



Swiss Federal Institute of Technology Zurich

Seminar for  
Statistics

1   **Department of Mathematics**

2

3

---

4

5   Master Thesis

Spring 2022

6

---

7

**Lukas Graz**

8

9                   **Interpolation and Correction**

10                  of

11                 **Multispectral Satellite Image Time Series**

11

---

12

Submission Date: September 18th 2022

13

---

14

Co-Adviser: Gregor Perich  
Adviser: Prof. Dr. Nicolai Meinshausen

# 15 Preface

## 16 Supplementary Material

- 17 Instructions and the relevant code needed to reproduce this thesis can be found in the  
18 GitHub repository:  
19 <https://github.com/LGraz/MasterThesis-Code>
- 20 To use our results we recommend the R-package:  
21 <https://github.com/LGraz/CorrectTimeSeries>
- 22 More information is given in the appendix A.

## 23 Acknowledgements

- 24 First, I wish to express my sincere gratitude to my supervisor Prof. Dr. Nicolai Mein-  
25 shausen who took the responsibility for my work and happily took the time to discuss  
26 conceptual and guiding questions and to inspire me with new ideas.
- 27 It is necessary to highlight that without Gregor Perich this project would not have been  
28 possible. His high personal commitment, reliability as well as the weekly instructive su-  
29 pervision meetings were, without question, essential for this work.
- 30 It was a real pleasure for me to be part of the *Crop Science* group for this time. Enjoying  
31 everyday company, a two-day excursion, and harvesting wheat together have made this  
32 time truly remarkable. In particular, I would like to thank Prof. Dr. Achim Walter, who  
33 supported this collaboration at its core.
- 34 Last but not least, I would like to express my gratitude to the *Seminar for Statistics*,  
35 which created the framework conditions for this work and did everything to help me with  
36 conceptional and administrative questions. I should also mention the computing resources  
37 provided by them, without which my computations would not have been feasible.

# 38 Abstract

39 Die Kern-Resultate müssen auch in den Abstract. Ebenso würde ich die vollständige  
Reproduzierbarkeit und die R-Package erwähnen.

- 40 Kurze problemerläuterung (NDVI-ts im Zentrum)
- 41 NDVI Interpolation gewinner
- 42 erforscht Robusification
- 43 NDVI Correction + yield-based evaluation

**44 Contents**

45	<b>Notation</b>	<b>vi</b>
46	<b>1 Introduction</b>	<b>1</b>
47	<b>2 Data and Methods</b>	<b>3</b>
48	2.1 Sentinel 2 Data . . . . .	3
49	2.2 Crop Yield Data . . . . .	3
50	2.3 Normalized Difference Vegetation Index (NDVI) . . . . .	5
51	2.4 Timescale Transformation . . . . .	6
52	2.5 The Concept of a ‘Pixel’ . . . . .	6
53	2.6 Challenges in S2 Data . . . . .	6
54	2.7 General Methods . . . . .	8
55	2.7.1 Root Mean Square Error (RMSE) . . . . .	8
56	2.7.2 Out-Of-Bag (OOB) and Leave-One-Out-Cross-Validation (LOOCV) . . . . .	8
57	<b>3 Interpolation Methods</b>	<b>9</b>
58	3.1 Interpolation Setup . . . . .	9
59	3.2 Parametric Regression . . . . .	9
60	3.2.1 Double Logistic (DL) . . . . .	11
61	3.2.2 Fourier Series (FS) . . . . .	11
62	3.2.3 Optimization Issues . . . . .	12
63	3.3 Non-Parametric Regression . . . . .	12
64	3.3.1 Kernel Regression: Nadaraya-Watson (NW) . . . . .	12
65	3.3.2 Universal Kriging (UK) . . . . .	13
66	3.3.3 Savitzky-Golay Filter (SG) . . . . .	14
67	3.3.4 Locally Weighted Regression (LOESS) . . . . .	16
68	3.3.5 B-Splines (BS) . . . . .	17
69	3.3.6 Smoothing Splines (SS) . . . . .	17
70	3.4 Tuning Parameter Estimation . . . . .	18
71	3.5 Robustification . . . . .	18
72	3.5.1 Our Adjustment: . . . . .	19
73	3.5.2 Examples and Conclusions . . . . .	20
74	3.5.3 Upper Envelope Approach - Penalty for Negative Residuals . . . . .	20
75	3.6 Performance Assessment . . . . .	20
76	<b>4 NDVI Correction</b>	<b>21</b>
77	4.1 Considering other SCL Classes . . . . .	21
78	4.2 Correction Models . . . . .	22
79	4.2.1 Ordinary Least Squares (OLS) . . . . .	22
80	4.2.2 Least Absolute Shrinkage and Selection Operator (LASSO) . . . . .	23
81	4.2.3 General Additive Model (GAM) . . . . .	24
82	4.2.4 Random Forest (RF) . . . . .	24
83	4.2.5 Multivariate Adaptive Regression Splines (MARS) . . . . .	25
84	4.3 Uncertainty Estimation . . . . .	26
85	4.4 Interpolation . . . . .	26
86	4.5 Resulting Interpolation Strategies . . . . .	26
87	4.6 Evaluation Method . . . . .	27

88	4.6.1 Yield Estimation . . . . .	27
89	<b>5 Results</b>	<b>30</b>
90	5.1 Goodness of Fit for Selected Interpolation Methods . . . . .	30
91	5.2 XXX (Robustification and) NDVI-Correction . . . . .	30
92	<b>6 Discussion</b>	<b>32</b>
93	6.1 Interpolation Methods . . . . .	32
94	6.1.1 Data Gaps in Time Series . . . . .	32
95	6.1.2 Preselection . . . . .	33
96	6.1.3 Candidate Selection . . . . .	33
97	6.2 NDVI Correction . . . . .	33
98	6.2.1 Bootstrap . . . . .	33
99	6.2.2 Using Additional Covariates . . . . .	33
100	6.2.3 Which Interpolation Strategy should we choose . . . . .	34
101	6.2.4 High RMSE in Yield Prediction . . . . .	34
102	<b>7 Conclusion</b>	<b>35</b>
103	7.1 Future Work . . . . .	37
104	7.1.1 Time Series Correction-Interpolation as a General Method . . . . .	37
105	7.1.2 Minor Improvements . . . . .	37
106	<b>Bibliography</b>	<b>38</b>
107	<b>A Reproducibility</b>	<b>40</b>
108	A.1 Reproduce Results . . . . .	40
109	A.2 R-Package . . . . .	40
110	<b>B Further Material</b>	<b>42</b>
111	B.1 Data and Methods . . . . .	42
112	B.1.1 GDD . . . . .	42
113	B.2 Interpolation . . . . .	43
114	B.3 NDVI correction . . . . .	44
115	B.3.1 OLS-SCL Model Outputs . . . . .	44

# 116 Todo list

117 Die Kern-Resultate müssen auch in den Abstract. Ebenso würde ich die vollständige Reproduzierbarkeit und die R-Package erwähnen. . . . .	iii
119 verdeutliche dem leser, dass ein auftrag das findne von interpolationmethoden war .	9
120 Paper zitieren wo eingeführt oder wo benutzt (falls einföhrung fast schon trivial) . .	9
121 figure / tabelle / pseudocode anstatt aufzählung . . . . .	15
122 consider naming the sub-plots . . . . .	20
123 defition of RYEA, it is not an accuracy but an error . . . . .	30
124 Here in the discussion, you should take up the points you mentioned in the introduction	32
125 wertend . . . . .	32
126 where does this section belong to? Chapter ‘NDVI Correction’ or ‘Further Work’? .	33
127 table mit OLS SCL als sieger diskutieren . . . . .	34
128 kurzer kontext von vergleichbaren values von gregor — diese sektion ist für dena uftraggebenr . . . . .	34
130 even in a perfect world the NDVI curve only holds a fraction of the information avialbe	34
131 You already capture the ”main” structure of your thesis with the interpolation and	
132 the NDVi correction sections. Can you combine them both in a ”synthesis”	
133 subsection at the end of the discussion? . . . . .	34
134 Frage: mehr details für die begründung der Interpolations-kandidaten? . . . . .	35
135 Anzahl von Beobachtungen, empfehlungen? – schwierig, weil regelmäßigkeit in	
136 ‘wichtigen’ zeiträumen (der veränderung) wichtiger ist. . . . .	36
137 page breaks . . . . .	44
138 replace space before ref by tilda . . . . .	45
139 check quantile definitions . . . . .	45
140 schwarz weiss färbung der IS tabelle korrigieren . . . . .	45
141 so wenig wie möglich abkürzungen in den fig und table captions . . . . .	45
142 refer to data aviability . . . . .	45
143 abkürzungen Fourier und in tabellen . . . . .	45

# 144 Notations

## 145 Variables

$c$	a (vector of) constant(s)
$\lambda \in \mathbb{R}$	a scalar
$n \in \mathbb{N}$	sample size
$i, j$	indices in $\{1, \dots, n\}$
$n \in \mathbb{R}^n$	time, usually in GDD
$w \in \mathbb{R}^n$	a vector of weights for each location $x$
$y \in \mathbb{R}^n$	response in 1-dim interpolation setting
$\hat{y} \in \mathbb{R}^n$	estimate of $y$
$\bar{y} \in \mathbb{R}$	sample mean of $y$
$r \in \mathbb{R}^n$	residuals given by $y - \hat{y}$
$X \in \mathbb{R}^{n \times p}$	the design matrix. Each row corresponds to one observation and each column to one covariate.
$X_{[:,j]}$	the $j$ -th column of $X$
$X_{[i,:]}$	the $i$ -th row of $X$

## 146 Abbreviations and Objects

Pixel	A pixel originates of an image pixel and describes a square of 10 x 10 meters in the field which coincides with the resolution (and location) of the Sentinel-2 pixels. Such pixels are illustrated in figure 2.1b. Additional information like yield is also attached.
$P_t$	the observed data (weather and spectral bands) at time $t$ and the location of one pixel.
$P$	a pixel. We see it as a collection of all the observations at the specified location within one season. More formally, $P := \{P_t   t \text{ is a valid sample time within a defined season}\}$
SCL	Scene Classification Layer provided by the European Space Agency (ESA) that gives an estimation of the land cover class of each pixel. It indicates what one can expect at a pixel at a sampled time. For an overview, c.f. table 2.2

---

$P_{SCL45}$	is similar to $P$ but we only consider observations which belong to the classes 4 and 5. This is used done to get a subset of observations which are less contaminated by clouds and shadows.
NDVI	Normalized Difference Vegetation Index ( <a href="#">Rouse, 1974</a> )
DAS	Days After Sowing
GDD	Growing Degree Days – cumulative sum of “ $\max(0, \text{temperature} - \text{threshold})$ ”
RYEA	Relative Yield-Estimation-Accuracy. Definition <a href="#">4.6.0.1</a>
OOB	Out Of the Box. Describes the procedure of estimating the value for a point by a model that has not seen this point before (c.f. section <a href="#">2.7.2</a> ).
LOOCV	Leave One Out Cross Validation. Describes the procedure of estimating the value for a point by a model that has seen all the points except the current one (c.f. section <a href="#">2.7.2</a> ).

## 147 Statistical Models

DL	Double Logistic (c.f. section <a href="#">3.2.1</a> )
FS	Fourier Series (c.f. section <a href="#">3.2.2</a> )
NW	Nadaraya-Watson (c.f. section <a href="#">3.3.1</a> )
UK	Universal Kriging (c.f. section <a href="#">3.3.2</a> )
SG	Savitzky-Golay Filter (c.f. section <a href="#">3.3.3</a> )
LOESS	Locally Weighted Regression (c.f. section <a href="#">3.3.4</a> )
BS	B-splines (c.f. section <a href="#">3.3.5</a> )
SS	Smoothing Splines (c.f. section <a href="#">3.3.6</a> )
OLS	Ordinary Least Squares (c.f. section <a href="#">4.2.1</a> )
OLS-SCL	OLS using only the observed NDVI and SCL classes (as factor variables)
OLS-all	OLS using the covariates OLS-SCL uses and the spectral bands
LASSO	Least Absolute Shrinkage and Selection Operator (c.f. section <a href="#">4.2.2</a> )
GAM	General Additive Model (c.f. section <a href="#">4.2.3</a> )
RF	Random Forest (c.f. section <a href="#">4.2.4</a> )
MARS	Multivariate Adaptive Regression Splines (c.f. section <a href="#">4.2.5</a> )

148 XXX only equations that are referenced are equipped with a number

149 **Chapter 1**

150 **Introduction**

151 Remote sensing zielt darauf ab, ziel-Größen effizient aus der Entfernung messen zu können.  
152 Hier finden Satellitenbilder Zeitreihen Verwendung, wie etwa die von der europäischen  
153 Raum Agentur (ESA) kostenlos veröffentlichten Bilder Zeitreihen der multi-spektralen  
154 Sentinel 2 (S2) Satelliten. Die Vegetationsentwicklung von Wäldern und landwirtschaftlich  
155 relevanten Flächen im großen stile zu überwachen, ist unter anderem für public angents,  
156 Versicherungen, Umwelt- und Klimaforscher von grossem Interesse. Mögliche Ziele sind  
157 hierbei eine crop Klassifizierung für das subventionieren von bauern oder das Erstellen  
158 von Pflanzenmodellen, um Ernteertrag oder Stickstoff-konzentrationen zu schätzen. Um  
159 die hochdimensionalen Satellitenbilder in leicht interpretierbare größen zu transformieren,  
160 werden spektrale Indizes, wie der Normalized Difference Vegetation Index (NDVI) benutzt.  
161 Dieser ist ein proxy für die Vegetationsdichte und die korrespondierende Zeitreihe spiegelt  
162 somit das pflanzenwachstum wieder. Der Informationsgehalt von einem Satellitenbild  
163 ist jedoch abhängig vom Zustand der Atmosphäre und so trägt der davon abgeleitete  
164 NDVI bei einer dichten Wolkendecke keine informationen über die Vegetation am Boden.  
165 Daher liefert die ESA zusätzlich eine Scene Classification Layer (SCL), welche Aufschluss  
166 gibt was beobachtet wird (z.B Schatten, Wolken, Vegetation, etc.). So können wir bei  
167 der extraktion der NDVI zeitreihe aus der S2 Satellitenbilder Zeitreihe, anhand der SCL  
168 Klassifizierung, die uninformativen Beobachtungen herausfiltern. Durch diese Filtration  
169 kann es jedoch leicht vorkommen, dass wir besonders im Winter über mehrere Wochen  
170 keine Observationen haben. Zudem kommt, dass manche Beobachtungen zu unrecht durch  
171 die SCL als informativ bewertet wird (z.B. als Vegetation) und somit in einem fehlerhaften  
172 NDVI resultiert. Diese beiden probleme versucht man gegenwärtig mit interpolation und  
173 smoothing zu lösen. Starke formannahmen über die NDVI kurve werden in ... getroffen.  
174 Flexiblere ansätze wurden von ... verwendet.

175 In dieser Thesis werden wir stärken und schwächen von solch gängigen Interpolations-  
176 methoden diskutieren und hinsichtlich der NDVI Interpolation bewerten. Dafür benutzen  
177 wir die S2 Satelitenbilder Zeitreihe und Ernteertragskarten von verschiedenen Feldern ver-  
178 schiedenen Weizenarten auf einer Farm in Witzwil in der Schweiz über die Jahre 2017-  
179 2021. Um die Interpolationmethoden zu verbessern, verallgemeinern und testen wir einen  
180 Ansatz, der Interpolationen robuster gegen Ausreisser machen soll. Zudem ermitteln wir,  
181 wie Datenlücken die verschiedenen Interpolationmethoden beeinflussen. Ausserdem stellen  
182 wir am Beispiel des NDVI eine generelle Interpolations-prozedur vor, welche anhand von  
183 zusätzlichen informationen die Zielvariable mit einer Unsicherheitsschätzung korrigiert und  
184 anschließend interpoliert. Somit müssen wir die Observationen nicht mehr a priori via der

185 SCL filtern , sondern korrigieren den beobachteten NDVI und filtern via der geschätzten  
186 Unsicherheiten. Schlussendlich benchmarken wir verschiedene interpolations strategien  
187 mit einem objektiven Qualitätsmaß, welches annimmt, dass je besser eine NDVI TS das  
188 Pflanzenwachstum modelliert, desto geeigneter ist sie, um den Ernteertrag zu schätzen.

189 Die Hauptfragestellungen, welchen wir in dieser Thesis nachgehen wollen lauten also:

- 190 i.) 1 review of interpolation methods
- 191 ii.) 2 erroruous observations — how to deal with them
- 192 iii.) 3 data gaps — influence itpl mehtods
- 193 iv.) 4 data gaps — how to deal with them
- 194 v.) 6 how to compare two NDVI interpolation strategies?

195 Roadmap ... 1 in 3 2 robustification 3.5 3 discussed in 6.1.1 4 utilize observations filterd  
196 before and estimating how reliable each of them are 4 6 4.6

197 “Similarly, smoothing the time series of satellite data is helpful to address inconsistency  
198 in observation frequency and timing due to clouds and other sensor artefacts Skakun,  
199 Vermote, Franch, Roger, Kussul, Ju, and Masek (2019)”

200 **Chapter 2**

201 **Data and Methods**

202 We will start by describing the available data and the challenges associated with it. Our  
203 study region is a farm of over 800ha, which is located in western Switzerland. From  
204 Perich, Turkoglu, Graf, Wegner, Aasen, Walter, and Liebisch (2022) we acquire satellite  
205 image data (section 2.1), yield maps of several cereals from 2017 to 2021 (section 2.2),  
206 and meteorological data (section 2.5). Afterwards, we will introduce general methods in  
207 section 2.7, which will be used in the remaining chapters.

208 **2.1 Sentinel 2 Data**

209 The European Space Agency (ESA)<sup>1</sup> freely distributes the high-quality images of the two  
210 Sentinel satellites (S2). Together, both satellites have a revisit time of 5 days at the  
211 Equator and 2-3 days at mid-latitudes. However, in our study region, we only receive an  
212 image every 5 days.

