



TECHNISCHE
UNIVERSITÄT
WIEN



AUTOMATION & CONTROL INSTITUTE
INSTITUT FÜR AUTOMATISIERUNGS-
& REGELUNGSTECHNIK

Hier Titel Einfügen

DIPLOMA THESIS

Conducted in partial fulfillment of the requirements for the degree of a
Diplom-Ingenieur (Dipl.-Ing.)

supervised by

Univ.-Prof. Dr. techn. A. Kugi
Dr. techn. T. Glück

submitted at the

TU Wien

Faculty of Electrical Engineering and Information Technology
Automation and Control Institute

by

Vorname Nachname

Matriculation number Matrikelnummer

Address line 1

Address line 2

Austria

Vienna, Month Day, YEAR

Complex Dynamical Systems Group

A-1040 Wien, Gusshausstr. 27, Internet: <http://www.acin.tuwien.ac.at>

Vorwort

Hier steht das Vorwort

Die Thematik dieser Diplomarbeit ...

Weiterer Text

Noch mehr Text

Vielen Dank!

Wien, XX. Monat, Jahr

Abstract

English

abstract.

More text with linebreak.

Kurzzusammenfassung

Deutsche

Kurzzusammenfassung.

Mehr Text mit linebreak dazwischen.

Contents

1	Introduction	1
1.1	Related Work	2
1.1.1	Parametric Models	2
1.1.2	Non-parametric Models	3
1.1.3	Artificial Neural Networks	5
1.1.4	Look-up Tables	5
1.2	Outline	6
2	Fundamentals	7
2.1	Linear Models	7
2.1.1	Parameter Estimation	9
	The Method of Least Squares	9
	Maximum Likelihood Estimation	10
2.1.2	Estimation of the Variance σ^2	10
	Maximum Likelihood Estimation	10
	Restricted Maximum Likelihood Estimation	11
2.1.3	The Hat Matrix	11
2.2	Model Selection	12
2.2.1	Model Assessment Criteria	12
	Sum of Expected Squared Prediction Error	12
	Corrected Coefficient of Determination R_{corr}^2	14
	Corrected Coefficient of Determination after McFadden $R_{McFadden}^2$	15
	Mallow's Cp	15
	Akaike Information Criterion	15
	Bayesian Information Criteria	16
	Cross-validation	16
2.2.2	Subset Selection Methods	17
2.2.3	Regularization	17
2.3	Splines	18
2.3.1	B-Splines	19
	Tensor-Product B-Splines	22
	Additive Regression	23
2.3.2	P-Splines	25
3	Solution Approach	29
3.1	Shape-constraint P-splines	30
3.1.1	Monotonic increasing constraint	30
3.1.2	Monotonic decreasing constraint	31

3.1.3	Convex constraint	31
3.1.4	Concave constraint	32
3.1.5	Peak constraint	32
3.1.6	Valley constraint	33
3.1.7	Jamming constraint	33
3.1.8	Boundedness constraint	33
3.2	Shape-constraint tensor-product P-splines	34
3.2.1	Monotonic increasing constraint in dimension 1	34
3.2.2	Monotonic decreasing constraint for dimension 1	35
3.2.3	Convex constraint for dimension 1	35
3.2.4	Concave constraint for dimension 1	36
3.2.5	Shape constraints in two dimensions	36
Appendices		37
A Appendix A		37
A.1	Kronecker product	37
A.2	Row-wise Kronecker product	37
B Appendix B		39

List of Figures

2.1	Bias-variance trade-off.	14
2.2	B-spline basis function of order $l = 0, 1, 2, 3$ for equidistant (left) and non-equidistant (right) knots.	19
2.3	Spatial neighborhood or adjacent parameters for a tensor-product B-spline.	27

List of Tables

3.1	Overview of the considered constraints	29
3.2	Overview of the constraints for shape-constraint tensor-product P-splines.	34

List of Abbreviations

PSM permanent-magnet synchronous motor

1 Introduction

The ongoing digitalization of complex, large-scale industrial plants is leading to a massive increase of process data. This data can be used to enhance the overall understanding of the characterizing physical process inside the plant. Modern observer and control concepts are used to enhance the efficiency and quality of the plant. They use mathematical models of the ongoing process. The exact physical description of the relevant quantities as mathematical model is nevertheless often not feasible because of the complexity of the process as well as computational and measurement limitations.

Data driven approaches are state-of-the-art in many fields, including e.g. image and speech recognition. The usage of data driven methods, e.g. artificial neural networks, parametric models, etc., to model process quantities for which measurements are expensive or not practical, gains more and more influence and acceptance in the field of process control and optimization. The application of the specific algorithm depends on issues like data quality, interpretability of the model and computational efficiency.

In most process optimization tasks massive amounts of domain specific knowledge in form of physical theories and a priori knowledge is available. The combination of the use of domain knowledge and data driven modeling techniques is called hybrid modeling or grey-box modeling. It lies between the two modeling extrema of white-box models, which are derived from first principles and physical models and black-box models, which are derived from data only[1]. The incorporation of this knowledge in state-of-the-art data driven approaches is not trivial and not solved for some algorithms. Nevertheless, its inclusion should improve the interpretability, which itself is of importance in the context of the emerging field of explainable artificial intelligence (XAI). XAI refers to modeling approaches and techniques in which the main goal is that the resulting model is understandable by humans[2].

In this thesis, we are going to develop an algorithm for efficient, static, multi-dimensional function approximation using a priori domain knowledge. The algorithm is based on structured additive regression using splines[3]. We use user-defined constraints to include the a priori domain knowledge in the fitting process[4]. It should produce interpretable and efficient models based on the given domain knowledge. The incorporation of domain specific knowledge should improve the model quality and robustness as well as interpolation and extrapolation behavior in situations where the measured data is sparse and/or noisy. Further, we are going to evaluate the algorithm using noisy samples from artificial test functions with known behavior as well as real world data collected in a heat treatment process.

1.1 Related Work

We will now discuss some of the most used data-driven algorithms. The discussion includes the following model approaches:

- Parametric models: Linear and polynomial regression
- Non-parametric models: Basis function models
- Artificial neural networks
- Look-up tables

and focuses on the interpretability, computational efficiency and the ability to include domain knowledge of the individual modeling approaches. The list given above is not intended to be complete.

The common starting point for the different data-driven modeling approaches is that we have some data $\{x_1^{(i)}, \dots, x_q^{(i)}; y^{(i)}\}$, $i = 1, \dots, n$. We restrict the following discussion to the single-input setting, i.e. $\{x^{(i)}, y^{(i)}\}$, $i = 1, \dots, n$. The generalization to multiple input dimensions is given in the respective literature. We then use the given data to estimate a model function $f(x)$ which is then used to predict the response or output variable y , i.e.

$$y = f(x). \quad (1.1)$$

Therefore we are in the setting of supervised learning.

1.1.1 Parametric Models

According to Nelles, parametric models are defined as models that can describe the true process behavior using a finite number of parameters[5]. An example is given by the linear regression model for one input variable x as

$$y = f(x) = \beta_0 + \beta_1 x. \quad (1.2)$$

Both parameters β_0 and β_1 allow for a direct interpretation as β_0 is the intercept, i.e. the output for the input $x = 0$, and β_1 is the slope, i.e. the constant defining the relationship between the increase of the output y with respect to the increase of the input x . The interpretability of linear regression models is therefore very high.

Linear regression models are widely used and part of standard software tools. Their parameters can be efficiently computed using the least squares algorithm. One major drawback is that they can only recover linear relationships between input and output variables. They are therefore quite restrictive and do not allow the incorporation of a priori domain knowledge except being increasing or decreasing by the sign of the slope β_1 .

An extension of the linear regression model is given by polynomial regression. Here, we try to model the output data y using a polynomial of degree p , i.e.

$$y = f(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \cdots + \beta_p x^p. \quad (1.3)$$

Polynomial regression introduces more flexibility in the fitting process, since the restriction of linear relationship is relaxed to a polynomial relationship of degree p . As for linear models, the interpretability of the parameters is given. We can use the least squares algorithm for parameter estimation. The incorporation of some a priori domain knowledge is possible, e.g. as the degree of the polynomial regression model. The major problem of polynomial regression is that the model function becomes quite wiggly for high polynomial degrees p .

Linear and polynomial regression models are so-called global models. Their parameters act on the complete input space. This property makes the incorporation of specific a priori domain knowledge, e.g. unimodal behavior, difficult and in most cases nearly impossible for parametric models.

1.1.2 Non-parametric Models

Nelles defines non-parametric models as models which require an infinite number of parameters to describe a process exactly[5]. In almost all practical applications this infinite series is approximated by a finite number of parameters using the basis function approach given by

$$y = f(x) = \sum_{i=1}^M \theta_i^{(l)} \Phi_i(x, \theta_i^{(nl)}) \quad (1.4)$$

with the linear parameters $\theta_i^{(l)}$, the basis functions $\Phi_i(\cdot)$, the input variable x and the non-linear parameters $\theta_i^{(nl)}$. The output y is therefore given by a linear combination of M basis functions $\Phi_i(\cdot)$. To model a non-linear relationship between y and x , the basis functions $\Phi(\cdot)$ need to be non-linear. Commonly used basis functions are e.g. the *hat function*, the *Gaussian*, *splines* or the *hinge function*.

A commonly used algorithm utilizing the basis function approach is called Multivariate Adaptive Regression Splines (MARS) [6]. MARS approximates data, as example again for a single input dimension, using the following model

$$y = \sum_{i=1}^M \theta_i \Phi_i(x) \quad (1.5)$$

using constant parameters θ_i . The basis functions $\Phi_i(\cdot)$ are one of the following three alternatives:

1. $\Phi_i(x) = 1$, representing the intercept.
2. $\Phi_i(x) = \max(0, x - c_i)$ or $= \max(0, c_i - x)$, representing the *hinge function* h_i using the constant value c_i .

3. $\Phi_i(x) = h_i h_j$, representing a product of two *hinge functions*.

MARS fits the model using a recursive splitting approach. More information can be found in [6] and [7]. MARS models are more flexible compared to the parametric linear and polynomial regression models. As only hinge functions and products of hinge functions are used, MARS models are efficient and in general simple to understand and interpret. To our knowledge, there is currently no possibility to include a priori domain knowledge in the fitting process when using MARS.

Another widely used methods using basis functions is the use of *splines*. Splines are defined as piece-wise polynomials on a sequence of knots. Further information can be found in Section 2.3 and [8]. Using k splines to model some data, we obtain the following model formulation

$$y = f(x) = \sum_{i=1}^k \beta_i b_i(x) \quad (1.6)$$

with the spline basis functions b_i and the parameters β_i . The parameters can be calculated using the least squares algorithm. The usage of splines allows a lot of flexibility and computational efficiency. A priori domain knowledge can be incorporated using the penalized least squares algorithm with a sophisticated choice of mapping and weighting matrices, see 3, and e.g. [4] or [27].

The basis function approach in (1.4) may be extended by changing the parameters $\theta_i^{(l)}$ to more complex forms. An example for this is the so-called local linear neuro-fuzzy model, for which each parameter $\theta_i^{(i)}$ is changed to be a *local linear model* and each basis function $\Phi_i(\cdot)$ is then called *validity function* determining the region of validity of the local linear model [5]. The validity functions are normalized for any model input x , i.e.

$$\sum_{i=1}^M \Phi_i(x) = 1. \quad (1.7)$$

and typically chosen to be *Gaussian* functions, i.e.

$$\Phi_i(x) = a_i \exp\left(-\frac{(x - \mu_i)^2}{\sigma_i^2}\right) \quad (1.8)$$

with the normalization constant a_i and the parameters μ_i and σ_i determining the location and scale of the Gaussian function. The output of the local linear neuro-fuzzy model using M local linear models is then given by

$$y = \sum_{i=1}^M (\beta_{i0} + \beta_{i1}x_1) \Phi_i(x). \quad (1.9)$$

The first term in the summation are the *local linear models*. The parameters β_{ij} for $i = 1, \dots, M$ and $j = 0, 1$ as well as the parameters μ_i and σ_i from the validity functions Φ_i need to be optimized. This is done using the LOLIMOT algorithm. Further information is given in [5].

Local linear models as extension of linear models possess more flexibility with regards to non-linear relationships in the data. They can also be efficiently evaluated after the iterative training process. The interpretability is high since each local linear model contributes to the prediction according to its validity function. The ability to include a priori domain knowledge in the fitting process is currently not available.

1.1.3 Artificial Neural Networks

Artificial neural networks are currently the state-of-the-art solution method for many problems ranging from computer vision over time-series prediction to regression tasks. They are constructed as coarse model of the human brain, consisting of neurons which are connected by some weights. These connections are adapted in the learning process using an algorithm called "backpropagation". They utilize a high number of parameters to model hidden, high-dimensional relationships in the data. Further information can be found in standard textbook about neural networks, e.g. [9] or [10].

In terms of modeling flexibility, artificial neural networks of sufficient size are proven to be able to represent a wide variety of functions by so-called universal approximation theorems, cf. [11] and [12]. The computational complexity of a neural network depends on its size, aka. the number of parameters. Large networks need many training samples to generate sufficiently accurate predictions. Artificial neural networks are an example of a black-box model. The inclusion of a priori domain knowledge into the learning process of neural networks is possible for specific types of knowledge using the concepts of hints, see [13] and [14].

1.1.4 Look-up Tables

A look-up table is an array of values, which allows to replace computational expensive computations with inexpensive array indexing operations. The values in the look-up table are most often computed and stored beforehand. To gain higher resolution, interpolation techniques such as linear or quadratic interpolation may be applied to look-up tables.

Look-up tables are a standard tool in many fields. They are extremely efficient in terms of computation time. One problem that occurs is the exponential increase in size with the number of dimensions for the look-up table. As example, a 2×2 -table needs to save 4 values, while a $2 \times 2 \times 2$ table already needs 8 values without gaining additional accuracy. Another problem is that the values in the look-up table may come from complex, computational or physical models.

