

# **Entwicklung eines Verfahrens zur Vorgabe der benötigten Sintereinstellungen von PTC- Bauteilen aus Daten der Granulatfreigabe und den Ergebnissen vorheriger Lose**

Masterarbeit  
von  
Katharina Mertens, BSc



eingereicht am  
Lehrstuhl Wirtschafts- und Betriebswissenschaften  
der  
Montanuniversität Leoben

Leoben, Mai 2017

# **Aufgabenstellung**

Nur bei Master- und Bachelorarbeiten!

Die Aufgabenstellung ist ein einseitiges, bei Masterarbeiten vom Lehrstuhlleiter, bei Bachelorarbeiten vom Betreuer unterschriebenes Dokument. Die Erstellung erfolgt vom Betreuer und wird dem Verfasser noch vor Abschluss der wissenschaftlichen Arbeit ausgehändigt. Die Aufgabenstellung muss direkt nach dem Titelblatt bei jedem Exemplar der Arbeit in Original (Institutsexemplar) bzw. in Kopie mit eingebunden werden. (Siehe Richtlinie zur Erstellung wissenschaftlicher Arbeiten)

## **Eidesstattliche Erklärung**

Ich erkläre an Eides statt, dass ich diese Arbeit selbständig verfasst, andere als die angegebenen Quellen und Hilfsmittel nicht benutzt und mich auch sonst keiner unerlaubten Hilfsmittel bedient habe.

I declare in lieu of oath, that I wrote this thesis and performed the associated research myself, using only literature cited in this volume.

Leoben, Mai 2017

---

(Katharina Mertens, BSc)

---

## **Gleichheitsgrundsatz**

Aus Gründen der Lesbarkeit wurde in dieser Arbeit darauf verzichtet, geschlechtsspezifische Formulierungen zu verwenden. Es wird ausdrücklich festgehalten, dass die bei Personen verwendeten maskulinen Formen für beide Geschlechter zu verstehen sind.

## Danksagung

An dieser Stelle möchte ich all jenen meinen Dank aussprechen, die durch ihre fachliche und persönliche Unterstützung zum Gelingen dieser Masterarbeit beigetragen haben.

Danken möchte ich o.Univ.-Prof. Dipl.-Ing. Dr. mont. Hubert Biedermann für die Möglichkeit, meine Masterarbeit am Lehrstuhl für Wirtschafts- und Betriebswissenschaften zu verfassen, und das Bereitstellen der Infrastruktur.

Ebenfalls möchte ich mich bei seinem Mitarbeiter Dipl.-Ing. Robert Bernerstätter bedanken. Vielen Dank für die Zeit, die moralische Unterstützung, die konstruktive Kritik und die Anregungen, die diese Arbeit bereichert haben.

Auch muss ich mich bei der EPCOS OHG und meinem Betreuer Dr. Johann Schmidt bedanken, die mit ihrer Problemstellung den Anstoß zu dieser Arbeit gegeben haben. Vielen Dank für die Unterstützung.

Zusätzlich möchte ich mich bei O.Univ.-Prof. Dipl.-Ing. Dr.techn Paul O'Leary und seinen Mitarbeitern vom Lehrstuhl Automation für ihr offenes Ohr, Anregungen und die Unterstützung beim Erlernen von Python und LaTeX bedanken.

Während meiner Zeit an dieser Universität habe ich viele unterschiedliche Menschen kennen und schätzen gelernt. Einige wurden Wegbegleiter, denen ich an dieser Stelle danken möchte.

Mein besonderer Dank gilt meinen Eltern und vor allem meiner Mutter. Ohne sie hätte ich dieses Studium nach dem Tod meines Vaters nicht weiterführen können. Ich danke ihr für ihre emotionale und finanzielle Unterstützung.

Ein großes Danke an meine Familie und Freunde, die immer für mich da waren und mich in meinem Weg bestärkt haben.