213 The S2 images contain 12 spectral bands with spatial resolutions up to 10 meters (see  
214 2.1). Bands with a lower resolution (20 and 60 meters) were upscaled to 10 meter reso-  
215 lution using cubic interpolation (Perich et al. (2022)). In order to decrease the effect of  
216 atmospheric conditions like reflections and scattering, bottom-of-atmosphere, radiometric  
217 corrected Level-2A data was used<sup>2</sup>. The ESA also supplies an algorithm<sup>3</sup> produces Scene  
218 Classification Layer (SCL) where for each location the observed subject is assigned to one  
219 of 11 SCL-classes (c.f. table 2.2). In this thesis, we will use this classification to filter out  
220 data points, which we believe to be less informative. That are all observations which SCL-  
221 class does not correspond to vegetation or bare soils (classes 4 and 5). For convenience,  
222 we define the set SCL45 as the observations which belong to SCL-class 4 or 5.

223 **2.2 Crop Yield Data**

224 The crop yield data were collected using a combine harvester. Equipped with GPS, the  
225 harvester drives over the fields and continuously estimates the dry crop yield density in

---

<sup>1</sup><https://sentinel.esa.int/web/sentinel/missions/sentinel-2>

<sup>2</sup>According to Perich et al. (2022): “Data prior to March 2018 was only available in the top-of-  
atmosphere L1C format and was downloaded as such [...] L1C data was processed to L2A product level  
using the ‘Sen2Cor’ processor provided by ESA”

<sup>3</sup><https://sentinels.copernicus.eu/web/sentinel/technical-guides/sentinel-2-msi/level-2a/>  
algorithm

Table 2.1: List of spectral bands of the S2-satellites. Each band has its center at the wavelength  $\lambda$  in nm with the spectral width  $\Delta\lambda$  in nm with a spatial resolution  $SR$  in m ([Jaramaz et al. \(2013\)](#)).

Band	$\lambda$	$\Delta\lambda$	$SR$	Purpose
1	443	20	60	Atmospheric correction (aerosol scattering)
2	490	65	10	Sensitive to vegetation senescing, carotenoid, browning and soil background; atmospheric correction (aerosol scattering)
3	560	35	10	Green peak, sensitive to total chlorophyll in vegetation
4	665	30	10	Maximum chlorophyll absorption
5	705	15	20	Position of red edge; consolidation of atmospheric corrections / fluorescence baseline.
6	740	15	20	Position of red edge, atmospheric correction, retrieval of aerosol load.
7	783	20	20	Leaf Area Index (LAI), edge of the Near-Infrared (NIR) plateau.
8	842	115	10	LAI
8a	865	20	20	NIR plateau, sensitive to total chlorophyll, biomass, LAI and protein; water vapor absorption reference; retrieval of aerosol load and type.
9	945	20	60	Water vapor absorption, atmospheric correction.
10	1375	30	60	Detection of thin cirrus for atmospheric correction.
11	1610	90	20	Sensitive to lignin, starch and forest above ground biomass. Snow/ice/-cloud separation.
12	2190	180	20	Assessment of Mediterranean vegetation conditions. Distinction of clay soils for the monitoring of soil erosion. Distinction between live biomass, dead biomass and soil, e.g. for burn scars mapping.

Table 2.2: Overview: Scene Classification Layers (SCL)

Color	No.	Class	Color	No.	Class
	0:	Missing Data		6:	Water
	1:	Saturated or defective pixel		7:	Cloud low probability
	2:	Dark features / Shadows		8:	Cloud medium probability
	3:	Cloud shadows		9:	Cloud high probability
	4:	Vegetation		10:	Thin cirrus cloud
	5:	Bare soils		11:	Snow or ice

226  $t/\text{ha}$  (see fig. [2.1a](#)). We take the data set derived in [Perich et al. \(2022\)](#), where error-prone measurement points (such as during a tight curve of the combine harvester) were removed and then the yield map was rasterized using linear interpolation (c.f. fig. [2.1b](#)).

229 We summarize the rasterized dry-yield values by the following statistics:

230 Minimum 1st Quartile Median Mean 3rd Quartile Maximum Variance  
0.107 6.186 7.560 7.359 8.756 13.35 4.035

231 Comparing the average per-field crop yield reported by the farmer with the yield estimated by the combine harvester shows that the latter overestimates crop yield by ca. 10% (c.f. [Perich et al. \(2022\)](#)). Since the relative estimation error is approximately constant and we do not aim for an accurate yield prediction, we will not consider this deviation.

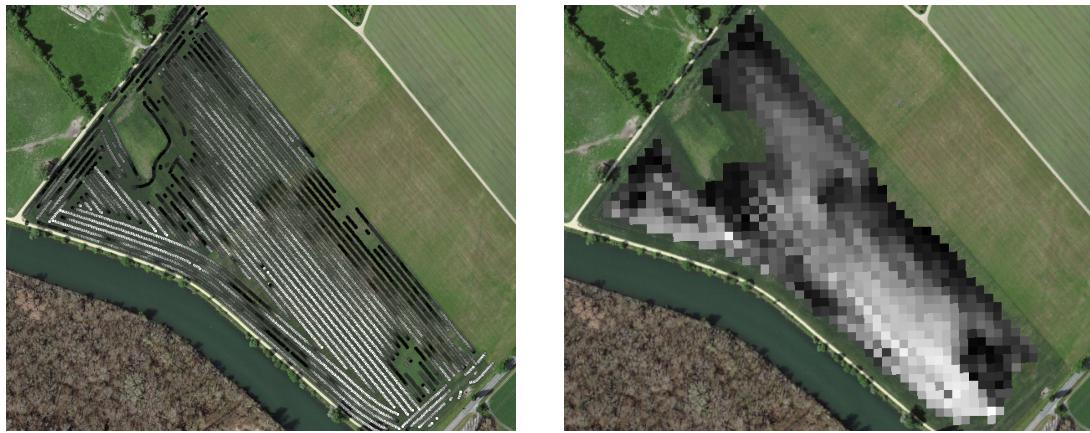


Figure 2.1: Crop yield density map of a field. Ranges from 0.1 t/ha (black) to 5.35 t/ha (white)

### 2.3 Normalized Difference Vegetation Index (NDVI)

The well-known (NDVI) introduced in [Rouse \(1974\)](#) is used to measure vegetation in remote sensing. It utilizes a large jump of reflectancy between red and infrared and can be calculated using the bands  $B4$  and  $B8$  (table 2.1) by:

$$NDVI = \frac{B8 - B4}{B8 + B4}$$

Since we measure the NDVI via the S2 satellites from space we can not expect to measure the true NDVI. This is especially true if we do not see the ground because of clouds or the ground signal is disturbed by cloud shadows. Even if we only use SCL45 observations we still encounter issues as will be described in section 2.6. Therefore, we call the calculated values merely the observed NDVI. In the following chapters, we will study the resulting NDVI time series (for one location and one season) extensively. Such a time series is shown in figure 2.2a.

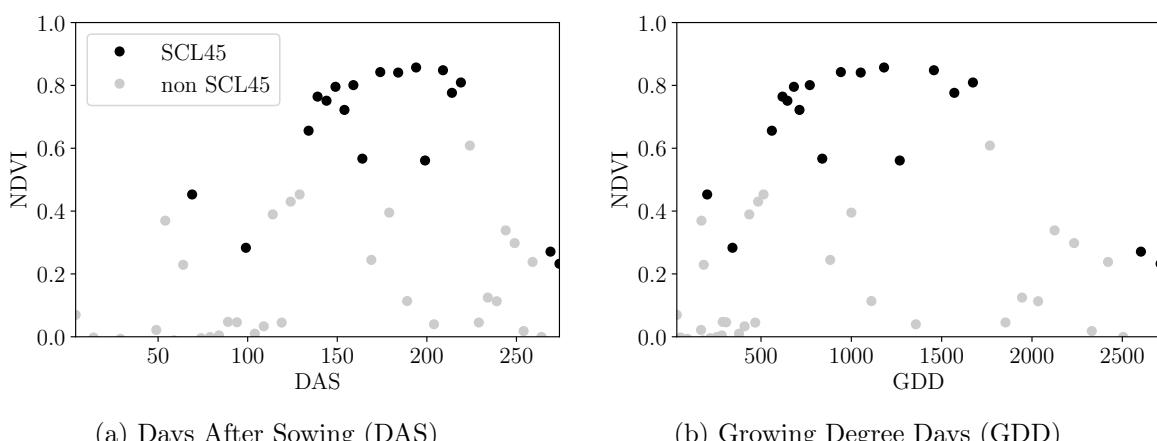


Figure 2.2: NDVI time series plotted against DAS and GDD. GDD are introduced in section 2.4.

## 246 2.4 Timescale Transformation

247 Regarding the Days After Sowing (DAS) time scale shown in fig. 2.2a, we detect two  
 248 drawbacks. First, this scale makes it difficult to compare two NDVI time series because  
 249 wheat is not always sown on the same day of the year and in some years plants begin  
 250 to emerge earlier. Second, because there are only few SCL45 observations in the winter,  
 251 we face significant data gaps in this period. The time scale transformation introduced in  
 252 McMaster and Wilhelm (1997) fixes both problems. The resulting Growing Degree Days  
 253 (GDD) are defined as the cumulative sum since sowing of temperature above a given base  
 254 temperature  $T_{base}$ . For cereals, we use  $T_{base} = 0$  (Perich et al. (2022)). Thus, the GGD  
 255 for  $n$  days after sowing will be equal to:

$$GDD_n := \sum_{i=0}^n \max(T_i - T_{base}, 0).$$

256 Important plant growth stages and their corresponding GDD values are tabultaed in B.1.1  
 257 In figure 2.2 we see an example for comparison of the DAS and GDD timescale. Here  
 258 we see that the first 120 DAS are compressed to just 500 GDD and hence the gap in  
 259 observations was succesfully compressed. Due to the reasons mentioned above, from now  
 260 on we will only consider GDD.

## 261 2.5 The Concept of a ‘Pixel’

262 Now we create a new data structure that we call Pixel. This originates from the pixels of  
 263 the S2 satellite images. It will contain all the information needed to confront the tasks in  
 264 the following chapters.

265 Consider a 10 by 10 meter square that coinsides with a S2 image pixel and  $T$  the GDD  
 266 values for which S2 images are avialable in a given season. For  $t \in T$  let  $P_t$  be a tupel of  
 267 all the spectral bands, the observed NDVI and the SCL class (at the considered location  
 268 at time  $t$ ). Then, define  $P$  as the collection of all the  $P_t$  and the estimated dry-yield for  
 269 this square. Analogously to  $P$ , define  $P^{SCL45}$  by only considering  $P_t$  with SCL-class 4 or  
 270 5 (vegetation and soil).

## 271 2.6 Challenges in S2 Data

272 Now, we shall illustrate with an example pixel the challenges, we will confront in the  
 273 coming chapters. The figure 2.3 shows a selection of 6 satellite images of a field, one  
 274 selected Pixel and the NDVI time series of that pixel. In February (image a), we see  
 275 no vegetation but bare soil and thus also a low NDVI. At the beginning of May (b), we  
 276 observe a cloudless dark green field with a high NDVI. In (c) heavy cloud cover (SCL class  
 277 9) leads to a complete loss of plant information in this S2 observation. Figure (d) shows  
 278 that the SCL classification is not reliable, since we evidently observe clouds which is also  
 279 reflected in a sudden NDVI drop. Even though SCL indicates that (e) are thin cirrus  
 280 clouds, we see a pale green and we also note a NDVI.

281 So in conclusion, we remark that some SCL45 observations are not accurate and even  
 282 though a few non-SCL45 observations contain useful information, most of them are too  
 283 unreliable (e.g. all SCL 9 observations). Thus, we aim to substitute the unreliable ones  
 284 with interpolated versions and correct corrupt ones.

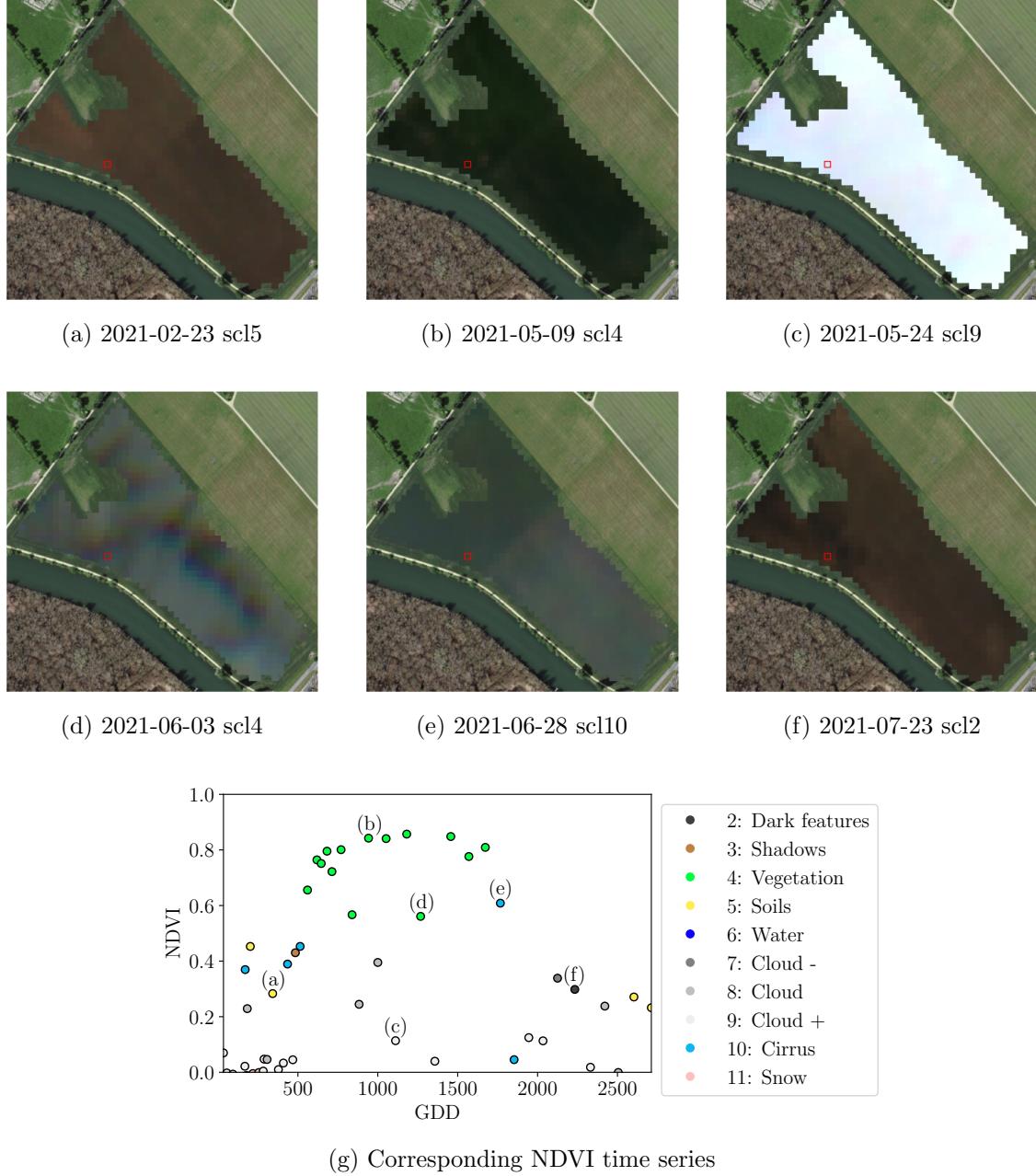


Figure 2.3: Satellite images of a field at selected times with a static background for orientation. Moreover, the NDVI time series of the red-highlighted pixel is shown in (g) colored by the SCL labels.

## 2.7 General Methods

Here we will only introduce Methods which will accure in several places. For interpolation methods we refer to sections 3.2 and 3.3, for a robust interpolation strategy to section 3.5. In section 3.4 we describe a method to objectively determine the quality of an interpolation, and in chapter 4 we present the NDVI correction together with an adapted interpolation strategy.

### 2.7.1 Root Mean Square Error (RMSE)

In this section we describe different criteria to evaluate models. Hence, given a vector  $y \in \mathbb{R}^n$  and its estimator  $\hat{y}$  (estimated using the model), we define the RMSE as:

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

### 2.7.2 Out-Of-Bag (OOB) and Leave-One-Out-Cross-Validation (LOOCV)

The rationale for OOB and LOOCV is that we intend to evaluate a model  $M$  with unseen data. That is, if  $D$  describes the entire dataset and we train a model on a subset of  $D$ , we can use the remaining data to evaluate the model.

To formally introduce this, let:

$$D = \{(X_{[j,:]}, y_j) \mid X \in \mathbb{R}^{n \times p}, y \in \mathbb{R}^n, j = 1, \dots, n\}$$

be a dataset,  $i \in \{1, \dots, n\}$  and  $M^{(-i)}$  a model fitted on a subset of  $D \setminus \{(X_{[i,:]}, y_i)\}$ . Then we call  $\hat{y}_i := M^{(-i)}(X_{[i,:]})$  an OOB estimator of  $y_i$ . If we do this for all  $i \in \{1, \dots, n\}$ , we obtain  $\hat{y} := (\hat{y}_1, \dots, \hat{y}_n)$  the OOB estimator for  $y \in \mathbb{R}^n$ .

In the bootstrap (e.g., random forest) framework, we define  $\hat{y}_i$  to be the average of all computed and admissible  $M^{(-i)}$ .

In the case that  $M^{(-i)}$  was fitted on the set  $D \setminus \{(X_i, y_i)\}$  (i.e., not a true subset), we call the corresponding  $\hat{y}_i$  also the LOOCV estimator.

If we optimize some parameter via OOB (or LOOCV) this means that we search for the parameter that minimizes some loss function which takes the OOB (or LOOCV) residuals.

Usually we approximate this parameter by searching on a grid.