Lattice regression tackles this problems by jointly estimating all lookup-table values by minimizing the regularized interpolation error on training data [15]. They state that using ensembles of lookup-tables which combine several *tiny* lattices enables linear scaling in the number of input dimension even for high dimensions [16]. They further state that

lattice regression may be used to incorporate a priori domain knowledge like monotonicity, shape or unimodality into the fitting process, see [17] or [18].

1.2 Outline

The base of this thesis is a literature study about function fitting using a priori domain knowledge focusing on non-parametric techniques and neural networks. We decided to use structured additive regression [3] utilizing B-splines and tensor product splines as non-linear basic functions. This approach enables flexible, multi-dimensional function fitting. We further expanded this method by applying a priori domain knowledge through the use of user-defined constraints. These constraints consist of mapping matrices determined by the type of constraint, and weighting matrices, determining whether the constraint is active, see [4] and [27].

We are able to incorporate the following a priori domain knowledge:

- Monotonicity, i.e. $f'(x) \geq 0$ or $f'(x) \leq 0$
- Curvature, i.e. $f''(x) \geq 0$ or $f''(x) \leq 0$
- Unimodality, i.e.
 $f'(x) > 0, x < m$ and $f'(x) < 0, x > m$ for $m = \arg \max_x f(x)$
- Boundedness, i.e. $f(x) \geq M$ or $f(x) \leq M$ for some value M
- Jamming, i.e. $f(x^{(p)}) \approx y^{(p)}$ for some point p with high fidelity

The thesis is divided into 5 chapters: Chapter 2 provides an overview of the fundamental mathematical concepts used. We focus on the description of linear models, splines and the topic of structured additive regression. In chapter 3, we develop the algorithm using the concepts given in chapter 2. In chapter 4, we test the algorithm using different artificial functions and a priori domain knowledge as well as real-world data. Chapter 5 gives a summary and outlines future, possible work.

2 Fundamentals

This chapter summarizes the fundamentals of regression. Excellent overviews can be found in the textbooks [19], [3], [7]. The shown fundamentals are strongly aligned with the presentation given in [3]. Section 2.1 gives an overview of the model assumptions used throughout this work. Furthermore, Section 2.2 outlines methods to evaluate and compare different models again each other in terms of complexity and accuracy. Section 2.3 is devoted to the spline definitions.

2.1 Linear Models

Given the data set $D = \{(x_1^{(i)}, \dots, x_q^{(i)}, y^{(i)})\}$, $i = 1, 2, \dots, n$, we aim to model the relation between the inputs x_1, \dots, x_q and the output y with a function $f(x_1, \dots, x_q)$. Since this relationship is not exact, there will be a random part ϵ . It is typically assumed that this part is additive and thus

$$y = f(x_1, \dots, x_q) + \epsilon. \quad (2.1)$$

We would like to estimate the unknown function f . For this, some assumptions on the model structure are made:

1. *The unknown function f is a linear combination of the inputs*

The function $f(x_1, \dots, x_q)$ is modeled as a linear combination of inputs, i.e. in the form

$$f(x_1, \dots, x_q) = \beta_0 + \beta_1 x_1 + \dots + \beta_q x_q, \quad (2.2)$$

with unknown parameters β_0, \dots, β_q . Note that the model (2.2) is linear in its parameters as well as in its inputs. The parameter β_0 is called intercept or bias in the machine learning community, see [9]. We introduce the input vectors $\mathbf{x}^T = [1, x_1, \dots, x_q] \in \mathbb{R}^{1 \times q+1}$ and the parameter vector $\boldsymbol{\beta}^T = [\beta_0, \beta_1, \dots, \beta_q] \in \mathbb{R}^{1 \times q+1}$ to obtain

$$y = f(\mathbf{x}^T) = \mathbf{x}^T \boldsymbol{\beta}. \quad (2.3)$$

2. *Additive errors*

An additional assumption of linear models is additivity of errors, which means that

$$y = \mathbf{x}^T \boldsymbol{\beta} + \epsilon. \quad (2.4)$$

This is reasonable for many practical applications, even though it appears quite restrictive.

To estimate the unknown parameters or coefficients β , we define the output vector $\mathbf{y}^T = [y^{(1)}, \dots, y^{(n)}] \in \mathbb{R}^{1 \times n}$ and error vector $\epsilon^T = [\epsilon^{(1)}, \dots, \epsilon^{(n)}] \in \mathbb{R}^{1 \times n}$ as well as the design matrix

$$\mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_q^{(1)} \\ \vdots & & & \vdots \\ 1 & x_1^{(n)} & \dots & x_q^{(n)} \end{bmatrix} \in \mathbb{R}^{n \times q+1} \quad (2.5)$$

and generate n equations like (2.4), which can be combined to

$$\mathbf{y} = \mathbf{X}\beta + \epsilon. \quad (2.6)$$

We assume that the design matrix \mathbf{X} has full column rank, i.e. $\text{rank}(\mathbf{X}) = q + 1 = p$, implying linear independence of the columns of \mathbf{X} , which is necessary to obtain a unique estimator for the regression coefficients β , see [3].

Another necessary requirement is that the number of data points n is larger or equal to the number of regression parameters p , which is equivalent to the statement that the linear system in (2.6) is not underdetermined.

In addition to the assumptions on the unknown function f , the necessary assumptions on the error vector ϵ are the following [3]:

1. *Expectation of the error*

The errors have a mean value of zero, i.e. $E(\epsilon^{(i)}) = 0$

2. *Variances and correlation structure of the errors*

We assume constant error variance with $\text{Var}(\epsilon^{(i)}) = \sigma^2$. This property is called homoscedasticity. Additionally, we assume that the errors are uncorrelated, which means $\text{Cov}(\epsilon^{(i)}, \epsilon^{(j)}) = 0$ for $i \neq j$. The combination of these assumptions lead to the covariance matrix $\text{Cov}(\epsilon) = E(\epsilon\epsilon^T) = \sigma^2\mathbf{I}$.

3. *Assumptions on the input and design matrix*

We have to distinguish two cases where the inputs are deterministic or stochastic. In most cases, the inputs and the output are stochastic and hence all model assumptions are conditioned on the design matrix (2.5). This means that the input $\mathbf{x}^{(i)}$ and the errors ϵ are not stochastically independent. For notational simplicity, we usually suppress the dependence on the design matrix.

4. *Gaussian errors*

The errors follow at least approximately a normal distribution. Together with Assumptions 1 and 2, we obtain that $\epsilon^{(i)} = \mathcal{N}(0, \sigma^2)$ holds.

Summarizing, we have the following model assumptions:

$$E(\mathbf{y}) = \mathbf{X}\beta \quad (2.7)$$

$$\text{Cov}(\mathbf{y}) = \sigma^2 \mathbf{I}, \quad (2.8)$$

yielding

$$\mathbf{y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}). \quad (2.9)$$

A linear model with multiple inputs can therefore be interpreted as a multi-variate normal distribution with its mean vector given by $\boldsymbol{\mu} = \mathbf{X}\boldsymbol{\beta}$ and its covariance matrix given by $\sigma^2 \mathbf{I}$, i.e. to specify the linear model given in (2.9), we need to estimate the regression coefficients $\boldsymbol{\beta}$ and the variance σ^2 .

2.1.1 Parameter Estimation

The linear model given in (2.9) features the unknown parameters $\boldsymbol{\beta}$ and σ , which need to be estimated using the given data. In the following, the estimators $\hat{\boldsymbol{\beta}}$ and $\hat{\sigma}$ are introduced, and their statistical properties are derived. The two methods to estimate the regression parameters in the context of linear models are the method of Least Squares (LS) and the method of Maximum Likelihood (ML).

The Method of Least Squares

The unknown regression parameters $\boldsymbol{\beta} \in \mathbb{R}^p$ are estimated by minimizing the sum of squared error

$$\text{LS}(\mathbf{y}, \boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 = \sum_{i=1}^n \epsilon_i^2 = \boldsymbol{\epsilon}^T \boldsymbol{\epsilon}, \quad (2.10)$$

with respect to $\boldsymbol{\beta}$, see, e.g. [7]. Rewriting (2.10) leads to the least squares criterion

$$\begin{aligned} \text{LS}(\mathbf{y}, \boldsymbol{\beta}) &= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\ &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta}. \end{aligned} \quad (2.11)$$

The first-order necessary condition for optimality, cf. [20], reads as

$$\frac{\partial \text{LS}(\mathbf{y}, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = 0. \quad (2.12)$$

The second-order condition requires the Hessian to be positive-definite, i.e.

$$\frac{\partial^2 \text{LS}(\mathbf{y}, \boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} = 2\mathbf{X}^T \mathbf{X} > 0. \quad (2.13)$$

Since the design matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ is assumed to have full rank, the matrix $\mathbf{X}^T \mathbf{X}$ is positive definite. The least squares estimate $\hat{\boldsymbol{\beta}}_{LS}$ is hence obtained, see (2.12), by solving the so-called *normal equations*

$$\mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^T \mathbf{y}. \quad (2.14)$$

Since $\mathbf{X}^T \mathbf{X}$ is positive definite, the *normal equations* (2.14) feature a unique solution given by the least squares estimator

$$\hat{\boldsymbol{\beta}}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}. \quad (2.15)$$

Maximum Likelihood Estimation

Under the normality assumption, the likelihood is defined as [19]

$$\mathcal{L}(\boldsymbol{\beta}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right). \quad (2.16)$$

The log-likelihood is then given by

$$l(\boldsymbol{\beta}, \sigma^2) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}). \quad (2.17)$$

Thus, maximizing the log-likelihood $l(\boldsymbol{\beta}, \sigma^2)$ with respect to $\boldsymbol{\beta}$ is equivalent to minimizing the least squares criterion given in (2.10). The maximum likelihood estimator $\hat{\boldsymbol{\beta}}_{ML}$ is therefore equivalent to the least squares estimator $\hat{\boldsymbol{\beta}}_{LS}$ in (2.15).

2.1.2 Estimation of the Variance σ^2

The estimation of the variance σ^2 is necessary for the construction of confidence intervals of the regression parameters and for the construction of prediction intervals. It is further used in all kinds of statistical tests as well as in model selection approaches and model assessment criteria [21].

Maximum Likelihood Estimation

The first-order necessary condition for optimality results in this case in

$$\frac{\partial l(\boldsymbol{\beta}, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = 0. \quad (2.18)$$

Substituting the maximum likelihood estimator $\hat{\boldsymbol{\beta}}_{LS}$, given in (2.15), for $\boldsymbol{\beta}$ results in the maximum likelihood estimator

$$\hat{\sigma}_{ML}^2 = \frac{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{LS})^T(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{LS})}{n} = \frac{\hat{\boldsymbol{\epsilon}}^T \hat{\boldsymbol{\epsilon}}}{n}. \quad (2.19)$$

This estimator for σ^2 is rarely used since it is biased, i.e. $E(\hat{\sigma}_{ML}^2) \neq \sigma^2$, see [3].

Restricted Maximum Likelihood Estimation

The mean value of the sum of squared residuals is $E(\hat{\boldsymbol{\epsilon}}^T \hat{\boldsymbol{\epsilon}}) = (n - p)\sigma^2$. Hence, another estimator for σ^2 is given by

$$\hat{\sigma}_{REML}^2 = \frac{1}{n - q} \hat{\boldsymbol{\epsilon}}^T \hat{\boldsymbol{\epsilon}}. \quad (2.20)$$

The restricted maximum likelihood estimator (REML) $\hat{\sigma}_{REML}^2$ is in general less biased [3]. Therefore, it is the commonly used estimator for the variance σ^2 .

2.1.3 The Hat Matrix

Using the least squares estimator (2.15), we can estimate the mean of \mathbf{y} by

$$\widehat{E(\mathbf{y})} = \hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}_{LS}. \quad (2.21)$$

Substituting (2.15) results in

$$\hat{\mathbf{y}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H}\mathbf{y}, \quad (2.22)$$

with the matrix $\mathbf{H} \in \mathbb{R}^{n \times n}$, which is called *hat matrix*. Using the *hat matrix*, we can express the residuals $\hat{\epsilon}^{(i)} = y^{(i)} - \hat{y}^{(i)}$ in matrix notation as

$$\hat{\boldsymbol{\epsilon}} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I} - \mathbf{H})\mathbf{y}. \quad (2.23)$$

The *hat matrix* \mathbf{H} has the following useful properties:

- \mathbf{H} is symmetric.
- \mathbf{H} is idempotent, i.e. $\mathbf{H}^2 = \mathbf{H}$.
- The rank of \mathbf{H} is equal to its trace.
- $\frac{1}{n} \leq h_{ii} \leq 1$, if all data points are different, i.e. $x^{(i)} \neq x^{(j)}$ for $i \neq j$. Here, h_{ii} are the diagonal elements of \mathbf{H} .
- The matrix $(\mathbf{I} - \mathbf{H})$ is also idempotent and symmetric, with $\text{rank}(\mathbf{I} - \mathbf{H}) = n - p$.

The hat matrix is used in model selection techniques like cross-validation since its trace acts as a measure for the degrees of freedom of the model, as well as in outlier detection and in diagnostic plots for linear models.

2.2 Model Selection

Linear models are widely exploited for regression problems on large data sets ($n \gg 0$), because the solution of the *normal equations* (2.14) can be computed efficiently even for large n . If these data sets also contain a large number of input variables ($q \gg 0$), the situation becomes more complicated since possible interaction effects or correlation of input variables may occur. These interaction terms limit the, otherwise perfect, interpretability of the linear model.

Therefore, we need techniques and criteria to select the *best possible model* out of the variety of different models for a given data set. Model assessment criteria, see Section 2.2.1, are used to compare different models while subset selection techniques, see Section 2.2.2, give an algorithmic approach to model generation. Further, we can influence the estimated coefficients β directly via regularization, see Section 2.2.3.

