



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**

12

13

14 Submission Date: September 18th 2022

15

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 Multispectral satellite imagery Time Series (TS) are utilized to model ground-based TS of
40 spectral indices. As such, the TS of the Normalized Difference Vegetation Index (NDVI)
41 — a proxy for vegetation density — is considered to contain information about vegetation
42 development. Due to atmospheric effects (e.g. clouds or shadows) satellite measurements
43 may not match the ground signal and therefore traditional approaches try to filter out such
44 corrupted observations before extracting and subsequently interpolating the NDVI. After
45 filtering, remaining corrupted observations and resulting data gaps are the two challenges
46 for interpolation that we address in this thesis.

47 For this purpose, we use crop yield maps from 2017-2021 of cereals from a farm in Switzerland
48 and corresponding Sentinel 2 satellite image TS published by the European Space
49 Agency, together with a Scene Classification Layer (SCL).

50 We give a benchmark-supported review of different interpolation methods and opt for
51 Smoothing Splines as a non-parametric (flexible) method and Double Logistic approximation
52 as a parametric method with implicit shape assumptions. In addition, we generalize
53 an iterative technique which robustifies interpolation methods against outliers by reducing
54 their weight. In most cases, this robustification successfully decreased the 50% and 75%
55 quantiles of the absolute out-of-bag residuals.

56 Moreover, using NDVI as an example, we present a general interpolation procedure that
57 utilizes additional information to correct the target variable with an uncertainty estimate
58 and then performs a weighted interpolation. Consequently, we do not filter using the SCL
59 but weight observations according to their reliability. The combination of different inter-
60 polation methods and correction models results in 28 interpolation strategies. We assume
61 that the better the resulting NDVI TS models crop growth, the more suitable it is to pre-
62 dict crop yield. Based on this, the best interpolation strategy uses Smoothing Splines and
63 corrects the NDVI with uncertainty estimation through a simple linear model considering
64 only of the observed NDVI and the associated SCL class. This strategy produces NDVI
65 TS that can explain 5% more variance in yield estimation than the NDVI TS obtained by
66 traditional Smoothing Splines.

67 Instructions and a codebase for reproducibility of the results, as well as an R package
68 making the presented general interpolation procedure accessible to the user, are supplied.

69 **Contents**

70	Notation	vi
71	1 Introduction	1
72	2 Data and Methods	3
73	2.1 Sentinel 2 Data	3
74	2.2 Crop Yield Data	3
75	2.3 Normalized Difference Vegetation Index (NDVI)	5
76	2.4 Timescale Transformation	6
77	2.5 The Concept of a ‘Pixel’	6
78	2.6 Challenges in S2 Data	6
79	2.7 General Methods	8
80	2.7.1 Root Mean Square Error (RMSE)	8
81	2.7.2 Out-Of-Bag (OOB) and Leave-One-Out-Cross-Validation (LOOCV)	8
82	3 Interpolation Methods	9
83	3.1 Interpolation Setup	9
84	3.2 Parametric Regression	9
85	3.2.1 Double Logistic (DL)	11
86	3.2.2 Fourier Series (FS)	11
87	3.2.3 Optimization Issues	12
88	3.3 Non-Parametric Regression	12
89	3.3.1 Kernel Regression: Nadaraya-Watson (NW)	12
90	3.3.2 Universal Kriging (UK)	13
91	3.3.3 Savitzky-Golay Filter (SG)	15
92	3.3.4 Locally Weighted Regression (LOESS)	16
93	3.3.5 B-Splines (BS)	17
94	3.3.6 Smoothing Splines (SS)	17
95	3.4 Tuning Parameter Estimation	18
96	3.5 Robustification	18
97	3.5.1 Our Adjustment:	19
98	3.5.2 Examples and Conclusions	20
99	3.5.3 Upper Envelope Approach - Penalty for Negative Residuals	20
100	3.6 Performance Assessment	20
101	4 NDVI Correction	21
102	4.1 Considering other SCL Classes	21
103	4.2 Correction Models	22
104	4.2.1 Ordinary Least Squares (OLS)	22
105	4.2.2 Least Absolute Shrinkage and Selection Operator (LASSO)	23
106	4.2.3 General Additive Model (GAM)	24
107	4.2.4 Random Forest (RF)	24
108	4.2.5 Multivariate Adaptive Regression Splines (MARS)	25
109	4.3 Weighted Interpolation	26
110	4.4 Resulting Interpolation Strategies	26
111	4.5 Evaluation via Yield Estimation Accuracy	27

112	5 Results	30
113	5.1 Goodness of Fit for Selected Interpolation Methods	30
114	5.2 XXX (Robustification and) NDVI-Correction	30
115	6 Discussion	32
116	6.1 Interpolation Methods	32
117	6.1.1 Data Gaps in Time Series	32
118	6.1.2 Preselection	33
119	6.1.3 Candidate Selection	33
120	6.2 NDVI Correction	33
121	6.2.1 Bootstrap	33
122	6.2.2 Using Additional Covariates	33
123	6.2.3 Choose Interpolation Strategy	34
124	6.2.4 High RMSE in Yield Prediction	34
125	7 Conclusion	35
126	7.1 Future Work	37
127	7.1.1 Time Series Correction-Interpolation as a General Method	37
128	7.1.2 Minor Improvements	37
129	Bibliography	38
130	A Reproducibility	40
131	A.1 Reproduce Results	40
132	A.2 R-Package	40
133	B Further Material	42
134	B.1 Data and Methods	42
135	B.1.1 GDD	42
136	B.2 Interpolation	43
137	B.3 NDVI correction	44
138	B.3.1 OLS-SCL Model Outputs	44

Todo list

139

140	verdeutliche dem leser, dass ein auftrag das findne von interpolationmethoden war	9
141	Paper zitieren wo eingeführt oder wo benutzt (falls einföhrung fast schon trivial)	9
142	figure / tabelle / pseudocode anstatt aufzählung	15
143	consider naming the sub-plots	20
144	defition of RYEA, it is not an accuracy but an error	30
145	Here in the discussion, you should take up the points you mentioned in the introduction .	32
146	where does this section belong to? Chapter ‘NDVI Correction’ or ‘Further Work’?	33
147	table mit OLS SCL als sieger diskutieren	34
148	kurzer kontext von vergleichbaren values von gregor — diese sektion ist für dena uftraggebeinr	34
149	even in a perfect world the NDVI curve only holds a fraction of the information avialbe .	34
150	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”	
151	subsection at the end of the discussion?	34
152	Frage: mehr details für die begründung der Interpolations-kandidaten?	35
153	Anzahl von Beobachtungen, empfehlungen? – schwierig, weil regelmäßigkeit in ‘wichtigen’ zeiträumen (der veränderung) wichtiger ist.	36
154	page breaks	44
155	replace space before ref by tilda	45
156	check quantile definitions	45
157	schwarz weiss färbung der IS tabelle korrigieren	45
158	so wenig wie möglich abkürzungen in den fig und table captions	45
159	refer to data aviability	45
160	abkürzungen Fourier und in tabellen	45
161	figure spacing (caption zu nah dran — manuell vspace einfügen wo nötig)	45

165 Notations

166 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

167 Abbreviations and Objects

Pixel	A pixel originates of an image pixel and describes a square of 10 x 10 meters in the field that 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 that 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 (Rouse, 1974)
DAS	Days After Sowing
GDD	Growing Degree Days – cumulative sum of “ $\max(0, \text{temperature} - \text{threshold})$ ”
RYEA	Relative Yield-Estimation-Accuracy. Definition 4.5.0.1
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 2.7.2).
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 2.7.2).

168 **Statistical Models**

DL	Double Logistic (c.f. section 3.2.1)
FS	Fourier Series (c.f. section 3.2.2)
NW	Nadaraya-Watson (c.f. section 3.3.1)
UK	Universal Kriging (c.f. section 3.3.2)
SG	Savitzky-Golay Filter (c.f. section 3.3.3)
LOESS	Locally Weighted Regression (c.f. section 3.3.4)
BS	B-splines (c.f. section 3.3.5)
SS	Smoothing Splines (c.f. section 3.3.6)
OLS	Ordinary Least Squares (c.f. section 4.2.1)
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 4.2.2)
GAM	General Additive Model (c.f. section 4.2.3)
RF	Random Forest (c.f. section 4.2.4)
MARS	Multivariate Adaptive Regression Splines (c.f. section 4.2.5)

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

170 **Chapter 1**

171 **Introduction**

172 Remote sensing zielt darauf ab, ziel-Grössen effizient aus der Entfernung messen zu können.
173 Hier finden Satellitenbilder Zeitreihen Verwendung, wie etwa die von der europäischen
174 Raum-Agentur (ESA) kostenlos veröffentlichten Bilder Zeitreihen der Multi-spektralen
175 Sentinel 2 (S2) Satelliten. Die Vegetationsentwicklung von Wäldern und landwirtschaftlich
176 relevanten Flächen im grossen Stile zu überwachen, ist unter anderem für public angents,
177 Versicherungen, Umwelt- und Klimaforscher von grossem Interesse. Mögliche Ziele sind
178 hierbei eine crop Klassifizierung für das Subventionieren von Bauern oder das Erstellen
179 von Pflanzenmodellen, um Ernteertrag oder Stickstoffkonzentration zu schätzen. Um die
180 hochdimensionalen Satellitenbilder in leicht interpretierbare Grössen zu transformieren,
181 werden spektrale Indizes, wie der Normalized Difference Vegetation Index (NDVI) benutzt.
182 Dieser ist ein Proxy für die Vegetationsdichte und die korrespondierende Zeitreihe spiegelt
183 somit das Pflanzenwachstum wider. Der Informationsgehalt von einem Satellitenbild ist
184 jedoch abhängig vom Zustand der Atmosphäre und so trägt der davon abgeleitete NDVI
185 bei einer dichten Wolkendecke keine Informationen über die Vegetation am Boden. Daher
186 liefert die ESA zusätzlich eine Scene Classification Layer (SCL), welche Aufschluss gibt,
187 was beobachtet wird (z.B Schatten, Wolken, Vegetation, etc.). So können wir bei der
188 Extraktion der NDVI Zeitreihe aus der S2 Satellitenbilder Zeitreihe, anhand der SCL
189 Klassifizierung, die uninformativen Beobachtungen herausfiltern. Durch diese Filtration
190 kann es jedoch leicht vorkommen, dass wir besonders im Winter über mehrere Wochen
191 keine Observationen haben. Zudem kommt, dass manche Beobachtungen zu Unrecht durch
192 die SCL als informativ bewertet wird (z.B. als Vegetation) und somit in einem fehlerhaften
193 NDVI resultiert. Diese beiden Probleme versucht man gegenwärtig mit Interpolation und
194 Smoothing zu lösen. Starke Formannahmen über die NDVI Kurve werden in ... getroffen.
195 Flexiblere Ansätze wurden von ... verwendet.