308

## Chapter 3

309

# Interpolation Methods

310

311 In section 2.6 we have established the need for interpolating the NDVI time series. In  
312 this chapter we first specify a setting for the interpolation and divide the interpolation  
313 methods into those that make fundamental shape assumptions (parametric) and those  
314 that are more flexible (non-parametric). We give an introduction for each method with  
315 an compact definition, highlight adjustments or give remarks where appropriate, and then  
316 point out strengths and weaknesses of each method. Additionally, a brief overview of  
317 the considered interpolation methods is provided in table 3.1. Afterwards, we extract an  
318 robustification strategy from the one interpolation method and generalize it so we can use  
319 it for all methods that allow for a priori weighted observations. Finally, using LOOCV,  
320 we tune the parameters (where necessary) and get a first idea of the performance of each  
321 method.

verdeutliche  
dem  
leser,  
dass ein  
auftrag  
das  
findne  
von  
interpo-  
lation-  
metho-  
den war

322

### 3.1 Interpolation Setup

In this chapter, we will only consider SCL45 observations, since they are more reliably. Hence, data in the form of  $(t_i, y_i)$  for  $i = 1, \dots, n$  is given, where  $t_i$  is the time in GDD and  $y_i$  denotes the NDVI at time  $t_i$ . Assume that it can be represented by

$$y_i = m(t_i) + \varepsilon_i,$$

where  $\varepsilon_i$  is some noise and  $m : \mathbb{R} \rightarrow \mathbb{R}$  is some (parametric or non-parametric) function. If we assume that  $\varepsilon_1, \dots, \varepsilon_n$  i.i.d. with  $\mathbb{E}[\varepsilon_i] = 0$  then

$$m(t) = \mathbb{E}[y | t]$$

323 We will introduce parametric and non-parametric approaches to estimate  $m$  in section 3.2  
324 and 3.3 Furthermore, in the subsequent, we denote  $w \in \mathbb{R}^n$  as the vector of weights such  
325 that  $w_i$  corresponds to the weight that  $(t_i, y_i)$  should have in the interpolation.

326 Paper zitieren wo eingeführt oder wo benutzt (falls einföhrung fast schon trivial)

327

### 3.2 Parametric Regression

328 Parametric Curve estimation tries to fit a parametric function, such as, for example, a  
329 Gaussian function with parameters  $\mu$  and  $\sigma$ , to a dataset. In the following, we introduce  
330 two parametric approaches.

Table 3.1: Summary of the studied interpolation methods containing important assumptions, advantages and disadvantages and whether the method supports weighted observations (w) and if the resulting interpolation is bounded w.r.t. a fixed interval (b).

	<b>Assumptions</b>	<b>Advantages</b>	<b>Disadvantages</b>	w	b
Double- Logistic	<ul style="list-style-type: none"> <li>- Function first increases then decreases</li> <li>- NDVI has a minimal value</li> </ul>	<ul style="list-style-type: none"> <li>- Good for evergreen plants (if snow masks NDVI)</li> <li>- Upper envelope</li> </ul>	<ul style="list-style-type: none"> <li>- Parameter estimation can be very difficult</li> <li>- Strange behavior for long data-gaps</li> </ul>	Yes	(Yes)
Fourier Series	<ul style="list-style-type: none"> <li>- NDVI can be approximated by a 2cd order Fourier series.</li> </ul>	<ul style="list-style-type: none"> <li>- Incorporates periodical growth-cycles</li> </ul>	<ul style="list-style-type: none"> <li>- Parameter estimation can be very difficult</li> <li>- Curve easily exceeds bounds of the NDVI</li> </ul>	Yes	No
Nadaraya- Watson (Kernel Smooth- ing)	<ul style="list-style-type: none"> <li>- Close points are related to each other via a kernel function</li> </ul>	<ul style="list-style-type: none"> <li>- Simple</li> <li>- Computationally very fast</li> </ul>	<ul style="list-style-type: none"> <li>- Biased, especially at ‘peaks’ and ‘valleys’</li> <li>- Bandwidth: fails if there are big data-gaps</li> </ul>	Yes	Yes
Universal Kriging	<ul style="list-style-type: none"> <li>- Function is a realization of a stationary Gaussian process</li> </ul>	<ul style="list-style-type: none"> <li>- Informative parameters</li> <li>- Flexible</li> </ul>	<ul style="list-style-type: none"> <li>- Regression to the mean</li> <li>- Assumptions clearly not met</li> </ul>	Yes	(Yes)
SG	<ul style="list-style-type: none"> <li>- High frequencies are noise (Low-Pass-Filter)</li> <li>- Equidistant points</li> <li>- Local polynomials</li> </ul>	<ul style="list-style-type: none"> <li>- Computationally very fast</li> </ul>	<ul style="list-style-type: none"> <li>- Cannot deal natively with missing data (need some interpolation)</li> </ul>	No	(Yes)
SG + NDVI	<ul style="list-style-type: none"> <li>- Upper envelope</li> <li>- Vegetation cannot grow faster than some slope</li> </ul>	<ul style="list-style-type: none"> <li>- Biological knowledge</li> </ul>	<ul style="list-style-type: none"> <li>- Bad “upper envelope” since weights are not used for the estimation itself</li> </ul>	(No)	(Yes)
LOESS	<ul style="list-style-type: none"> <li>- Local polynomial with points closer to the estimated point are more important</li> </ul>	<ul style="list-style-type: none"> <li>- Flexible</li> <li>- Generalization of SG</li> <li>- Weighting function makes intuitive sense</li> </ul>	<ul style="list-style-type: none"> <li>- Computationally expensive</li> </ul>	Yes	(Yes)
B-Splines (Smoothed)	<ul style="list-style-type: none"> <li>- Function can be approximated by a linear combination of B-splines basis functions</li> </ul>	<ul style="list-style-type: none"> <li>- General assumption</li> <li>- Flexible shape</li> </ul>	<ul style="list-style-type: none"> <li>- Unbounded</li> <li>- No intuitive meaning for smoothing</li> </ul>	Yes	No
Smoothing splines	<ul style="list-style-type: none"> <li>- 2cd derivative of function is integrable</li> </ul>	<ul style="list-style-type: none"> <li>- Intuitive meaning of penalty</li> <li>- General assumptions</li> <li>- Flexible shape</li> </ul>	<ul style="list-style-type: none"> <li>- Choice of smoothing parameter</li> </ul>	Yes	No

331 **3.2.1 Double Logistic (DL)**

The Double Logistic smoothing as described in [Beck, Atzberger, Høgda, Johansen, and Skidmore \(2006\)](#)REF heavily relies on shape assumptions of the fitted curve (i.e. the NDVI time series). First, we assume that there is a minimum NDVI level  $y_{\min}$  in the winter (e.g. due to evergreen plants), which might be masked by snow. This can be estimated beforehand, taking several years into account. Second, we assume that the growth cycle can be divided into an increase and a decrease period, where the time series follows a logistic function. The maximum increase (or decrease) is observed at  $t_0$  (or  $t_1$ ) with a slope of  $d_0$  (or  $d_1$ ). The equation of the double-logistic fit is given by:

$$y(t) = y_{\min} + (y_{\max} - y_{\min}) \left( \frac{1}{1 + e^{-d_0(t-t_0)}} + \frac{1}{1 + e^{-d_1(t-t_1)}} - 1 \right)$$

- 332 Where the five free parameters:  $y_{\max}$ ,  $d_0$ ,  $d_1$ ,  $t_0$ ,  $t_1$  are initially estimated by least squares.  
 333 Such fit can be seen in figure [3.1](#).

334 **Robustification**

- 335 Similar as for the SG (c.f. section [3.3.3](#)) one can reestimate (only once) the parameters by  
 336 giving less weight to the overestimated observations and more weight to the underestimated  
 337 observations. For the details on the choice of the weights we refer to [Beck et al. \(2006\)](#). We  
 338 will not apply this reestimation but rather the robustification introduced later in section  
 339 [3.5](#).

Advantages	Disadvantages
<ul style="list-style-type: none"> <li>— Incorporates subject specific knowledge in the case of evergreen plants covered in snow.</li> <li>— Optimized parameters have an intuitive meaning.</li> </ul>	<ul style="list-style-type: none"> <li>— Strong shape assumptions on the NDVI curve.</li> <li>— Parameter optimization might go wrong. This can be mitigated to some extent to provide bounds for the parameters</li> <li>— Strange behavior in regions with little observations. (c.f. figure <a href="#">3.1</a>)</li> </ul>

340 **3.2.2 Fourier Series (FS)**

Analogous to section [3.2.1](#) we fit a parametric curve to the data by least squares. Here we take the second order FS approximation:

$$\text{NDVI}(t) = \sum_{j=0}^2 a_j \times \cos(j \times \Phi_t) + b_j \times \sin(j \times \Phi_t)$$

- 341 where  $\Phi = 2\pi \times (t - 1)/n$ . Thus, we periodical behavior. If we would set the period to  
 342 match one year this would coinced with the nothion that plans grow every year. Example  
 343 fits can be seen in figure [3.1](#)

Advantages	Disadvantages
<ul style="list-style-type: none"> <li>— Assumption of periodicity can be helpful if we are modelling multiyear grow cycles</li> <li>— Flexible curve shape</li> </ul>	<ul style="list-style-type: none"> <li>— Bad behavior in regions with little data (c.f. figure 3.1)</li> <li>— Hard to interpret estimated parameters</li> <li>— Parameter estimation can go wrong. Introducing bounds can help.</li> </ul>

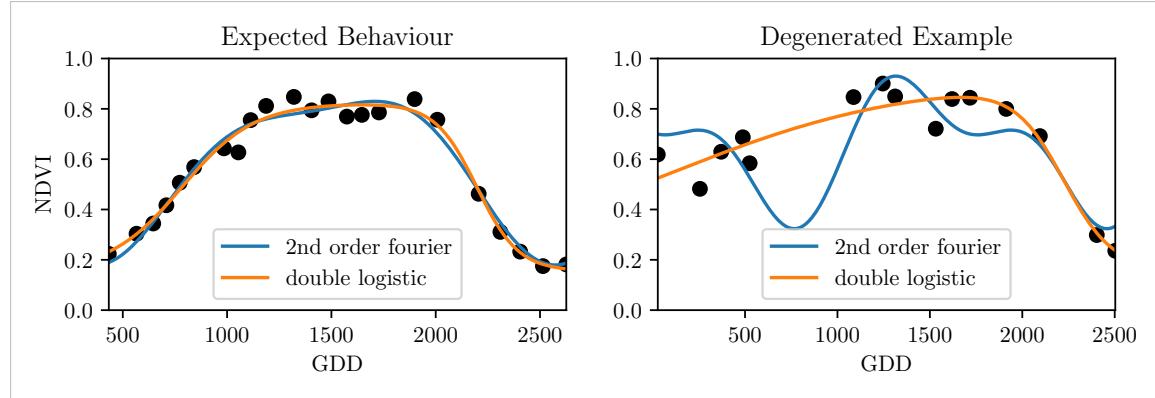


Figure 3.1: Here we observe the possibilities of a precise fit for the two parametric methods but notice also some misbehavior

### 344 3.2.3 Optimization Issues

345 We shall mention some optimization issues we countered during implementation. Since we  
 346 aim to minimize the residual sum of squares over 5 (or 6) parameters, we try to solve a  
 347 non-convex optimization problem. Thus, the algorithm<sup>1</sup> either struggles to find the global  
 348 minimum or fails to converge. This was fixed by providing for each parameter reasonable  
 349 initial values and generous bounds (which match our experience).

## 350 3.3 Non-Parametric Regression

351 In non-parametric curve estimation, the curve does no longer have to be fully determined  
 352 by parameters, but we allow it to flexibly approximate the data. Note, that we do not  
 353 exclude the use of tuning-parameters.

### 354 3.3.1 Kernel Regression: Nadaraya-Watson (NW)

355 As described in section 3.1, we aim to estimate

$$\mathbb{E}[Y \mid T = t] = \int_{\mathbb{R}} y f_{Y|T}(y \mid t) dy = \frac{\int_{\mathbb{R}} y f_{T,Y}(t,y) dy}{f_T(t)}, \quad (3.3.1.1)$$

356 where  $f_{Y|T}$ ,  $f_{T,Y}$ ,  $f_T$  denote the conditional, joint and marginal densities. This can be done  
 357 with a kernel  $K$ :

$$\hat{f}_T(t) = \frac{\sum_{i=1}^n K\left(\frac{t-t_i}{h}\right)}{nh}, \quad \hat{f}_{T,Y}(t,y) = \frac{\sum_{i=1}^n K\left(\frac{t-t_i}{h}\right) K\left(\frac{y-Y_i}{h}\right)}{nh^2}, \quad (3.3.1.2)$$

<sup>1</sup>We used the python function `scipy.optimize.curve_fit`.

where  $h$ , the bandwidth, symbolizes the windowsize of to consider. By using the above function in equation (3.3.1.1) we arrive at the NW kernel estimator:

$$\hat{m}(t) = \frac{\sum_{i=1}^n K((t - t_i)/h) Y_i}{\sum_{i=1}^n K((t - t_i)/h)}$$

358 Common choices for the kernel are the normal function or a uniform function (also called  
 359 ‘bot’ function).

### 360 Choose Bandwidth

361 Note that we still need to choose the bandwidth of the function. This can be done with  
 362 the help of LOOCV while optimizing the RMSE. For non-equidistant data we refere to  
 363 [Brockmann, Gasser, and Herrmann \(1993\)](#) where a local adaptive bandwidth selection is  
 364 presented.

Advantages	Disadvantages
— flexible due to different possible kernels	— if the $t \mapsto K(t)$ is not continuous, $\hat{m}$ isn't either
— can be assigned degrees of freedom (trace of the hat-matrix)	— choice of bandwidth, especially if $t_i$ are not equidistant.
— estimation of the noise variance $\hat{\sigma}_\varepsilon^2$ (REF c.f. CompStat 3.2.2)	

### 365 3.3.2 Universal Kriging (UK)

366 UK as described in [Diggle and Ribeiro \(2007\)](#) was developed in geostatistics to deal with  
 367 autocorrelation of the response variable at locations which are spatially close. By applying  
 368 the notion that two spectral indices which are timewise close should also take similar values,  
 369 we justify the application of UK. In the end, we would like to fit a smooth Gaussian process  
 370 to the data.

371 A Gaussian Process  $\{S(t) : t \in \mathbb{R}\}$  is a stochastic process if  $(S(t_1), \dots, S(t_k))$  has a multi-  
 372 variate Gaussian distribution for every collection of times  $t_1, \dots, t_k$ .  $S$  can be fully charac-  
 373 terized by the mean  $\mu(t) := E[S(t)]$  and its covariance function  $\gamma(t, t') := \text{Cov}(S(t), S(t'))$ .  
 374 Furthermore, we will assume the Gaussian process to be stationary. That is for  $\mu(t)$  to be  
 375 constant in  $t$  and  $\gamma(t, t')$  to depend only on  $h = t - t'$ . Thus, we will write in the following  
 376 only  $\gamma(h)$ .<sup>2</sup>

Now, we need to make some assumption on the covariance function. For this we introduce the variogram of a Gaussian process as

$$V(h) := V(t, t + h) := \frac{1}{2} \text{Var}(S(t) - S(t + h)) = \gamma(0) + \gamma(t)$$

and define  $\gamma$  via the above equation by choosing the Gaussian Variogram defined by

$$V(h) = p \cdot \left( 1 - e^{-\frac{h^2}{(\frac{4}{7}r)^2}} \right) + n.$$

---

<sup>2</sup>Note that the process is also *isotropic* (i.e.  $\gamma(h) = \gamma(\|h\|)$ ) since we are in a one-dimensional setting and the covariance is symmetric.

377 Here  $h$  denotes the distance,  $n$  is the nugget,  $r$  is the range and  $p$  is the partial sill. The  
378 influence of the parameters is visualized in figure 3.2.<sup>3</sup>

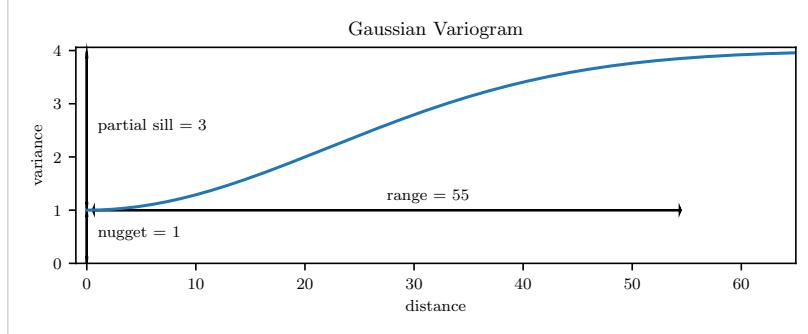


Figure 3.2: Gaussian Variogram with nugget=1, partial sill=3, range=55

379 Finally, we consider a one-dimensional Gaussian process  $G_\gamma$  with variogram  $\gamma$  and tune the  
380 variogram parameters using maximum likelihood<sup>4</sup>. Let  $z$  be a vector with the new values  
381 to extrapolate, then we can determine the values  $m(z) = \mathbb{E}[G_\gamma(z)|(t, y)]$  using Bayes rule<sup>5</sup>.  
382 For an example fit, we refer to figure 3.3.

### 383 Violated Assumption

384 Since we observe a clear pattern of a growth period in spring and harvest in the end  
385 of summer, we have to admit that our stationarity assumption with the constant mean  
386 is structurally violated. This is also the reason why we observe (for every variogram  
387 parameter) a tendency to the mean, as indicated in figure 3.3.

Advantages	Disadvantages
<ul style="list-style-type: none"> <li>— It is a well-studied method.</li> <li>— Variogram parameters have an intuitive meaning.</li> <li>— Flexible covariance structure.</li> </ul>	<ul style="list-style-type: none"> <li>— Regression to the mean.</li> <li>— Violated assumption of constant mean and constant variance. Thus, the NDVI is not a stationary process.</li> <li>— Pure maximum likelihood can result in overfitting.</li> </ul>

### 388 3.3.3 Savitzky-Golay Filter (SG)

389 The SG, introduced in [Savitzky and Golay \(1964\)](#) is a technique in signal processing and  
390 can be used to filter out high frequencies (low-pass filter) ([Schafer, 2011](#)). Furthermore,  
391 it can also be used for smoothing by filtering high frequency noise while keeping the low  
392 frequency signal.

First, we choose a window size  $m$ . Then, for each point,  $j \in \{m, m+1, \dots, n-m\}$  we fit

<sup>3</sup>Strictly speaking we use a scaled version of the variogram. Thus, only the ratio of  $p/n$  matters.

<sup>4</sup>As illustrated in figure 3.3 maximum likelihood estimation can lead to overfitting. Thus, we will in practice sample several such optimized parameters and use their median in the end.