## Kurzfassung

Diese Masterarbeit behandelt Möglichkeiten der produktionswirtschaftlichen Verbesserungen durch Data Mining. Anstoß hierfür ist die Problemstellung der EPCOS OHG, die in einem mehrstufigen Produktionsprozess PTC-Bauteile (Positive Temperature Coefficient, Kaltleiter) produziert. Das als Produktionsengpass identifizierte Sintern beeinflusst den elektrischen Widerstand des fertigen Bauteils. Um den gewünschten Widerstand zu erhalten, ist das wiederholte Anpassen der Sintereinstellungen notwendig. Dies führt zu einer schlechten Auslastung der Fertigungskapazität und Verlängerung der Durchlaufzeit. Um den Anpassungsaufwand zu reduzieren, wird ein Vorhersagemodell für den Widerstandswert auf Basis der vom Unternehmen bereitgestellten Daten erstellt. Die Daten beziehen sich auf verschiedene Produktionsabschnitte und Freigabemessungen innerhalb des Produktionsprozesses, die in zwei unterschiedlichen Standorten stattfinden.

Der Prozessstandard CRISP-DM (Cross Industry Standard Process for Data Mining) bietet den Rahmen für die Durchführung der Arbeit. Für die Datenanalyse werden Methoden wie die Berechnung statistischer Größen, die Ausreißeranalyse, die Analyse der fehlenden Daten, die Gruppierung und die Dimensionalitätsreduktion, verwendet. Zur Vorhersage des Widerstandswerts finden die Methoden der Klassifikation Anwendung, zu denen Entscheidungs bäume, Diskriminanzanalyse, Support Vector Machines, k-Nearest Neighbor und Ensemblemethoden zählen. Das Ergebnis des praktischen Teils sind vier aufbereitete und analysierte Datensets, je eines mit und ohne Ausreißer je Standort. Die Evaluation erfolgt durch die Anwendung der Modelle auf unabhängige Testdaten und dem Vergleich der Resultate. Die nicht zufriedenstellenden Ergebnisse beruhen auf den zu kleinen Verhältnissen zwischen der Anzahl der Datensätze und der Anzahl der Attribute. Somit kann keine Implementierung innerhalb des Unternehmens erfolgen. Allerdings würde ein funktionierendes Modell zu 13 % freier Kapazität im betrachteten Zeitraum führen und hätte eine Verkürzung der Durchlaufzeit zur Folge. Um diese in der Zukunft mit Hilfe des Data Minings auszunutzen, werden am Ende der Arbeit Möglichkeiten zur Verbesserung der Datenbasis und der Datenaufbereitung präsentiert.

## Abstract

This master thesis deals with production-related improvements by data mining. The problem is posed by EPCOS OHG, which produces PTC (Positive Temperature Coefficient) components in a multi-stage production process. Sintering is the bottleneck of production and influences the electrical properties (resistance value) of the finished component. In order to obtain the desired resistance, repeated adjustments of the sintering settings are necessary. These setting times lead to a poor utilization of the production capacity and an increase in the lead time. For this reason, a prediction model for the resistance value is established based on the data provided by the company. The data comes from two different production sites which include different production stages and release measurements.

The process standard CRISP-DM (Cross Industry Standard Process for Data Mining) provides the framework for this thesis. To analyse data methods like calculating statistical parameters, outlier analysis and analysis of missing data, grouping and dimensionality reduction are applied to the data. For the prediction of the resistance value, methods of classification are used, including decision trees, discriminant analysis, support vector machines, k-nearest neighbor and ensemble methods. The results of the practical part are four prepared and analysed data sets, one each with and without outliers per location. The models are applied to independent test data and their results are compared. The unsatisfactory output is the result of small ratios between the number of attribute values and the number of attributes. Therefore, the implementation within the company won't be promising. A working model results in 13 % free capacity during the period in question and shortens the lead time. In order to utilise this in the future while using data mining, possibilities to improve the data and the data processing are presented at the end of the thesis.