196 In dieser Thesis werden wir stärken und schwächen von solch gängigen Interpolations-
197 methoden diskutieren und hinsichtlich der NDVI Interpolation bewerten. Dafür benutzen
198 wir die S2 Satellitenbilder Zeitreihe und Ernteertragskarten von verschiedenen Feldern
199 verschiedenen Weizenarten auf einer Farm in Witzwil in der Schweiz über die Jahre 2017-
200 2021. Um die Interpolationmethoden zu verbessern, verallgemeinern und testen wir eine
201 iterativen Technik, die Interpolationen robuster gegen Ausreisser machen soll, indem sie
202 weniger Gewicht bekommen. Zudem ermitteln wir, wie Datenlücken die verschiedenen
203 Interpolationmethoden beeinflussen. Ausserdem stellen wir am Beispiel des NDVI eine
204 generelle Interpolations-prozedur vor, welche anhand von zusätzlichen Informationen die
205 Zielvariable mit einer Unsicherheitsschätzung korrigiert und anschliessend interpoliert.

206 Somit müssen wir die Observationen nicht mehr a priori via der SCL filtern , sondern
207 korrigieren den beobachteten NDVI und gewichten die Beobachtungen via der geschätzten
208 Unsicherheiten. Durch die Kombination von Interpolationsmethoden mit den NDVI kor-
209 rigierenden Modellen ergeben sich somit 28 Intepolationsstrategien. Diese benchmarken wir
210 mit einem objektiven Qualitätsmass, welches annimmt, dass je besser eine NDVI TS das
211 Pflanzenwachstum modelliert, desto geeigneter ist sie, um den Ernteertrag zu schätzen.

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

- 213 i.) 1 review of interpolation methods
- 214 ii.) 2 erroruous observations — how to deal with them
- 215 iii.) 3 data gaps — influence itpl mehtods
- 216 iv.) 4 data gaps — how to deal with them
- 217 v.) 5 how to compare two NDVI interpolation strategies?

218 Roadmap ...

219 **Chapter 2**

220 **Data and Methods**

221 We will start by describing the available data and the challenges associated with it. Our
222 study region is a farm of over 800ha, which is located in western Switzerland. From Perich
223 et al. (2022) we acquire satellite image data (section 2.1), yield maps of several cereals
224 from 2017 to 2021 (section 2.2), and meteorological data (section 2.5). Afterwards, we will
225 introduce general methods in section 2.7, which will be used in the remaining chapters.

226 **2.1 Sentinel 2 Data**

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

231 The S2 images contain 12 spectral bands with spatial resolutions up to 10 meters (see
232 2.1). Bands with a lower resolution (20 and 60 meters) were upscaled to 10 meter reso-
233 lution using cubic interpolation (Perich et al. (2022)). In order to decrease the effect of
234 atmospheric conditions like reflections and scattering, bottom-of-atmosphere, radiometric
235 corrected Level-2A data was used². The ESA also supplies an algorithm³ produces Scene
236 Classification Layer (SCL) where for each location the observed subject is assigned to one
237 of 11 SCL-classes (c.f. table 2.2). In this thesis, we will use this classification to filter out
238 data points, that we believe to be less informative. That are all observations which SCL-
239 class does not correspond to vegetation or bare soils (classes 4 and 5). For convenience,
240 we define the set SCL45 as the observations that belong to SCL-class 4 or 5.

241 **2.2 Crop Yield Data**

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

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

²According 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”

³[https://sentinels.copernicus.eu/web/sentinel/technical-guides/sentinel-2-msi/level-2a/
algorithm](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 λ in nm with the spectral width $\Delta\lambda$ in nm with a spatial resolution SR in m ([Jaramaz et al. \(2013\)](#)).

Band	λ	$\Delta\lambda$	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.

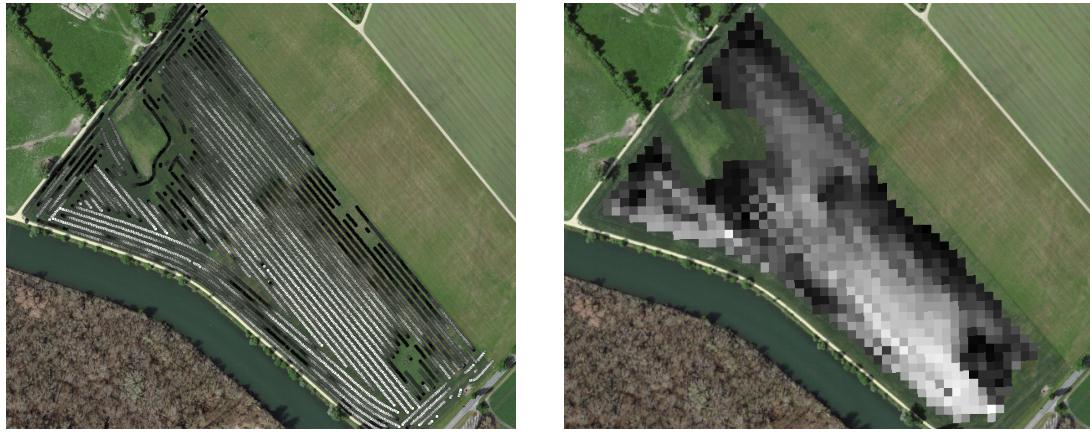
Table 2.2: Overview: Scene Classification Layers (SCL)

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

244 t/ha (see fig. [2.1a](#)). We take the data set derived in [Perich et al. \(2022\)](#), where error-
 245 prone measurement points (such as during a tight curve of the combine harvester) were
 246 removed and then the yield map was rasterized using linear interpolation (c.f. fig. [2.1b](#)).
 247 We summarize the rasterized dry-yield values by the following statistics:

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

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



(a) Raw combine harvester data (cleaned)

(b) rasterized to Sentinel 2 resolution.

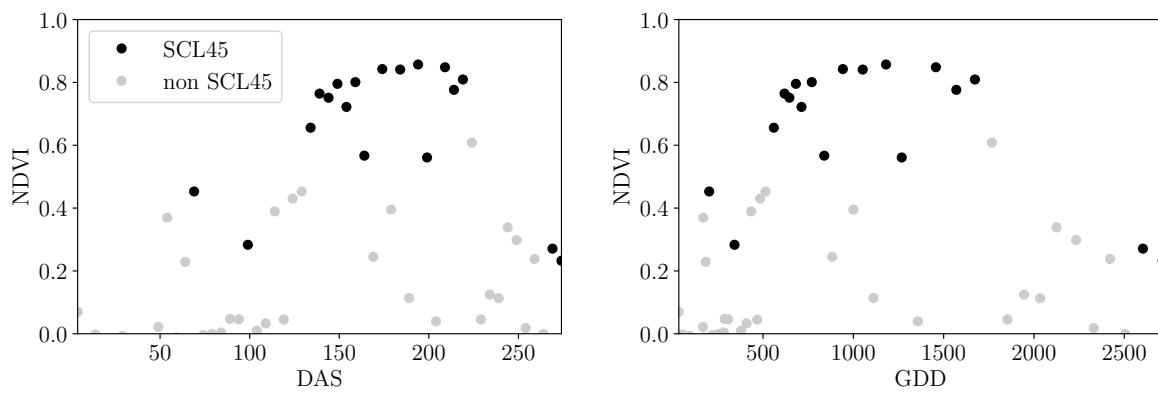
Figure 2.1: Crop yield density map of a field. Ranges from 0.1 t/ha (black) to 5.35 t/ha (white)

253 2.3 Normalized Difference Vegetation Index (NDVI)

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

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

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



(a) Days After Sowing (DAS)

(b) Growing Degree Days (GDD)

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

264 2.4 Timescale Transformation

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

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

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

279 2.5 The Concept of a ‘Pixel’

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

283 Consider a 10 by 10 meter square that coinsides with a S2 image pixel and T the GDD
 284 values for which S2 images are avialable in a given season. For $t \in T$ let P_t be a tupel of
 285 all the spectral bands, the observed NDVI and the SCL class (at the considered location
 286 at time t). Then, define P as the collection of all the P_t and the estimated dry-yield for
 287 this square. Analogously to P , define P^{SCL45} by only considering P_t with SCL-class 4 or
 288 5 (vegetation and soil).