<sup>5</sup>Bayes rule generally claims, that for two random variables  $A$  and  $B$  we have that  $P(A|B) = P(B|A)/P(B)$

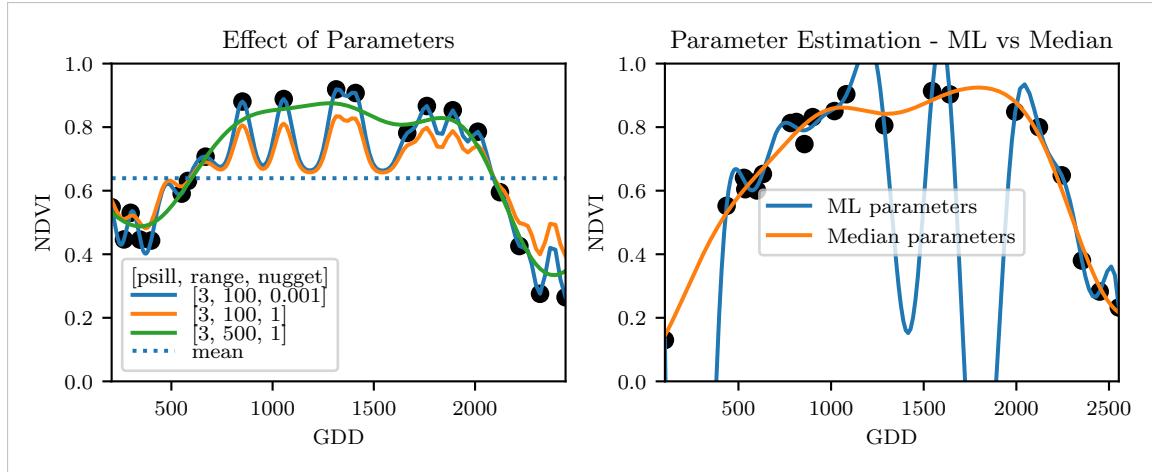


Figure 3.3: On the left, we see how the interpolation change if we increase the nugget and the range parameter. On the right, we compare two UK interpolations, where one takes parameters by numerically matimizing the (which results in a very small nugget) and the other takes the median of many such numerical optimizations.

a polynomial of degree  $k$  by:

$$\hat{y}_j = \min_{p \in P_k} \sum_{i=-m}^m (p(t_{j+i}) - y_{i+j})^2,$$

where  $P_k$  denotes the Polynomials of degree  $k$  over  $\mathbb{R}$ . For equidistant points this can efficiently be calculated by

$$\hat{y}_j = \sum_{i=-m}^m c_i y_{j+i},$$

393 where the  $c_i$  are only dependent on the  $m$  and  $k$  and are tabulated in the original paper.

394 Chen, Jönsson, Tamura, Gu, Matsushita, and Eklundh (2004) developed a ‘robust’ inter-  
395 polation method for the NDVI based on the SG. The method is based on the assumption  
396 that due to atmospheric effects the observed NDVI tends to be underestimated and that  
397 it cannot increase too quickly. The latter is argued by the biological impossibility of such  
398 fast vegetation changes. Their proposed algorithm is:

- 399 i.) Remove non-SCL45 points.
- 400 ii.) Remove points which would indicate an increase greater than 0.4 within 20 days.
- 401 iii.) Linearly interpolate to obtain an equidistant time series  $X^0$ .
- 402 iv.) Apply the SG to obtain a new time series  $X^1$ .
- 403 v.) Update  $X^1$  by applying again a SG. Repeat this until  $w^T |X^1 - X^0|$  stops decreasing,  
404 where  $w$  is a weight vector with  $w_i = \min\left(1, 1 - \frac{X_i^1 - X_i^0}{\max_i \|X_i^1 - X_i^0\|}\right)$ . This reduces the  
405 penalty introduced by outliers<sup>6</sup> and by repeating this step we approach the “upper  
406 NDVI envelope”.

figure /  
tabelle /  
pseu-  
doode  
anstatt  
aufzäh-  
lung

<sup>6</sup>Here we call a point  $i$  an outlier if  $X_i^0 < X_i^1$ .

407 **Extension: Spatial-Temporal SG**

408 One notable adaptation of the SG is the presented by [Cao, Chen, Shen, Chen, Zhou, Wang, and Yang \(2018\)](#). The key difference is the additional assumption of the cloud cover  
 409 being discontinuous and that we can improve by looking at adjacent pixels<sup>7</sup>. Because we  
 410 are working with rather high resolution satellite data, and we need the variance in the  
 411 predictors, we will waive this extension.

Advantages	Disadvantages
— Popular technique in signal processing.	— No natural way of how to estimate points which are not in the data.
— Efficient calculation for equidistant points.	— Not generalizable to other spectral indices.
— Upper envelope matches intuition for the NDVI. Therefore, it is robust against outliers with small values.	— Linear interpolation to account for missing data might be not appropriate.
	— No smooth interpolation between two measurements.

413 **3.3.4 Locally Weighted Regression (LOESS)**

414 The LOESS introduced by [Cleveland \(1979\)](#) can be understood as a generalization of the  
 415 SG (c.f. sec. [3.3.3](#)).

Given a proportion  $\alpha \in (0, 1]$ , we estimate each  $y_i$  separately by fitting a polynomial of order  $d$  by weighted least squares. The weights are (usually) defined by

$$w_i(t_j) = \begin{cases} \left(1 - \left(\frac{|t_j - t_i|}{h_i}\right)^3\right)^3, & \text{for } |t_j - t_i| < h_i \\ 0, & \text{for } |t_j - t_i| \geq h_i \end{cases}$$

416 where  $h_i$  is the minimal distance such that  $\lceil \alpha n \rceil$  observations are in the ball  $B_{h_i}(t_i)$ .<sup>8</sup> So  
 417 for each  $y_i$  we only consider a proportion  $\alpha$  of the observations.

418 **Differences between the Robust LOESS and the SG?**

419 The LOESS smoother takes a fraction of points instead of a fixed number and therefore  
 420 automatically adapts to the size of the data we wish to interpolate. However, we run  
 421 into the danger of considering too little observations, since the estimation breaks down if  
 422  $\lceil \alpha n \rceil < d + 1$ .<sup>8</sup> Furthermore, LOESS gives less weight to points further away. This yields  
 423 a "smoother" estimate, since when we slide the window (e.g. for estimating the next value)  
 424 an influential point at the border does not suddenly get zero weight from being weighted  
 425 equally before. Finally, the LOESS also can be used for non-equidistant data and allows  
 426 for arbitrary interpolation.

<sup>7</sup>Here, we say that a pixel is adjacent if it is the same pixel but from a different year (keeping the same day of the year) or (if not enough of such temporal-adjacent pixel are found) it is spatially adjacent

<sup>8</sup>If too many weights are set to zero, we might end up considering not enough observations and thus get a singular design-matrit (for the least squares estimation). Therefore, we substitute  $h_i$  with  $1.01h_i$ , so that the observation on the boundary of  $B_{h_i}(t_i)$  does not get completely ignored. But we also have to assure that  $\alpha$  is big enough.

Advantages	Disadvantages
<ul style="list-style-type: none"> <li>— Flexible generalization of SG</li> <li>— arbitrary interpolation possible</li> <li>— Intuitive parameters</li> </ul>	<ul style="list-style-type: none"> <li>— The nature of local regression might lead to surprising estimates (no smoothness guarantees for the second derivative)</li> </ul>

427 **3.3.5 B-Splines (BS)**

BS as discussed in [Lyche and Mørken \(2005\)](#) are piecewise cubic polynomials defined by

$$S(t) = \sum_{j=0}^{n-1} c_j B_{j,k;t}(t),$$

428 where  $B$  are basis functions and recursively defined by:

429

$$\begin{aligned} B_{i,0}(z) &= 1, \text{ if } t_i \leq z < t_{i+1}, \text{ otherwise } 0 \\ B_{i,k}(z) &= \frac{z - t_i}{t_{i+k} - t_i} B_{i,k-1}(z) + \frac{t_{i+k+1} - z}{t_{i+k+1} - t_{i+1}} B_{i+1,k-1}(z). \end{aligned}$$

Assuming that all  $t_i$  are distinct, this yields an interpolation which fits the data perfectly. To reduce the amount of overfitting and increase the smoothness, we relax the constraint that we have to perfectly interpolate. Thus, we use the minimum number of basis functions<sup>9</sup> such that:

$$\sum_{i=1}^n (w_i(y_i - \hat{y}_i))^2 \leq s$$

Advantages	Disadvantages
<ul style="list-style-type: none"> <li>— can be assigned degrees of freedom</li> <li>— extendable to "smooth" version</li> <li>— performs also well if points are not equidistant</li> </ul>	<ul style="list-style-type: none"> <li>— smoothing process does not translate well to a interpretation (unlike SS)</li> <li>— choice of smoothing parameter <math>s</math></li> </ul>

430 **3.3.6 Smoothing Splines (SS)**

431 Let  $\mathcal{F}$  be the Sobolev space (the space of functions of which the second derivative is  
432 integrable). Then the unique<sup>10</sup> minimizer

$$\hat{m} := \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n w_i (y_i - f(t_i))^2 + \lambda \int f''(t)^2 dt \quad (3.3.6.1)$$

433 is a cubic spline (i.e. a piecewise cubic polynomial function). The objective function  
434 ensures that we decrease the curvature while keeping the RMSE low.

<sup>9</sup>So we do not require one basis function for each neighboring pair of knots. SciPy uses FITPACK and DFITPACK, the documentation suggests that smoothness is achieved by reducing the number of knots used

<sup>10</sup>Strictly speaking it is only unique for  $\lambda > 0$

Advantages	Disadvantages
<ul style="list-style-type: none"> <li>— Can be assigned degrees of freedom (trace of the hat-matrix).</li> <li>— Efficient estimation (closed form solution).</li> <li>— Intuitive penalty (we don't want the function to be too "wobbly" — change slopes).</li> <li>— Also performs well if points are not equidistant.</li> <li>— Fixes the Runge's phenomenon (fluctuation of high degree polynomial interpolation).</li> </ul>	<ul style="list-style-type: none"> <li>— The tuning parameter <math>\lambda</math> must be chosen. This can be done via cross validation and optimizing a score function (e.g. the RMSE).</li> </ul>

### 3.4 Tuning Parameter Estimation

Many of the interpolation methods introduced in section 3.2 and 3.3 include a free parameter. To determine this parameter for a specific interpolation method, we will estimate the absolute residuals using OOB estimation and then optimize the parameter using a score function. We clarify the procedure step by step:

- i.) Construct a set  $\Lambda$  of candidate parameters that generously covers the parameter space.
- ii.) Consider  $\mathcal{P}$ , a set of Pixels.
- iii.) For each parameter  $\lambda \in \Lambda$  consider the individual pixels and compute the LOOCV<sup>11</sup> for the absolute residuals of the specific NDVI interpolation method for all Pixels in  $\mathcal{P}$  and store them in the set  $R_\lambda$ .
- iv.) Determine  $\lambda_{optimal} = \arg \min_{\lambda \in \Lambda} q_{90}(R_\lambda)$ , where we describe the 90% quantile with  $q_{90}$ .

We choose quantile(90) as our optimization function because we want to allow 10% of outliers (corrupt points) but also aim for an accurate fit in 90% of the cases.

Figure 3.4 exemplifies the effect of the optimization function (different quantiles). To summarize, we may say that the higher the quantile, the stronger the smoothing.

### 3.5 Robustification

Now we discuss a general approach of how to make an interpolation more robust against outliers. The main idea is to give less weight to observations that have high residuals after the initial (or if we reiterate, the previous) fit.

Even though the procedure is taken from the robust version of the LOESS smoother (c.f. section 3.3.4 and Cleveland (1979)), we can apply it to every interpolation method that allows for prior weighting of observations.

<sup>11</sup>For a definition of the leave-one-out-cross-validation we refer to section 2.7.2

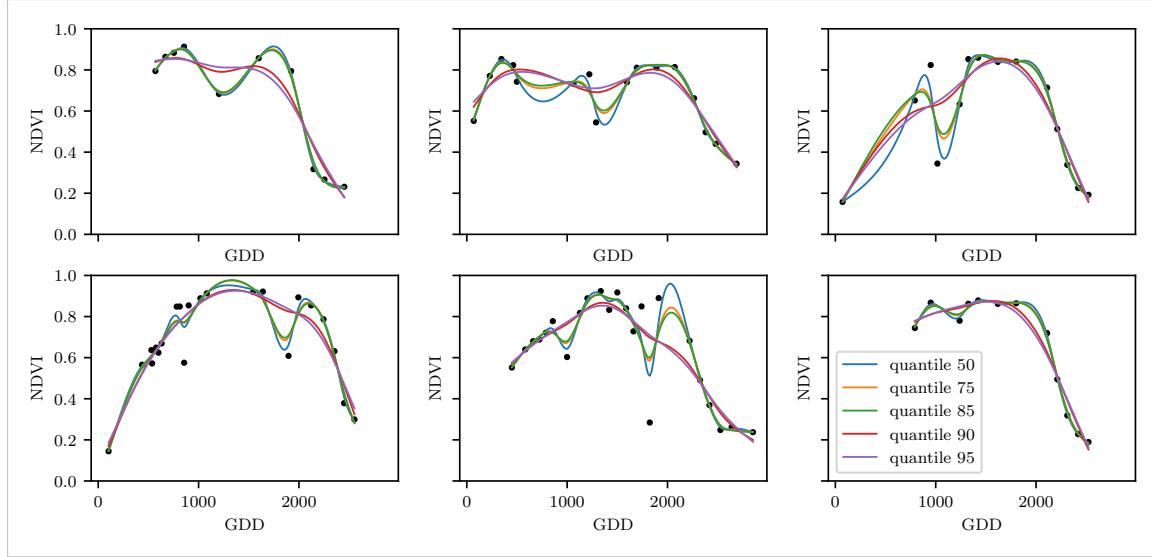


Figure 3.4: Smoothing splines fit with smoothing parameter optimized by minimizing the given quantile of the absolute leave-one-out residuals. Note that the larger the considered quantile is, the smoother the resulting curve becomes.

459 After an initial fit we calculate the residuals  $r_i := y_i - \hat{y}_i$  and obtain  $\tilde{r}_i$  by scaling with the  
460 median of the absolute residuals:

$$\tilde{r}_i := \frac{r_i}{6 \text{ med}(|r_1|, \dots, |r_n|)}$$

461 Next, we compute new weights by

$$w_i^{\text{new}} := w_i^{\text{old}} \begin{cases} (1 - \tilde{r}_i^2)^2, & \text{if } |\tilde{r}_i| < 1 \\ 0, & \text{else} \end{cases}; \quad (3.5.0.1)$$

462 Using the new weights, we can re-interpolate. This reweighting can be iterated for several  
463 steps or till the change of the values is smaller than some tolerance.

464 Note that this procedure is indeed robust since we use the median for the normalization  
465 which has a breakdown point<sup>12</sup> of 50%.<sup>13</sup>

### 466 3.5.1 Our Adjustment:

During the iterations or when supplying prior weights, low-weighted observations can corrupt our estimation of scale (the median of absolute residuals). Thus, we introduce the weighted median as

$$\text{med}_{\text{weighted}}(r, w) := \arg \min_{\lambda \in \mathbb{R}} \sum_{i=1}^n |r_i w_i - \lambda|$$

467 for  $r, w \in \mathbb{R}^n$ .

<sup>12</sup>Intuitively, the breakdown point denotes the fraction of observations a “vicious” player can replace without breaking the estimator. For example, the median has a breakdown point of 50%.

<sup>13</sup>The breakdown point relates only to outliers in the  $y$  values. Note that we do not require the interpolation methods to be robust, since the residual for an outlier will still be larger than for non-outliers and thus will be down weighted more and more in each iteration (because for the next iteration the residual of the outlier will be even larger, since we gave less weight to it).

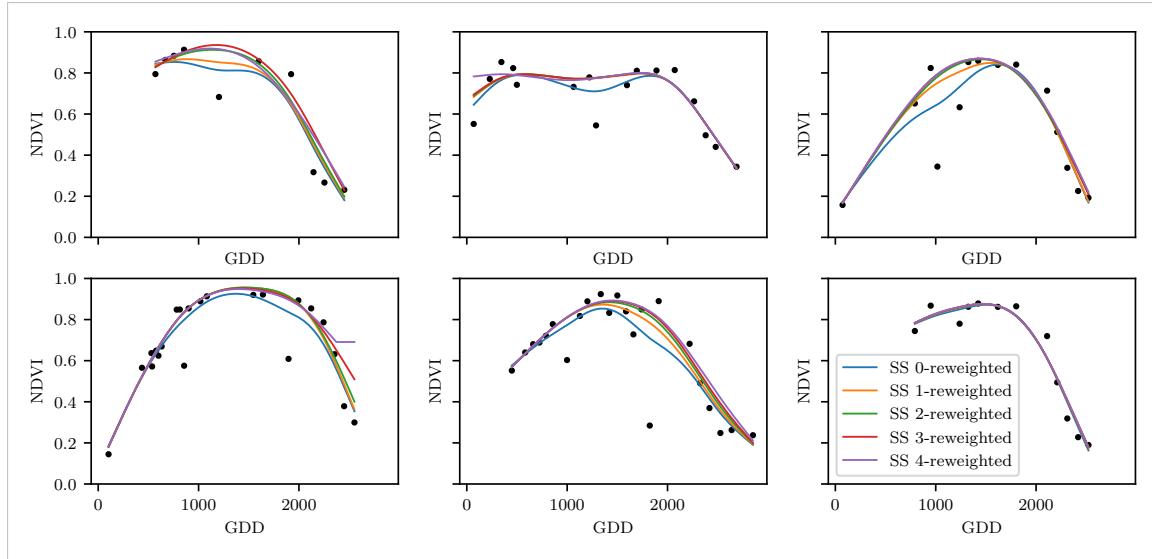
468 **3.5.2 Examples and Conclusions**

Figure 3.5: Smoothing splines fitted to different (SCL45) NDVI time series. Iterations of a robustifying refit (as indicated in section 3.5) are also displayed

469 Examples of the first four iterative fits using SS are shown in figure 3.5 for six pixels. For  
 470 the analogous figures of the other interpolation methods c.f. figures B.1, B.2, B.3 and B.1.  
 471 Indeed, we observe how the interpolated time series is less affected by outliers after each  
 472 iteration. We notice the biggest difference in the first iteration. Furthermore, in the plot  
 473 at the bottom left we see how the interpolation ‘escapes’ from the right endpoint with  
 474 each successive iteration, even though our intuition does not necessarily identify this point  
 475 as an outlier. Therefore, in the following, we will always stop after one iteration.

consider  
naming  
the sub-  
plots

476 **3.5.3 Upper Envelope Approach - Penalty for Negative Residuals**

477 If we artificially increase the negative residuals in 3.5.0.1 by multiplying (e.g. factor 2),  
 478 the corresponding points will get less weight in the next iteration. This allows us to create  
 479 an interpolation that resembles an upper envelope. Intuitively, this upper envelope can be  
 480 thought of as a sheet that is laid on top of the points.