# Inhaltsverzeichnis

<b>1</b>	<b>Einleitung</b> .....	<b>1</b>
1.1	Ausgangssituation und Problemstellung .....	1
1.2	Zielsetzung und Forschungsfrage .....	3
1.3	Vorgehensweise und Aufbau der Arbeit .....	3
<b>2</b>	<b>Produktionswirtschaftliche Theorie</b> .....	<b>5</b>
2.1	Das Dilemma der Ablaufplanung .....	6
2.2	Durchlaufzeit .....	8
2.3	Kapazität .....	12
2.3.1	Einführung .....	12
2.3.2	Anpassungen der Kapazität .....	15
<b>3</b>	<b>Datenanalyse</b> .....	<b>20</b>
3.1	Einführung .....	20
3.2	Cross-Industry Standard Process for Data Mining .....	25
3.3	Datenverständnis .....	27
3.3.1	Attributtypen .....	27
3.3.2	Univariate Deskription .....	29
3.3.3	Multivariate Deskription .....	38
3.4	Datenaufbereitung .....	42
3.4.1	Ausreißeranalyse .....	42
3.4.2	Analyse der Fehlenden Daten .....	43
3.4.3	Gruppierung .....	47
3.4.4	Dimensionalitätsreduktion .....	48
3.5	Modellierung .....	51
3.5.1	Einführung .....	52
3.5.2	Methoden .....	53
<b>4</b>	<b>Praktische Fallstudie</b> .....	<b>57</b>
4.1	Einführung .....	57
4.2	Standort Deutschlandsberg .....	60
4.2.1	Analyse der Fehlenden Daten .....	61
4.2.2	Datenzusammenführung .....	63
4.2.3	Auswertung .....	64
4.2.4	Ausreißeranalyse .....	66
4.2.5	Gruppierung .....	69

---

4.2.6	Eignung für die PCA.....	70
4.2.7	Modellierung .....	71
4.3	Standort Šumperk .....	82
4.3.1	Aufbereitung des Zielattributs.....	83
4.3.2	Datenzusammenführung .....	84
4.3.3	Auswertung .....	86
4.3.4	Eignung für PCA .....	90
4.3.5	Modellierung .....	90
4.4	Produktionswirtschaftliche Betrachtung.....	97
<b>5</b>	<b>Zusammenfassung und Ausblick.....</b>	<b>99</b>

## Abbildungsverzeichnis

Abbildung 1: Produktionsprozess der PTC-Bauteile .....	2
Abbildung 2: Konzept zur Lösung der Forschungsfragen .....	3
Abbildung 3: Aufbau der Arbeit.....	4
Abbildung 4: Produktion als Kombinations- und Transformationsprozess.....	5
Abbildung 5: Logistische Erfolgsfaktoren von Produktionsunternehmen.....	6
Abbildung 6: Gewichtung von logistischen Zielgrößen.....	7
Abbildung 7: Durchlaufzeitanteile und Durchlaufelement.....	8
Abbildung 8: Berechnung des OEE .....	14
Abbildung 9: Verlauf von Nutz- und Leerkosten.....	15
Abbildung 10: Kostenremanenz: (1) Hysteresis-Schleife, (2) remanenter Sprung .....	16
Abbildung 11: Gesamtkostenverlauf beim Übergang.....	18
Abbildung 12: Lohnkostenverlauf bei zeitlicher Anpassung .....	19
Abbildung 13: Komponenten von Data Science.....	21
Abbildung 14: Qualitative Unterscheidung verschiedener Ansätze zur Datenanalyse..	22
Abbildung 15: Die evolutionstechnischen Stufen der Analytik.....	23
Abbildung 16: Die sechs Phasen des CRISP-DM.....	25
Abbildung 17: Aufgaben und Output des CRISP-DM Modells .....	26
Abbildung 18: Beispiel eines Histogramms.....	30
Abbildung 19: Eine linkssteile (a), symmetrische (b) und rechtssteile Verteilung (c) ....	31
Abbildung 20: Beispiel eines Boxplots .....	33
Abbildung 21: Beispiel für eine Dichtefunktion in einem Histogramm .....	35
Abbildung 22: Zwei Normalverteilungsdichtekurven mit kleinem und großem $\sigma$ .....	36
Abbildung 23: Beispiel eines Streudiagramms.....	38
Abbildung 24: Punktekonfiguration und Korrelationskoeffizienten.....	39
Abbildung 25: Extremfälle für Spearmans Korrelationskoeffizienten .....	41
Abbildung 26: Prozess zum Identifizieren und Bereinigen fehlender Daten .....	44
Abbildung 27: Beispiel eines Entscheidungsbaums.....	54
Abbildung 28: Mögliche Platzierungen der Trennebene .....	55
Abbildung 29: Produktionsprozess der PTC-Bauteile .....	58
Abbildung 30: Anzahl einfließender Chargen je Produktionsabschnitts .....	58
Abbildung 31: Aufbau der Ausgangsdaten .....	59
Abbildung 32: R( $\vartheta$ )-Kennlinie eines Kaltleiters (schematisch).....	59
Abbildung 33: Beispiele für Histogramme für DL mit Ausreißer .....	64
Abbildung 34: Histogramm für R25 Vergl.% mit Ausreißern .....	65