289 2.6 Challenges in S2 Data

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

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

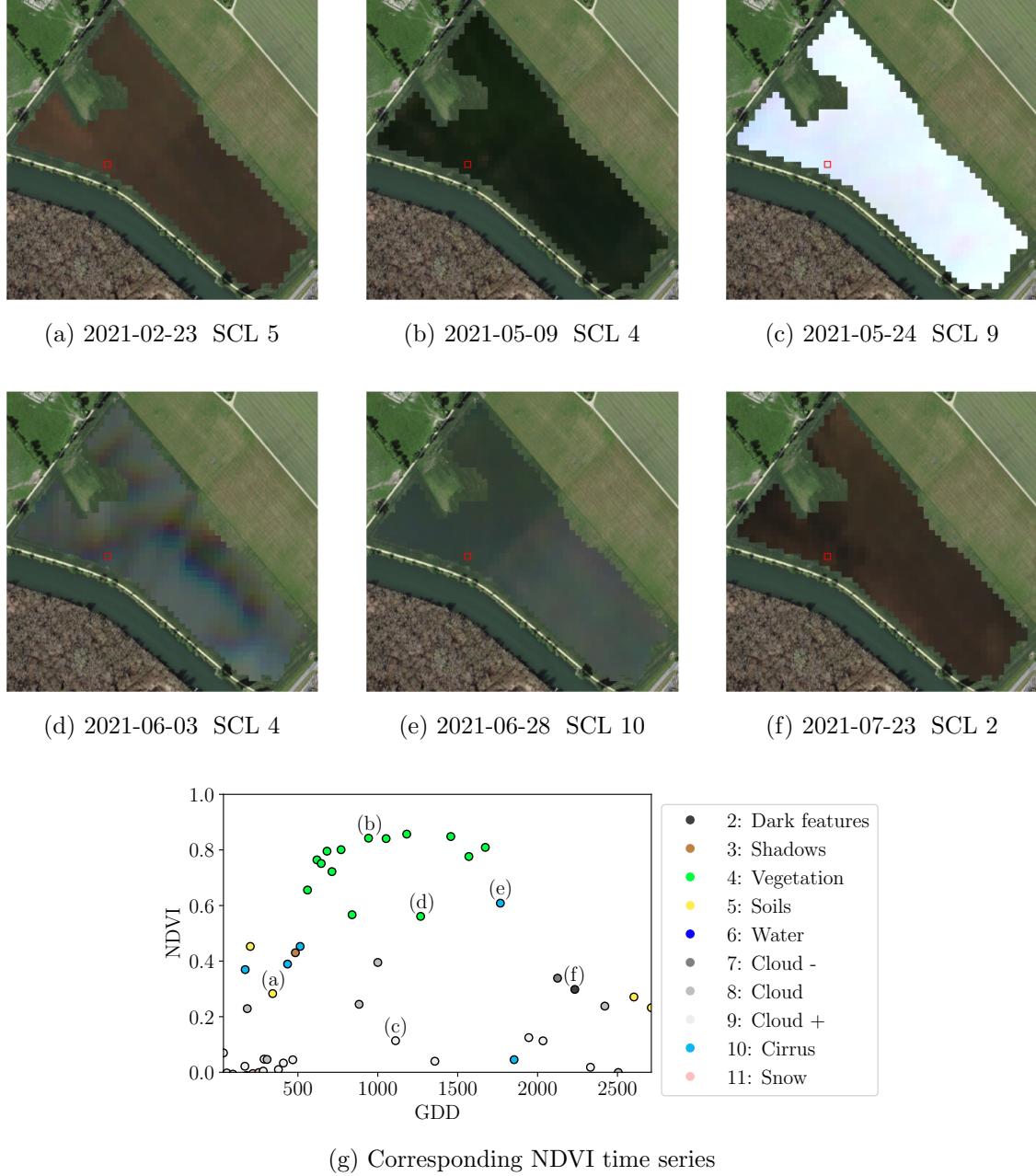


Figure 2.3: Satellite images of a field at selected times with a static background for orientation. Moreover, the NDVI time series of the red-highlighted pixel is shown in (g) colored by the SCL labels.

303 **2.7 General Methods**

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

309 **2.7.1 Root Mean Square Error (RMSE)**

310 In this section we describe different criteria to evaluate models. Hence, given a vector
 311 $y \in \mathbb{R}^n$ and its estimator \hat{y} (estimated using the model), we define the RMSE as:

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

312 **2.7.2 Out-Of-Bag (OOB) and Leave-One-Out-Cross-Validation (LOOCV)**

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

To formally introduce this, let:

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

316 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
 317 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
 318 obtain $\hat{y} := (\hat{y}_1, \dots, \hat{y}_n)$ the OOB estimator for $y \in \mathbb{R}^n$.

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

321 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
 322 the corresponding \hat{y}_i also the LOOCV estimator.

323 If we optimize some parameter via OOB (or LOOCV) this means that we search for the
 324 parameter that minimizes some loss function which takes the OOB (or LOOCV) residuals.
 325 Usually we approximate this parameter by searching on a grid.

326 **Chapter 3**

327 **Interpolation Methods**

328

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

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

340 **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 ε_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

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

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

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

345 **3.2 Parametric Regression**

346 Parametric Curve estimation tries to fit a parametric function, such as, for example, a
347 Gaussian function with parameters μ and σ , to a dataset. In the following, we introduce
348 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).

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

349 **3.2.1 Double Logistic (DL)**

The Double Logistic smoothing as described in Beck et al. (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:

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

- 350 Where the five free parameters: y_{\max} , d_0 , d_1 , t_0 , t_1 are initially estimated by least squares.
 351 Such fit can be seen in figure 3.1.

352 **Robustification**

- 353 Similar as for the SG (c.f. section 3.3.3) one can reestimate (only once) the parameters by
 354 giving less weight to the overestimated observations and more weight to the underestimated
 355 observations. For the details on the choice of the weights we refer to Beck et al. (2006). We
 356 will not apply this reestimation but rather the robustification introduced later in section
 357 3.5.

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

358 **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:

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

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

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

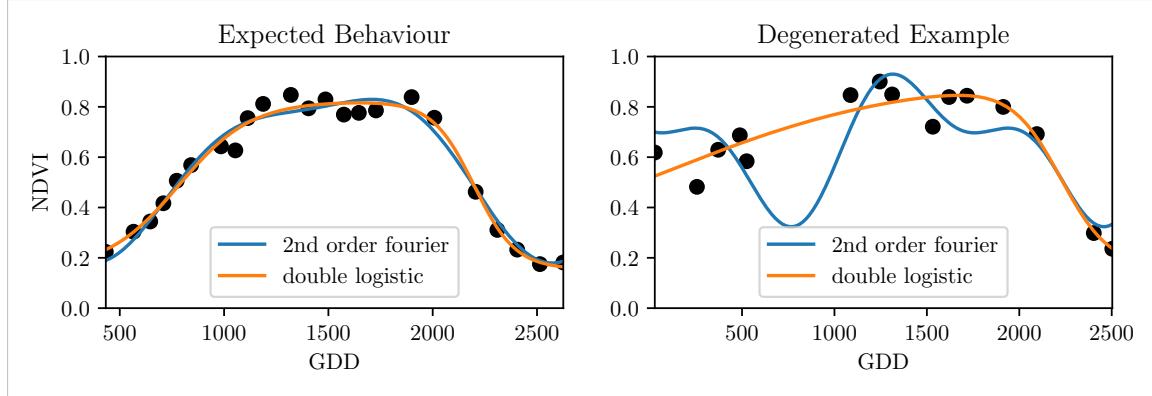


Figure 3.1: Here we observe the possibilities of a precise fit for the two parametric methods but notice also some misbehavior

362 3.2.3 Optimization Issues

363 We shall mention some optimization issues we countered during implementation. Since we
 364 aim to minimize the residual sum of squares over 5 (or 6) parameters, we try to solve a
 365 non-convex optimization problem. Thus, the algorithm¹ either struggles to find the global
 366 minimum or fails to converge. This was fixed by providing for each parameter reasonable
 367 initial values and generous bounds (that match our experience).

368 3.3 Non-Parametric Regression

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

372 3.3.1 Kernel Regression: Nadaraya-Watson (NW)

373 As described in section 3.1, we aim to estimate

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

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

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

¹We 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:

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

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

378 Choose Bandwidth

379 Note that we still need to choose the bandwidth of the function. This can be done with
 380 the help of LOOCV while optimizing the RMSE. For non-equidistant data we refere to
 381 Brockmann et al. (1993) where a local adaptive bandwidth selection is presented.

Advantages	Disadvantages
— fletible 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-matrit)	— 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)	

382 3.3.2 Universal Kriging (UK)

383 UK as described in dig (2007) was developed in geostatistics to deal with autocorrelation
 384 of the response variable at locations that are spatially close. By applying the notion that
 385 two spectral indices that are timewise close should also take similar values, we justify the
 386 application of UK. In the end, we would like to fit a smooth Gaussian process to the data.

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

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 γ via the above equation by choosing the Gaussian Variogram defined by

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

393 Here h denotes the distance, n is the nugget, r is the range and p is the partial sill. The
 394 influence of the parameters is visualized in figure 3.2.³

²Note 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.

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

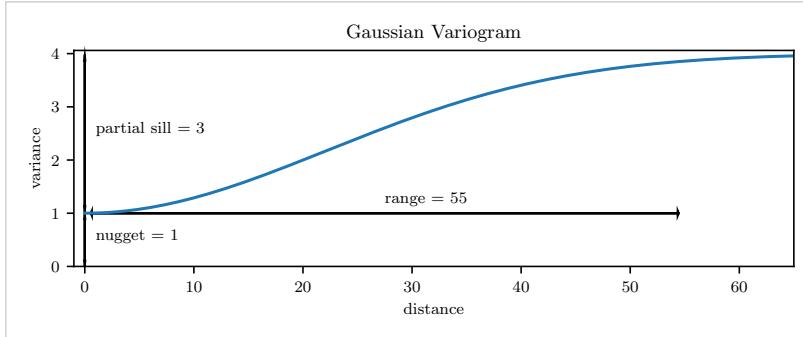


Figure 3.2: Gaussian Variogram with nugget=1, partial sill=3, range=55

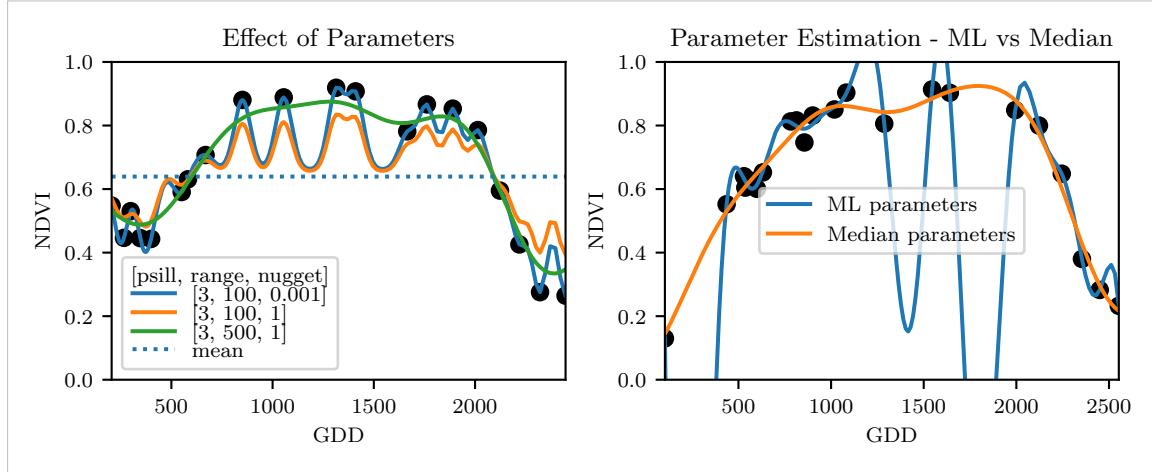


Figure 3.3: On the left, we see how the interpolation change if we increase the nugget and the range parameter. On the right, we compare two UK interpolations, where one takes parameters by numerically maximizing the (which results in a very small nugget) and the other takes the median of many such numerical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

395 Finally, we consider a one-dimensional Gaussian process G_γ with variogram γ and tune the
 396 variogram parameters using maximum likelihood⁴. Let z be a vector with the new values
 397 to extrapolate, then we can determine the values $m(z) = \mathbb{E}[G_\gamma(z)|(t, y)]$ using Bayes rule⁵.
 398 For an example fit, we refer to figure 3.3.