481 This approach is based on the premise that we tend to underestimate the NDVI (as argued  
 482 in Cao et al. (2018)). Since we want to develop a general method that is in principle not  
 483 related to the NDVI, we will not pursue this approach further.

484 **3.6 Performance Assessment**

485 Next, we will benchmark the in section 6.1.2 preselected interpolation methods with and  
 486 without robustification. For this, we will use the same technique as we did for the param-  
 487 eter determination in section 3.4. On  $B_\lambda$  we apply the RMSE and different quantiles.

488 The results are presented in section 5.1 and are discussed in section 6.1. The double logistic  
 489 turns out to be the best convincing parametric method and from the non-parametric  
 490 methods we choose the SS.

491 **Chapter 4**

492 **NDVI Correction**

493 Let's remind ourselves that the data from the S2 satellites is distributed with an SCL and  
494 we therefore have some evidence about what is observed at each pixel for each sampled  
495 time (c.f. table 2.2). So far, we have only considered points, labeled as cloud- and shadow-  
496 free (SCL45). However, we remind ourselves of the satellite images in figure 2.3d, where  
497 we had cloudy images despite the 'vegetation' label and see vegetation in figure 2.3e even  
498 though we are supposed to observe 'cirrus clouds'.

499 In this chapter, we will try to improve our NDVI interpolation by not relying only on the  
500 observed NDVI, but by training our own model to correct the NDVI using all S2 bands.  
501 For this, we introduce several statistical modelling approaches and discuss the strengths  
502 and weaknesses for each of them. After correcting the observed NDVI, we will assess the  
503 uncertainties of our corrections and translate them into weights. These will be used for  
504 the subsequent interpolation. This step-by-step procedure is illustrated by the figure B.4  
505 in the appendix. Finally, we will evaluate which combination of interpolation methods  
506 and correction model performs the best.

507 **4.1 Considering other SCL Classes**

508 In figure 4.1 we plot the observed NDVI and notice that some blue points which correspond  
509 to the SCL-class 10 (thin cirrus clouds) follow the interpolated line closely. Hence, they  
510 might be useful in improving an interpolation fit.

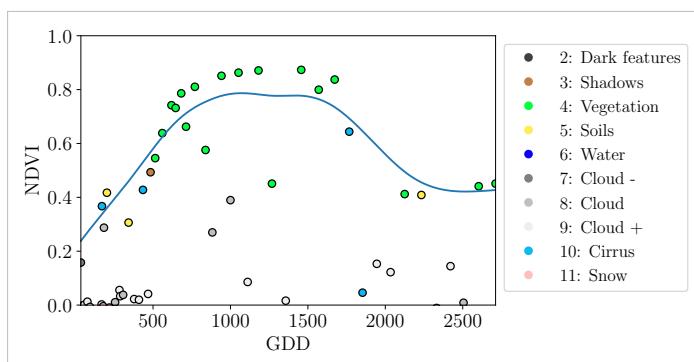


Figure 4.1: A smoothing splines fit considering green and yellow points (SCL45)

511 To get an impression of whether there is some useful information contained in non-SCL45

512 observations, we would like to compare the observed NDVI with the true NDVI. But since,  
 513 we do not have any ground truth data, we will make the following assumption:

514 **Assumption 1.** The “true” NDVI value at time  $t$  can be successfully estimated by robustified  
 515 LOOCV interpolation using high-quality observations. That is, the interpolated value  
 516 (using a robustified interpolation method from chapter 3) considering the points  $P^{SCL45} \setminus$   
 517  $P_t$ . In the following, we will call this estimate the “true”-NDVI.

518 We would like to get an idea if there is any information that can be recovered from non-  
 519 SCL45 observations. For that, we will check for the other SCL-classes if there is a relation  
 520 between the “true” NDVI (derived with robustified SS) and the observed NDVI. Thus, we  
 521 pair each “true” NDVI with its observed one, collect all pairs, and create a scatter plot  
 522 for each SCL-class in fig 4.2. As expected, the “true” and the observed NDVI seem to be  
 523 highly correlated for SCL45. But we can also detect some patterns of correlation in the  
 524 SCL-classes 2, 3, 7, 8 and 10.

525 It might be tempting to just include some of the mentioned SCL classes for interpolation.  
 526 But on the one hand, the choice would not be objective and on the other hand, the  
 527 correlation seems to be weaker than for SCL45. Therefore, in the following section, we  
 528 will correct the observed NDVI and estimate the uncertainty of each correction.

## 529 4.2 Correction Models

530 For training an NDVI correction model, we require ground-truth data which we will aim to  
 531 model using informative covariates. Since ground-truth NDVI data is not available, we will  
 532 again use the assumption 1 and use the “true” NDVI instead. There is no canonical answer  
 533 to the question of which covariates we should use. It is a tradeoff between simplicity,  
 534 generalizability and performance (with the danger of overfitting). Our desire with the  
 535 NDVI correction is to develop a product that is simple to use and understand. Therefore,  
 536 in the subsequent, we will only take the spectral data of the satellite (i.e. all the bands)  
 537 and the observed NDVI derived from it as covariates. We organize the chosen covariates  
 538 in the design matrix  $X^1$ , where each row corresponds to a  $P_t$  (i.e., a pixel at a time  $t$ ) and  
 539 each column to one covariate.

540 In the following, we will introduce different approaches, to model the relationship between  
 541 the response  $y := \text{NDVI}^{\text{true}} \in \mathbb{R}^n$  and the design matrix  $X \in \mathbb{R}^{n \times p}$ . First, we will  
 542 study the basic OLS. Second, we look at the LASSO, an penalized adaptation of the  
 543 OLS which is known to successfully deal with highly correlated covariates. Afterwards,  
 544 GAMs are introduced which model the response similar to OLS but allow for non-linear  
 545 relations. Last but not least, we discuss RF and MARS, which are both flexible modelling  
 546 approaches.

547 Note that in order to reduce computation time, only 10% of the data has been used to fit  
 548 the subsequent models, which are still more than 120'000 observations.

### 549 4.2.1 Ordinary Least Squares (OLS)

550 The OLS is a linear model which aims to minimize the sum of the squared residuals. We  
 551 assume a linear relationship between  $y$  and  $X$  and allow for Gaussian noise. That is:

$$y = X\beta + \epsilon \quad \text{where } \epsilon \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

---

<sup>1</sup>Strictly speaking, we include also the intercept and introduce one dummy variable for each SCL-class

552 Assuming that  $(X^T X)$  is regular, we can estimate the regression coefficients  $\beta$  by

$$\hat{\beta} = (X^T X)^{-1} X^T y = \arg \min_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2$$

553 We will train two models, one using all covariates discussed above and one using only the  
554 SCL-classes and the observed NDVI.

Advantages	Disadvantages
<ul style="list-style-type: none"> <li>— Simple method with good interpretability of coefficients.</li> <li>— Computationally cheap.</li> </ul>	<ul style="list-style-type: none"> <li>— Catches only linear relationships.</li> <li>— No integrated variable selection.<sup>2</sup></li> </ul>

555 **4.2.2 Least Absolute Shrinkage and Selection Operator (LASSO)**

556 The LASSO can be similarly expressed than the OLS but adds a penalty to the minimization  
557 problem:

$$\hat{\beta}_\lambda = \arg \min_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 = \arg \min_{\beta \in \mathbb{R}^p \text{ and } \|\beta\|_1 < \lambda} \|y - X\beta\|_2^2. \quad (4.2.2.1)$$

558 Even though we do not have a closed form solution for equation (4.2.2.1) we can solve  
559 it easily via optimization, since the function  $\beta \in \{\beta \in \mathbb{R}^p | \|\beta\|_1 < \lambda\} \mapsto \|y - X\beta\|_2^2$  is  
560 continuous and convex.

561 Tibshirani (2011) shows that the LASSO solution tends to be sparse. That is  $\beta_i = 0$  for  
562 most  $i = 1, \dots, p$ . The larger  $\lambda$ , the more  $\beta_i = 0$  and hence the simpler the resulting  
563 model.

564 In order to know which  $\lambda$  to choose, we try a huge range of possible values. For each  
565  $\beta_\lambda$ , we calculate the cross-validated  $RMSE_\lambda$ <sup>4</sup> (and its standard deviation  $\sigma_\lambda$  using the  $k$   
566 folds) and define the  $\lambda$  with the smallest corresponding  $RMSE_\lambda$  as  $\lambda_{min}$ . From here we  
567 choose the largest  $\lambda$  for which the  $RMSE_\lambda$  is smaller than  $RMSE_{\lambda_{min}} + \sigma_\lambda$ . This yields  
568 a simpler model while keeping the  $RMSE$  reasonable model.

569 We will apply the LASSO using the selected covariates in section 4.2 and their second  
570 degree of interactions.<sup>5</sup>

Advantages	Disadvantages
<ul style="list-style-type: none"> <li>— Usually yields a sparse solution. This tends to give better generalizability (prediction performance on unseen data).</li> <li>— Successfully deals with correlation in covariates.</li> <li>— Interpretable results.</li> </ul>	<ul style="list-style-type: none"> <li>— Estimate is biased.</li> <li>— Computationally expensive.</li> </ul>

<sup>3</sup>The last two terms are equivalent by lagrangian optimization

<sup>4</sup>The cross validated Root Mean Square Error is the mean of the RMSE's obtained for each fold using the model trained on the remaining folds.

<sup>5</sup>This is if our covariates are  $\{1, a, b\}$ , then we will now use  $\{1, a, b, ab, a^2, b^2\}$ .

571 **4.2.3 General Additive Model (GAM)**

572 GAMs as described in [Hastie and Tibshirani \(1987\)](#) are a special case of Projection Pursuit  
 573 Regression, where only the  $p$  directions parallel to the coordinate axes are considered. The  
 574 result is different to a linear model since the coordinate functions are not restricted to be  
 575 linear but are assumed to be non-parametric functions. The model can be written as:

$$g_{add}(x) = \mu + \sum_{i=1}^p g_j(x_j).^6$$

576 To estimate the non-parametric functions, we can use SS (ref sec. [3.3.6](#)). For this let  $\mathcal{S}_j$   
 577 be the function which takes some  $z \in \mathbb{R}^n$  and returns the SS fitted to  $(X_{:,j}, z)$  where the  
 578 smoothing parameter is optimized by LOOCV<sup>7</sup>. Since we cannot fit all  $g_j$  simultaneously,  
 579 we will use a strategy named Backfitting. We basically cycle through the indices  $1, \dots, p$   
 580 and refit  $\hat{g}_j$  each time. The following illustrates the procedure:

- 1)  $\hat{g}_1 = \mathcal{S}_1(y - \mu)$
  - 2)  $\hat{g}_j = \mathcal{S}_j(y - \mu - \hat{g}_1(X_{:,1}) - \dots - \hat{g}_{j-1}(X_{:,j-1}))$  for  $j = 2, \dots, p$
  - 3)  $\hat{g}_1 = \mathcal{S}_1(y - \mu - \hat{g}_2(X_{:,2}) - \dots - \hat{g}_p(X_{:,p}))$
  - 4)  $\hat{g}_j = \mathcal{S}_j(y - \mu - \sum_{k \neq j} \hat{g}_k(X_{:,k}))$  for  $j = 2, \dots, p$
- $\vdots$

581 We repeat step 3) and 4) until the change falls below some tolerance.

Advantages	Disadvantages
— Captures non-linearity.	— No automatic variable selection.
— Good interpretability.	— Computationally expensive.

582 **4.2.4 Random Forest (RF)**

583 To define a random Forest introduced by [Breiman \(2001\)](#) we will first define what a Tree  
 584 is. A (*decision*) *Tree* is a graph  $(V, E)$  without circles, a distinct root node, every node  
 585 has at most two children and every leaf has a value assigned to it. At each node there  
 586 is a boolean condition testing if one variable is greater than some value and a pointer to  
 587 one child depending on the boolean value. To evaluate a tree we start at the root node,  
 588 test the boolean expression and go to the node indicated by the resulting pointer. This  
 589 we repeat until we end up at a leaf-node, where we return the value assigned to it.

590 To build such a Tree, we will recursively partition the covariate space using greedy splits<sup>8</sup>  
 591 decreasing the RMSE<sup>9</sup> each time. If the set we want to split contains less than a certain  
 592 amount of training points, we stop.

<sup>6</sup>where  $g_j$  is a real-valued function. For identifiability we also demand  $\mathbb{E}[g_j(X_{:,j})] = 0$  for  $j = 1, \dots, p$ .

<sup>7</sup>For efficiency an proxy of the LOOCV is used called generalized cross validation.

<sup>8</sup>For computational reasons, we will only use splits along one covariate. So we ‘cut’ our covariate space into rectangles.

<sup>9</sup>To calculate the RMSE, we need a prediction. Let  $P$  be the current partition, then the predicted value for some  $x \in A \in P$  is the mean of the responses of all the points in  $A$  (included in the training data).

593 To build a Random Forest we will bootstrap-aggregate<sup>10</sup> many such Trees<sup>11</sup>. The prediction  
 594 of the Random Forest for a new point  $x$  is then the mean of the predictions from all  
 595 the Trees.

Advantages	Disadvantages
— Captures non-linear relationships.	— The resulting (prediction) function is not continuous but locally constant.
— Captures all interactions and performs automatic variable selection.	— Computationally expensive.
— Can deal with missing data.	— No interpretability.

596 **4.2.5 Multivariate Adaptive Regression Splines (MARS)**

597 A MARS model as introduced in [Friedman \(1991\)](#) can be described by

$$g(x) = \sum_{m=0}^M \beta_m h_m(x),$$

598 where the  $h_m$  are simple functions (explained later) and the  $\beta_m$  are estimated via Least  
 599 Squares.

600 In the building procedure of a MARS model, we first select many of those simple functions  
 601 and later drop some of them to avoid overfitting. For the construction of those simple  
 602 functions, define  $\mathcal{B}$  be the set of pairs of ‘hockystick functions’

$$\mathcal{B} := \left\{ (b_1, b_2) \mid (b_1(x), b_2(x)) = ((x_j - d)_+, (d - x_j)_+), d = X_{1,j}, \dots, X_{n,j}, j = 1, \dots, p \right\}$$

603 and the set  $\mathcal{M} = \{1\}$  of all functions currently in the model. Now, consider  $\mathcal{C}$  the set of  
 604 candidate functions-pairs

$$\mathcal{C} := \{(h(\cdot)b_1(\cdot), h(\cdot)b_2(\cdot)) \mid h \in \mathcal{M}, (b_1, b_2) \in \mathcal{B}\} \quad (4.2.5.1)$$

605 and select the pair (which when added to  $\mathcal{M}$  and the coefficients refitted) reduces the  
 606 RMSE the most. Add the selected pair to  $\mathcal{M}$  and repeat until the RMSE reduction  
 607 becomes insignificant.

608 Finally, to avoid overfitting, we prune the set  $\mathcal{M}$  by optimizing a LOOCV score.<sup>12</sup>

609 To reduce computational complexity, we follow the recommendation from [Stephen \(2021\)](#)  
 610 and restrict  $h$  in equation (4.2.5.1) to be of degree one (so it is also in a pair of  $\mathcal{B}$ ).  
 611 Consequently,  $\mathcal{C}$  contains functions with a degree of at most 2.

<sup>10</sup>That is we will sample (with replacement) several times  $n$  observations from our original data and fit a Tree to each such sample.

<sup>11</sup>Building the Tree, this time we will not test every covariate at each node (for the RMSE minimization) but a node-specific subsample of the covariates. Thus, also the “second best split” can be selected.

<sup>12</sup>This means that we perform an iterative procedure to reduce the number of functions in  $\mathcal{M}$ . For every function  $h$  in  $\mathcal{M}$ , we compute the model using  $\mathcal{M} \setminus \{h\}$ . We discard the function which – when excluding from  $\mathcal{M}$  – leads to the best LOOCV score.

Advantages	Disadvantages
<ul style="list-style-type: none"> <li>— Catches non-linear relationships.</li> <li>— Interpretability via functions in <math>\mathcal{M}</math> and their coefficients.</li> <li>— Allows for interactions with variable selection.</li> </ul>	<ul style="list-style-type: none"> <li>— Computationally expensive (can be reduced by restricting the degree of interactions).</li> </ul>

## 612 4.3 Uncertainty Estimation

613 Once we corrected the NDVI using the models described in the previous section, we are left  
 614 with the problem that not every correction is equally reliable.<sup>13</sup> Hence, we are interested  
 615 in a measure of how uncertain an estimate is.

616 We achieve this analogously as we corrected the NDVI, by replacing the response (NDVI<sup>“true”</sup>)  
 617 with the absolute residuals  $v := |y - \hat{y}|$  and modeling their relationship with the covariates  
 618 defined by  $X$ . In this way, we obtain a model for the absolute residuals  $v$  and the estimator  
 619  $\hat{v}$ .

## 620 4.4 Interpolation

621 Consider now a pixel  $P$ ,  $\hat{y}^{(P)}$  its corrected NDVI and  $\hat{v}^{(P)}$  the estimated uncertainties of  
 622  $\hat{y}^{(P)}$ . In order to interpolate  $\hat{y}^{(P)}$ , we will give less weight to unreliable observations. Thus,  
 623 we define the weight function:

$$w_\tau^{(P)} := \frac{1}{R} \frac{1}{\hat{v}_\tau^{(P)}}, \quad \text{for } \tau = 1, \dots, n_P \quad (4.4.0.1)$$

624 where  $\tau$  is an index over the satellite images and  $R := \frac{\sum_i^{n_P} \hat{v}_i^{(P)}}{n_P}$  a normalization constant.  
 625 The normalization is needed since for some interpolation methods, inflating the sum of  
 626 weights would decrease the effect of the smoothing.