Abbildung 35: Beispiele für Streudiagramme für DL mit Ausreißer .....	65
Abbildung 36: Korrelationsmatrix (graphisch) für DL mit Ausreißern.....	66
Abbildung 37: Beispiele für Histogramme für DL ohne Ausreißer .....	67
Abbildung 38: Histogramm für R25 Vergl.% ohne Ausreißer .....	67
Abbildung 39: Beispiele für Streudiagramme für DL ohne Ausreißer .....	68
Abbildung 40: Korrelationsmatrizen (graphisch) für DL mit und ohne Ausreißer .....	68
Abbildung 41: Korrelationsmatrix der VM+NM.....	74
Abbildung 42: Vergleich verschiedener Kombinationsmöglichkeiten .....	77
Abbildung 43: Ergebnisse des CL für VM+NM mit Ausreißern.....	79
Abbildung 44: Simple Tree für Datenset VM+NM mit Ausreißern .....	80
Abbildung 45: Ergebnisse des CL für VM+NM ohne Ausreißer .....	81
Abbildung 46: Simple Tree für Datenset VM+NM ohne Ausreißer .....	82
Abbildung 47: Histogramm für gemittelten gesinterten Widerstand (MW) .....	83
Abbildung 48: Streudiagramm für R25 Vergl.% und MW .....	83
Abbildung 49: Aufbau der Ausgangsdaten für Šumperk .....	85
Abbildung 50: Beispiele für Histogramme für Šumperk mit Ausreißern.....	86
Abbildung 51: Beispiele für Streudiagramme für Šumperk mit Ausreißern.....	87
Abbildung 52: Korrelationsmatrix (graphisch) für Šumperk mit Ausreißern .....	87
Abbildung 53: Beispiele für Histogramme für Šumperk ohne Ausreißer.....	88
Abbildung 54: Beispiele für Streudiagramme für Šumperk ohne Ausreißer .....	89
Abbildung 55: Korrelationsmatrizen (graphisch) für Šumperk mit und ohne Ausreißer .....	89
Abbildung 56: Korrelationsmatrizen für Šumperk.....	91
Abbildung 57: Vergleich verschiedener Kombinationsmöglichkeiten für Šumperk .....	93
Abbildung 58: Ergebnisse des CL für Šumperk mit Ausreißern .....	94
Abbildung 59: Simple Tree für Šumperk mit Ausreißern .....	95
Abbildung 60: Ergebnisse des CL für Šumperk ohne Ausreißer .....	96
Abbildung 61: Simple Tree für Šumperk ohne Ausreißer .....	97
Abbildung 62: CT: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern).....	q
Abbildung 63: MT: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern) .....	q
Abbildung 64: ST: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern).....	r
Abbildung 65: LSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern).....	r
Abbildung 66: QSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern).....	s
Abbildung 67: CSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern).....	s
Abbildung 68: FGSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern).....	t
Abbildung 69: MGSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern).....	t
Abbildung 70: CGSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern) .....	u
Abbildung 71: BoT: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern).....	u
Abbildung 72: BaT: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern).....	v