2.2.1 Model Assessment Criteria

One way of comparing various models, i.e. models using different sets of inputs, is the use of model assessment criteria. Generally, model assessment criteria can be split in two components. The first one measures the goodness of fit, e.g. using the sum of squared errors, while the second measures the complexity of the model. Most model assessment criteria are based on the sum of the expected squared prediction error (SPSE), which is also known as *generalization error*. Therefore, the derivation of the SPSE is given next.

Sum of Expected Squared Prediction Error

Given independent observations y_i , $i = 1, 2, \dots, n$ and a subset of inputs $\{x_0 = 1, x_1, \dots, x_q\}$, we want to measure the prediction quality. The specific models are defined by numbers $M \subset \{0, 1, \dots, q\}$ of used inputs with corresponding design matrix \mathbf{X}_M . Moreover, $|M|$ is the cardinal number of M , i.e. the number of inputs included in the model. The least squares estimator for β , cf. (2.15), is then given by

$$\hat{\beta}_M = (\mathbf{X}_M^T \mathbf{X}_M)^{-1} \mathbf{X}_M^T \mathbf{y}.$$

The data \mathbf{y} can be interpreted as a random variable. We can then define an estimator $\hat{\mathbf{y}}_M$ for the vector $\boldsymbol{\mu}$ of expectations $\mu^{(i)} = E(y^{(i)})$ by

$$\hat{\mathbf{y}}_M = \mathbf{X}_M \hat{\beta}_M. \quad (2.24)$$

Moreover, it is easy to show that the following properties hold true:

- (i) $E(\hat{\mathbf{y}}_M) = \mathbf{X}_M (\mathbf{X}_M^T \mathbf{X}_M)^{-1} \mathbf{X}_M^T E(\mathbf{y})$
- (ii) $\text{Cov}(\hat{\mathbf{y}}_M) = \sigma^2 \mathbf{X}_M (\mathbf{X}_M^T \mathbf{X}_M)^{-1} \mathbf{X}_M^T$
- (iii) $\sum_{i=1}^n \text{Var}(\hat{y}_M^{(i)}) = \sigma^2 \text{tr}(\mathbf{X}_M (\mathbf{X}_M^T \mathbf{X}_M)^{-1} \mathbf{X}_M^T) = \sigma^2 |M|$

(iv) Sum of Mean Squared Error

$$\begin{aligned}
 \text{SMSE} &= \sum_{i=1}^n \mathbb{E}(\hat{y}_M^{(i)} - \mu_M^{(i)})^2 \\
 &= \sum_{i=1}^n \mathbb{E}\left((\hat{y}_M^{(i)} - \mathbb{E}(\hat{y}_M^{(i)})) + (\mathbb{E}(\hat{y}_M^{(i)}) - \mu_M^{(i)})\right)^2 \\
 &= |M|\sigma^2 + \sum_{i=1}^n (\mathbb{E}(\hat{y}_M^{(i)}) - \mu_M^{(i)})^2.
 \end{aligned} \tag{2.25}$$

Note that the estimator (2.24) can be regarded as a prediction of future observations of the form

$$y^{(n+i)} = \mu^{(i)} + \epsilon^{(n+i)} \tag{2.26}$$

for new input data $\{(x_1^{(i)}, \dots, x_q^{(i)}), i = 1, 2, \dots, n\}$. Thus, we can derive the SPSE as

$$\begin{aligned}
 \text{SPSE} &= \sum_{i=1}^n \mathbb{E}(y^{(n+i)} - \hat{y}_M^{(i)})^2 \\
 &= \sum_{i=1}^n \mathbb{E}((y^{(n+i)} - \mu_M^{(i)}) - (\hat{y}_M^{(i)} - \mu_M^{(i)}))^2 \\
 &= \sum_{i=1}^n \mathbb{E}(y^{(n+i)} - \mu_M^{(i)})^2 + 2\mathbb{E}((y^{(n+i)} - \mu_M^{(i)})(\hat{y}_M^{(i)} - \mu_M^{(i)})) + \mathbb{E}(\hat{y}_M^{(i)} - \mu_M^{(i)})^2 \\
 &= \sum_{i=1}^n \mathbb{E}(y^{(n+i)} - \mu_M^{(i)})^2 + \sum_{i=1}^n (\mathbb{E}(\hat{y}_M^{(i)}) - \mu_M^{(i)})^2 \\
 &= n\sigma^2 + \text{SMSE} \\
 &= n\sigma^2 + |M|\sigma^2 + \sum_{i=1}^n (\mathbb{E}(\hat{y}_M^{(i)}) - \mu_M^{(i)})^2.
 \end{aligned} \tag{2.27}$$

The SPSE can be split into three parts:

1. *Irreducible Prediction Error Term:* $n\sigma^2$

This term cannot be reduced through model selection techniques since it only contains the number of data points n and the variance σ^2 .

2. *Variance Error Term:* $|M|\sigma^2$

The second term contains the number of used variables $|M|$ as well as the variance σ^2 . It can therefore be reduced by reducing the model complexity, i.e. by using a smaller number of inputs.

3. *Squared Bias Error Term:* $\sum_{i=1}^n (\mathbb{E}(\hat{y}_M^{(i)}) - \mu_M^{(i)})^2$

The last term can be interpreted as bias. It can be reduced by increasing the model complexity, i.e. by using additional inputs.

The SPSE acts as an example of the bias-variance trade-off, which is characteristic for all statistical models. It states that by increasing the model complexity, the bias is reduced but instead the variance is increased. On the other hand, by decreasing model complexity, the variance of the model is reduced, but the bias is increased, see Figure 2.1 and [9].

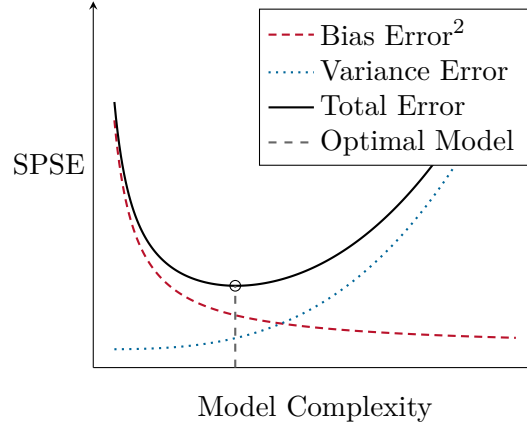


Figure 2.1: Bias-variance trade-off.

In practice, the true value for the SPSE is not accessible since $\mu^{(i)}$ and σ^2 are unknown. Therefore, we need to estimate the SPSE. This can be done by using one of the following two strategies:

1. *Estimate SPSE using new and independent data*

If new observations are available, the SPSE can be estimated by

$$\widehat{\text{SPSE}} = \sum_{i=1}^n (y^{(n+i)} - \hat{y}_M^{(i)})^2. \quad (2.28)$$

These new observations can also be some held-out validation data from a train-validation split of the given data.

2. *Estimate SPSE using existing data*

When using existing data, the estimate for the SPSE is given by the squared error and an additional term depending on the estimated variance and the model complexity.

The estimate is thus given by

$$\widehat{\text{SPSE}} = \sum_{i=1}^n (y^{(i)} - \hat{y}_M^{(i)})^2 + 2|M|\hat{\sigma}^2. \quad (2.29)$$

Typically used model assessment criteria follow the basic idea of the SPSE, see [3].

Corrected Coefficient of Determination R_{corr}^2

The corrected coefficient of determination R_{corr}^2 is an improvement over the coefficient of determination R^2 , which is defined as

$$R^2 = 1 - \frac{\sum_{i=1}^n (y^{(i)} - \hat{y}_M^{(i)})^2}{\sum_{i=1}^n (y^{(i)} - \bar{y})^2}, \quad (2.30)$$

where $\bar{y} = E(y^{(i)})$ is the mean value of \mathbf{y} . The major drawback of R^2 is that it will never decrease when further inputs are included in the model, e.g. the R^2 of a model using $\{x_1, x_2, x_3\}$ is always larger or equal the R^2 of a model using $\{x_1, x_2\}$, even if the variable does not enhance the prediction quality.

The corrected coefficient of determination R_{corr}^2 reduces this problem by an correction term depending on the number of parameters and is given by

$$R_{corr}^2 = 1 - \frac{n-1}{n-p}(1 - R^2). \quad (2.31)$$

The corrected coefficient of determination is a standard output parameter in many statistical programs and may be used to compare even models with different number of used variables [3].

Corrected Coefficient of Determination after McFadden $R_{McFadden}^2$

The corrected coefficient of determination after McFadden is defined as

$$R_{McFadden}^2 = 1 - \frac{\ln(L_M) - M}{\ln(L_0)} \quad (2.32)$$

using the likelihood of the model M given by L_M and the likelihood of the zero model L_0 . A standard zero model is given by the mean value $\bar{y} = E(y^{(i)})$. Higher values of $R_{McFadden}^2$ correspond to better fits.

Mallow's Cp

Mallow's complexity parameter is based directly on the ideas specified for the estimation of the SPSE and is given by

$$C_p = \frac{\sum_{i=1}^n (y^{(i)} - \hat{y}_M^{(i)})^2}{\hat{\sigma}^2} - n + 2|M|. \quad (2.33)$$

A lower value of Mallow's C_p corresponds to a better model fit [3].

Akaike Information Criterion

The AIC is among the most used model assessment criteria and defined by

$$\text{AIC} = -2\ln(\hat{\beta}_{ML}, \hat{\sigma}_{ML}^2) + 2(|M| + 1), \quad (2.34)$$

where $l(\hat{\beta}_{ML}, \hat{\sigma}_{ML}^2)$ is the value of the log-likelihood (2.17) at its maximum, i.e. at $\hat{\beta}_{ML}$ and $\hat{\sigma}_{ML}$. It is worth noting that the total number of parameters is $|M| + 1$ because the variance is also counted as parameter. The log-likelihood for a linear model assuming Gaussian errors is given by, cf. (2.17),

$$-2l(\hat{\beta}_{ML}, \hat{\sigma}_{ML}^2) = n \log(\hat{\sigma}_{ML}^2) + n. \quad (2.35)$$

Therefore, neglecting the constant n , the AIC evaluates to

$$\text{AIC} = n \log(\hat{\sigma}_{ML}^2) + 2(|M| + 1). \quad (2.36)$$

A lower value of the AIC means a to a better model fit [3].

Bayesian Information Criteria

The BIC is similar to the AIC, but it penalizes more complex models much harder than the AIC. In its general form, it is given as

$$\text{BIC} = -2l(\hat{\beta}_{ML}, \hat{\sigma}_{ML}^2) + \log(n)(|M| + 1). \quad (2.37)$$

Again, assuming Gaussian errors for a linear model and neglecting the constant term n , the BIC evaluates to

$$\text{BIC} = n \log(\hat{\sigma}_{ML}^2) + \log(n)(|M| + 1). \quad (2.38)$$

A lower value of the BIC correspond to a better model fit [3].

Cross-validation

The basic idea of cross-validation (CV) is to split the given data set into a training set to estimate the parameters and a validation set to assess the prediction quality. A special case of cross-validation is the "leave-one-out" cross-validation, where all but one data point are used for training and the model is then evaluated on this held-out data point. This seems to be quite expensive, since one needs to estimate one model per data point. However, in the context of linear models, it can be shown that the cross-validation score can be computed using one model trained on all data \mathbf{y} and the *hat matrix* $\mathbf{H}_M = \mathbf{X}_M(\mathbf{X}_M^T \mathbf{X}_M)^{-1} \mathbf{X}_M^T$. The cross-validation score is then given by

$$\text{CV} = \frac{1}{n} \sum_{i=1}^n \left(\frac{y^{(i)} - \hat{y}_M^{(i)}}{1 - h_{ii,M}} \right)^2, \quad (2.39)$$

where $h_{ii,M}$ denote the diagonal elements of the *hat matrix* and $\hat{y}_M^{(i)}$ is defined as the prediction of the model for the input $\{x_1^{(i)}, \dots, x_q^{(i)}\}$. A lower cross-validation score corresponds to a better model fit [22].

An approximation to the cross-validation score is given by the so-called generalized cross-validation (GCV) score. It is mainly used in the context of non-parametric regression or when the hat matrix \mathbf{H}_M is numerically expensive to compute. In the GCV score, the diagonal elements of the hat matrix $h_{ii,M}$ are replaced by the mean of the trace of \mathbf{H}_M . The GCV score is then given by

$$\text{GCV} = \frac{1}{n} \sum_{i=1}^n \left(\frac{y^{(i)} - \hat{y}_M^{(i)}}{1 - \text{trace}(\mathbf{H}_M)/n} \right)^2. \quad (2.40)$$

The numerical advantage comes from the fact that the trace of a product of matrices is invariant to cyclical permutations, i.e.

$$\text{trace}(\mathbf{H}_M) = \text{trace}(\mathbf{X}_M(\mathbf{X}_M^T \mathbf{X}_M)^{-1} \mathbf{X}_M^T) = \text{trace}(\mathbf{X}_M^T \mathbf{X}_M (\mathbf{X}_M^T \mathbf{X}_M)^{-1}). \quad (2.41)$$

The trace can therefore be computed from the product of two matrices of shape $p \times p$, see [3].

2.2.2 Subset Selection Methods

To make use of the various model assessment criteria, some algorithmic approach to model selection needs to be given. The most commonly used approaches are forward, backward and stepwise selection [3].

In forward selection, we start with a candidate model, which includes a small number of variables. In each iteration of forward selection, an additional variable is added to the candidate model. The added variable is the one which leads to the largest reduction of a predefined model assessment criteria. The algorithm stops, if no further reduction is achieved.

In backward selection, we start with a candidate model, which includes all variables. In each iteration of backward selection, we eliminate the variable from the model which provides the largest reduction of a predefined model assessment criteria. The algorithm stops, if no further reduction is possible.

In step-wise selection, forward and backward selection are combined to enable the inclusion and deletion of a variable in every operation. The algorithm stops, if no further reduction is possible.