399 Violated Assumption

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

⁴ 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.

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

Advantages	Disadvantages
— It is a well-studied method.	— Regression to the mean.
— Variogram parameters have an intuitive meaning.	— Violated assumption of constant mean and constant variance. Thus, the NDVI is not a stationary process.
— Flexible covariance structure.	— Pure maximum likelihood can result in overfitting.

404 **3.3.3 Savitzky-Golay Filter (SG)**

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

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

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

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

410 [Chen et al. \(2004\)](#) developed a ‘robust’ interpolation method for the NDVI based on the
 411 SG. The method is based on the assumption that due to atmospheric effects the observed
 412 NDVI tends to be underestimated and that it cannot increase too quickly. The latter is
 413 argued by the biological impossibility of such fast vegetation changes. Their proposed
 414 algorithm is:

- 415 i.) Remove non-SCL45 points.
- 416 ii.) Remove points that would indicate an increase greater than 0.4 within 20 days.
- 417 iii.) Linearly interpolate to obtain an equidistant time series X^0 .
- 418 iv.) Apply the SG to obtain a new time series X^1 .
- 419 v.) Update X^1 by applying again a SG. Repeat this until $w^T |X^1 - X^0|$ stops decreasing,
 420 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
 421 penalty introduced by outliers⁶ and by repeating this step we approach the “upper
 422 NDVI envelope”.

figure /
 tabelle /
 pseudocode
 anstatt
 aufzählung

423 **Extension: Spatial-Temporal SG**

424 One notable adaptation of the SG is the presented by [Cao et al. \(2018\)](#). The key difference
 425 is the additional assumption of the cloud cover being discontinuous and that we can

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

426 improve by looking at adjacent pixels⁷. Because we are working with rather high resolution
 427 satellite data, and we need the variance in the predictors, we will waive this extension.

Advantages	Disadvantages
<ul style="list-style-type: none"> — Popular technique in signal processing. — Efficient calculation for equidistant points. — Upper envelope matches intuition for the NDVI. Therefore, it is robust against outliers with small values. 	<ul style="list-style-type: none"> — No natural way of how to estimate points that are not in the data. — Not generalizable to other spectral indices. — Linear interpolation to account for missing data might be not appropriate. — No smooth interpolation between two measurements.

428 3.3.4 Locally Weighted Regression (LOESS)

429 The LOESS introduced by Cleveland (1979) can be understood as a generalization of the
 430 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

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

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

433 Differences between the Robust LOESS and the SG?

434 The LOESS smoother takes a fraction of points instead of a fixed number and therefore
 435 automatically adapts to the size of the data we wish to interpolate. However, we run
 436 into the danger of considering too little observations, since the estimation breaks down if
 437 $\lceil \alpha n \rceil < d + 1$.⁸ Furthermore, LOESS gives less weight to points further away. This yields
 438 a "smoother" estimate, since when we slide the window (e.g. for estimating the next value)
 439 an influential point at the border does not suddenly get zero weight from being weighted
 440 equally before. Finally, the LOESS also can be used for non-equidistant data and allows
 441 for arbitrary interpolation.

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

⁷Here, 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

⁸If 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 α is big enough.

442 **3.3.5 B-Splines (BS)**

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

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

443 where B are basis functions and recursively defined by:

444

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

Assuming that all t_i are distinct, this yields an interpolation that 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 functions⁹ such that:

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

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

- can be assigned degrees of freedom
 - extendable to "smooth" version
 - performs also well if points are not equidistant
- smoothing process does not translate well to a interpretation (unlike SS)
 - choice of smoothing parameter s

445 **3.3.6 Smoothing Splines (SS)**

446 Let \mathcal{F} be the Sobolev space (the space of functions of which the second derivative is
447 integrable). Then the unique¹⁰ minimizer

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

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

⁹So 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

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

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

450 3.4 Tuning Parameter Estimation

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

- 455 i.) Construct a set Λ of candidate parameters that generously covers the parameter
 456 space.
- 457 ii.) Consider \mathcal{P} , a set of Pixels.
- 458 iii.) For each parameter $\lambda \in \Lambda$ consider the individual pixels and compute the LOOCV¹¹
 459 for the absolute residuals of the specific NDVI interpolation method for all Pixels in
 460 \mathcal{P} and store them in the set R_λ .
- 461 iv.) Determine $\lambda_{optimal} = \arg \min_{\lambda \in \Lambda} q_{90}(R_\lambda)$, where we describe the 90% quantile with
 462 q_{90} .

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

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

467 3.5 Robustification

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

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

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

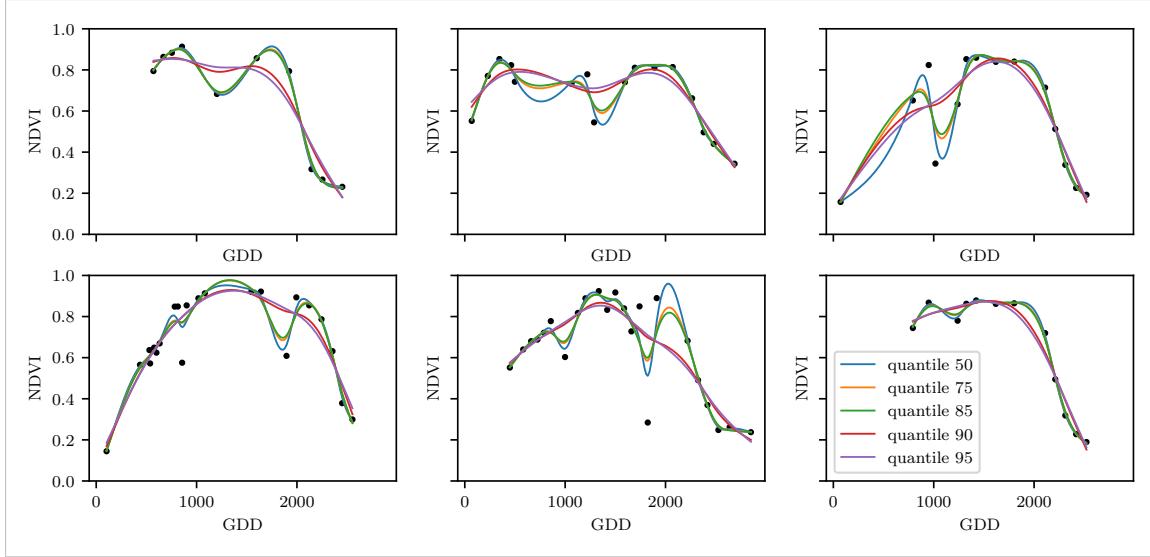


Figure 3.4: Smoothing splines fit with smoothing parameter optimized by minimizing the given quantile of the absolute leave-one-out residuals. Note that the larger the considered quantile is, the smoother the resulting curve becomes.

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

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

476 Next, we compute new weights by

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

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

479 Note that this procedure is indeed robust since we use the median for the normalization
480 which has a breakdown point¹² of 50%.¹³

481 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

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

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

¹²Intuitively, 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%.

¹³The 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).

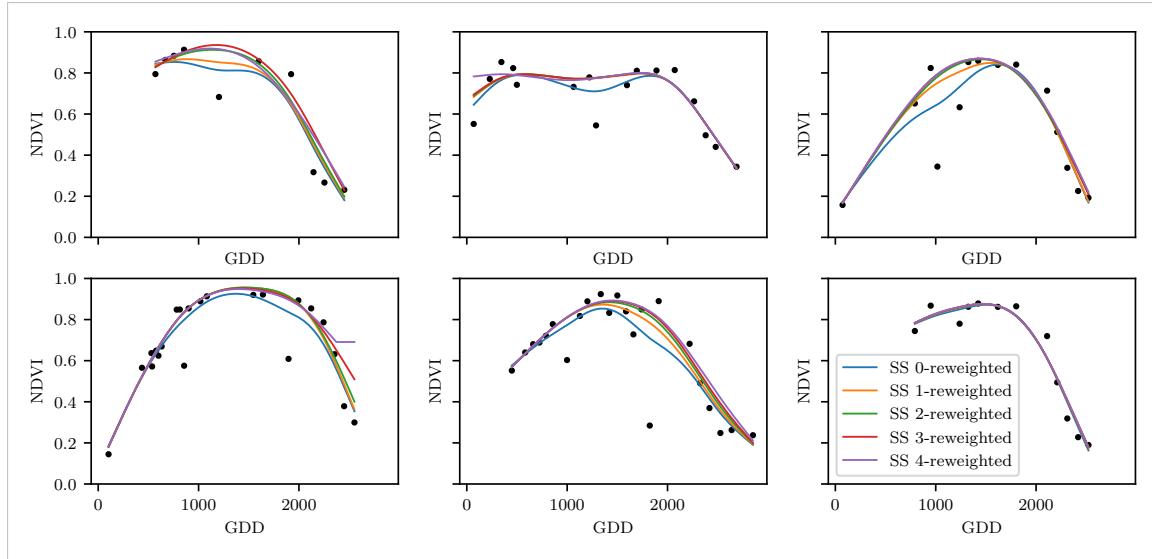
483 **3.5.2 Examples and Conclusions**

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

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

consider
naming
the sub-
plots

491 **3.5.3 Upper Envelope Approach - Penalty for Negative Residuals**

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

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

499 **3.6 Performance Assessment**

500 Next, we will benchmark the in section 6.1.2 preselected interpolation methods with and
 501 without robustification. For this, we will use the same technique as we did for the param-
 502 eter determination in section 3.4. On B_λ we apply the RMSE and different quantiles.