Abbildung 73: RUSBT: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern) .....	v
Abbildung 74: CT: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer) .....	aa
Abbildung 75: MT: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer) .....	aa
Abbildung 76: ST: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer) .....	bb
Abbildung 77: LSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer) .....	bb
Abbildung 78: QSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer) .....	cc
Abbildung 79: CSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer) .....	cc
Abbildung 80: FGSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer) .....	dd
Abbildung 81: MGSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer) .....	dd
Abbildung 82: CGSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer) .....	ee
Abbildung 83: BoT: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer) .....	ee
Abbildung 84: BaT: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer) .....	ff
Abbildung 85: RUSBT: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer) .....	ff
Abbildung 86: CT: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern) .....	qq
Abbildung 87: MT: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern) .....	qq
Abbildung 88: ST: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern) .....	rr
Abbildung 89: LSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern) .....	rr
Abbildung 90: QSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern) .....	ss
Abbildung 91: CSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern) .....	ss
Abbildung 92: FGSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern) .....	tt
Abbildung 93: MGSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern) .....	tt
Abbildung 94: CGSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern) .....	uu
Abbildung 95: BoT: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern) .....	uu
Abbildung 96: BaT: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern) .....	vv
Abbildung 97: RUSBT: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern) .....	vv
Abbildung 98: CT: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer) .....	yy
Abbildung 99: MT: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer) .....	zz
Abbildung 100: ST: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer) .....	zz
Abbildung 101: LSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer) .....	aaa
Abbildung 102: QSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer) .....	aaa
Abbildung 103: CSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer) .....	bbb
Abbildung 104: FGSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer) .....	bbb
Abbildung 105: MGSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer) .....	ccc
Abbildung 106: CGSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer) .....	ccc
Abbildung 107: BoT: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer) .....	ddd
Abbildung 108: BaT: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer) .....	ddd

Abbildung 109: RUSBT: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer).....eee

## Tabellenverzeichnis

Tabelle 1: Kenngrößen der Produktionslogistik .....	11
Tabelle 2: Ausprägungen des Kapazitätsbegriffs.....	13
Tabelle 3: Beispiele für die Aufgaben der Datenanalyse .....	24
Tabelle 4: Vergleich der Schritte von KDD, SEMMA und CRISP-DM .....	26
Tabelle 5: Skalenniveau .....	28
Tabelle 6: Einstufung des MSA-Kriterium hinsichtlich der Eignung für die PCA.....	51
Tabelle 7: Anwendbarkeit auf Attribute unterschiedlicher Skalenniveaus .....	53
Tabelle 8: Veränderung der Anzahl und Prozente der fehlenden Werte (VM).....	61
Tabelle 9: Veränderung der Anzahl und Prozente der fehlenden Werte (NM) .....	62
Tabelle 10: Einteilung von R25 Vergl.% in Klassen .....	69
Tabelle 11: Entwicklung der Anzahl der Attribute (VM+NM mit Ausreißern) .....	70
Tabelle 12: MSA für Datenset (VM+NM mit Ausreißern) .....	70
Tabelle 13: Entwicklung der Anzahl der Attribute (VM+NM ohne Ausreißer) .....	71
Tabelle 14: MSA für Datenset (VM+NM ohne Ausreißer) .....	71
Tabelle 15: Übersicht über Methoden in Matlab (1) .....	72
Tabelle 16: Übersicht über Methoden in Matlab (2) .....	73
Tabelle 17: Auswertung der Testdaten (Anzahl je Klasse).....	75
Tabelle 18: Kombinationen für VM+NM.....	76
Tabelle 19: Einteilung von MW in Klassen.....	84
Tabelle 20: Entwicklung der Anzahl der Attribute (Šumperk mit Ausreißern) .....	90
Tabelle 21: MSA für Datenset (Šumperk mit Ausreißern) .....	90
Tabelle 22: Auswertung der Testdaten (Anzahl je Klasse).....	91
Tabelle 23: Kombinationen für Šumperk.....	92
Tabelle 24: Auswertung der Bestückung des Ofens .....	98
Tabelle 25: Verhältnisse Datensätze zu Attributen für VM+NM .....	100
Tabelle 26: Verhältnisse Datensätze zu Attribute für Šumperk .....	100
Tabelle 27: VM: Einteilung der Attribute in MAR und MCAR (fehlenden Daten in %).....	b
Tabelle 28: NM: Einteilung der Attribute in MAR und MCAR (fehlenden Daten in %) .....	c
Tabelle 29: Attribute nach der Zusammenführung (VM+NM).....	d
Tabelle 30: Statistische Maße der VM+NM (mit Ausreißern) .....	e
Tabelle 31: Ausreißer in Anzahl und Prozent je Attribut.....	h
Tabelle 32: MSA für einzelne Attribute (VM+NM mit Ausreißern) .....	k
Tabelle 33: MSA für einzelne Attribute (VM+NM ohne Ausreißer) .....	n
Tabelle 34: Ergebnis CL und Testdaten: VM+NM mit Ausreißern.....	p

