa^\circ)} \{\|\boldsymbol{\xi}_{\varkappa, \varkappa^\circ}^b\|^2 > \mathbf{C}_0(|\varkappa| - |\varkappa^\circ|)\} \right) \leq e^{-x}. \quad (9.6)$$

Here \mathbb{P}^b is to be understood as the empirical bootstrap measure.

To be done: An analytic bound on C_0

9.2 Sparse penalty

This section briefly discusses the use of a sparse penalty based on the ℓ_1 -norm of the vector $\boldsymbol{\theta}$. Below for a vector $\boldsymbol{\theta} \in \mathbb{R}^p$ we denote

$$\|\boldsymbol{\theta}\|_1 = \sum_{j=1}^p |\theta_j|, \quad \|\boldsymbol{\theta}\|^2 = \sum_{j=1}^p \theta_j^2, \quad \|\boldsymbol{\theta}\|_\infty = \max_{j \leq p} |\theta_j|.$$

We consider the LASSO type procedure which is based on minimization of the empirical risk $\|\mathbf{Y} - \Psi^\top \boldsymbol{\theta}\|^2$ penalized by $\lambda \|\boldsymbol{\theta}\|_1$:

$$\hat{\boldsymbol{\theta}} \stackrel{\text{def}}{=} \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \mathcal{J}_\lambda(\boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \{ \|\mathbf{Y} - \Psi^\top \boldsymbol{\theta}\|^2 + \lambda \|\boldsymbol{\theta}\|_1 \}. \quad (9.7)$$

Note that the formulation of the problem implicitly assumes that all components θ_j of the vector $\boldsymbol{\theta}$ have the same impact. This can be translated into the scaling condition on the matrix Ψ . Usually this matrix and thus, the coefficients θ_j are rescaled in a way that the diagonal elements of the matrix $\Psi\Psi^\top$ are equal to one:

$$\sum_{i=1}^n \psi_{i,j}^2 = 1.$$

The problem (9.7) has a closed form solution only in very special situation. One of them is when the matrix Ψ is orthogonal.

Theorem 9.2.1. *Let $n \geq p$ and $\Psi\Psi^\top = I_p$. Then $\hat{\boldsymbol{\theta}}$ is obtained by soft-thresholding of $\tilde{\boldsymbol{\theta}} = \Psi\mathbf{Y}$:*

$$\hat{\theta}_j = \begin{cases} (\tilde{\theta}_j - \lambda)_+ & \tilde{\theta}_j \geq 0, \\ -(|\tilde{\theta}_j| - \lambda)_+ & \tilde{\theta}_j < 0 \end{cases}$$

However, compared to the complexity penalization, the problem (9.7) can be solved numerically because the objective function is convex.

Now we briefly discuss some properties of the solution. The underlying structural assumption is that the true vector $\boldsymbol{\theta}^*$ is sparse, that is, most of its entries vanish. By \varkappa^* we denote the corresponding oracle support. We first consider the case when Ψ_{\varkappa^*} is orthogonal to the rest of Ψ . Our first result shows that a proper choice of the parameter λ ensures a sparse solution: the non-zero coefficients of $\hat{\boldsymbol{\theta}}$ are all located within \varkappa^* .

Theorem 9.2.2. Let $\boldsymbol{\theta}^*$ be supported on \varkappa_0 , \varkappa_0^c be the complement of \varkappa_0 , and