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

506 **Chapter 4**

507 **NDVI Correction**

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

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

522 **4.1 Considering other SCL Classes**

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

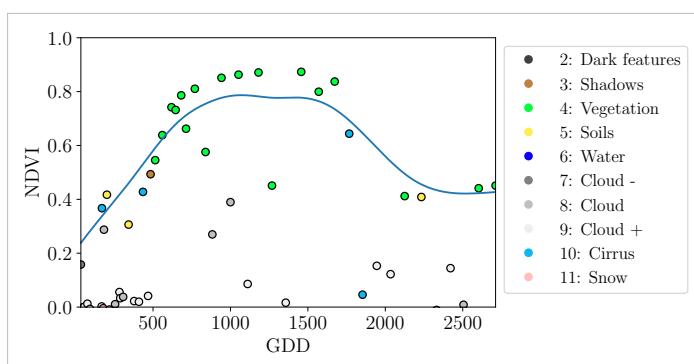


Figure 4.1: A smoothing splines fit considering green and yellow points (SCL45)

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

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

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

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

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

544 4.2 Correction Models

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

555 In the following, we will introduce different approaches, to model the relationship between
 556 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
 557 study the basic OLS. Second, we look at the LASSO, an penalized adaptation of the
 558 OLS which is known to successfully deal with highly correlated covariates. Afterwards,
 559 GAMs are introduced which model the response similar to OLS but allow for non-linear
 560 relations. Last but not least, we discuss RF and MARS, which are both flexible modelling
 561 approaches.

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

564 4.2.1 Ordinary Least Squares (OLS)

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

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

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

567 Assuming that $(X^T X)$ is regular, we can estimate the regression coefficients β by

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

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

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

570 4.2.2 Least Absolute Shrinkage and Selection Operator (LASSO)

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

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

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

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

579 In order to know which λ to choose, we try a huge range of possible values. For each
580 β_λ , we calculate the cross-validated $RMSE_\lambda$ ⁴ (and its standard deviation σ_λ using the k
581 folds) and define the λ with the smallest corresponding $RMSE_\lambda$ as λ_{min} . From here we
582 choose the largest λ for which the $RMSE_\lambda$ is smaller than $RMSE_{\lambda_{min}} + \sigma_\lambda$. This yields
583 a simpler model while keeping the $RMSE$ reasonable model.

584 We will apply the LASSO using the selected covariates in section 4.2 and their second
585 degree of interactions.⁵

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

³The last two terms are equivalent by lagrangian optimization

⁴The 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.

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

586 **4.2.3 General Additive Model (GAM)**

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

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

591 To estimate the non-parametric functions, we can use SS (ref sec. [3.3.6](#)). For this let \mathcal{S}_j
 592 be the function that takes some $z \in \mathbb{R}^n$ and returns the SS fitted to $(X_{:,j}, z)$ where the
 593 smoothing parameter is optimized by LOOCV⁷. Since we cannot fit all g_j simultaneously,
 594 we will use a strategy named Backfitting. We basically cycle through the indices $1, \dots, p$
 595 and refit \hat{g}_j each time. The following illustrates the procedure:

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

596 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.

597 **4.2.4 Random Forest (RF)**

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

605 To build such a Tree, we will recursively partition the covariate space using greedy splits⁸
 606 decreasing the RMSE⁹ each time. If the set we want to split contains less than a certain
 607 amount of training points, we stop.

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

⁷For efficiency an proxy of the LOOCV is used called generalized cross validation.

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

⁹To 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).

608 To build a Random Forest we will bootstrap-aggregate¹⁰ many such Trees¹¹. The prediction
 609 of the Random Forest for a new point x is then the mean of the predictions from all
 610 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.

611 4.2.5 Multivariate Adaptive Regression Splines (MARS)

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

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

613 where the h_m are simple functions (explained later) and the β_m are estimated via Least
 614 Squares.

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

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

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

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

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

623 Finally, to avoid overfitting, we prune the set \mathcal{M} by optimizing a LOOCV score.¹²

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

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

¹¹Building 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.

¹²This 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 that – when excluding from \mathcal{M} – leads to the best LOOCV score.

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

627 4.3 Weighted Interpolation

628 Once we corrected the NDVI using the models described in the previous section, we are left
 629 with the problem that not every correction is equally reliable.¹³. Hence, we are interested
 630 in a measure of how uncertain an estimate is. We achieve this analogously as we corrected
 631 the NDVI, by replacing the response (NDVI-“true”) with the absolute residuals $v := |y - \hat{y}|$
 632 and modeling their relationship with the covariates defined by X . In this way, we obtain
 633 a model for the absolute residuals v and the estimator \hat{v} .

634 In the following we will convert our uncertainty estimate into weights that can be used for
 635 interpolation. For this, consider a pixel P , $\hat{y}^{(P)}$ its corrected NDVI and $\hat{v}^{(P)}$ the estimated
 636 uncertainties of $\hat{y}^{(P)}$. In order to interpolate $\hat{y}^{(P)}$, we will give less weight to unreliable
 637 observations. Thus, we define the weight function:

$$w_{\tau}^{(P)} := \frac{1}{R} \frac{1}{\hat{v}_{\tau}^{(P)}}, \quad \text{for } \tau = 1, \dots, n_P \quad (4.3.0.1)$$

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

641 4.4 Resulting Interpolation Strategies

642 We have developed the following procedure to obtain a new interpolation (keyword-wise):
 643 i.) LOOCV Interpolation (+ robustify?) to get “true” NDVI
 644 ii.) Correction
 645 iii.) Uncertainty estimation
 646 iv.) Interpolation (+ robustify?)

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

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

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

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

- 654 — We use the same interpolation method each time.
 655 — Either we robustify both times, or we do not robustify at all.
 656 — We use the same underlying method for correction and uncertainty estimation.
- 657 In this fashion, we obtain 28 distinct interpolation strategies, which we will benchmark in
 658 the next section.

659 4.5 Evaluation via Yield Estimation Accuracy

660 In this section, we introduce the relative yield-estimation-accuracy (RYEA) and utilize it
 661 to evaluate the 28 interpolation strategies from section 4.4. The fundamental assumption is
 662 that the closer the interpolated NDVI time series is to the true one, the better it can be used
 663 to determine crop yield. Implicitly, we believe that an NDVI time series that better models
 664 yield will incorporate more true information about the underlying vegetation. Therefore,
 665 we want to determine a comparable RYEA for each interpolation strategy and choose it
 666 as a benchmark criterion. This is an objective measure, since we have not considered crop
 667 yield in any of our previous steps. Moreover, this criterion is justified by the fact that
 668 yield estimation has been a motivation for the interpolation.

669 **Definition 4.5.0.1.** (RYEA) Let $y \in \mathbb{R}^n$ be the yield, M be a model for estimating y , and
 670 $\hat{y} = M(X)$ where X describes the data¹⁴. We define the RYEA as the relative RMSE in
 671 yield estimation. Formally expressed:

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

672 where \bar{y} denotes the sample mean.

673 We would like to estimate the yield from the NDVI time series produced by all the in-
 674 terpolation strategies for all pixels. However, given the high dimensionality and different
 675 lengths of the interpolation (not every time series has the same start and end point), we
 676 must first map each NDVI time series into a low-dimensional vector space of covariates.
 677 For this, we will use the following statistics:

- Maximum slope
- Minimum slope
- Integral¹⁵ over all
- Peak (i.e. maximal NDVI)
- GDD for the Peak
- Integral¹⁵ up to the peak
- Integral¹⁵ after peak
- Integral¹⁵ from 0-685 GDD
- Integral¹⁵ from 685-1075 GDD

678 For the choice we were inspired by (c.f. table 2 in Kamir et al. (2020)). However, we
 679 deliberately omit any statistic that involves the minimum (e.g. the NDVI-range), since we
 680 regard the minimum as a very error-prone measure due to the large influence of clouds in
 681 the time series.

682 As a result, for each interpolation strategy, a matrix is obtained in which each row corre-
 683 sponds to a pixel and both the yield and the covariates (computed by applying the above

¹⁴We will use the matrixes derived in section 4.5

¹⁵We 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

684 statistics) are contained. Using this matrix, we train a random forest for yield estimation,
685 and compute the integrated OOB estimates¹⁶ \hat{y} . Note that the choice of the modeling
686 approach does not matter much, as long as it is general enough (i.e. able to approximate
687 any function) and we use the same one for each interpolation strategy. Finally, for each
688 interpolation strategy, we calculate the RYEA and describe the results in section 5.2.

¹⁶By 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 $\frac{1}{e}$).