## 627 4.5 Resulting Interpolation Strategies

628 We have developed the following procedure to obtain a new interpolation (keyword-wise):

- 629 i.) LOOCV Interpolation (+ robustify?) to get “true” NDVI
- 630 ii.) Correction
- 631 iii.) Uncertainty estimation
- 632 iv.) Interpolation (+ robustify?)

633 At each step we have a choice, more precisely:

- 634 — Interpolation: Smoothing Splines / Double Logistic
- 635 — Robustify: Yes / No
- 636 — Correction & uncertainty estimation: RF / OLS – considering only SCL-classes /  
 637 OLS – considering all selected covariates / MARS / GAM / LASSO / no correction.

---

<sup>13</sup>One correction is illustrated in the figure B.4f. In this figure, the outer points (labeled as clouds) have a large scatter.

638 As it is not feasible to try every possible combination, we make the following restrictions  
 639 on which combinations we will consider:

- 640 — We use the same interpolation method each time.  
 641 — Either we robustify both times, or we do not robustify at all.  
 642 — We use the same underlying method for correction and uncertainty estimation.

643 In this fashion, we obtain 28 distinct interpolation strategies, which we will benchmark in  
 644 the next section.

## 645 4.6 Evaluation Method

646 In this section, we introduce the relative yield-estimation-accuracy (RYEA) and utilize it  
 647 to evaluate the 28 interpolation strategies from section 4.5. The fundamental assumption  
 648 is that the closer the interpolated NDVI time series is to the true one, the better it  
 649 can be used to determine crop yield. Implicitly, we believe that an NDVI time series  
 650 which better models yield will incorporate more true information about the underlying  
 651 vegetation. Therefore, we want to determine a comparable RYEA for each interpolation  
 652 strategy and choose it as a benchmark criterion. This is an objective measure, since we  
 653 have not considered crop yield in any of our previous steps. Moreover, this criterion is  
 654 justified by the fact that yield estimation has been a motivation for the interpolation.

655 **Definition 4.6.0.1.** (RYEA) Let  $y \in \mathbb{R}^n$  be the yield,  $M$  be a model for estimating  $y$ , and  
 656  $\hat{y} = M(X)$  where  $X$  describes the data<sup>14</sup>. We define the RYEA as the relative RMSE in  
 657 yield estimation. Formally expressed:

$$\text{RYEA} = \frac{\sqrt{\sum_{i=1}^n (y_i - \hat{y}_i)^2}}{\bar{y}},$$

658 where  $\bar{y}$  denotes the sample mean.

### 659 4.6.1 Yield Estimation

660 For all the pixels, we will interpolate the NDVI time series with every interpolation strat-  
 661 egy. From the interpolated NDVI time series, we would like to estimate the yield. However,  
 662 given the high dimensionality and different lengths of the interpolation (not every time  
 663 series has the same start and end point), we must first map each NDVI time series into a  
 664 low-dimensional vector space of covariates. For this, we will use the following statistics:

- Maximum slope
- Minimum slope
- Integral<sup>15</sup> over all
- Peak (i.e. maximal NDVI)
- GDD for the Peak
- Integral<sup>15</sup> up to the peak
- Integral<sup>15</sup> after peak
- Integral<sup>15</sup> from 0-685 GDD
- Integral<sup>15</sup> from 685-1075 GDD

---

<sup>14</sup>We will use the matrixes derived in section 4.6.1

<sup>15</sup>We will only consider the integral of the function  $\max(0, NDVI - 0.3)$ , where 0.3 is assumed to be a minimal NDVI value. REF

665 For the choice we were inspired by (c.f. table 2 in Kamir, Waldner, and Hochman (2020)).  
666 However, we deliberately omit any statistic that involves the minimum (e.g. the NDVI-  
667 range), since we regard the minimum as a very error-prone measure due to the large  
668 influence of clouds in the time series.

669 As a result, for each interpolation strategy, a matrix is obtained in which each row corre-  
670 sponds to a pixel and both the yield and the covariates (computed by applying the above  
671 statistics) are contained. Using this matrix, we train a random forest for yield estimation,  
672 and compute the integrated OOB estimates<sup>16</sup>  $\hat{y}$ . Note that the choice of the modeling  
673 approach does not matter much, as long as it is general enough (i.e. able to approximate  
674 any function) and we use the same one for each interpolation strategy. Finally, for each  
675 interpolation strategy, we calculate the RYEA and describe the results in section 5.2.

---

<sup>16</sup>By the integrated OOB estimates, we denote the predictions for each pixel where only trees are used, where the pixel has not been used (as  $n_{tree}$ , the number of Trees, grows the fraction of trees which do not contain a certain pixel converges to  $\frac{1}{e}$ ).

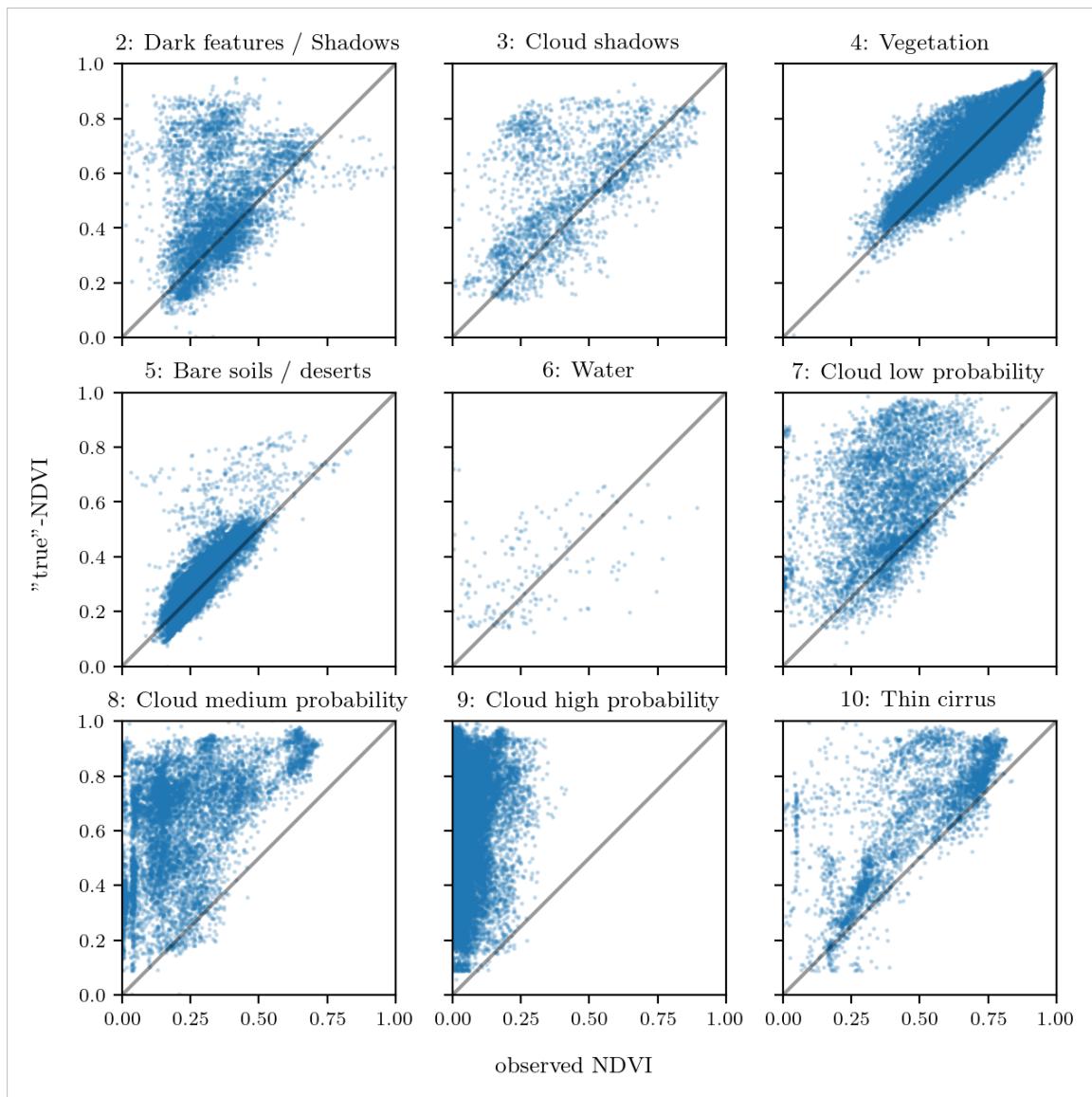


Figure 4.2: For each SCL class, we compare the true NDVI with the observed NDVI. (The true NDVI was estimated with LOOCV smoothing splines, and we used all observations of 10% of the total pixels.)

676 **Chapter 5**

677 **Results**

678 **5.1 Goodness of Fit for Selected Interpolation Methods**

679 Table 5.1 benchmarks the selected<sup>1</sup> interpolation methods (on  $P^{SCL45}$ ) with respect to  
680 various score functions. The score functions take the absolute values of the LOOCV  
681 residuals and summarize them in a number (the smaller, the better). For each of the 5  
682 selected interpolation methods, we consider the basic and the robustified (see section 3.5)  
683 version.

Table 5.1: Comparing the goodness of fit for selected interpolation methods (on  $P^{SCL45}$ ) measured with the score functions (which take the LOOCV residuals as input) listed in the left column.  $q_X$  denotes here the  $X\%$  quantile.

	SS	LOESS	DL	BSPL	FR	$SS^{\text{rob}}$	$\text{LOESS}^{\text{rob}}$	$DL^{\text{rob}}$	$BSPL^{\text{rob}}$	$FR^{\text{rob}}$
RMSE	0.063	0.061	0.061	0.074	0.075	0.070	0.065	0.065	0.079	0.208
qtile50	0.036	0.034	0.027	0.043	0.031	0.032	0.031	0.022	0.037	0.049
qtile75	0.063	0.061	0.051	0.077	0.058	0.061	0.057	0.044	0.070	0.099
qtile85	0.080	0.079	0.070	0.098	0.083	0.081	0.076	0.063	0.094	0.158
qtile90	0.092	0.092	0.088	0.112	0.108	0.097	0.090	0.082	0.113	0.226
qtile95	0.119	0.115	0.122	0.142	0.161	0.132	0.115	0.124	0.157	0.375

684 DL is the best among both robustified and non-robustified with respect to most of the  
685 score functions used (all except q95) and is especially superior to the other parametric  
686 approach, which is FS. Especially the robust FS performs poorly. The LOESS dominates  
687 (i.e. is superior on every score function) all other non-parametric methods, but is closely  
688 followed by the SS. The BSPL, on the other hand, is the worst non-parametric method  
689 tested here.

690 **5.2 XXX (Robustification and) NDVI-Correction**

691 defition of RYEA, it is not an accuracy but an error

692 The RYEA for the 28 (in section 4.5) chosen interpolation strategies is given in table 5.2.  
693 Robustification in the interpolation strategies, does not improve the quality of the fit

<sup>1</sup> For the discussion which methods have been selected c.f. section 6.1.2.

Table 5.2: RYEAs. For the non-relative RMSE and the coefficient of determination ( $R^2$ ) see table B.1 and B.2.

	RF	OLS-SCL	OLS-all	MARS	GAM	LASSO	no-correction
ss	0.155	0.140	0.143	0.142	0.142	0.142	0.149
dl	0.156	0.151	0.152	0.152	0.149	0.149	0.158
ss-rob	0.155	0.143	0.147	0.149	0.146	0.145	0.148
dl-rob	0.157	0.153	0.152	0.145	0.148	0.150	0.157

694 (measured via the RYEAs) in most cases. In addition, SS (rob) are better than DL(rob)  
 695 in terms of RYEAs, with one exception.

696 The interpolation strategy that leads to the lowest RYEAs is the OLS-SCL with SS. Given  
 697 that the OLS-SCL models have very good interpretability, we also present the regression  
 698 equations below. The corrected NDVI is calculated using

$$\begin{aligned} \text{NDVI}_{\text{corr}} = & 0.711 \text{NDVI}_{\text{observed}} + \mathbb{1}_{SCL=2}0.215 + \mathbb{1}_{SCL=3}0.237 + \mathbb{1}_{SCL=4}0.210 \\ & + \mathbb{1}_{SCL=5}0.116 + \mathbb{1}_{SCL=6}0.162 + \mathbb{1}_{SCL=7}0.327 + \mathbb{1}_{SCL=8}0.474 \quad (5.2.0.1) \\ & + \mathbb{1}_{SCL=9}0.575 + \mathbb{1}_{SCL=10}0.306 + \mathbb{1}_{SCL=11}0.512 \end{aligned}$$

699 where  $\mathbb{1}_{SCL=2}$  is equal to one if the current observation corresponds to SCL class 2 and  
 700 zero otherwise.<sup>2</sup>. Whereas, we obtain the estimated absolute residuals by:

$$\begin{aligned} \widehat{\text{abs}}(\text{NDVI}^{\text{"true"}} - \text{NDVI}_{\text{corr}}) = & -0.133 \text{NDVI}_{\text{observed}} + \mathbb{1}_{SCL=2}0.186 + \mathbb{1}_{SCL=3}0.185 \\ & + \mathbb{1}_{SCL=4}0.146 + \mathbb{1}_{SCL=5}0.089 + \mathbb{1}_{SCL=6}0.167 \\ & + \mathbb{1}_{SCL=7}0.203 + \mathbb{1}_{SCL=8}0.181 + \mathbb{1}_{SCL=9}0.173 \\ & + \mathbb{1}_{SCL=10}0.180 + \mathbb{1}_{SCL=11}0.172 \quad (5.2.0.2) \end{aligned}$$

701 In the equation 5.2.0.1, we notice the strongest upwards correction for SCL classes 8, 9 and  
 702 11 (correspond to ‘medium probability clouds’, ‘high probability clouds’ and ‘thin cirrus  
 703 clouds’). The estimated absolute residuals, however, are the smallest for SCL classes 4 and  
 704 5 (‘vegetation’ and ‘bare soil’). Furthermore, the higher the observed NDVI the lower are  
 705 the estimated absolute residuals.

706 For the R-output of the `summary` function of the two models, we refer to the appendix  
 707 B.3.1.

<sup>2</sup>  $\mathbb{1}$  is also called an indicator function or characteristic function in mathematics.

708 **Chapter 6**

709 **Discussion**

710 Here in the discussion, you should take up the points you mentioned in the introduction

711 **6.1 Interpolation Methods**

712 **6.1.1 Data Gaps in Time Series**

713 NW estimates the value for  $t$  by relating to the points near  $t$ . To determine what “near”  
714 means, a bandwidth  $h$  is used (c.f. equation 3.3.1.2). This gets problematic as soon as the  
715 data gaps become larger than  $h$ , since in this case no points are left that are considered  
716 to be close to  $t$ .

717 Regarding the GK, we expect that because of the stationarity assumption, the interpolation  
718 will tend to the mean if data gaps are present (c.f. figure 3.3).

719 Since the SG requires equidistant points, it is clear that data gaps will break it. The linear  $\text{F}_{\text{wertend}}$   
720 interpolation, which is supposed to recover this, we consider as not being a satisfying  
721 solution.

722 We do not trust the FR interpolation if there are noticeable data gaps. On the one hand,  
723 it corresponds to our experience that the curve can escape strongly there (c.f. figure  
724 3.1). On the other hand, the unreliability is illustrated by the poor values in table 5.1 for  
725 the robustified variant. These are meaningful in describing the ability to cope with data  
726 gaps, since more data points are ignored during the robustification and thus data gaps are  
727 simulated.

728 Similarly, for SS, LOESS, DL and BS we compare the values in table 5.1 between the  
729 robustified and non-robust variant. We find that the robust variant is not very different  
730 from the non-robust variant (unlike FR). Thus, we conclude that these methods do not  
731 have systematic failures.

732 Regarding the LOESS, we observe in the figure B.1 in plot (c) a strange peak between  
733 the first and second observation. This peak is due to the local weighting. In case of data  
734 gaps, the weights can attain non-intuitive values. For instance, the first data point in the  
735 plot, although adjacent to the peak, is given a low weight compared to the points to the  
736 right of the peak (for estimating the value at this peak).

737 In our experience, the DL handles data gaps well, but it may happen that the model  
 738 describes the NDVI increase as abrupt. This however was fixed, by bounding the first  
 739 derivative (c.f. section 3.2.3).

### 740 6.1.2 Preselection

741 We shall now justify our preselection of the interpolation methods tested in section 3.6.  
 742 We decided against NW because it has systematic errors at peaks and valleys. Moreover,  
 743 this method handles data gaps poorly (c.f. 6.1.1). Moreover, we will not consider UK since  
 744 the underlying assumptions are not met and therefore a systematic bias is introduced. On  
 745 top of that, ML parameter finding occasionally fails. Also, we do not include the SG in  
 746 the next selection, since we think of it as a special case of LOESS.

### 747 6.1.3 Candidate Selection

748 Given that DL convinces regarding most of the selected score functions in table 5.1 we will  
 749 certainly investigate this method in chapter 4. Moreover, we see that the robustification  
 750 mostly improved the score regarding the 50, 75, 85, and 90 % Quantiles. Only for the  
 751 outlier-sensitive score functions (RMSE and q95)<sup>1</sup> we notice significant worsening (we  
 752 consider the robust FS separately in section 6.1.1). Consequently, we will also use the  
 753 robustification in section 4. Not wanting to rely on the form assumptions of the DL, we  
 754 further choose a non-parametric method for further consideration. Despite the LOESS  
 755 slightly dominating the SS in table 5.1, we choose the SS. This is due to the strange  
 756 behavior of the LOESS in case of data gaps (see section 6.1.1) and the good interpretability  
 757 of the SS using the minimization function 3.3.6.1.

758 XXX discuss results from table B.1

## 759 6.2 NDVI Correction

### 760 6.2.1 Bootstrap

761 The question arises if we can build the correction model on the same year as we want to  
 762 apply it on. Usually, a similar approach might carry the danger of overfitting. However, we  
 763 have not used any ground truth at any point (until the evaluation). Instead, we estimated  
 764 the “true” NDVI with the assumption 1 via OOB. Thus, we have bootstrapped our way  
 765 out of the problem. Consequently, we reason that we can apply our method to a new  
 766 (comparable) dataset and solve the correction again via this bootstrap.