$$\Psi_{\varkappa_0} \Psi_{\varkappa_0^c}^\top = 0. \quad (9.8)$$

If the coefficient λ fulfills

$$\|\Psi \boldsymbol{\varepsilon}\|_\infty \leq \lambda,$$

then

$$\widehat{\varkappa} \leq \varkappa_0 \quad (9.9)$$

Proof. It suffices to check for each candidate $\boldsymbol{\theta}$ that the criteria in the optimization problem (9.7) only improves if we kill all its entries which do not enter in the oracle set \varkappa^* . Let \varkappa be the support of $\boldsymbol{\theta}$. Define $\boldsymbol{\theta}_{\varkappa_0} = \Pi_{\varkappa_0} \boldsymbol{\theta}$ as the restriction of the parameter vector $\boldsymbol{\theta}$ to \varkappa_0 , and similarly $\boldsymbol{\theta}_{\varkappa_0^c} = \Pi_{\varkappa_0^c} \boldsymbol{\theta}$ is the projection on the complement set \varkappa_0^c . Obviously $\boldsymbol{\theta}_{\varkappa_0^c} = \boldsymbol{\theta} - \boldsymbol{\theta}_{\varkappa_0}$. Then the model equation $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$ implies

$$\begin{aligned} \|\mathbf{Y} - \Psi^\top \boldsymbol{\theta}\|^2 - \|\mathbf{Y} - \Psi^\top \boldsymbol{\theta}_{\varkappa_0}\|^2 &= \|\boldsymbol{\varepsilon} - \Psi^\top (\boldsymbol{\theta}_{\varkappa_0} + \boldsymbol{\theta}_{\varkappa_0^c} - \boldsymbol{\theta}^*)\|^2 - \|\boldsymbol{\varepsilon} - \Psi^\top (\boldsymbol{\theta}_{\varkappa_0} - \boldsymbol{\theta}^*)\|^2 \\ &= \|\Psi^\top \boldsymbol{\theta}_{\varkappa_0^c}\|^2 - 2\{\boldsymbol{\varepsilon} - \Psi^\top (\boldsymbol{\theta}_{\varkappa_0} - \boldsymbol{\theta}^*)\}^\top \Psi^\top \boldsymbol{\theta}_{\varkappa_0^c}. \end{aligned}$$

Exercise 9.2.1. Check that

$$\{\Psi^\top (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)\}^\top \Psi^\top \boldsymbol{\theta}_{\varkappa_0^c} = 0.$$

Hint: use that $\boldsymbol{\theta}_{\varkappa_0} - \boldsymbol{\theta}^*$ is supported on \varkappa^* , and $\boldsymbol{\theta}_{\varkappa_0^c} = \boldsymbol{\theta} - \boldsymbol{\theta}_{\varkappa_0}$ on \varkappa_0^c . Then the result follows from (9.8).

Now $\|\boldsymbol{\theta}\|_1 - \|\boldsymbol{\theta}_{\varkappa_0}\|_1 = \|\boldsymbol{\theta}_{\varkappa_0^c}\|_1$ and

$$\begin{aligned} \|\mathbf{Y} - \Psi^\top \boldsymbol{\theta}\|^2 - \|\mathbf{Y} - \Psi^\top \boldsymbol{\theta}_{\varkappa_0}\|^2 + \lambda(\|\boldsymbol{\theta}\|_1 - \|\boldsymbol{\theta}_{\varkappa_0}\|_1) \\ = \|\Psi^\top \boldsymbol{\theta}_{\varkappa_0^c}\|^2 + \lambda\|\boldsymbol{\theta}_{\varkappa_0^c}\|_1 - 2(\Psi \boldsymbol{\varepsilon})^\top \boldsymbol{\theta}_{\varkappa_0^c}. \end{aligned}$$

It remains to check that the condition $2\|\Psi \boldsymbol{\varepsilon}\|_\infty \leq \lambda$ ensures that

$$\lambda\|\boldsymbol{\theta}_{\varkappa_0^c}\|_1 - 2(\Psi \boldsymbol{\varepsilon})^\top \boldsymbol{\theta}_{\varkappa_0^c} \geq 0.$$

Therefore, reduction of $\boldsymbol{\theta}$ to $\boldsymbol{\theta}_{\varkappa_0}$ only improves the objective function, and the result follows.