2.2.3 Regularization

Model selection can also be achieved using regularization techniques by directly influencing the parameters β , which need to be estimated given a data set. In general, regularization restricts the parameter space by adding some penalty term depending on the complexity of the model to the least squares objective function according to (2.10). This leads to the penalized least squares (PLS) criterion

$$\text{PLS}(\mathbf{y}, \beta; \lambda) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \cdot \text{pen}(\beta), \quad (2.42)$$

where λ is the so-called smoothing parameter and $\text{pen}(\beta)$ is the penalty term describing the regularization technique.

In Ridge regression, the penalty term in the penalized least squares criterion in (2.42) is given by the squared weighted L_2 -norm of the parameter vector β , i.e. $\text{pen}(\beta) = \|\beta\|_{\mathbf{K}}^2 = \beta^T \mathbf{K} \beta$ with penalty matrix $\mathbf{K} \in \mathbb{R}^{p \times p}$. The closed form solution reads as

$$\hat{\beta}_{PLS} = \arg \min_{\beta} (\text{PLS}(\mathbf{y}, \beta; \lambda)) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{K})^{-1} \mathbf{X}^T \mathbf{y}. \quad (2.43)$$

The additional penalty term in Ridge regression leads to smaller parameter estimates $\hat{\beta}_{PLS}$ compared to the unpenalized estimate $\hat{\beta}_{LS}$. For large values of the smoothing parameter λ , the parameter estimates will converge towards, but never reach, zero.

Ridge regression is commonly used when the input dimension q is high, i.e. the number of parameters β_i is large, and also known as Tikhonov regularization [23]. Note that it is also possible to use a nonlinear penalty matrix $\mathbf{K}(\beta)$ resulting in

$$\text{pen}(\beta) = \|\beta\|_{\mathbf{K}(\beta)}^2 = \beta^T \mathbf{K}(\beta) \beta. \quad (2.44)$$

However, the resulting penalized least squares problem has no closed form solution and must be solved by an iterative approach. We start with an initial guess $\beta^{[0]}$ and iterate for $k = 0, 1, 2, \dots$ the iteration

$$\beta^{[k+1]} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{K}(\beta^{[k]}))^{-1} \mathbf{X}^T \mathbf{y}, \quad (2.45)$$

until $|\beta^{[k+1]} - \beta^{[k]}| \leq \text{Tol}$ with Tol being a given tolerance.

Note that by the introduction of the penalty matrix \mathbf{K} , we reduce the degrees of freedom of the model. This can be seen by comparing the trace of the hat matrix of the regularized model, i.e.

$$\mathbf{H}_{pen} = \mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{K})^{-1} \mathbf{X}^T, \quad (2.46)$$

with the trace of the hat matrix of the unpenalized model, see (2.22).

2.3 Splines

A spline is a piecewise polynomial defined on a sequence of knots. This definition is quite general. Therefore, a large variety of splines exists, ranging from regression splines [24], over B-splines [8] to natural cubic splines and many more. We will focus on the definition of B-splines in Section 2.3.1, tensor-product B-splines as the multi-dimensional expansion of B-splines in Section 2.3.1, and P-splines in Section 2.3.2, see for more information [3] and [25].

2.3.1 B-Splines

We put the focus on the definition and use of B-splines $s(x)$, which are constructed using the d B-spline basis functions $B_j^l(x)$ of order l as

$$s(x) = \sum_{j=1}^d B_j^l(x) \beta_j \quad (2.47)$$

given the knot sequence

$$K = \{\kappa_{1-l}, \kappa_{1-l+1}, \dots, \kappa_{d+1}\}. \quad (2.48)$$

The B-spline basis function $B_j^l(x)$ of order l is defined by means of the Cox-de Boor recursion formula as

$$B_j^0(x) = \begin{cases} 1 & \text{for } \kappa_j \leq x < \kappa_{j+1} \\ 0 & \text{otherwise} \end{cases} \quad (2.49)$$

$$B_j^l(x) = \frac{x - \kappa_{j-l}}{\kappa_j - \kappa_{j-l}} B_{j-1}^{l-1}(x) + \frac{\kappa_{j+1} - x}{\kappa_{j+1} - \kappa_{j+1-l}} B_j^{l-1}(x) \quad (2.50)$$

using the knot sequence (2.48). Hence it is composed of $(l + 1)$ -polynomial pieces of degree l , see [3]. An example of a B-spline basis function of order $l = 0, 1, 2, 3$ is given in Figure 2.2.

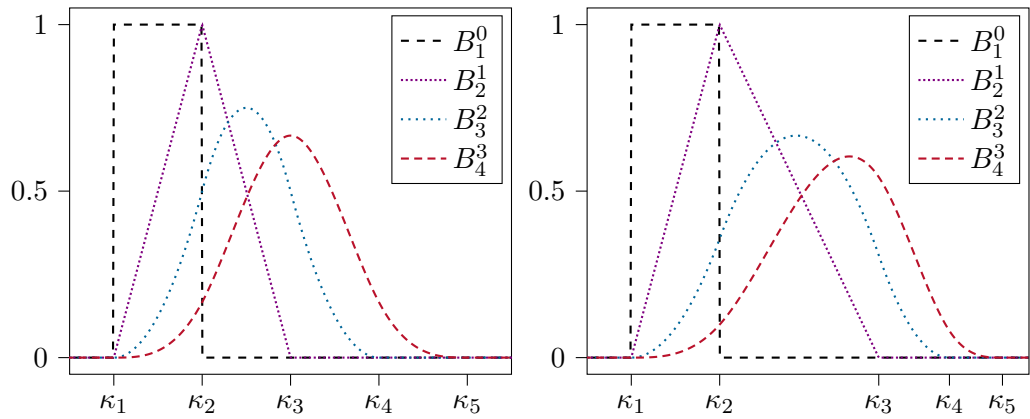


Figure 2.2: B-spline basis function of order $l = 0, 1, 2, 3$ for equidistant (left) and non-equidistant (right) knots.

The left plot shows the B-spline basis functions based on an equidistant sequence of knots. The B-spline basis function B_1^0 is the zero function, except for $x \in [\kappa_1, \kappa_2]$ where it

is equal to 1, see (2.49). The B-spline basis function B_2^1 is the well known *hat function*, being zero except for $x \in [\kappa_1, \kappa_3]$. It consists of two linear pieces, one defined from κ_1 to κ_2 , the other from κ_2 to κ_3 . This can be seen by expanding the recursive definition (2.50) as

$$B_2^1(x) = \frac{x - \kappa_1}{\kappa_2 - \kappa_1} B_1^0(x) + \frac{\kappa_3 - x}{\kappa_3 - \kappa_2} B_2^0(x). \quad (2.51)$$

Everywhere else, B_2^1 is equal to zero. At the joining points, the values of the linear pieces are equal. The B-spline basis function B_3^2 consists of three quadratic pieces, joining at the knots κ_2 and κ_3 . Expanding the recursive definition, using equidistant knots with the knot spacing h , shows this as

$$\begin{aligned} B_3^2(x) &= \frac{x - \kappa_1}{\kappa_3 - \kappa_1} B_2^1(x) + \frac{\kappa_4 - x}{\kappa_4 - \kappa_2} B_3^1(x) \\ &= \dots \\ &= \frac{1}{2h^2} \left[(x - \kappa_1)^2 B_1^0(x) \right. \\ &\quad + [(x - \kappa_1)(\kappa_3 - x) + (\kappa_4 - x)(x - \kappa_2)] B_2^0(x) \\ &\quad \left. + (\kappa_4 - x)^2 B_3^0(x) \right]. \end{aligned} \quad (2.52)$$

At κ_2 and κ_3 , the values of the quadratic pieces, as well as their first derivatives are equal. Finally, the B-spline basis function B_4^3 consists of 4 cubic pieces with the joining points at κ_2 , κ_3 and κ_4 at which respective cubic polynomials possess equal values as well as equal first and second-order derivatives.

The right plot in Figure 2.2 shows the B-spline basis functions of the same order $l = 0, 1, 2, 3$ defined on a non-equidistant knot sequence. The basic properties of these are the same as for B-spline basis function based on equidistant knots. The shown locality, i.e. being nonzero only over a sequence of $l + 2$ knots, is a very attractive feature leading to an enhanced numerical stability compared to other types of splines. Some general properties of a B-spline basis function of order l are summarized in the following list:

- It consists of $l + 1$ polynomial pieces of degree l , e.g. a cubic B-spline basis function ($l = 3$) consists of 4 cubic pieces.
- The pieces join at l inner knots.
- At these knots, the derivatives up to order $l - 1$ are continuous.
- The B-spline basis function is positive on the domain spanned by $l + 2$ knots, everywhere else it is zero, e.g. for $l = 2$, a sequence of 4 knots is necessary.
- At every given x , only $l + 1$ B-spline basis functions are nonzero.

Using the definition of B-spline basis functions, see (2.49) and (2.50), the first-order derivative of a B-spline basis function of order l can be given as

$$\frac{\partial}{\partial x} B_j^l(x) = l \left[\frac{1}{\kappa_j - \kappa_{j-l}} B_{j-1}^{l-1}(x) - \frac{1}{\kappa_{j+1} - \kappa_{j+1-l}} B_j^{l-1}(x) \right] \quad (2.53)$$

using B-spline basis functions of order $l - 1$. Higher order derivatives are obtained by using lower order B-spline basis functions, see [8].

As shown in Figure 2.2, the knots can either be an equidistant sequence, which facilitates the construction and estimation of the coefficients, or a non-equidistant sequence. For equidistant knots, we split the domain $[a, b]$ into $m - 1$ intervals, i.e.

$$h = \frac{b - a}{m - 1}, \quad (2.54)$$

where m is given by $m = d - l + 1$ and obtain the sequence

$$\kappa_j = a + h(j - 1), \quad j = 1, \dots, m. \quad (2.55)$$

Non-equidistant knot placement can be obtained using quantile-based knots, i.e. by using the $(j - 1)/(m - 1)$ -quantiles for $j = 1, \dots, m$ of the observed inputs $x^{(1)}, \dots, x^{(n)}$ as knots. Using this approach, more knots are placed in the areas where lots of data is present. The boundary knots, i.e. $\{\kappa_{1-l}, \dots, \kappa_0\}$ on the left side and $\{\kappa_{d-l+2}, \dots, \kappa_{d+1}\}$ on the right side, are usually set to be apart from each other by at least the minimal knot distance [3].

The collection of d B-spline basis functions of order l over a sequence $K = \{\kappa_{1-l}, \kappa_{1-l+1}, \dots, \kappa_{d+1}\}$ knots is called B-spline basis. The basis is created such that it covers the domain $[a, b]$, i.e.

$$\sum_{j=1}^d B_j^l(x) = 1 \text{ for } x \in [a, b]. \quad (2.56)$$

A function $f(x)$ can then be represented by

$$f(x) = \sum_{j=1}^d B_j^l(x) \beta_j = \mathbf{b}^T \boldsymbol{\beta}, \quad (2.57)$$

using the B-spline basis functions $B_j^l(x)$ of appropriate order l and the parameter vector $\boldsymbol{\beta}^T = [\beta_1, \dots, \beta_d] \in \mathbb{R}^{1 \times d}$. The basis functions can be given in vector notation as $\mathbf{b}^T = [B_1^l(x), \dots, B_d^l(x)] \in \mathbb{R}^{1 \times d}$. Using the set of data $D = \{(x^{(i)}, y^{(i)}), i = 1, 2, \dots, n\}$, the B-spline basis matrix for d splines of order l is given by the matrix \mathbf{B} as

$$\mathbf{B} = \begin{bmatrix} B_1^l(x^{(1)}) & \dots & B_d^l(x^{(1)}) \\ \vdots & & \vdots \\ B_1^l(x^{(n)}) & \dots & B_d^l(x^{(n)}) \end{bmatrix} \in \mathbb{R}^{n \times d}. \quad (2.58)$$

The n equations (2.57) can then be arranged as a linear model in the form

$$\mathbf{y} = \mathbf{B}\boldsymbol{\beta} + \boldsymbol{\epsilon}. \quad (2.59)$$

Once the basis matrix in (2.58) is given, the parameters $\boldsymbol{\beta}$ can be estimated using the Least Squares algorithm given in Section 2.1.1 by optimizing the objective function

$$\text{LS}(\mathbf{y}) = \|\mathbf{y} - \mathbf{B}\boldsymbol{\beta}\|_2^2. \quad (2.60)$$

Therefore, the estimation is computationally efficient and easy to implement since closed-form solutions exists. Further, the advanced theoretical framework of linear models can be applied to use model selection and regularization approaches as well as to calculate e.g. confidence intervals for the regression coefficients and the prediction.

The derivative of the function $f(x)$ in (2.57) can be calculated by summing over all d basis functions and including the estimated parameters $\boldsymbol{\beta}$ into the B-spline basis function derivative (2.53) as

$$\frac{\partial f(x)}{\partial x} = \frac{\partial}{\partial x} \sum_{j=1}^d B_j^l(x) \beta_j = l \sum_{j=2}^d \frac{\Delta \beta_j}{\kappa_j - \kappa_{j-l}} B_{j-1}^{l-1}(x). \quad (2.61)$$

The second-order derivative of $f(x)$ can be given as

$$\frac{\partial^2 f(x)}{\partial x^2} = \frac{\partial^2}{\partial x^2} \sum_{j=1}^d B_j^l(x) = l^2 \sum_{j=3}^d \frac{\Delta^2 \beta_j}{(\kappa_j - \kappa_{j-l})(\kappa_{j+1} - \kappa_{j+1-l})} B_{j-2}^{l-2}(x). \quad (2.62)$$

Here, the finite difference operators $\Delta \beta_j = \beta_j - \beta_{j-1}$ and $\Delta_j^\beta = \Delta(\Delta \beta_j) = \beta_j - 2\beta_{j-1} + \beta_{j-2}$ are used. Therefore, by estimating the B-spline parameters $\boldsymbol{\beta}$, we also generate an estimate for the derivatives of the function $f(x)$.