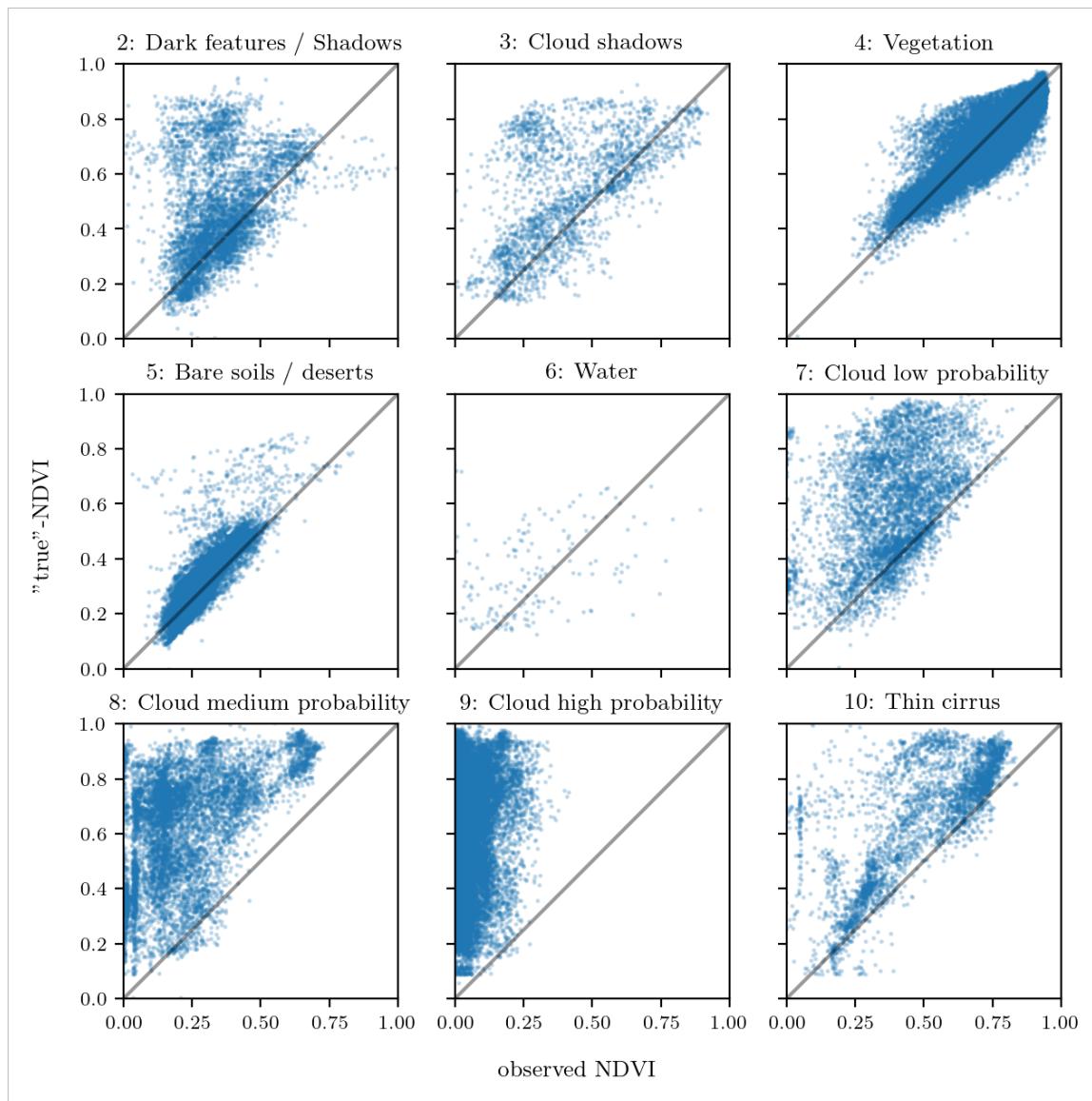


Figure 4.2: For each SCL class, we compare the true NDVI with the observed NDVI. (The true NDVI was estimated with LOOCV smoothing splines, and we used all observations of 10% of the total pixels.)

689 **Chapter 5**

690 **Results**

691 **5.1 Goodness of Fit for Selected Interpolation Methods**

692 Table 5.1 benchmarks the selected¹ interpolation methods (on P^{SCL45}) with respect to
693 various score functions. The score functions take the absolute values of the LOOCV
694 residuals and summarize them in a number (the smaller, the better). For each of the 5
695 selected interpolation methods, we consider the basic and the robustified (see section 3.5)
696 version.

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

	SS	LOESS	DL	BSPL	FR	SS^{rob}	$\text{LOESS}^{\text{rob}}$	DL^{rob}	$BSPL^{\text{rob}}$	FR^{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

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

703 **5.2 XXX (Robustification and) NDVI-Correction**

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

705 The RYEA for the 28 (in section 4.4) chosen interpolation strategies is given in table 5.2.
706 Robustification in the interpolation strategies, does not improve the quality of the fit

¹ 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

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

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

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

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

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

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

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

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

721 **Chapter 6**

722 **Discussion**

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

724 **6.1 Interpolation Methods**

725 **6.1.1 Data Gaps in Time Series**

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

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

732 Since the SG requires equidistant points, it follows that data gaps will break it. The
733 linear interpolation, that is supposed to recover this, we consider as not being a satisfying
734 solution.

735 We do not trust the FR interpolation if there are noticeable data gaps. On the one hand,
736 it corresponds to our experience that the curve can escape strongly there (c.f. figure
737 3.1). On the other hand, the unreliability is illustrated by the poor values in table 5.1 for
738 the robustified variant. These are meaningful in describing the ability to cope with data
739 gaps, since more data points are ignored during the robustification and thus data gaps are
740 simulated.

741 Similarly, for SS, LOESS, DL and BS we compare the values in table 5.1 between the
742 robustified and non-robust variant. We find that the robust variant does not differ strongly
743 from the non-robust variant (unlike FR). Thus, we conclude that these methods do not
744 have systematic failures.

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

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

753 6.1.2 Preselection

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

760 6.1.3 Candidate Selection

761 Given that DL convinces regarding most of the selected score functions in table 5.1 we will
 762 certainly investigate this method in chapter 4. Moreover, we see that the robustification
 763 mostly improved the score regarding the 50, 75, 85, and 90 % Quantiles. Only for the
 764 outlier-sensitive score functions (RMSE and q95)¹ we notice significant worsening (we
 765 consider the robust FS separately in section 6.1.1). Consequently, we will also use the
 766 robustification in section 4. Not wanting to rely on the form assumptions of the DL, we
 767 further choose a non-parametric method for further consideration. Despite the LOESS
 768 slightly dominating the SS in table 5.1, we choose the SS. This is due to the strange
 769 behavior of the LOESS in case of data gaps (see section 6.1.1) and the good interpretability
 770 of the SS using the minimization function 3.3.6.1.

771 XXX discuss results from table B.1

772 6.2 NDVI Correction

773 6.2.1 Bootstrap

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

780 6.2.2 Using Additional Covariates

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

786 We decided against using this data, because on the one hand we have the problem that
 787 we have practically too few observations (we observe only 5 years) and we expect the
 788 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’?

¹For 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.

789 that the weather data published by Meteoswiss are for a grid with a resolution of 1 km.
 790 On the other hand, we want the underlying model not to learn improper relationships.
 791 For example, the model might automatically predict a high NDVI for a day in summer
 792 (detected by high GDD / many sunshine hours / high temperature) just because it is
 793 “used” to observing a lot of vegetation in summer. Including temporally (e.g., P_{t-1} and
 794 P_{t+1}) and geographically adjacent pixels would likely improve performance. However, for
 795 simplicity, we omit it here².

796 6.2.3 Choose Interpolation Strategy

797 table mit OLS SCL als sieger diskutieren

798 if we use no-correctionXss-rob instead of OLS-SCLXss we loose $(0.148 - 0.14)/0.148 =$
 799 5,4% of the information.

800 6.2.4 High RMSE in Yield Prediction

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

- 803 i.) Uncertainty in Yield data collected by the combine harvester
- 804 ii.) Uncertainty in Yield data through rasterization
- 805 iii.) Uncertainty in satellite images through “measurement errors” introduced via clouds
806 and other atmospheric effects
- 807 iv.) Uncertainty introduced by interpolating (especially when long data-gaps are present)

808 even in a perfect world the NDVI curve only holds a fraction of the information
avialbe

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

809 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?

²This 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).

810 **Chapter 7**

811 **Conclusion**

812 In dieser Thesis haben wir studiert, wie wir mit aus Satellitenbildern das Pflanzenwach-
813 stum via NDVI-Zeitreihen modellieren können. Die grösste Herausforderungen waren hi-
814 erbei die fragen, wie man mit (durch Wolken oder Schatten) verfälschten Beobachtungen
815 umgehen soll und wie man die einzelnen Beobachtungen zu interpolieren hat. Für eine
816 zusammenfassung der betrachteten interpolationsmethoden verweisen wir auf die Tabelle
817 **3.**

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

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

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

844 traditionelle Antwort, dass wir nur Beobachtungen beachten, welche als Vegetation oder
845 als bare soil gelabelt sind (SCL45). Dies wird mit der von der European Space Agency
846 gelieferten ‘Scene Classification Layer’ (SCL) bewerkstelligt. In figure 2.3 wird jedoch die
847 Unzuverlässigkeit dieses Labelings illustriert. Zudem haben wir die festgestellt, dass auch
848 nicht-SCL45 Beobachtungen wertvolle Informationen enthalten seien können (vgl. Sektion
849 4.1). Wir haben uns entschieden, nicht an der traditionellen (SCL-)Filtration festzuhalten.
850 Stattdessen betrachten wir alle Beobachtungen und korrigieren den beobachteten NDVI.
851 Dafür benutzen wir statistische Modelle, die zusätzliche Informationen wie die verbleiben-
852 den Spektralbänder in Betracht nehmen. Bevor wir aber die korrigierten NDVI Werte
853 interpolieren, weisen wir jeder Beobachtung ein Gewicht zu, korrespondieren zu ihrer Un-
854 sicherheit. Die Unsicherheit wird analog wie die NDVI Korrektur geschätzt. Durch die
855 Wahl verschiedener Interpolationsmethoden (mit und ohne robustifizierung) und statis-
856 tischer Modelle, erhalten wir somit 28 verschiedene Interpolationsstrategien (vgl. Sektion
857 4.4). Um zu beurteilen, welche dieser Interpolationsstrategie am besten ist, machen wir
858 die folgende Annahme: “je besser die Interpolationsstrategie, desto besser kann damit in-
859 terpolierte NDVI Zeitreihe den Ertrag voraussagen”. Überraschender Weise ist die beste
860 Strategie, die mit nicht-robustifizierten SS und dem einfachsten betrachteten statischen
861 Modell, welches nur den beobachteten NDVI und die SCL Klassifizierung benutzt. Let
862 us recapitulate the best interpolation strategy: First, we estimate the “true” NDVI using
863 SS via LOOCV, then obtain the corrected NDVI using the OLS-SCL model (c.f. equa-
864 tion 5.2.0.1). Subsequently, we estimate the absolute error with the OLS-SCL model (c.f.
865 equation 5.2.0.1) and thereby obtain weights which are supposed to reflect the reliability
866 of the corrected NDVI (c.f. equation 4.3.0.1). Finally, we perform a weighted interpolation
867 with SS.