9.2.1 Basic inequality

If we drop the orthogonality condition (9.8), the wonderful oracle result (9.9) does not hold any more. However, one can establish an oracle bound on the quadratic risk in terms of the sparsity value $\|\boldsymbol{\theta}^*\|_1$. We again check the condition that $\boldsymbol{\theta}$ is better than $\boldsymbol{\theta}^*$ w.r.t. the criterion $\mathcal{J}_\lambda(\boldsymbol{\theta})$ from (9.7). It holds due to the model equation $\mathbf{Y} = \boldsymbol{\Psi}^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$

$$\begin{aligned} \|\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta}\|^2 - \|\mathbf{Y} - \boldsymbol{\Psi}^\top \boldsymbol{\theta}^*\|^2 &= \|\boldsymbol{\varepsilon} - \boldsymbol{\Psi}^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 - \|\boldsymbol{\varepsilon}\|^2 \\ &= -2(\boldsymbol{\Psi}\boldsymbol{\varepsilon})^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \|\boldsymbol{\Psi}^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2. \end{aligned}$$

The event $\widehat{\boldsymbol{\theta}} = \boldsymbol{\theta}$ is only possible if $\mathcal{J}_\lambda(\boldsymbol{\theta}) \leq \mathcal{J}_\lambda(\boldsymbol{\theta}^*)$. This yields

$$\|\boldsymbol{\Psi}^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 + \lambda \|\boldsymbol{\theta}\|_1 \leq 2(\boldsymbol{\Psi}\boldsymbol{\varepsilon})^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \lambda \|\boldsymbol{\theta}^*\|_1 \quad (9.10)$$

Moreover, if the score $\nabla = \boldsymbol{\Psi}\boldsymbol{\varepsilon}$ fulfills $\|\nabla\|_\infty \leq \mathbf{C}_0$, then

$$2(\boldsymbol{\Psi}\boldsymbol{\varepsilon})^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \leq 2\mathbf{C}_0 \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_1 \quad (9.11)$$

By the triangle inequality $\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_1 \geq \|\boldsymbol{\theta}\|_1 - \|\boldsymbol{\theta}^*\|_1$, and (9.10) and (9.11) imply

$$\|\boldsymbol{\Psi}^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 + (\lambda - 2\mathbf{C}_0) \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_1 \leq 2\lambda \|\boldsymbol{\theta}^*\|_1$$

If $\lambda > 2\mathbf{C}_0$, this inequality provides a number of informative messages. The prediction loss $\|\boldsymbol{\Psi}^\top (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2$ and the estimation loss $\|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_1$ are can be bounded as

$$\begin{aligned} \|\boldsymbol{\Psi}^\top (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2 &\leq 2\lambda \|\boldsymbol{\theta}^*\|_1, \\ \|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_1 &\leq \frac{2\lambda}{\lambda - 2\mathbf{C}_0} \|\boldsymbol{\theta}^*\|_1. \end{aligned}$$

The last inequality can be made more exact if we consider for any $\boldsymbol{\theta}$ the decomposition $\boldsymbol{\theta} = \boldsymbol{\theta}_{\mathcal{X}_0} + \boldsymbol{\theta}_{\mathcal{X}_0^c}$, where $\boldsymbol{\theta}_{\mathcal{X}_0}$ is supported on \mathcal{X}^* and $\boldsymbol{\theta}_{\mathcal{X}_0^c}$ on its complement \mathcal{X}_0^c . Obviously $\|\boldsymbol{\theta}\|_1 = \|\boldsymbol{\theta}_{\mathcal{X}_0}\|_1 + \|\boldsymbol{\theta}_{\mathcal{X}_0^c}\|_1$. Denote also $\mathbf{u} = \boldsymbol{\theta}_{\mathcal{X}_0} - \boldsymbol{\theta}^*$. By construction, \mathbf{u} is supported on \mathcal{X}^* . It holds

$$\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_1 = \|\mathbf{u}\| + \|\boldsymbol{\theta}_{\mathcal{X}_0^c}\|_1.$$