B-splines of appropriate order $l > 2$ produce smooths curves, i.e. first- and second-order derivatives are continuous, where the smoothness is mostly determined by the number of splines used. By using a low number d , the curve will be quite smooth, but possess a large data error. When using a high number of splines d , the data error will be small but the variance of the curve will be large. This is an example of the bias-variance trade-off, a classical problem of regression and supervised learning, see Section 2.2.1 and [8].

Tensor-Product B-Splines

Tensor-product B-splines can be seen as the multi-dimensional extension of B-splines. We examine an example for two input dimensions x_1 and x_2 . Note that tensor-product B-splines can be constructed for arbitrary dimensions using the approach given below. Assume that we have B-spline bases B_1 and B_2 available using d_1 and d_2 basis functions for the respective dimension and order $l = 3$. For readability, we may omit the order

l in the further description. The tensor-product B-spline basis is then constructed by considering all pairwise products of the two B-spline basis functions, i.e.

$$B_{j,r}(x_1, x_2) = B_j^l(x_1)B_r^l(x_2) \quad (2.63)$$

for $j = 1, 2, \dots, d_1$ and $r = 1, 2, \dots, d_2$. We then obtain the basis function representation for the tensor-product B-spline $t(x_1, x_2)$ as

$$t(x_1, x_2) = \sum_{j=1}^{d_1} \sum_{r=1}^{d_2} B_{j,r}(x_1, x_2) \beta_{j,r}. \quad (2.64)$$

We can therefore combine the individual basis functions to generated new basis functions for the tensor-product B-spline. For any set of data

$$D = \{(x_1^{(i)}, x_2^{(i)}, y^{(i)}), i = 1, 2, \dots, n\}, \quad (2.65)$$

the relationship between the basis matrix \mathbf{X} of the tensor-product B-spline and the basis matrices \mathbf{B}_1 and \mathbf{B}_2 is given as

$$\mathbf{X} = \mathbf{B}_2 \odot \mathbf{B}_1, \quad (2.66)$$

where \odot indicates the use of the row-wise Kronecker product, see Appendix A, $\mathbf{X} \in \mathbb{R}^{n \times d_1 d_2}$ denotes the tensor-product B-spline basis matrix, $\mathbf{B}_1 \in \mathbb{R}^{n \times d_1}$ denotes the B-spline basis matrix for dimension x_1 and $\mathbf{B}_2 \in \mathbb{R}^{n \times d_2}$ denotes the B-spline basis matrix for dimension x_2 [3].

We can then model a two dimensional function using the data D similar to (2.59) as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\gamma} + \boldsymbol{\epsilon}, \quad (2.67)$$

with the tensor-product B-spline basis matrix $\mathbf{X} \in \mathbb{R}^{n \times d_1 d_2}$ and the parameter vector $\boldsymbol{\gamma}^T = [\gamma_1, \dots, \gamma_{d_1 d_2}] \in \mathbb{R}^{1 \times d_1 d_2}$.

This approach can in theory be repeated for as many input dimensions as required. In practice, modeling more than two input dimensions using tensor-product B-splines becomes infeasible because of the exponential increase of basis functions and therefore parameters to estimate.

Additive Regression

To circumvent the latter problem, we now assume the restrictive structure of additive models, see [3], given by

$$f_{add} = f(x_1, \dots, x_q) = f_1(x_1) + \dots + f_q(x_q). \quad (2.68)$$

Hence, we use one function $f_i(x_i)$ per input dimension. For some given data $D = \{(x_1^{(i)}, \dots, x_q^{(i)}, y^{(i)}), i = 1, 2, \dots, n\}$, by using a B-spline $s_i(x_i)$ for each function $f_i(x_i)$ we obtain a linear model

$$f_i(\mathbf{x}_i) = \mathbf{X}_{s_i} \boldsymbol{\beta}_{s_i}, \quad (2.69)$$

where $\mathbf{X}_{s_i} \in \mathbb{R}^{n \times d_i}$ is the B-spline basis matrix using d_i B-spline basis functions for $i = 1, 2, \dots, q$, $\mathbf{x}_i^T = [x_i^{(1)}, \dots, x_i^{(n)}] \in \mathbb{R}^{1 \times n}$ is the data vector of input dimension i and $\boldsymbol{\beta}_{s_i} \in \mathbb{R}^{d_i \times 1}$ are the parameters to estimate. This leads to the model structure

$$\mathbf{y} = \mathbf{X}_{s_1} \boldsymbol{\beta}_{s_1} + \dots + \mathbf{X}_{s_q} \boldsymbol{\beta}_{s_q} + \boldsymbol{\epsilon}, \quad (2.70)$$

which can be written as linear model by concatenation of the B-spline basis matrices and parameter vectors as

$$\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\epsilon} = \begin{bmatrix} \mathbf{X}_{s_1} & \dots & \mathbf{X}_{s_q} \end{bmatrix} \begin{bmatrix} \boldsymbol{\beta}_{s_1} \\ \vdots \\ \boldsymbol{\beta}_{s_q} \end{bmatrix} + \boldsymbol{\epsilon}, \quad (2.71)$$

with the matrix $\mathbf{X} \in \mathbb{R}^{n \times \sum_{i=1}^q d_i}$ and parameter vector $\boldsymbol{\beta} \in \mathbb{R}^{\sum_{i=1}^q d_i \times 1}$. The model (2.71) does not contain interaction terms between variables. Nevertheless, these can be easily introduced for two dimensions using tensor-product B-splines without an overflowing increase in the number of coefficients. We can then write the additive model with interaction terms as

$$f(x_1, \dots, x_q) = f_{add} + f_{1,2}(x_1, x_2) + \dots + f_{q-1,q}(x_{q-1}, x_q). \quad (2.72)$$

Hence, we use one function $f_i(x_i)$ per input dimension and per interaction term. Using one tensor-product B-spline $t_{j,r}(x_j, x_r)$ for each interaction term, we obtain the model

$$\mathbf{y} = \mathbf{X}_{s_1} \boldsymbol{\beta}_{s_1} + \dots + \mathbf{X}_{s_q} \boldsymbol{\beta}_{s_q} + \sum_{j=1}^{q-1} \sum_{r>j}^q \mathbf{X}_{t_{j,r}} \boldsymbol{\beta}_{t_{j,r}} + \boldsymbol{\epsilon}, \quad (2.73)$$

using the tensor-product B-spline basis matrices $\mathbf{X}_{t_{j,r}} \in \mathbb{R}^{n \times d_j d_r}$ and the parameter $\boldsymbol{\beta}_{t_{j,r}} \in \mathbb{R}^{d_j d_r \times 1}$. Using the notation in (2.73), the theoretical framework of linear models can be applied to the additive regression model, since (2.73) can be formulated as linear model yielding

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} = \begin{bmatrix} \mathbf{X}_{s_1} & \dots & \mathbf{X}_{s_q} & \mathbf{X}_{t_{1,2}} & \dots & \mathbf{X}_{t_{q-1,q}} \end{bmatrix} \begin{bmatrix} \beta_{s_1} \\ \vdots \\ \beta_{s_q} \\ \beta_{t_{1,2}} \\ \vdots \\ \beta_{t_{q-1,q}} \end{bmatrix} + \boldsymbol{\epsilon}, \quad (2.74)$$

with $\mathbf{X} \in \mathbb{R}^{n \times d_{total}}$ as design matrix, $\boldsymbol{\beta} \in \mathbb{R}^{d_{total} \times 1}$ as parameter vector and $d_{total} = \sum_{i=1}^q d_i + \sum_{j=1}^{q-1} \sum_{r>j}^q d_j d_r$ as total number of parameters in the model. Therefore, the parameters can be calculated efficiently using the Least Squares (LS) algorithm, see Section 2.1.1. Further, the assumptions given in Section 2.1 on the error term, as well as on the model function are used.

2.3.2 P-Splines

P-splines combine the concepts of B-spline basis functions and regularization to produce smooth function estimations. A function is said to be smooth if its second-order derivative is continuous and does not vary much. Therefore, a penalty of the form

$$\lambda \int (f''(x))^2 dx \quad (2.75)$$

is typically introduced to penalize the curvature of a function which is measured by its second-order derivative [26]. The penalty is weighted by the so-called *smoothing parameter* λ . This yields the penalized least squares objective function, cf. Section 2.2.3, as

$$\text{PLS}(\mathbf{y}, \boldsymbol{\beta}; \lambda) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \int (f''(x))^2 dx, \quad (2.76)$$

using the B-spline basis matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ for some data $D = \{(x^{(i)}, y^{(i)}), i = 1, 2, \dots, n\}$. Inserting the B-spline basis function formulation (2.57), using order l and d basis functions, into (2.75) results in

$$\begin{aligned} \int (f''(x))^2 dx &= \int \left(\sum_{j=1}^d B_j''(x) \beta_j \right)^2 dx \\ &= \int \sum_{j=1}^d \sum_{r=1}^d \beta_j \beta_r B_j''(x) B_r''(x) dx \\ &= \boldsymbol{\beta}^T \mathbf{K} \boldsymbol{\beta}, \end{aligned} \quad (2.77)$$

with the penalty matrix $\mathbf{K}[j, r] = \int B_j''(x) B_r''(x) dx$ as a matrix of dimension $\mathbb{R}^{d \times d}$. The entries of \mathbf{K} are given by the integrated products of the second-order derivatives of the

B-spline basis functions $B_j(x)$ and $B_r(x)$. For readability, the order l is omitted. These second-order derivatives can be obtained by using the derivative properties of B-spline basis functions given in (2.62), see [3].

Eilers and Marx proposed to base the penalty on finite-differences of higher order of the parameters of adjacent B-spline basis functions which circumvents the direct calculation of the derivative and the integral and therefore, reduces the complexity from n , the number of data points to evaluate the integral on, to d , the number of parameters [25]. The squared second-order finite difference gives a good discrete approximation of the integral of the squared second-order derivative in (2.75), i.e.

$$\lambda \sum_{j=3}^d (\Delta^2 \beta_j)^2 \approx \lambda \int (f''(x))^2 dx, \quad (2.78)$$

where $\Delta^2 \beta_j$ is defined as

$$\begin{aligned} \Delta^2 \beta_j &= \Delta(\Delta \beta_j) \\ &= \Delta(\beta_j - \beta_{j-1}) \\ &= \beta_j - 2\beta_{j-1} + \beta_{j-2}. \end{aligned} \quad (2.79)$$

In matrix form, (2.79) may be given as

$$\mathbf{D}_2 \boldsymbol{\beta} = \begin{bmatrix} 1 & -2 & 1 & & & \\ & 1 & -2 & 1 & & \\ & & \ddots & \ddots & \ddots & \\ & & & 1 & -2 & 1 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_d \end{bmatrix}, \quad (2.80)$$

with $\mathbf{D}_2 \in \mathbb{R}^{(d-2) \times d}$. Inserting the matrix form of the second-order finite difference operator (2.80) into (2.78) yields

$$\lambda \sum_{j=3}^d (\Delta^2 \beta_j)^2 = \lambda \boldsymbol{\beta}^T \mathbf{D}_2^T \mathbf{D}_2 \boldsymbol{\beta} = \lambda \boldsymbol{\beta}^T \mathbf{K} \boldsymbol{\beta}. \quad (2.81)$$

We then obtain the objective function to minimize as

$$\text{PLS}(\mathbf{y}, \boldsymbol{\beta}; \lambda) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \boldsymbol{\beta}^T \mathbf{K} \boldsymbol{\beta}, \quad (2.82)$$

which is equivalent to the objective function of Ridge regression, cf. Section 2.2.3, for the special choice of \mathbf{K} given by $\mathbf{K} = \mathbf{D}_2^T \mathbf{D}_2 \in \mathbb{R}^{d \times d}$. As in Ridge regression, the smoothness parameter λ plays a critical role. For $\lambda \rightarrow 0$, the P-spline approaches the underlying B-spline since the penalty term in (2.82) goes to 0. For $\lambda \rightarrow \infty$, the P-spline approaches a polynomial model. The order of the polynomial is given by the order of the finite

difference penalty, e.g. for second-order finite difference penalty, we penalize the discrete approximation of the second-order derivative leading to a linear function, because for these, the second-order derivative is equal to zero. Note that higher-order difference penalties are also possible.

The main advantage of P-splines is their easy set up by replacing the integral of the squared second-order derivative of the B-spline basis functions with the squared second-order finite differences of parameters of adjacent B-spline basis functions. This reduces the computational complexity and allows faster training and evaluation. Hence, P-splines are widely used in practice [25].

A similar penalty term for tensor-product B-splines can be constructed using the Kronecker product. Recall the definition of a tensor-product B-spline given in (2.64) and (2.67).

The spatial alignment of the B-spline basis functions and the corresponding parameters of the two-dimensional tensor-product B-spline needs to be incorporated by the definition of the term *adjacent parameters*. An example for these adjacent parameters, also called spatial neighborhood, is taken from [3] and given in Figure 2.3. Here, we choose the parameters left and right, i.e. $\beta_{j,r-1}$ and $\beta_{j,r+1}$, as well as above and below, i.e. $\beta_{j-1,r}$ and $\beta_{j+1,r}$, as spatial neighborhood for $\beta_{j,r}$.

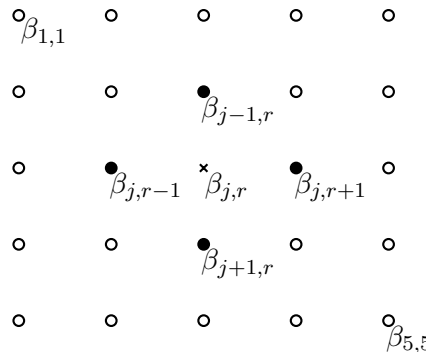


Figure 2.3: Spatial neighborhood or adjacent parameters for a tensor-product B-spline.