### 768 6.2.2 Using Additional Covariates

769 In section 4.2 we have only used the spectral data (and the observational NDVI calculated  
 770 from them) as covariates. Since we have the weather data available (c.f. REF-SEC), it  
 771 would be a small effort to incorporate it, together with statistics collected from it (i.e.  
 772 GDD or ‘rainfall in the last 30 days’).

773 We decided against using this data, because on the one hand we have the problem that  
 774 we have practically too few observations (we observe only 5 years) and we expect the  
 775 weather in our study region to be rather homogeneous which is suggested by the fact

where  
does  
this sec-  
tion be-  
long to?  
Chapter  
‘NDVI  
Correc-  
tion’ or  
‘Further  
Work’?

<sup>1</sup>For the RMSE one outlier is enough to take away the usefulness of the statics, in the case of q95 it is enough if 5% of the data are corrupt to break the statics.

776 that the weather data published by Meteoswiss are for a grid with a resolution of 1 km.  
 777 On the other hand, we want the underlying model not to learn improper relationships.  
 778 For example, the model might automatically predict a high NDVI for a day in summer  
 779 (detected by high GDD / many sunshine hours / high temperature) just because it is  
 780 “used” to observing a lot of vegetation in summer. Including temporally (e.g.,  $P_{t-1}$  and  
 781  $P_{t+1}$ ) and geographically adjacent pixels would likely improve performance. However, for  
 782 simplicity, we omit it here<sup>2</sup>.

### 783 6.2.3 Which Interpolation Strategy should we choose

784 table mit OLS SCL als sieger diskutieren

785 if we use no-correctionXss-rob instead of OLS-SCLXss we loose  $(0.148 - 0.14)/0.148 =$   
 786 5,4% of the information.

### 788 6.2.4 High RMSE in Yield Prediction

789 How much can we expect to get? We have multiple sources of uncertainty in the data:

- 790 i.) Uncertainty in Yield data collected by the combine harvester
- 791 ii.) Uncertainty in Yield data through rasterization
- 792 iii.) Uncertainty in satellite images through “measurement errors” introduced via clouds  
 793 and other atmospheric effects
- 794 iv.) Uncertainty introduced by interpolating (especially when long data-gaps are present)

795 even in a perfect world the NDVI curve only holds a fraction of the information  
 available

796 You already capture the ”main” structure of your thesis with the interpolation and the  
 NDVi correction sections. Can you combine them both in a ”synthesis” subsection at  
 the end of the discussion?

kurzer  
 kontext  
 von  
 vergle-  
 ichbaren  
 values  
 von  
 gregor  
 — diese  
 sektion  
 ist für  
 dena uf-  
 traggeber

---

<sup>2</sup>This is done for simplicity of understanding and using the model, since one would need to adapt to some convention of how to supply the data of adjacent pixels without redundancy (i.e. supplying  $P_t$  multiple times).

797 **Chapter 7**

798 **Conclusion**

799 In dieser Thesis haben wir studiert, wie wir mit aus Satellitenbildern das Pflanzenwach-  
800 stum via NDVI-Zeitreihen modellieren können. Die grösste Herausforderungen waren hi-  
801 erbei die fragen, wie man mit (durch Wolken oder Schatten) verfälschten Beobachtungen  
802 umgehen soll und wie man die einzelnen Beobachtungen zu interpolieren hat. Für eine  
803 zusammenfassung der betrachteten interpolationsmethoden verweisen wir auf die Tabelle  
804 **3.**

805 Durch Wolken und Schatten manipulierte Beobachtungen führen dazu, dass wir fehlerhafte  
806 NDVI werte erhalten. Zwar können wir diese bis zu einem gewissen grad filtern, haben aber  
807 trotzdem noch fehlerhafte Beobachtungen. Um mit diesen Ausreißern umzugehen haben  
808 wir eine Technik verallgemeinert, welche die Interpolation robuster gegen Ausreisser en-  
809 twickelt macht. Durch die Filtration von fehlerhaften Beobachtungen, erhalten wir beson-  
810 ders im Winter Datenlücken. Daher ist es ein Kriterium für unsere gewählten interpo-  
811 lationsmethode, dass sie gut mit solchen Datenlücken umgehen können. Der Nadaraya-  
812 Watson kernel schätzer, Universal Kriging, 2cd order Fourier Series und Savitzky-Golay  
813 Filter konnten hier nicht überzeugen (vgl. sektion [6.1.1](#)). Vereinzelt hat hier auch eine  
814 Generalisierung des Savitzky-Golay Filters – der LOESS — überraschendes verhalten  
815 aufgezeigt. Dieser konnte hingegen bei der Leave-One-Out-Cross-Validation (LOOCV)  
816 überzeugen (c.f. table [5.1](#)), jedoch bevorzugen wir die Smoothing Splines (SS), da sie  
817 dort nur wenig schlechter abscheiden, aber eine deutlich glattere kurve produzieren (vgl.  
818 Abbildung [3.5](#) und [B.1](#)). Die SS approximieren flexibel die Daten, halten aber gleichzeitig  
819 die Krümmung gering (c.f. equation [3.3.6.1](#)). B-Splines hingegen waren hinsichtlich jeder  
820 getesteten Score Funktion schlechter als SS und ihr smoothing Mechanismus ist auch  
821 schlechter Interpretierbar. Am besten schneiden hier jedoch die Approximation durch eine  
822 Double logistic (DL) ab, welche starke annahmen über die Form der NDVI kurve macht.  
823 Probleme für die Parameterschätzung des DL (und der Fourierreihe) haben wir behoben,  
824 indem wir den parameterraum durch großzügige aber realistische werte beschränkt haben.  
825 Probleme mit overfitting beim Universal Kriging haben wir behoben, indem wir die pa-  
826 rameter für ein subsample an NDVI zeitreihen bestimmt haben und schlussendlich den  
827 median jeweiliger parameter benutzt haben. Schlussendlich wählen wir DL und SS als  
828 unsere Favoriten der Interpolationsmethoden.

829 Frage: mehr details für die begründung der Interpolations-kandidaten?

830 Auf die Frage, wie wir mit den verfälschten Beobachtungen umgehen sollen, lautet die

831 traditionelle Antwort, dass wir nur Beobachtungen beachten, welche als Vegetation oder  
832 als bare soil gelabelt sind (SCL45). Dies wird mit der von der European Space Agency  
833 gelieferten ‘Scene Classification Layer’ (SCL) bewerkstelligt. In figure 2.3 wird jedoch die  
834 Unzuverlässigkeit dieses Labelings illustriert. Zudem haben wir die festgestellt, dass auch  
835 nicht-SCL45 Beobachtungen wertvolle Informationen enthalten seien können (vgl. Sektion  
836 4.1). Wir haben uns entschieden, nicht an der traditionellen (SCL-)Filtration festzuhalten.  
837 Stattdessen betrachten wir alle Beobachtungen und korrigieren den beobachteten NDVI.  
838 Dafür benutzen wir statistische Modelle, die zusätzliche Informationen wie die verbleibenden  
839 Spektralbänder in Betracht nehmen. Bevor wir aber die korrigierten NDVI Werte  
840 interpolieren, weisen wir jeder Beobachtung ein Gewicht zu, korrespondieren zu ihrer Un-  
841 sicherheit. Die Unsicherheit wird analog wie die NDVI Korrektur geschätzt. Durch die  
842 Wahl verschiedener Interpolationsmethoden (mit und ohne robustifizierung) und statis-  
843 tischer Modelle, erhalten wir somit 28 verschiedene Interpolationsstrategien (vgl. Sektion  
844 4.5). Um zu beurteilen, welche dieser Interpolationsstrategie am besten ist, machen wir  
845 die folgende Annahme: “je besser die Interpolationsstrategie, desto besser kann damit in-  
846 terpolierte NDVI Zeitreihe den Ertrag voraussagen”. Überraschender Weise ist die beste  
847 Strategie, die mit nicht-robustifizierten SS und dem einfachsten betrachteten statischen  
848 Modell, welches nur den beobachteten NDVI und die SCL Klassifizierung benutzt. Let  
849 us recapitulate the best interpolation strategy: First, we estimate the “true” NDVI using  
850 SS via LOOCV, then obtain the corrected NDVI using the OLS-SCL model (c.f. equa-  
851 tion 5.2.0.1). Subsequently, we estimate the absolute error with the OLS-SCL model (c.f.  
852 equation 5.2.0.1) and thereby obtain weights that are supposed to reflect the reliability of  
853 the corrected NDVI (c.f. equation 4.4.0.1). Finally, we perform a weighted interpolation  
854 with SS.

855 Zwar ist die die robustifizierung nicht teil der besten Interpolationsstrategie, verfehlt dieses  
856 Ziel aber nur knapp. Hingegen sehen wir in tabelle 5.1, dass die robustifizierung in den  
857 meisten Fällen zu kleineren LOOCV Residuen führt (mit ausnahme von der Fourier Ap-  
858 proxmiation). Daher empfehlen wir die robustifizierung durchzuführen, wenn wir mit  
859 Fehlerhaften beobachtungen rechnen.

860 Auf die Frage welche interpolationsmethode wir schlussendlich empfehlen, wollen wir zwei  
861 Fälle betrachten. Wenn es nur darum geht möglichst präzise eine Kurve den daten anzu-  
862 passen, empfehlen wir die robustifizierten DL, da diese die LOOCV residuals in den meis-  
863 ten fällen minimieren (vgl. tabelle 5.1). Falls wir eine interpolation erhalten wollen  
864 die möglichst viele informationen über die pflanze enthält empfehlen wir die SS. Diese  
865 empfehlung gilt besonders, falls wir traditionell nur SCL45 beobachtungen betrachten  
866 wollen ohne die vorgeschlagene NDVI zu korrigieren. Jedoch empfehlen wir die oben  
867 aufgeführte interpolationsstrategie, da uns ansonsten über 5% der informationen aus der  
868 NDVI zeitreihe abhanden kommmen (vgl. sektion 6.2.3). Im anbetracht aller Fehlerquellen  
869 (c.f. section 6.2.4) und der tatsache dass wir nur die NDVI Zeitreihe betrachten wir die  
870 5% als eine solide verbesserung.

Anzahl von Beobachtungen, empfehlungen? – schwierig, weil regelmäßigkeit in ‘wichtigen’ zeiträumen (der veränderung) wichtiger ist.

872 

## 7.1 Future Work

873 

### 7.1.1 Time Series Correction-Interpolation as a General Method

874 Throughout this thesis, we developed a correction and interpolation method for the NDVI.  
875 However, we never used features of the NDVI. Only the parameter estimated via cross-  
876 validation in chapter 3.4 depends on the scale of the time series. For simplicity, we could  
877 thus determine the parameter using Generalized Cross Validation (as Ripley and Maechler  
878 suggest). Therefore, our approach of interpolation and correction of time series can be  
879 applied to arbitrary time series as long as additional information is available. However,  
880 further research is required, to demonstrate the general usefulness of this approach.

881 

#### Example: Cloud Correction with Uncertainty Estimation and Interpolation

882 This generalization can be used in particular for cloud correction. In the same manner as  
883 we corrected the NDVI time series in chapter 4, we can correct each spectral band and  
884 reunite the corrected bands with the uncertainties. If desired, the time series can also be  
885 interpolated before merging as in chapter 4.4. The resulting question would be how well  
886 this approach performs.

887 

### 7.1.2 Minor Improvements

888 During this project, we also noticed some minor issues that we would have liked to investi-  
889 giate further if more resources were available. The most relevant of these are:

- 890 — **Data:** Method how combine harvester point data has been extrapolated to the grid  
891 could possibly be improved.
- 892 — **Data:** For computational reasons, we mostly considered all years and split the data  
893 (on the pixel level) randomly into a train/test set. A leave one year out cross  
894 validation might yield more accurate results.
- 895 — **Data:** We have not included the spectral bands which have a resolution of 60 m. But  
896 precisely these seem to be promising for cloud correction, since they are a proxy of  
897 the water (content and form) in the atmosphere.
- 898 — **Data:** Raiyani, Gonçalves, Rato, Salgueiro, and Marques da Silva (2021) presents  
899 an Machine Learing approach that supposedly improves the SCL and thus could  
900 improve our results which are based on the SCL.
- 901 — **NDVI Correction:** Explore the effect of different link and normalizing functions in  
902 section 4.4. Currently we run into the danger of some outer points getting nearly  
903 ignored just because one estimated absolute residual for some interior point is very  
904 small.
- 905 — **NDVI Correction:** Yield is not the only target variable of interest. Other variables  
906 like protein content could also be used in section 4.6 for the method evaluation.

# 907 Bibliography

- 908 (2007). Gaussian models for geostatistical data. In P. J. Diggle and P. J. Ribeiro (Eds.),  
909 *Model-Based Geostatistics*, pp. 46–78. New York, NY: Springer.
- 910 Bailey, S. J. (2018, July). Using Growing Degree Days to Predict Plant Stages. pp. 8.
- 911 Beck, P. S. A., C. Atzberger, K. A. Høgda, B. Johansen, and A. K. Skidmore (2006,  
912 February). Improved monitoring of vegetation dynamics at very high latitudes: A new  
913 method using MODIS NDVI. *Remote Sensing of Environment* 100(3), 321–334.
- 914 Breiman, L. (2001, October). Random Forests. *Machine Learning* 45(1), 5–32.
- 915 Brockmann, M., T. Gasser, and E. Herrmann (1993, December). Locally Adaptive Band-  
916 width Choice for Kernel Regression Estimators. *Journal of the American Statistical  
917 Association* 88(424), 1302–1309.
- 918 Cao, R., Y. Chen, M. Shen, J. Chen, J. Zhou, C. Wang, and W. Yang (2018, November). A simple method to improve the quality of NDVI time-series data by integrating  
919 spatiotemporal information with the Savitzky-Golay filter. *Remote Sensing of Environ-  
920 ment* 217, 244–257.
- 921 Chen, J., P. Jönsson, M. Tamura, Z. Gu, B. Matsushita, and L. Eklundh (2004, June). A simple method for reconstructing a high-quality NDVI time-series data set based on the  
922 Savitzky–Golay filter. *Remote Sensing of Environment* 91(3), 332–344.
- 923 Cleveland, W. S. (1979, December). Robust Locally Weighted Regression and Smoothing  
924 Scatterplots. *Journal of the American Statistical Association* 74(368), 829–836.
- 925 Friedman, J. H. (1991, March). Multivariate Adaptive Regression Splines. *The Annals of  
926 Statistics* 19(1), 1–67.
- 927 Hastie, T. and R. Tibshirani (1987, June). Generalized Additive Models: Some Applica-  
928 tions. *Journal of the American Statistical Association* 82(398), 371–386.
- 929 Jaramaz, D., V. Perović, S. Belanovic Simic, E. Salnikov, D. Cakmak, V. Mrvić, and  
930 L. Zivotic (2013, May). The ESA Sentinel-2 mission Vegetation variables for Remote  
931 sensing of Plant monitoring.
- 932 Kamir, E., F. Waldner, and Z. Hochman (2020, February). Estimating wheat yields  
933 in Australia using climate records, satellite image time series and machine learning  
934 methods. *ISPRS Journal of Photogrammetry and Remote Sensing* 160, 124–135.
- 935 Lyche, T. and K. Mørken (2005, January). Spline Methods.
- 936 McMaster, G. S. and W. W. Wilhelm (1997, December). Growing degree-days: One  
937 equation, two interpretations. *Agricultural and Forest Meteorology* 87(4), 291–300.

- 940 Perich, G., M. O. Turkoglu, L. V. Graf, J. D. Wegner, H. Aasen, A. Walter, and F. Liebisch  
941 (2022, July). Pixel-based yield mapping and prediction from Sentinel-2 using spectral  
942 indices and neural networks.
- 943 Raiyani, K., T. Gonçalves, L. Rato, P. Salgueiro, and J. R. Marques da Silva (2021,  
944 January). Sentinel-2 Image Scene Classification: A Comparison between Sen2Cor and  
945 a Machine Learning Approach. *Remote Sensing* 13(2), 300.
- 946 Ripley, B. D. and M. Maechler. R: Fit a Smoothing Spline. [https://stat.ethz.ch/R-  
947 manual/R-patched/library/stats/html/smooth.spline.html](https://stat.ethz.ch/R-manual/R-patched/library/stats/html/smooth.spline.html).
- 948 Rouse, J. W. (1974, May). Monitoring the vernal advancement and retrogradation (green  
949 wave effect) of natural vegetation. Technical Report NASA-CR-139243.
- 950 Savitzky, A. and M. J. E. Golay (1964, July). Smoothing and Differentiation of Data by  
951 Simplified Least Squares Procedures. *Analytical Chemistry* 36(8), 1627–1639.
- 952 Schafer, R. W. (2011, July). What Is a Savitzky-Golay Filter? [Lecture Notes]. *IEEE  
953 Signal Processing Magazine* 28(4), 111–117.
- 954 Skakun, S., E. Vermote, B. Franch, J.-C. Roger, N. Kussul, J. Ju, and J. Masek (2019,  
955 July). Winter Wheat Yield Assessment from Landsat 8 and Sentinel-2 Data: Incorporating  
956 Surface Reflectance, Through Phenological Fitting, into Regression Yield Models.  
957 *Remote Sensing* 11(15), 1768.
- 958 Stephen, M. (2021, July). Earth: Multivariate Adaptive Regression Splines.
- 959 Tibshirani, R. (2011). Regression shrinkage and selection via the lasso: A retrospective.  
960 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 73(3), 273–  
961 282.

962 **Appendix A**

963 **Reproducibility**

964 **A.1 Reproduce Results**

965 For reproducibility of the whole computations, we refer to our codebase at:

966 <https://github.com/LGraz/MasterThesis-Code>

967 In order to reproduce our computations and results, set up the directory as described  
968 in the README and execute the computations via `./shell_scripts/reproduce.sh`  
969 and do not execute the python and R scripts by hand (unless you follow the order in  
970 `./shell_scripts/reproduce.sh`).