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

873 Auf die Frage welche interpolationsmethode wir schlussendlich empfehlen, wollen wir zwei
874 Fälle betrachten. Wenn es nur darum geht möglichst präzise eine Kurve den daten anzu-
875 passen, empfehlen wir die robustifizierten DL, da diese die LOOCV residuals in den meis-
876 ten fällen minimieren (vgl. tabelle 5.1). Falls wir eine interpolation erhalten wollen
877 die möglichst viele informationen über die pflanze enthält empfehlen wir die SS. Diese
878 empfehlung gilt besonders, falls wir traditionell nur SCL45 beobachtungen betrachten
879 wollen ohne die vorgeschlagene NDVI zu korrigieren. Jedoch empfehlen wir die oben
880 aufgeführte interpolationsstrategie, da uns ansonsten über 5% der informationen aus der
881 NDVI zeitreihe abhanden kommmen (vgl. sektion 6.2.3). Im anbetracht aller Fehlerquellen
882 (c.f. section 6.2.4) und der tatsache dass wir nur die NDVI Zeitreihe betrachten wir die
883 5% als eine solide verbesserung.

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

885 **7.1 Future Work**

886 **7.1.1 Time Series Correction-Interpolation as a General Method**

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

894 **Example: Cloud Correction with Uncertainty Estimation and Interpolation**

895 This generalization can be used in particular for cloud correction. In the same manner as
896 we corrected the NDVI time series in chapter 4, we can correct each spectral band and
897 reunite the corrected bands with the uncertainties. If desired, the time series can also be
898 interpolated before merging as in chapter 4.3. The resulting question would be how well
899 this approach performs.

900 **7.1.2 Minor Improvements**

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

- 903 — **Data:** Method how combine harvester point data has been extrapolated to the grid
904 could possibly be improved.
- 905 — **Data:** For computational reasons, we mostly considered all years and split the data
906 (on the pixel level) randomly into a train/test set. A leave one year out cross
907 validation might yield more accurate results.
- 908 — **Data:** We have not included the spectral bands that have a resolution of 60 m. But
909 precisely these seem to be promising for cloud correction, since they are a proxy of
910 the water (content and form) in the atmosphere.
- 911 — **Data:** Raiyani et al. (2021) presents an Machine Learing approach that supposedly
912 improves the SCL and thus could improve our results that are based on the SCL.
- 913 — **NDVI Correction:** Explore the effect of different link and normalizing functions in
914 section 4.3. Currently we run into the danger of some outer points getting nearly
915 ignored just because one estimated absolute residual for some interior point is close
916 to zero.
- 917 — **NDVI Correction:** Yield is not the only target variable of interest. Other variables
918 like protein content could also be used in section 4.5 for the method evaluation.

919 Bibliography

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

- 952 Perich, G., M. O. Turkoglu, L. V. Graf, J. D. Wegner, H. Aasen, A. Walter, and F. Liebisch
953 (2022, July). Pixel-based crop yield mapping and prediction using spectral indices and
954 neural networks on {Sentinel}-2 time series data.
- 955 Raiyani, K., T. Gonçalves, L. Rato, P. Salgueiro, and J. R. Marques da Silva (2021,
956 January). Sentinel-2 Image Scene Classification: A Comparison between Sen2Cor and
957 a Machine Learning Approach. *Remote Sensing* 13(2), 300.
- 958 Ripley, B. D. and M. Maechler. R: Fit a Smoothing Spline. [https://stat.ethz.ch/R-
959 manual/R-patched/library/stats/html/smooth.spline.html](https://stat.ethz.ch/R-manual/R-patched/library/stats/html/smooth.spline.html).
- 960 Rouse, J. W. (1974, May). Monitoring the vernal advancement and retrogradation (green
961 wave effect) of natural vegetation. Technical Report NASA-CR-139243.
- 962 Savitzky, A. and M. J. E. Golay (1964, July). Smoothing and Differentiation of Data by
963 Simplified Least Squares Procedures. *Analytical Chemistry* 36(8), 1627–1639.
- 964 Schafer, R. W. (2011, July). What Is a Savitzky-Golay Filter? [Lecture Notes]. *IEEE
965 Signal Processing Magazine* 28(4), 111–117.
- 966 Stephen, M. (2021, July). Earth: Multivariate Adaptive Regression Splines.
- 967 Tibshirani, R. (2011). Regression shrinkage and selection via the lasso: A retrospective.
968 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 73(3), 273–
969 282.

970 **Appendix A**

971 **Reproducibility**

972 **A.1 Reproduce Results**

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

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

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

979 **A.2 R-Package**

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

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

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

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

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

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

```
988 1 library(CorrectTimeSeries)
989 2
990 3 # load a list of dataframes, each one describes one pixel with the covariates and
991 # the response
992 4 data(timeseries_list)
993 5 str(timeseries_list[[1]])
994 6
995 7 # Train/Load RF
996 8 train_model_myself <- TRUE
997 9 if (train_model_myself){
998 10   # Add "true" NDVI (or generally the response), by Out-Of-Bag estimation
999 11   timeseries_list <- lapply(timeseries_list, function(df) {
1000 12     df$oob_ndvi <- OOB_est(df$gdd, df$ndvi_observed) # gdd is the time-axis
1001 13     df
1002 14   })
1003 15   # Train correction model
1004 16   formula <- "oob_ndvi ~ B02+B03+B04+B05+B06+B07+B08+B8A+B11+B12+scl_class"
1005 17   RF <- train_RF_with_fromula(formula, timeseries_list, robustify=TRUE)
1006 18 } else {
```

```
1008 19  data(RF_for_NDVI)
1009 20  RF <- RF_for_NDVI
1010 21 }
1011 22
1012 23 # ADD CORRECTION
1013 24 timeseries_list <- lapply(timeseries_list, function(df) {
1014 25   df$corrected_ndvi <- randomForest:::predict.randomForest(RF, df)
1015 26   df
1016 27 })
1017 28
1018 29 # Get interpolation for each timeseries
1019 30 newx <- 1:1000
1020 31 lapply(timeseries_list, function(df){
1021 32   ss <- smoothing_spline(df$gdd, df$corrected_ndvi)
1022 33   predict(ss, newx)$y
1023 34 })
```

Example of how to use the `CorrectTimeSeries` package

1025 **Appendix B**

1026 **Further Material**

1027 **B.1 Data and Methods**

1028 **B.1.1 GDD**

1029 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

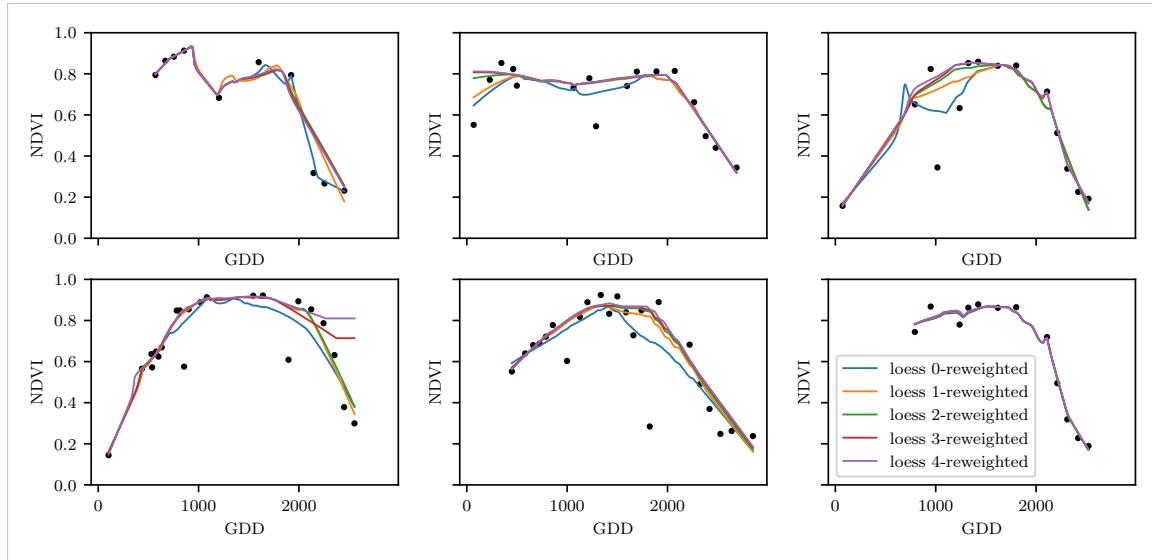
1030 **B.2 Interpolation**

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

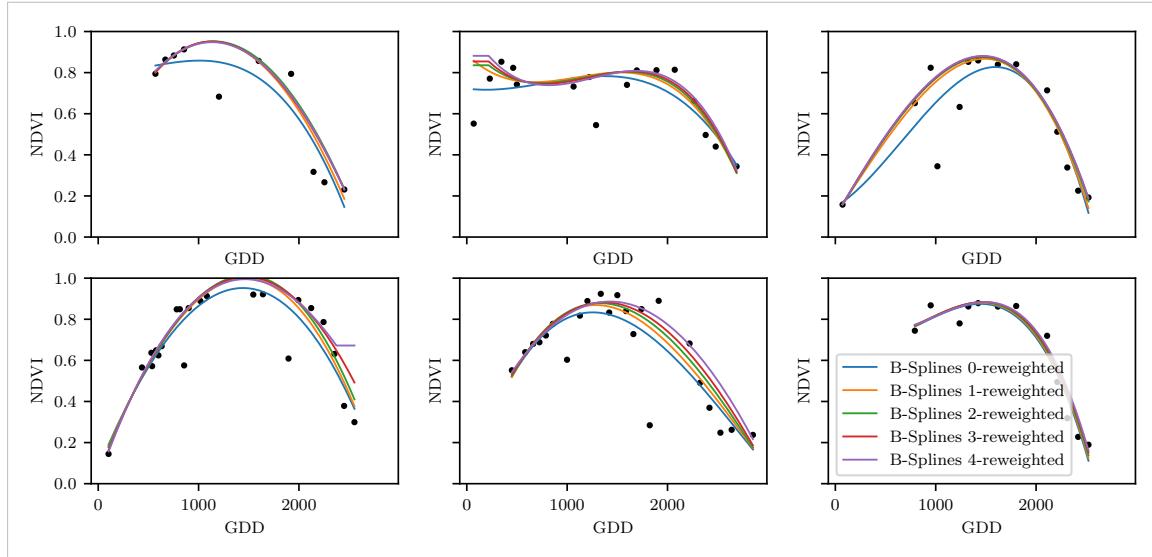


Figure B.2: B-splines fitted to different (SCL45) NDVI time series. Iterations of a robustifying refit (as indicated in section 3.5) are also displayed