We now use the concepts of P-splines in both dimensions and penalize the integral of the squared Hessian of the tensor-product B-spline, see (2.75), by a higher-order finite difference approximation in both dimensions. Using second-order finite differences as in (2.78) leads to the following definition of the penalty term

$$\lambda \left[\sum_{j=3}^{d_1} \sum_{r=1}^{d_2} (\Delta_1^2 \beta_{j,r})^2 + \sum_{j=1}^{d_1} \sum_{r=3}^{d_2} (\Delta_2^2 \beta_{j,r})^2 \right] \approx \iint (f''(x_1, x_2))^2 dx_1 dx_2, \quad (2.83)$$

as discrete approximation of the integral of the squared Hessian of the tensor-product B-spline where the first term on the left side calculates the "row-wise" squared second-order differences using

$$\Delta_1^2 \beta_{j,r} = \beta_{j,r} - 2\beta_{j-1,r} + \beta_{j-2,r} \quad (2.84)$$

and the second term on the left side calculates the "column-wise" squared second-order differences using

$$\Delta_2^2 \beta_{j,r} = \beta_{j,r} - 2\beta_{j,r-1} + \beta_{j,r-2}. \quad (2.85)$$

The subscript for Δ in (2.84) and (2.85) indicates the direction of the finite differences. Using the matrix form of the second-order finite difference operator and the Kronecker product, see Appendix A, as well as the parameter vector $\gamma \in \mathbb{R}^{d_1 d_2 \times 1}$, we can then write the "row-wise" penalty as

$$\gamma^T (\mathbf{I}_{d_2} \otimes \mathbf{D}_{1,2})^T (\mathbf{I}_{d_2} \otimes \mathbf{D}_{1,2}) \gamma = \sum_{j=3}^{d_1} \sum_{r=1}^{d_2} (\Delta_1^2 \beta_{j,r})^2 \quad (2.86)$$

and the "column-wise" penalty as

$$\gamma^T (\mathbf{D}_{2,2} \otimes \mathbf{I}_{d_1})^T (\mathbf{D}_{2,2} \otimes \mathbf{I}_{d_1}) \gamma = \sum_{j=1}^{d_1} \sum_{r=3}^{d_2} (\Delta_2^2 \beta_{j,r})^2 \quad (2.87)$$

using the identity matrices $\mathbf{I}_{d_1} \in \mathbb{R}^{d_1 \times d_1}$ and $\mathbf{I}_{d_2} \in \mathbb{R}^{d_2 \times d_2}$ and the second-order difference matrices $\mathbf{D}_{1,2} \in \mathbb{R}^{(d_1-2) \times d_1}$ and $\mathbf{D}_{2,2} \in \mathbb{R}^{(d_2-2) \times d_2}$. The first subscript of \mathbf{D} indicates the direction of the finite differences and the second subscript indicates the use of the second-order finite differences. Summing up both penalties leads to a formulation similar to (2.82) given by

$$\text{PLS}(\mathbf{y}, \gamma; \lambda) = \|\mathbf{y} - \mathbf{X}\gamma\|_2^2 + \lambda \gamma^T \mathbf{K} \gamma \quad (2.88)$$

with the tensor-product B-spline basis matrix $\mathbf{X} \in \mathbb{R}^{n \times d_1 d_2}$, the smoothing parameter λ and the penalty matrix \mathbf{K} given by

$$\mathbf{K} = [(\mathbf{I}_{d_2} \otimes \mathbf{D}_{1,2})^T (\mathbf{I}_{d_2} \otimes \mathbf{D}_{1,2}) + (\mathbf{D}_{2,2} \otimes \mathbf{I}_{d_1})^T (\mathbf{D}_{2,2} \otimes \mathbf{I}_{d_1})] \in \mathbb{R}^{d_1 d_2 \times d_1 d_2}. \quad (2.89)$$

3 Solution Approach

The main goal of this thesis is the development of an algorithm to include a priori domain knowledge like monotonicity, curvature, unimodality, etc. into the data fitting process. Using this additional information should improve the generalization capabilities of the model in situations of sparse, noisy or even partly wrong data. In this chapter, we use the theory discussed in Chapter 2, i.e. B-splines in Section 2.3.1 and P-splines in Section 2.3.2, and extend it by adding a shape-constraint penalty term of the form

$$\lambda_c \cdot \text{con}(\beta) \quad (3.1)$$

depending on the user-defined a priori domain knowledge leading to the so-called *shape-constraint P-splines* (SCP-splines). The various types of a priori domain knowledge that can be included are listed in Table 3.1.

Constraint	Description		Section
Jamming	$f(x^{(p)}) \approx y^{(p)}$		3.1.7
Boundedness	lower	$f(x) \geq M$	3.1.8
	upper	$f(x) \leq M$	3.1.8
Monotonicity	increasing	$f'(x) \geq 0$	3.1.1
	decreasing	$f'(x) \leq 0$	3.1.2
Curvature	convex	$f''(x) \geq 0$	3.1.3
	concave	$f''(x) \leq 0$	3.1.4
Unimodality	peak	$m = \arg \max_x f(x)$ $f'(x) \geq 0 \quad \text{if } x < m$ $f'(x) \leq 0 \quad \text{if } x > m$	3.1.5
	valley	$m = \arg \min_x f(x)$ $f'(x) \leq 0 \quad \text{if } x < m$ $f'(x) \geq 0 \quad \text{if } x > m$	3.1.6

Table 3.1: Overview of the considered constraints

The focus of this chapter is the definition of shape-constraint P-splines, see Section 3.1, which are characterized by their parameters β given by solving the optimization problem

$$\arg \min_{\beta} \text{PLS}_{\text{SC}}(\mathbf{y}, \beta; \lambda, \lambda_c) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \cdot \text{pen}(\beta) + \lambda_c \cdot \text{con}(\beta), \quad (3.2)$$

where $\text{pen}(\beta)$ is the smoothness penalty term for P-splines, see Section 2.3.2, and $\text{con}(\beta)$ specifies the user-defined shape-constraint penalty term to incorporate a priori domain

knowledge with, see [4] and [27]. Further, we extend the concept of shape-constraint P-splines to two dimensions and discuss shape-constraint tensor-product P-splines, see Section 3.2.

3.1 Shape-constraint P-splines

In Section 2.3.2, we enforced smoothness by penalizing the second-order derivative of the underlying B-spline using finite differences of adjacent parameters β over the whole input space to create the so-called P-splines. We will now utilize the same idea to create the shape-constrained penalty term $\text{con}(\beta)$ in (3.2). We motivate the approach using the example of monotonic increasing behavior. The descriptions for the other constraints listed in Table 3.1 follow afterwards. In the following discussion, we use d B-spline basis functions of order l and the data set $D = \{(x^{(i)}, y^{(i)}), i = 1, 2, \dots, n\}$ resulting in the B-spline basis matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$.

3.1.1 Monotonic increasing constraint

A function is monotonic increasing if its first-order derivative is larger than or equal to zero for the whole input space. We therefore introduce a penalty of the form

$$\lambda_c \int (f(x)')^2 dx \quad \text{if } f'(x) < 0, \quad (3.3)$$

to penalize negative first-order derivatives of the estimated function. The penalty is weighted by the *constraint parameter* λ_c . Making use of the finite difference approximation, see Section 2.3.2, this time for the first-order derivative, leads to a penalty of the form

$$\lambda_c \cdot \text{con}(\beta) = \lambda_c \beta^T \mathbf{K}_c \beta, \quad (3.4)$$

with the shape-constraint penalty matrix $\mathbf{K}_c = \mathbf{D}_1^T \mathbf{V}_c \mathbf{D}_1 \in \mathbb{R}^{d \times d}$. The first-order derivative is approximated using the matrix form of the first-order finite difference operator $\Delta^1 \beta_j = \beta_j - \beta_{j-1}$ given by

$$\mathbf{D}_1 \beta = \begin{bmatrix} -1 & 1 & & \\ & \ddots & \ddots & \\ & & -1 & 1 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_d \end{bmatrix}, \quad (3.5)$$

with the mapping matrix $\mathbf{D}_1 \in \mathbb{R}^{(d-1) \times d}$. The weighting matrix $\mathbf{V}_c \in \mathbb{R}^{(d-1) \times (d-1)}$ handles the if-condition in (3.3). It is a diagonal matrix with the diagonal elements v_j defined as

$$v_{j-1}(\beta) = \begin{cases} 0, & \text{if } \Delta^1 \beta_j \geq 0 \\ 1, & \text{if } \Delta^1 \beta_j < 0 \end{cases} \quad \text{for } j = 2, 3, \dots, d. \quad (3.6)$$

Therefore, the weighting matrix $\mathbf{V}_c := \mathbf{V}_c(\boldsymbol{\beta})$ depends on the parameters $\boldsymbol{\beta}$ and we arrive at a formulation similar to Ridge regression with a parameter dependent, non-linear penalty matrix $\mathbf{K} := \mathbf{K}(\boldsymbol{\beta})$, see Section 2.2.3. The objective function is then of the form

$$\text{PLS}_{\text{SCP}}(\mathbf{y}, \boldsymbol{\beta}; \lambda, \lambda_c) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda\boldsymbol{\beta}^T\mathbf{K}\boldsymbol{\beta} + \lambda_c\boldsymbol{\beta}^T\mathbf{K}_c\boldsymbol{\beta}. \quad (3.7)$$

We use the iterative approach given in Algorithm 1 to estimate the optimal parameters $\hat{\boldsymbol{\beta}}_{SC}$ under the user-defined shape constraint of monotonic increasing behavior.

Algorithm 1: Estimation of the shape-constraint P-spline coefficients.

Result: $\hat{\boldsymbol{\beta}}_{SC}$
 $i \leftarrow 1$;
 $\hat{\boldsymbol{\beta}}_i \leftarrow$ Solution from (3.2) for $\lambda = \lambda_c = 0$;
 $\mathbf{V}_c^0 \leftarrow \mathbf{0}$;
 $\mathbf{V}_c^1 \leftarrow \mathbf{V}_c(\hat{\boldsymbol{\beta}}_i)$;
while $\mathbf{V}_c^i \neq \mathbf{V}_c^{i-1}$ **do**
 $\hat{\boldsymbol{\beta}}_{i+1} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{D}_2^T\mathbf{D}_2 + \lambda_c\mathbf{D}_c^T\mathbf{V}_c^i\mathbf{D}_c)^{-1}\mathbf{X}^T\mathbf{y}$;
 $\mathbf{V}_c^{i+1} \leftarrow \mathbf{V}_c(\hat{\boldsymbol{\beta}}_{i+1})$;
 $i \leftarrow i + 1$;
end
 $\hat{\boldsymbol{\beta}}_{SC} \leftarrow \hat{\boldsymbol{\beta}}_i$;

In Algorithm 1, we use a Newton-Raphson scheme for the estimation of $\hat{\boldsymbol{\beta}}_{i+1}$. For more information, see Appendix B and [27]. The approach described above incorporates the shape-constraint as soft constraint depending on the constraint parameter λ_c with the limits of no constraint for $\lambda_c \rightarrow 0$ and hard constraint for $\lambda_c \rightarrow \infty$. Therefore, λ_c should be set by hand reflecting the user's confidence in the a priori domain knowledge. To incorporate the other constraints listed in Table 3.1, we need to specify the shape-constraint penalty matrix \mathbf{K}_c depending on the respective constraint, i.e. we define the constraint specific mapping matrix \mathbf{D}_c and weighting matrix \mathbf{V}_c .

3.1.2 Monotonic decreasing constraint

Monotonic decreasing behavior can be introduced by penalizing positive first-order derivatives. Therefore, we use the matrix form of the first-order finite difference operator given in (3.5) as mapping matrix $\mathbf{D}_1 \in \mathbb{R}^{(d-1) \times d}$ and define the diagonal elements of the weighting matrix $\mathbf{V} \in \mathbb{R}^{(d-1) \times (d-1)}$ as

$$v_{j-1}(\boldsymbol{\beta}) = \begin{cases} 0, & \text{if } \Delta^1\beta_j \leq 0 \\ 1, & \text{if } \Delta^1\beta_j > 0 \end{cases} \quad \text{for } j = 2, 3, \dots, d. \quad (3.8)$$

3.1.3 Convex constraint

Convex behavior can be introduced by penalizing negative second-order derivatives. Therefore, we use the matrix form of the second-order finite difference operator $\Delta^2\beta_j =$

$\beta_j - 2\beta_{j-1} + \beta_{j-2}$ given by

$$\mathbf{D}_2 \boldsymbol{\beta} = \begin{bmatrix} 1 & -2 & 1 & & \\ & \ddots & \ddots & \ddots & \\ & & 1 & -2 & 1 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_d \end{bmatrix}, \quad (3.9)$$

with the mapping matrix $\mathbf{D}_2 \in \mathbb{R}^{(d-2) \times d}$ and define the diagonal elements of the weighting matrix $\mathbf{V}_c \in \mathbb{R}^{(d-2) \times (d-2)}$ as

$$v_{j-2}(\boldsymbol{\beta}) = \begin{cases} 0, & \text{if } \Delta^2 \beta_j \geq 0 \\ 1, & \text{if } \Delta^2 \beta_j < 0 \end{cases} \quad \text{for } j = 3, 4, \dots, d. \quad (3.10)$$

3.1.4 Concave constraint

Concave behavior can be introduced by penalizing positive second-order derivatives. Therefore, we use the matrix form of the second-order finite difference operator, see (3.9), as mapping matrix $\mathbf{D}_2 \in \mathbb{R}^{(d-2) \times d}$ and define the diagonal elements of the weighting matrix $\mathbf{V}_c \in \mathbb{R}^{(d-2) \times (d-2)}$ as

$$v_{j-2}(\boldsymbol{\beta}) = \begin{cases} 0, & \text{if } \Delta^2 \beta_j \leq 0 \\ 1, & \text{if } \Delta^2 \beta_j > 0 \end{cases} \quad \text{for } j = 3, 4, \dots, d. \quad (3.11)$$