Now (9.10) and (9.11) imply

$$\|\boldsymbol{\Psi}^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 + \lambda (\|\boldsymbol{\theta}_{\mathcal{X}_0}\|_1 + \|\boldsymbol{\theta}_{\mathcal{X}_0^c}\|_1) \leq 2\mathbf{C}_0 (\|\mathbf{u}\| + \|\boldsymbol{\theta}_{\mathcal{X}_0^c}\|_1) + \lambda \|\boldsymbol{\theta}^*\|_1$$

and therefore, the component $\boldsymbol{\theta}_{\mathcal{X}_0^c}$ of $\boldsymbol{\theta}$ fulfills

$$\|\boldsymbol{\Psi}^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 + (\lambda - 2\mathbf{C}_0) \|\boldsymbol{\theta}_{\mathcal{X}_0^c}\|_1 \leq 2\mathbf{C}_0 \|\mathbf{u}\|_1 + \lambda \|\boldsymbol{\theta}^*\|_1 - \lambda \|\boldsymbol{\theta}_{\mathcal{X}_0}\|_1 \leq (\lambda + 2\mathbf{C}_0) \|\mathbf{u}\|.$$

Theorem 9.2.3. *Let $\mathbf{Y} = \Psi^\top \boldsymbol{\theta}^* + \boldsymbol{\varepsilon}$ and $\nabla = \Psi^\top \boldsymbol{\varepsilon}$ fulfill $\|\nabla\|_\infty \leq \mathbf{C}_0$. If $\lambda > 2\mathbf{C}_0$, then the estimate $\widehat{\boldsymbol{\theta}}$ fulfills*

$$\|\Psi^\top(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2 + (\lambda - 2\mathbf{C}_0)\|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_1 \leq 2\lambda\|\boldsymbol{\theta}^*\|_1$$

In addition, if $\boldsymbol{\theta}^$ is supported on \varkappa_0 , then the projections $\widehat{\boldsymbol{\theta}}_{\varkappa_0}$ and $\widehat{\boldsymbol{\theta}}_{\varkappa_0^c}$ of $\widehat{\boldsymbol{\theta}}$ to \varkappa^* and \varkappa_0^c can be related as*

$$\|\Psi^\top(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2 + (\lambda - 2\mathbf{C}_0)\|\widehat{\boldsymbol{\theta}}_{\varkappa_0^c}\|_1 \leq (\lambda + 2\mathbf{C}_0)\|\widehat{\boldsymbol{\theta}}_{\varkappa^*} - \boldsymbol{\theta}^*\|_1.$$

To be done: Compatibility condition

To be done: Restricted isometry and oracle bound

9.2.2 Dual problem and Danzig selector

The LASSO optimization can be viewed as minimizing the fit $\|\mathbf{Y} - \Psi^\top \boldsymbol{\theta}\|^2$ under the constraint on the ℓ_1 -norm of $\boldsymbol{\theta}$. Then the objective (9.7) is obtained by Lagrange multiplier method. One can also consider the dual problem: minimizing ℓ_1 -norm of $\boldsymbol{\theta}$ under the fit constraints. The dual problem is known as Danzig selector and reads as

$$\inf \|\boldsymbol{\theta}\|_1 \quad \text{s.t.} \quad \|\Psi(\mathbf{Y} - \Psi^\top \boldsymbol{\theta})\|_\infty \leq 2\lambda$$

The true value $\boldsymbol{\theta}^*$ is a natural candidate. Then $\|\Psi(\mathbf{Y} - \Psi^\top \boldsymbol{\theta})\|_\infty = \|\Psi\boldsymbol{\varepsilon}\|_\infty$. If λ is selected properly to ensure $\|\Psi\boldsymbol{\varepsilon}\|_\infty \leq 2\lambda$, then the constraints meet, and the objective functional is equal to $\|\boldsymbol{\theta}^*\|_1$. Therefore, the solution $\widehat{\boldsymbol{\theta}}$ cannot be less sparse than $\boldsymbol{\theta}^*$ in the sense $\|\widehat{\boldsymbol{\theta}}\|_1 \leq \|\boldsymbol{\theta}^*\|_1$.