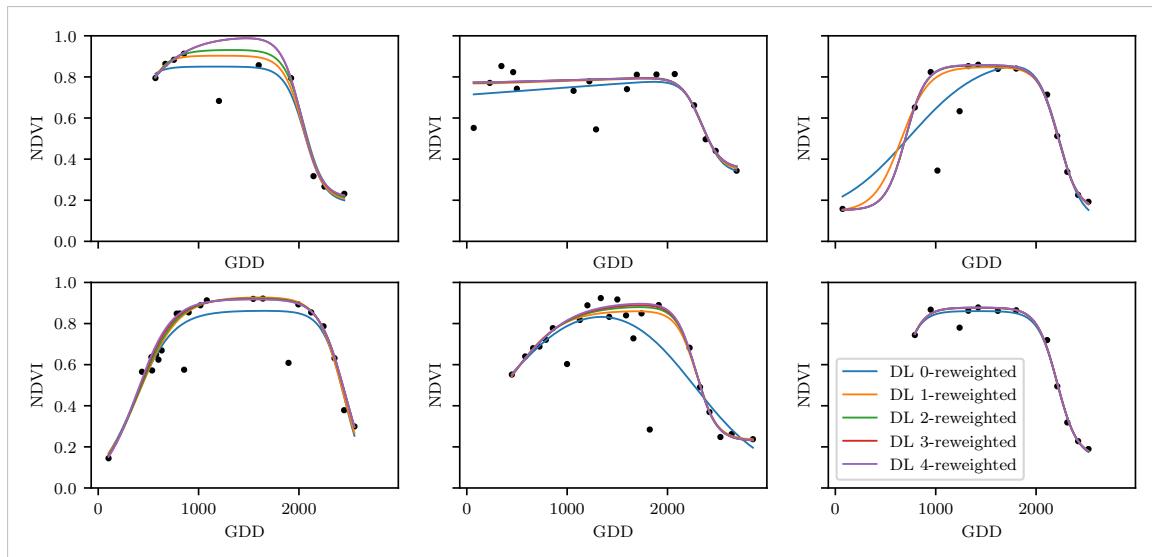


Figure B.3: A Double Logistic curve fitted to different (SCL45) NDVI time series. Iterations of a robustifying refit (as indicated in section 3.5) are also displayed

1031 B.3 NDVI correction

1032 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

1033 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:

```

```

1044      Estimate Std. Error t value Pr(>|t|)
1045 10 (Intercept) 0.21465 0.00230 93.46 < 2e-16 ***
1046 12 ndvi_observed 0.71116 0.00346 205.65 < 2e-16 ***
1047 13 scl_class3 0.02205 0.00356 6.20 5.8e-10 ***
1048 14 scl_class4 -0.00431 0.00251 -1.72 0.085 .
1049 15 scl_class5 -0.09875 0.00234 -42.15 < 2e-16 ***
1050 16 scl_class6 -0.05301 0.01104 -4.80 1.6e-06 ***
1051 17 scl_class7 0.11245 0.00274 41.09 < 2e-16 ***
1052 18 scl_class8 0.25963 0.00253 102.57 < 2e-16 ***
1053 19 scl_class9 0.35994 0.00236 152.47 < 2e-16 ***
1054 20 scl_class10 0.09091 0.00308 29.54 < 2e-16 ***
1055 21 scl_class11 0.29784 0.00392 76.06 < 2e-16 ***
1056 22 ---
1057 23 Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ',' 1
1058 24
1059 25 Residual standard error: 0.146 on 124978 degrees of freedom
1060 26 Multiple R-squared: 0.532, Adjusted R-squared: 0.532
1061 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)

```

1063
1064 1 Call:
1065 2 lm(formula = (paste(get_res(), " ~ ", "ndvi_observed + scl_class")),
1066 3   data = ndvi_df)
1067 4
1068 5 Residuals:
1069 6   Min     1Q   Median     3Q    Max
1070 7 -0.2051 -0.0427 -0.0074  0.0329  0.6589
1071 8
1072 9 Coefficients:
1073      Estimate Std. Error t value Pr(>|t|)
1074 10 (Intercept) 0.18647 0.00126 147.74 < 2e-16 ***
1075 12 ndvi_observed -0.13265 0.00190 -69.80 < 2e-16 ***
1076 13 scl_class3 -0.00180 0.00196 -0.92 0.3587
1077 14 scl_class4 -0.04069 0.00138 -29.55 < 2e-16 ***
1078 15 scl_class5 -0.09698 0.00129 -75.32 < 2e-16 ***
1079 16 scl_class6 -0.01906 0.00606 -3.14 0.0017 **
1080 17 scl_class7 0.01641 0.00150 10.91 < 2e-16 ***
1081 18 scl_class8 -0.00560 0.00139 -4.02 5.7e-05 ***
1082 19 scl_class9 -0.01384 0.00130 -10.67 < 2e-16 ***
1083 20 scl_class10 -0.00690 0.00169 -4.08 4.5e-05 ***
1084 21 scl_class11 -0.01446 0.00215 -6.72 1.8e-11 ***
1085 22 ---
1086 23 Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ',' 1
1087 24
1088 25 Residual standard error: 0.08 on 124978 degrees of freedom
1089 26 Multiple R-squared: 0.352, Adjusted R-squared: 0.352
1090 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)

replace space before ref by tilda
check quantile definitions
schwarz weiss färbung der IS tabelle korrigieren
so wenig wie möglich abkürzungen in den fig und table captions
refer to data aviability
abkürzungen Fourier und in tabellen
figure spacing (caption zu nah dran — manuell vspace einfügen wo nötig)

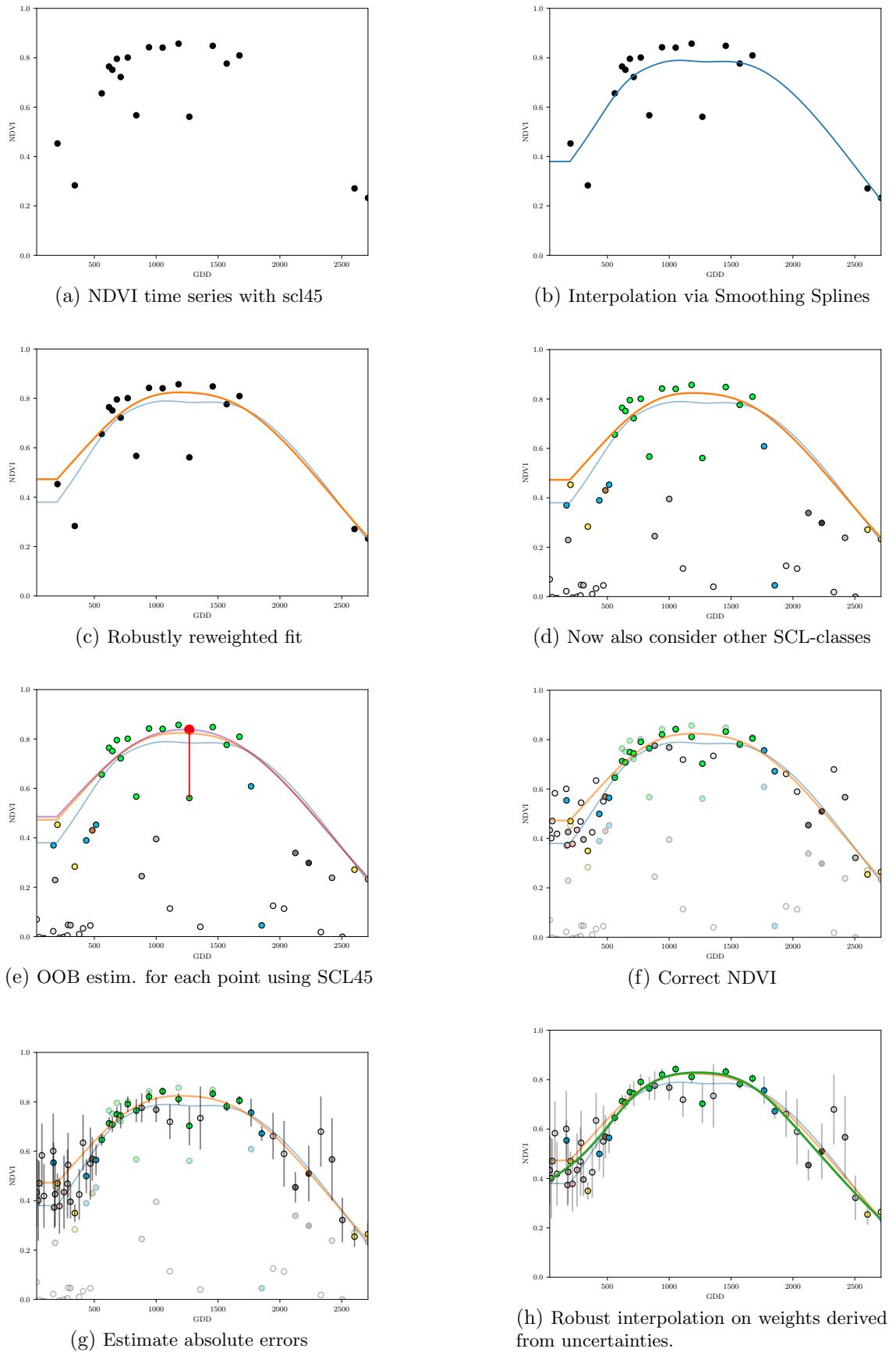


Figure B.4: Stepwise illustration of robust NDVI-Correction. For the color encoding of the SCL classes we refer to table 2.2.