3.1.5 Peak constraint

Peak behavior can be introduced by penalizing negative first-order derivatives for the increasing part and positive first-order derivatives for the decreasing part of the function. Therefore, we use the matrix form of the first-order finite difference operator, see (3.5), as mapping matrix $\mathbf{D}_1 \in \mathbb{R}^{(d-1) \times d}$. The weighting matrix $\mathbf{V}_c \in \mathbb{R}^{(d-1) \times (d-1)}$ now has a special structure. First, we find the data point $x^{(\max)}$ corresponding to the peak value in the data, i.e. $\max = \max_i y^{(i)}$ for $i = 1, 2, \dots, n$. Next, we identify the dominant B-spline basis function $B_p^l(x)$ around $x^{(\max)}$, i.e. the B-spline basis function with the maximal value at $x^{(\max)}$, such that

$$B_p^l(x^{(\max)}) > B_j^l(x^{(\max)}), \quad \text{for } j = 1, 2, \dots, p-1, p+1, \dots, d. \quad (3.12)$$

We now use the index p of the dominant B-spline basis function in the definition of the diagonal elements of the weighting matrix $\mathbf{V}_c \in \mathbb{R}^{(d-1) \times (d-1)}$ as

$$v_{j-1}(\boldsymbol{\beta}) = \begin{cases} 0, & \text{if } \Delta^1 \beta_j \geq 0 \\ 1, & \text{if } \Delta^1 \beta_j < 0 \end{cases}, \quad \text{for } j = 2, 3, \dots, p \quad (3.13)$$

and

$$v_{j-1}(\beta) = \begin{cases} 0, & \text{if } \Delta^1 \beta_j \leq 0 \\ 1, & \text{if } \Delta^1 \beta_j > 0 \end{cases}, \quad \text{for } j = p+1, p+2, \dots, d. \quad (3.14)$$

3.1.6 Valley constraint

Valley behavior can be introduced by the same approach as above by multiplying the data with -1 or by always doing the inverse operation. Therefore, we use the matrix form of the first-order finite difference operator, see (3.5), as mapping matrix $\mathbf{D}_1 \in \mathbb{R}^{(d-1) \times d}$ and define the diagonal elements of the weighting matrix $\mathbf{V}_c \in \mathbb{R}^{(d-1) \times (d-1)}$ as

$$v_{j-1}(\beta) = \begin{cases} 0, & \text{if } \Delta^1 \beta_j \leq 0 \\ 1, & \text{if } \Delta^1 \beta_j > 0 \end{cases}, \quad \text{for } j = 2, 3, \dots, p \quad (3.15)$$

and

$$v_{j-1}(\beta) = \begin{cases} 0, & \text{if } \Delta^1 \beta_j \geq 0 \\ 1, & \text{if } \Delta^1 \beta_j < 0 \end{cases}, \quad \text{for } j = p+1, p+2, \dots, d, \quad (3.16)$$

for p being the index of the B-spline basis function $B_p^l(x)$ with the maximal value at $x^{(\min)}$ with $\min = \min_i y^{(i)}$ for $i = 1, 2, \dots, n$.

3.1.7 Jamming constraint

Jamming the function $f(x)$ by some point $p = \{x^{(p)}, y^{(p)}\}$ means that the estimated function $f(x^{(p)}) \approx y^{(p)}$. This can be incorporated using the B-spline basis matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ as mapping matrix $\mathbf{D}_c \in \mathbb{R}^{n \times d}$ and a weighting matrix $\mathbf{V}_c \in \mathbb{R}^{n \times n}$ with diagonal elements v_j given by

$$v_j(\beta) = \begin{cases} 0, & \text{if } x^{(j)} \neq x^{(p)} \\ 1, & \text{if } x^{(j)} = x^{(p)} \end{cases} \quad \text{for } j = 1, 2, \dots, n. \quad (3.17)$$

3.1.8 Boundedness constraint

The user-defined constraint for boundedness from below by e.g. $M = 0$ uses as mapping matrix $\mathbf{D}_c \in \mathbb{R}^{n \times d}$ the B-spline basis matrix $\mathbf{X} \in \mathbb{R}^{n \times k}$. For the weighting matrix $\mathbf{V}_c \in \mathbb{R}^{n \times n}$, the diagonal weights v_j are defined as

$$v_j(\beta) = \begin{cases} 0, & \text{if } f(x^{(j)}) \geq M \\ 1, & \text{if } f(x^{(j)}) < M \end{cases} \quad \text{for } j = 1, 2, \dots, n. \quad (3.18)$$

Using different values of M allows us to bound from below from any number M . Switching the comparison operators in (3.18) enables us to bound functions from above.

3.2 Shape-constraint tensor-product P-splines

In Section 2.3.1, we extended the univariate approach of B-splines to bivariate tensor-product B-splines by using the row-wise Kronecker product, see Appendix A. We will now extend the tensor-product P-splines, see Section 2.3.2, to incorporate a priori domain knowledge as in Section 3.1 into the fitting process. We start by showing how to include a priori domain knowledge in one dimension, see Section 3.2.1, and describe the various constraints listed in Table 3.2. Afterwards, we show how to include a constraint based on both dimensions, see Section 3.2.5. In the following discussion, we use d_1 and d_2 B-spline basis functions of order l for the dimensions 1 and 2 respectively and the data set $D = \{(x_1^{(i)}, x_2^{(i)}, y^{(i)}), i = 1, 2, \dots, n\}$ resulting in the tensor-product B-spline basis matrix $\mathbf{X} \in \mathbb{R}^{d_1 d_2 \times n}$.

Constraint	Description		Section
Monotonicity	increasing	$\frac{\partial f(x_1, x_2)}{\partial x_1} \geq 0$	3.2.1
	decreasing	$\frac{\partial f(x_1, x_2)}{\partial x_1} \leq 0$	3.2.2
Curvature	convex	$\frac{\partial^2 f(x_1, x_2)}{\partial x_1^2} \geq 0$	3.2.3
	concave	$\frac{\partial^2 f(x_1, x_2)}{\partial x_1^2} \leq 0$	3.2.4
Multiple			3.2.5

Table 3.2: Overview of the constraints for shape-constraint tensor-product P-splines.

Note, that the constraints *Jamming* and *Boundedness* in Table 3.1 are dimension independent and therefore not listed in Table 3.2. Nevertheless, they can also be applied to tensor-product B-splines using the descriptions in Section 3.1.7 and 3.1.8.

3.2.1 Monotonic increasing constraint in dimension 1

A function is monotonic increasing in dimension 1 if its first-order derivative in dimension 1 is larger than or equal to zero for the whole input space. Therefore, we introduce a penalty of the form

$$\lambda_c \iint \left(\frac{\partial f(x_1, x_2)}{\partial x_1} \right)^2 dx_1 dx_2 \quad \text{if } \frac{\partial f(x_1, x_2)}{\partial x_1} < 0, \quad (3.19)$$

to penalize negative first-order partial derivatives of the estimated function. The penalty term is again weighted by the *constraint parameter* λ_c . Approximating the first-order derivative by finite-differences of order 1, similar to (2.86) and (2.87), leads to a penalty of the known form

$$\lambda_c \cdot \text{con}(\gamma) = \lambda_c \gamma^T \mathbf{K}_c \gamma, \quad (3.20)$$

with the shape-constraint penalty matrix $\mathbf{K}_c = \mathbf{D}_c^T \mathbf{V}_c \mathbf{D}_c \in \mathbb{R}^{d_1 d_2 \times d_1 d_2}$. The first-order derivative is approximated by the matrix form of the first-order finite difference operator, see (3.5), and by using the Kronecker product as

$$\mathbf{D}_c \boldsymbol{\gamma} = (\mathbf{I}_{d_2} \otimes \mathbf{D}_1) \boldsymbol{\gamma}, \quad (3.21)$$

with the mapping matrix $\mathbf{D}_c \in \mathbb{R}^{(d_1-1)d_2 \times d_1 d_2}$ using the identity matrix $\mathbf{I}_{d_2} \in \mathbb{R}^{d_2 \times d_2}$ and the finite-difference matrix for the first dimension $\mathbf{D}_1 \in \mathbb{R}^{(d_1-1) \times d_1}$. The diagonal elements v_j of the weighting matrix $\mathbf{V}_c \in \mathbb{R}^{(d_1-1)d_2 \times (d_1-1)d_2}$ are defined as

$$v_{j+(i-1)(d_1-1)}(\boldsymbol{\gamma}) = \begin{cases} 0, & \text{if } \Delta_{d_1}^1 \gamma_{j+(i-1)d_1+1} \geq 0 \\ 1, & \text{if } \Delta_{d_1}^1 \gamma_{j+(i-1)d_1+1} < 0 \end{cases} \quad \text{for } j = 1, 2, \dots, d_1 - 1 \text{ and } i = 1, 2, \dots, d_2, \quad (3.22)$$

using the first-order finite-difference operator for dimension 1 $\Delta_{d_1}^1$ defined as

$$\Delta_{d_1}^1 \gamma_j = \gamma_j - \gamma_{j-1}. \quad (3.23)$$

Hence, we obtain an objective function similar to (3.2) as

$$\arg \min_{\boldsymbol{\gamma}} \text{PLS}_{\text{SC-TP}}(\mathbf{y}, \boldsymbol{\gamma}; \lambda, \lambda_c) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\gamma}\|_2^2 + \lambda \boldsymbol{\gamma}^T \mathbf{K} \boldsymbol{\gamma} + \lambda_c \boldsymbol{\gamma}^T \mathbf{K}_c \boldsymbol{\gamma}, \quad (3.24)$$

with the smoothness penalty matrix $\mathbf{K} \in \mathbb{R}^{d_1 d_2 \times d_1 d_2}$, see Section 2.3.2. The optimal parameters under the shape constraint $\hat{\boldsymbol{\gamma}}_{\text{SC-TP}}$ can be estimated using the iterative scheme given in Algorithm 1.

3.2.2 Monotonic decreasing constraint for dimension 1

Monotonic decreasing behavior in dimension 1 can be introduced by penalizing positive first-order partial derivatives. Therefore, we use the matrix form of the first-order finite difference operator given in (3.5), i.e. $\mathbf{D}_1 \in \mathbb{R}^{(d_1-1) \times d_1}$, in the set up of the mapping matrix $\mathbf{D}_c \in \mathbb{R}^{(d_1-1)d_2 \times d_1 d_2}$, see (3.21), and define the diagonal elements of the weighting matrix $\mathbf{V} \in \mathbb{R}^{(d_1-1)d_2 \times (d_1-1)d_2}$ as

$$v_{j+(i-1)(d_1-1)}(\boldsymbol{\gamma}) = \begin{cases} 0, & \text{if } \Delta_{d_1}^1 \gamma_{j+(i-1)d_1+1} \leq 0 \\ 1, & \text{if } \Delta_{d_1}^1 \gamma_{j+(i-1)d_1+1} > 0 \end{cases} \quad \text{for } j = 1, 2, \dots, d_1 - 1 \text{ and } i = 1, 2, \dots, d_2, \quad (3.25)$$

using the first-order finite difference operator for dimension 1 in (3.23).

3.2.3 Convex constraint for dimension 1

Convex behavior in dimension 1 can be introduced by penalizing negative second-order partial derivatives. Therefore, we use the matrix form of the second-order finite difference operator in (3.9), i.e. $\mathbf{D}_2 \in \mathbb{R}^{(d_1-2) \times d_1}$, in the set up of the mapping matrix $\mathbf{D}_c \in$

$\mathbb{R}^{(d_1-2)d_2 \times d_1 d_2}$, see (3.21), and define the diagonal elements of the weighting matrix $\mathbf{V}_c \in \mathbb{R}^{(d_1-2)d_2 \times (d_1-2)d_2}$ as

$$v_{j+(i-1)(d_1-2)}(\gamma) = \begin{cases} 0, & \text{if } \Delta_{d_1}^2 \gamma_{j+(i-1)d_1+2} \geq 0 \\ 1, & \text{if } \Delta_{d_1}^2 \gamma_{j+(i-1)d_1+2} < 0 \end{cases} \quad \text{for } j = 1, 2, \dots, d_1 - 2 \text{ and } i = 1, 2, \dots, d_2, \quad (3.26)$$

using the second-order finite difference operator for dimension 1 defined as

$$\Delta_{d_1}^2 \gamma_j = \gamma_j - 2\gamma_{j-1} + \gamma_{j-2}. \quad (3.27)$$

3.2.4 Concave constraint for dimension 1

Concave behavior in dimension 1 can be introduced by penalizing positive second-order partial derivatives. Therefore, we use the matrix form of the second-order finite difference operator in (3.9), i.e. $\mathbf{D}_2 \in \mathbb{R}^{(d_1-2) \times d_1}$, in the set up of the mapping matrix $\mathbf{D}_c \in \mathbb{R}^{(d_1-2)d_2 \times d_1 d_2}$, see (3.21), and define the diagonal elements of the weighting matrix $\mathbf{V}_c \in \mathbb{R}^{(d_1-2)d_2 \times (d_1-2)d_2}$ as

$$v_{j+(i-1)(d_1-2)}(\gamma) = \begin{cases} 0, & \text{if } \Delta_{d_1}^2 \gamma_{j+(i-1)d_1+2} \leq 0 \\ 1, & \text{if } \Delta_{d_1}^2 \gamma_{j+(i-1)d_1+2} > 0 \end{cases} \quad \text{for } j = 1, 2, \dots, d_1 - 2 \text{ and } i = 1, 2, \dots, d_2, \quad (3.28)$$

using the second-order finite difference operator for dimension 1 defined in (3.27).