9.2.3 Data-driven choice of λ

The necessary requirement to the choice of λ is that the score vector $\nabla = \Psi\boldsymbol{\varepsilon}$

$$\|\nabla\|_\infty \leq 2\lambda.$$

This condition can be assessed by replacing the true noise by its bootstrap counterpart ∇^b . With a pilot $\widetilde{\boldsymbol{\theta}}$, define $\boldsymbol{\varepsilon}^b = \mathcal{W}^b(\mathbf{Y} - \Psi^\top \widetilde{\boldsymbol{\theta}})$ and fix λ^b by the condition

$$P^b(\|\nabla^b\|_\infty > \lambda^b) \leq e^{-x}.$$

The original procedure has to be applied with $\lambda = \lambda^b + \beta$.

References

- Birgé, L. and Massart, P. (1993). Rates of convergence for minimum contrast estimators. *Probab. Theory Relat. Fields*, 97(1-2):113–150.
- Chernozhukov, V., Chetverikov, D., and Kato, K. (2013). Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors. *The Annals of Statistics*, 41(6):2786–2819.
- Ghosal, S. (1999). Asymptotic normality of posterior distributions in high-dimensional linear models. *Bernoulli*, 5(2):315–331.
- Ghosal, S. (2000). Asymptotic normality of posterior distributions for exponential families when the number of parameters tends to infinity. *J. Multivariate Anal.*, 74(1):49–68.
- Ibragimov, I. and Khas'minskij, R. (1981). *Statistical estimation. Asymptotic theory. Transl. from the Russian by Samuel Kotz*. New York - Heidelberg -Berlin: Springer-Verlag .
- Kim, Y. (2006). The Bernstein-von Mises theorem for the proportional hazard model. *Ann. Statist.*, 34(4):1678–1700.
- Kleijn, B. J. K. and van der Vaart, A. W. (2012). The Bernstein-von-Mises theorem under misspecification. *Electronic J. Statist.*, 6:354–381.
- Mammen, E. (1993). Bootstrap and wild bootstrap for high dimensional linear models. *Ann. Stat.*, 21(1):255–285.
- Mammen, E. (1996). Empirical process of residuals for high-dimensional linear models. *Ann. Stat.*, 24(1):307–335.
- Portnoy, S. (1984). Asymptotic behavior of M-estimators of p regression parameters when p^2/n is large. I. Consistency. *Ann. Stat.*, 12:1298–1309.
- Portnoy, S. (1985). Asymptotic behavior of M estimators of p regression parameters when p^2/n is large. II: Normal approximation. *Ann. Stat.*, 13:1403–1417.
- Portnoy, S. (1986). Asymptotic behavior of the empiric distribution of M-estimated residuals from a regression model with many parameters. *Ann. Stat.*, 14:1152–1170.

- Portnoy, S. (1988). Asymptotic behavior of likelihood methods for exponential families when the number of parameters tends to infinity. *Ann. Statist.*, 16(1):356–366.
- Shen, X. (1997). On methods of sieves and penalization. *Ann. Stat.*, 25(6):2555–2591.
- Shen, X. and Wong, W. H. (1994). Convergence rate of sieve estimates. *Ann. Stat.*, 22(2):580–615.
- Spokoiny, V. (2012). In [Spokoiny \(2012\)](#).
- Spokoiny, V., Wang, W., and Härdle, W. (2013). Local quantile regression (with rejoinder). *J. of Statistical Planning and Inference*, 143(7):1109–1129. ArXiv:1208.5384.
- Van de Geer, S. (1993). Hellinger-consistency of certain nonparametric maximum likelihood estimators. *Ann. Stat.*, 21(1):14–44.
- van de Geer, S. (2002). M-estimation using penalties or sieves. *J. Stat. Plann. Inference*, 108(1-2):55–69.
- Zaitsev, A., Burnaev, E., and Spokoiny, V. (2013). Properties of the posterior distribution of a regression model based on gaussian random fields. *Automation and Remote Control*, 74(10):1645–1655.