971 **A.2 R-Package**

972 We also provide an R package for a general time series correction and interpolation if  
973 additional data is available at:

974 <https://github.com/LGraz/CorrectTimeSeries>

975 In our case we consider the NDVI time series and the additional data consists of the unused  
976 spectral bands.

977 We recommend installing it via the `devtools` package by:

978 `devtools::install_github("LGraz/CorrectTimeSeries")`

979 In the following, we shall give a stand-alone example of how the R package can be used:

```
980
981 1 library(CorrectTimeSeries)
982 2
983 3 # load a list of dataframes, each one describes one pixel with the covariates and
984 4 # the response
985 5 data(timeseries_list)
986 6 str(timeseries_list[[1]])
987 7
988 8 # Train/Load RF
989 9 train_model_myself <- TRUE
990 10 if (train_model_myself){
991 11   # Add "true" NDVI (or generally the response), by Out-Of-Bag estimation
992 12   timeseries_list <- lapply(timeseries_list, function(df) {
993 13     df$oob_ndvi <- OOB_est(df$gdd, df$ndvi_observed) # gdd is the time-axis
994 14     df
995 15   })
996 16   # Train correction model
997 17   formula <- "oob_ndvi ~ B02+B03+B04+B05+B06+B07+B08+B8A+B11+B12+scl_class"
998 18   RF <- train_RF_with_fromula(formula, timeseries_list, robustify=TRUE)
999 } else {
```

```
1000 19  data(RF_for_NDVI)
1001 20  RF <- RF_for_NDVI
1002 21 }
1003 22
1004 23 # ADD CORRECTION
1005 24 timeseries_list <- lapply(timeseries_list, function(df) {
1006 25   df$corrected_ndvi <- randomForest:::predict.randomForest(RF, df)
1007 26   df
1008 27 })
1009 28
1010 29 # Get interpolation for each timeseries
1011 30 newx <- 1:1000
1012 31 lapply(timeseries_list, function(df){
1013 32   ss <- smoothing_spline(df$gdd, df$corrected_ndvi)
1014 33   predict(ss, newx)$y
1015 34 })
```

Example of how to use the `CorrectTimeSeries` package

1017 **Appendix B**

1018 **Further Material**

1019 **B.1 Data and Methods**

1020 **B.1.1 GDD**

1021 Bailey (2018) tabulates the corresponding GDD for each stage of wheat.

Stage	Description	GDD
Emergence	Leaf tip just emerging from above-ground coleoptile.	125 – 160
Leaf development	Two leaves unfolded.	169 – 208
Tillering	First tiller visible	369 – 421
Stem elongation	First node detectable.	592 – 659
Anthesis	Flowering commences; first anthers of cereals are visible.	807 – 901
Seed fill	Seed fill begins. Caryopsis of cereals watery ripe (first grains have reached half of their final size).	1068 – 1174
Dough stage	Soft dough stage, grain contents soft but dry, fingernail impression does not hold.	1434 – 1556
Maturity complete	Grain is fully mature and drydown begins. Ready for harvest when dry.	1538 – 1665

## 1022 B.2 Interpolation

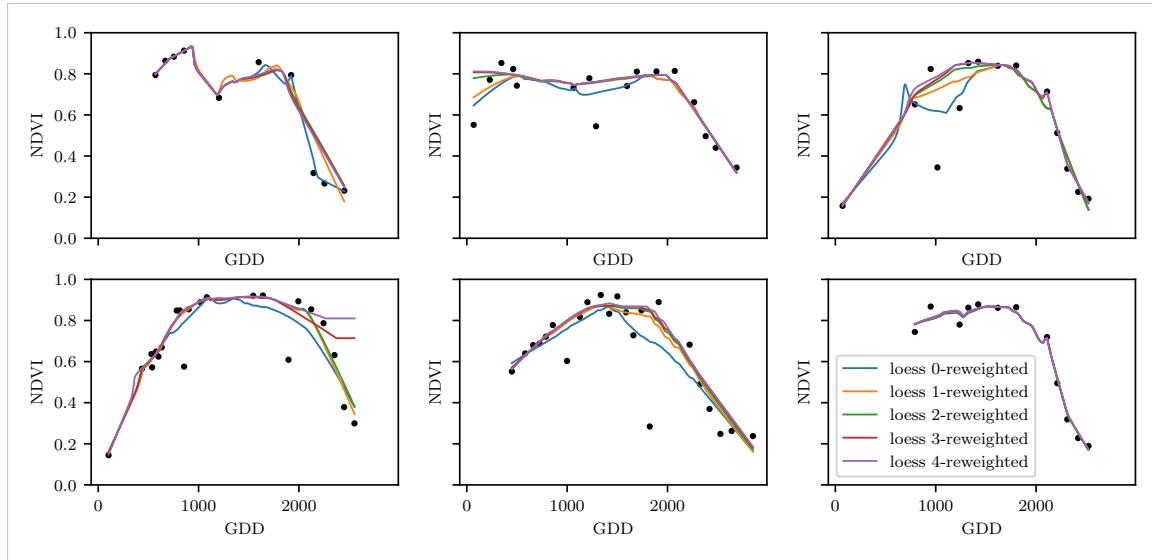


Figure B.1: The LOESS smoother fitted to different (SCL45) NDVI time series. Iterations of a robustifying refit (as indicated in section 3.5) are also displayed

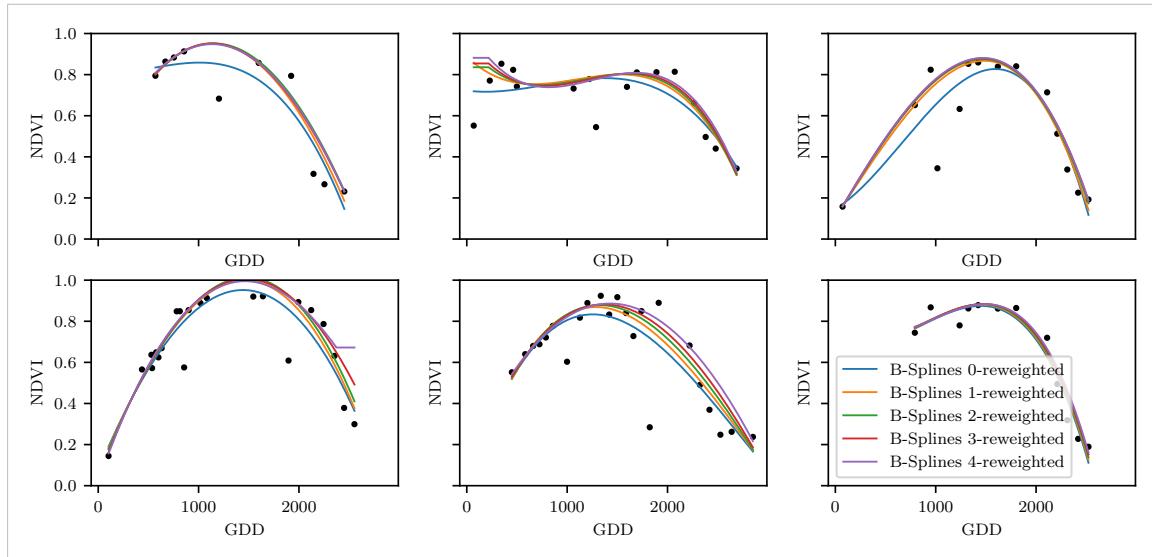


Figure B.2: B-splines fitted to different (SCL45) NDVI time series. Iterations of a robustifying refit (as indicated in section 3.5) are also displayed

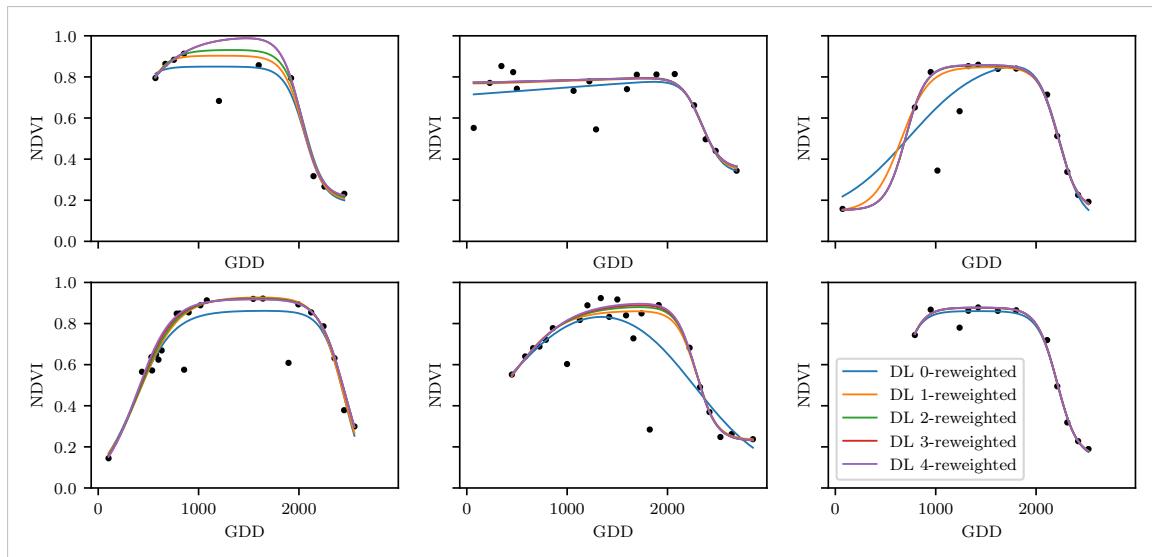


Figure B.3: A Double Logistic curve fitted to different (SCL45) NDVI time series. Iterations of a robustifying refit (as indicated in section 3.5) are also displayed

### 1023 B.3 NDVI correction

1024 page breaks

Table B.1: Non-relative RMSE for yield prediction

	RF	OLS-SCL	OLS-all	MARS	GAM	LASSO	no-correction
ss	1.144	1.033	1.051	1.042	1.046	1.042	1.095
dl	1.150	1.115	1.116	1.116	1.097	1.098	1.159
ss-rob	1.144	1.054	1.084	1.094	1.072	1.071	1.091
dl-rob	1.159	1.128	1.117	1.064	1.093	1.105	1.156

Table B.2: Coefficient of determination ( $R^2$ ) of yield prediction

	RF	OLS-SCL	OLS-all	MARS	GAM	LASSO	no-correction
ss	0.431	0.486	0.477	0.481	0.479	0.481	0.455
dl	0.427	0.445	0.444	0.444	0.454	0.453	0.423
ss-rob	0.431	0.475	0.461	0.456	0.467	0.467	0.457
dl-rob	0.423	0.439	0.444	0.470	0.456	0.450	0.424

#### 1025 B.3.1 OLS-SCL Model Outputs

```

1 Call:
2 lm(formula = (paste(response, " ~ ", "ndvi_observed + scl_class"))),
3   data = ndvi_df)
4
5 Residuals:
6   Min     1Q Median     3Q    Max
7 -0.7997 -0.0717  0.0039  0.0695  0.6632
8
9 Coefficients:

```

```

1036      Estimate Std. Error t value Pr(>|t|)
1037 10 (Intercept) 0.21465 0.00230 93.46 < 2e-16 ***
1038 12 ndvi_observed 0.71116 0.00346 205.65 < 2e-16 ***
1039 13 scl_class3 0.02205 0.00356 6.20 5.8e-10 ***
1040 14 scl_class4 -0.00431 0.00251 -1.72 0.085 .
1041 15 scl_class5 -0.09875 0.00234 -42.15 < 2e-16 ***
1042 16 scl_class6 -0.05301 0.01104 -4.80 1.6e-06 ***
1043 17 scl_class7 0.11245 0.00274 41.09 < 2e-16 ***
1044 18 scl_class8 0.25963 0.00253 102.57 < 2e-16 ***
1045 19 scl_class9 0.35994 0.00236 152.47 < 2e-16 ***
1046 20 scl_class10 0.09091 0.00308 29.54 < 2e-16 ***
1047 21 scl_class11 0.29784 0.00392 76.06 < 2e-16 ***
1048 ---
1049 23 Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ',' 1
1050 24
1051 25 Residual standard error: 0.146 on 124978 degrees of freedom
1052 26 Multiple R-squared: 0.532, Adjusted R-squared: 0.532
1053 27 F-statistic: 1.42e+04 on 10 and 124978 DF, p-value: <2e-16

```

R Summary of the NDVI correction model (c.f. equation 5.2.0.1)

```

1055
1056 1 Call:
1057 2 lm(formula = (paste(get_res(), " ~ ", "ndvi_observed + scl_class")),
1058 3   data = ndvi_df)
1059
1060 5 Residuals:
1061 6   Min     1Q   Median     3Q    Max
1062 7 -0.2051 -0.0427 -0.0074  0.0329  0.6589
1063
1064 9 Coefficients:
1065      Estimate Std. Error t value Pr(>|t|)
1066 11 (Intercept) 0.18647 0.00126 147.74 < 2e-16 ***
1067 12 ndvi_observed -0.13265 0.00190 -69.80 < 2e-16 ***
1068 13 scl_class3 -0.00180 0.00196 -0.92 0.3587
1069 14 scl_class4 -0.04069 0.00138 -29.55 < 2e-16 ***
1070 15 scl_class5 -0.09698 0.00129 -75.32 < 2e-16 ***
1071 16 scl_class6 -0.01906 0.00606 -3.14 0.0017 **
1072 17 scl_class7 0.01641 0.00150 10.91 < 2e-16 ***
1073 18 scl_class8 -0.00560 0.00139 -4.02 5.7e-05 ***
1074 19 scl_class9 -0.01384 0.00130 -10.67 < 2e-16 ***
1075 20 scl_class10 -0.00690 0.00169 -4.08 4.5e-05 ***
1076 21 scl_class11 -0.01446 0.00215 -6.72 1.8e-11 ***
1077 ---
1078 23 Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ',' 1
1079 24
1080 25 Residual standard error: 0.08 on 124978 degrees of freedom
1081 26 Multiple R-squared: 0.352, Adjusted R-squared: 0.352
1082 27 F-statistic: 6.8e+03 on 10 and 124978 DF, p-value: <2e-16

```

R Summary of the NDVI correction model (c.f. equation 5.2.0.2)

```

1084 replace space before ref by tilda
1085 check quantile definitions
1086 schwarz weiss färbung der IS tabelle korrigieren
1087 so wenig wie möglich abkürzungen in den fig und table captions
1088 refer to data availability
1089 abkürzungen Fourier und in tabellen

```

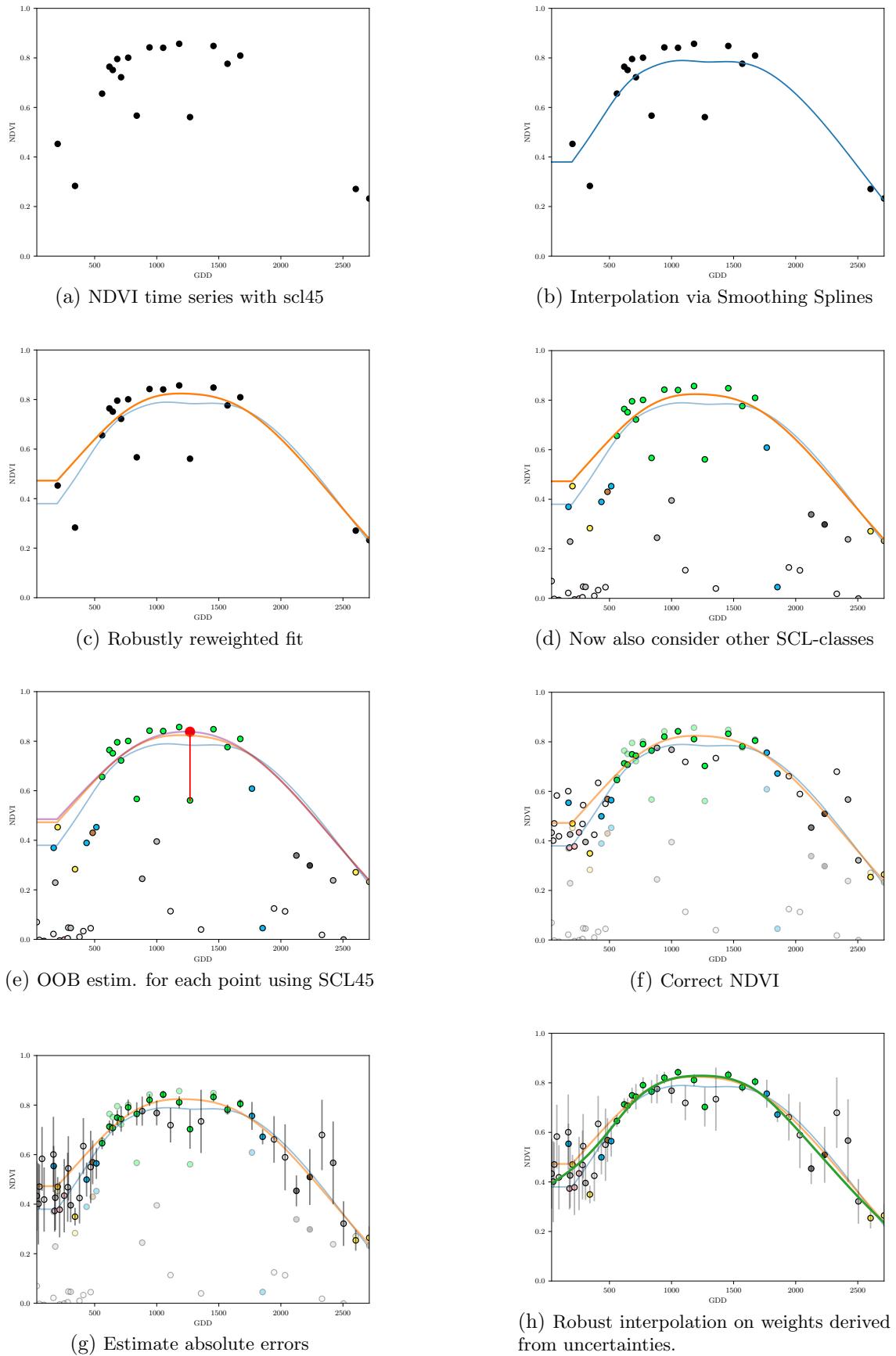


Figure B.4: Stepwise illustration of robust NDVI-Correction. For the color encoding of the SCL classes we refer to table 2.2.