3.2.5 Shape constraints in two dimensions

To enforce shape constraints in two dimensions, e.g monotonic increasing behavior in dimension 1 (c_1) and monotonic decreasing behavior in dimension 2 (c_2), we use the same idea as given in 3.2.1 for both dimensions. Hence, we obtain an objective function similar to (3.2) as

$$\arg \min_{\gamma} \text{PLS}_{\text{SC-TP}}(\mathbf{y}, \gamma; \lambda, \lambda_{c_1}, \lambda_{c_2}) = \|\mathbf{y} - \mathbf{X}\gamma\|_2^2 + \lambda \gamma^T \mathbf{K} \gamma + \lambda_{c_1} \gamma^T \mathbf{K}_{c_1} \gamma + \lambda_{c_2} \gamma^T \mathbf{K}_{c_2} \gamma, \quad (3.29)$$

with the smoothness penalty matrix $\mathbf{K} \in \mathbb{R}^{d_1 d_2 \times d_1 d_2}$, see Section 2.3.2, and two *constraint parameter* λ_{c_1} and λ_{c_2} reflecting the users trust in the a priori domain knowledge. The optimal parameters under the shape constraint $\hat{\gamma}_{\text{SC-TP}}$ can then be estimated again using the iterative scheme given in Algorithm 1.

A Appendix A

A.1 Kronecker product

The Kronecker product $\mathbf{A} \otimes \mathbf{B}$ of the $n \times p$ -matrix \mathbf{A} and the $r \times q$ -matrix \mathbf{B} is defined as $nr \times pq$ -matrix

$$\mathbf{C} = \mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & \dots & a_{1p}\mathbf{B} \\ \vdots & & \vdots \\ a_{n1}\mathbf{B} & \dots & a_{np}\mathbf{B} \end{bmatrix}. \quad (\text{A.1})$$

As example, we define $\mathbf{A} \in \mathbb{R}^{3 \times 2}$ and $\mathbf{B} \in \mathbb{R}^{2 \times 2}$ as

$$\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} 1 & 2 \\ 1 & 2 \end{bmatrix}, \quad (\text{A.2})$$

resulting in the matrix $\mathbf{C} \in \mathbb{R}^{6 \times 4}$ given by

$$\mathbf{C} = \begin{bmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} \end{bmatrix} = \begin{bmatrix} 1 & 2 & 2 & 4 \\ 1 & 2 & 2 & 4 \\ 3 & 6 & 4 & 8 \\ 3 & 6 & 4 & 8 \\ 5 & 10 & 6 & 12 \\ 5 & 10 & 6 & 12 \end{bmatrix}. \quad (\text{A.3})$$

A.2 Row-wise Kronecker product

Given the $n \times p$ -matrix \mathbf{A} and the $n \times q$ -matrix \mathbf{B} , the row-wise Kronecker product is defined as $n \times pq$ -matrix \mathbf{C} given by

$$\mathbf{C} = \mathbf{A} \odot \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B}_1 & \dots & a_{1p}\mathbf{B}_1 \\ \vdots & & \vdots \\ a_{n1}\mathbf{B}_n & \dots & a_{np}\mathbf{B}_n \end{bmatrix}. \quad (\text{A.4})$$

As example, we define $\mathbf{A} \in \mathbb{R}^{3 \times 3}$ and $\mathbf{B} \in \mathbb{R}^{3 \times 2}$ as

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \\ \mathbf{A}_3 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \\ \mathbf{B}_3 \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 1 & 2 \\ 1 & 2 \end{bmatrix}, \quad (\text{A.5})$$

resulting in the matrix $\mathbf{C} \in \mathbb{R}^{3 \times 6}$ given by

$$\mathbf{C} = \begin{bmatrix} \mathbf{A}_1 \otimes \mathbf{B}_1 \\ \mathbf{A}_2 \otimes \mathbf{B}_2 \\ \mathbf{A}_3 \otimes \mathbf{B}_3 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 2 & 4 & 3 & 6 \\ 4 & 8 & 5 & 10 & 6 & 12 \\ 7 & 14 & 8 & 16 & 9 & 19 \end{bmatrix}. \quad (\text{A.6})$$

The row-wise Kronecker product is also known as *face-splitting product* or as *transposed Khatri-Rao product* [28].

B Appendix B

The following proof is take from [27]. To find an optimal solution for the penalized least squares objective function for shape-constraint P-splines given by

$$L(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda\boldsymbol{\beta}^T\mathbf{K}\boldsymbol{\beta} + \lambda_c\boldsymbol{\beta}^T\mathbf{K}_c\boldsymbol{\beta}, \quad (\text{B.1})$$

we use a Newton-Raphson scheme. At each iteration i , the new estimate $\hat{\boldsymbol{\beta}}_{i+1}$ is computed such that

$$\mathbf{g}(\boldsymbol{\beta}_i) + \mathbf{H}(\boldsymbol{\beta}_i)(\boldsymbol{\beta}_{i+1} - \boldsymbol{\beta}_i) = 0, \quad (\text{B.2})$$

with \mathbf{g} being the gradient of and \mathbf{H} being the Hessian of $L(\boldsymbol{\beta})$. The gradient of $L(\boldsymbol{\beta})$ equals to

$$\mathbf{g}(\boldsymbol{\beta}) = -2\mathbf{X}^T\mathbf{y} + 2\left[\mathbf{X}^T\mathbf{X} + \lambda\mathbf{D}_2^T\mathbf{D}_2 + \lambda_c\mathbf{D}_c^T\mathbf{V}_c(\boldsymbol{\beta})\mathbf{D}_c\right]\boldsymbol{\beta} + \lambda_c\boldsymbol{\beta}^T\mathbf{D}_c^T\frac{\partial\mathbf{V}_c(\boldsymbol{\beta})}{\partial\boldsymbol{\beta}}\mathbf{D}_c\boldsymbol{\beta}, \quad (\text{B.3})$$

which can be simplified, see [27], to

$$\mathbf{g}(\boldsymbol{\beta}) = -2\mathbf{X}^T\mathbf{y} + 2\left[\mathbf{X}^T\mathbf{X} + \lambda\mathbf{D}_2^T\mathbf{D}_2 + \lambda_c\mathbf{D}_c^T\mathbf{V}_c(\boldsymbol{\beta})\mathbf{D}_c\right]\boldsymbol{\beta}. \quad (\text{B.4})$$

The gradient is therefore a piecewise linear function of $\boldsymbol{\beta}$ since $\mathbf{V}(\boldsymbol{\beta})$ is a diagonal matrix with 0s and 1s as diagonal elements, see (B.4). This implies that the Hessian is a step function of $\boldsymbol{\beta}$ and equal to

$$\mathbf{H}(\boldsymbol{\beta}) = 2\mathbf{X}^T\mathbf{X} + 2\lambda\mathbf{D}_2^T\mathbf{D}_2 + 2\lambda_c\mathbf{D}_c^T\mathbf{V}_c(\boldsymbol{\beta})\mathbf{D}_c + 2\lambda_c\mathbf{D}_c^T\frac{\partial\mathbf{V}_c(\boldsymbol{\beta})}{\partial\boldsymbol{\beta}}\mathbf{D}_c. \quad (\text{B.5})$$

The last part can be neglected again, see [27]. Substituting (B.4) and (B.5) into (B.2) leads to

$$\begin{aligned} 0 = & -2\mathbf{X}^T\mathbf{y} \\ & + 2\left[\mathbf{X}^T\mathbf{X} + \lambda\mathbf{D}_2^T\mathbf{D}_2 + \lambda_c\mathbf{D}_c^T\mathbf{V}_c^i(\boldsymbol{\beta})\mathbf{D}_c\right]\boldsymbol{\beta}_i \\ & + 2\left[\mathbf{X}^T\mathbf{X} + \lambda\mathbf{D}_2^T\mathbf{D}_2 + \lambda_c\mathbf{D}_c^T\mathbf{V}_c^i(\boldsymbol{\beta})\mathbf{D}_c\right](\boldsymbol{\beta}_{i+1} - \boldsymbol{\beta}_i), \end{aligned} \quad (\text{B.6})$$

which can be reformulated as

$$-\mathbf{X}^T \mathbf{y} + \left[\mathbf{X}^T \mathbf{X} + \lambda \mathbf{D}_2^T \mathbf{D}_2 + \lambda_c \mathbf{D}_c^T \mathbf{V}_c^i(\boldsymbol{\beta}) \mathbf{D}_c \right] \boldsymbol{\beta}_{i+1} = 0, \quad (\text{B.7})$$

and, hence, we obtain

$$\boldsymbol{\beta}_{i+1} = \left[\mathbf{X}^T \mathbf{X} + \lambda \mathbf{D}_2^T \mathbf{D}_2 + \lambda_c \mathbf{D}_c^T \mathbf{V}_c^i(\boldsymbol{\beta}) \mathbf{D}_c \right]^{-1} \mathbf{X}^T \mathbf{y}. \quad (\text{B.8})$$

Bibliography

- [1] W. R. Ashby, *An Introduction to Cybernetics*. Chapman & Hall Ltd, 1961.
- [2] F. K. Došilović, M. Brčić, and N. Hlupić, “Explainable artificial intelligence: A survey,” in *International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO)*, IEEE, 2018, pp. 0210–0215.
- [3] L. Fahrmeir, T. Kneib, S. Lang, and B. Marx, *Regression*. Springer, 2007.
- [4] B. Hofner, J. Müller, and T. Hothorn, “Monotonicity-constrained species distribution models,” *Ecology*, vol. 92, no. 10, pp. 1895–1901, 2011.
- [5] O. Nelles, *Nonlinear System Identification*. Springer, 2001.
- [6] J. H. Friedman, “Multivariate adaptive regression splines,” *The Annals of Statistics*, pp. 1–67, 1991.
- [7] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*. Springer, 2001.
- [8] C. De Boor, *A practical Guide to Splines*. Springer, 1978, vol. 27.
- [9] C. Bishop, *Pattern Recognition and Machine Learning*. Springer, 2006.
- [10] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio, *Deep Learning*. MIT Press, 2016, <http://www.deeplearningbook.org>.
- [11] G. Cybenko, “Approximation by superpositions of a sigmoidal function,” *Mathematics of Control, Signals and Systems*, vol. 2, no. 4, pp. 303–314, 1989.
- [12] K. Hornik, “Approximation capabilities of multilayer feedforward networks,” *Neural Networks*, vol. 4, no. 2, pp. 251–257, 1991.
- [13] Y. S. Abu-Mostafa, “Learning from hints in neural networks,” *Journal of Complexity*, vol. 6, no. 2, pp. 192–198, 1990.
- [14] J. Sill and Y. S. Abu-Mostafa, “Monotonicity hints,” in *Advances in Neural Information Processing Systems*, 1997, pp. 634–640.
- [15] E. Garcia and M. Gupta, “Lattice regression,” in *Advances in Neural Information Processing Systems*, 2009, pp. 594–602.
- [16] M. M. Fard, K. Canini, A. Cotter, J. Pfeifer, and M. Gupta, “Fast and flexible monotonic functions with ensembles of lattices,” in *Advances in Neural Information Processing Systems*, 2016, pp. 2919–2927.
- [17] M. Gupta, A. Cotter, J. Pfeifer, K. Voevodski, K. Canini, A. Mangylov, W. Moczydlowski, and A. Van Esbroeck, “Monotonic calibrated interpolated look-up tables,” *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 3790–3836, 2016.

- [18] S. You, D. Ding, K. Canini, J. Pfeifer, and M. Gupta, "Deep lattice networks and partial monotonic functions," in *Advances in Neural Information Processing Systems*, 2017, pp. 2981–2989.
- [19] S. N. Wood, *Generalized Additive Models*. CRC Press, 2017.
- [20] D. G. Luenberger, Y. Ye, *et al.*, *Linear and Nonlinear Programming*. Springer, 2016, vol. 2.
- [21] V. Blobel and E. Lohrmann, *Statistische und numerische Methoden der Datenanalyse*. Springer, 1998.
- [22] G. H. Golub, M. Heath, and G. Wahba, "Generalized cross-validation as a method for choosing a good ridge parameter," *Technometrics*, vol. 21, no. 2, pp. 215–223, 1979.
- [23] A. E. Hoerl and R. W. Kennard, "Ridge regression: Biased estimation for nonorthogonal problems," *Technometrics*, vol. 12, no. 1, pp. 55–67, 1970.
- [24] R. L. Eubank and C. H. Spiegelman, "Testing the goodness of fit of a linear model via nonparametric regression techniques," *Journal of the American Statistical Association*, vol. 85, no. 410, pp. 387–392, 1990.
- [25] P. H. Eilers and B. D. Marx, "Flexible smoothing with *b*-splines and penalties," *Statistical Science*, pp. 89–102, 1996.
- [26] F. O'Sullivan, "A statistical perspective on ill-posed inverse problems," *Statistical Science*, pp. 502–518, 1986.
- [27] K. Bollaerts, P. H. C. Eilers, and I. Van Mechelen, "Simple and multiple *p*-splines regression with shape constraints," *British Journal of Mathematical and Statistical Psychology*, vol. 59, no. 2, pp. 451–469, 2006.
- [28] V. I. Slyusar, "Analytical model of the digital antenna array on a basis of face-splitting matrixs product," in *Proceedings of International Conference on Antenna Theory and Techniques*, 1997, pp. 108–109.

Eidesstattliche Erklärung

Hiermit erkläre ich, dass die vorliegende Arbeit gemäß dem Code of Conduct – Regeln zur Sicherung guter wissenschaftlicher Praxis (in der aktuellen Fassung des jeweiligen Mitteilungsblattes der TU Wien), insbesondere ohne unzulässige Hilfe Dritter und ohne Benutzung anderer als der angegebenen Hilfsmittel, angefertigt wurde. Die aus anderen Quellen direkt oder indirekt übernommenen Daten und Konzepte sind unter Angabe der Quelle gekennzeichnet. Die Arbeit wurde bisher weder im In- noch im Ausland in gleicher oder in ähnlicher Form in anderen Prüfungsverfahren vorgelegt.

Wien, XX. Monat, Jahr

Vorname Nachname