Tabelle 35: Ergebnis CL und Testdaten: VM+NM ohne Ausreißer .....	z
Tabelle 36: Statistische Maße für Šumperk (mit Ausreißern).....	kk
Tabelle 37: MSA für einzelne Attribute (Šumperk mit Ausreißern).....	nn
Tabelle 38: Ergebnis CL und Testdaten: Šumperk mit Ausreißern .....	pp
Tabelle 39: Ergebnis CL und Testdaten: Šumperk ohne Ausreißer .....	yy

## Abkürzungsverzeichnis

AVG	Arbeitsvorgang
BaT	Bagged Trees
BG	Beschäftigungsgrad
BoT	Boosted Trees
bzw.	beziehungsweise
CGSVM	Coarse Gaussian Support Vector Machine
CL	Classification Learner
CRISP-DM	Cross-Industry Standard Process for Data Mining
CSVM	Cubic Support Vector Machine
CT	Complex Tree
d. h.	das heißt
DIN	Deutsches Institut für Normung
DL	Deutschlandsberg
ebda	ebenda
engl.	englisch
etc.	et cetera
exkl.	exklusive
FGSVM	Fine Gaussian Support Vector Machine
ggf.	gegebenenfalls
IDC	International Data Corporation
inkl.	inklusive
IQR	Interquartilsabstand
IT	Informationstechnik
KDD	Knowledge Discovery in Databases
KMO	Kaiser-Meyer-Olkin
KNN	$k$ -Nächste-Nachbarn
LSVM	Linear Support Vector Machine
MAR	Missing At Random
MCAR	Missing At Completely At Random
MGSVM	Medium Gaussian Support Vector Machine
MSA	measure of sampling adequacy
MT	Medium Tree
MW	mittlerer gesinterter Widerstand
NM	Nachmahlung

PCA	Hauptkomponentenanalyse (Principal Component Analysis)
PCs	Hauptkomponenten (Principal Components)
PTC	Positive Temperature Coefficient
QSVM	Quadratic Support Vector Machine
RUSBT	RUSBoosted Trees
SEMMA	Sample, Explore, Modify, Model, Assess
ST	Simple Tree
SVD	Singulärwertzerlegung (Singular Value Decomposition)
SVM	Support Vector Machine
TBE	Bearbeitungsende
TBEV	Bearbeitungsende Vorgänger
TRA	Rüstanfang
u. a.	unter anderem
UD	Ofen
VM	Vormahlung
z. B.	zum Beispiel
ZDF	Durchführungszeit
ZDL	Durchlaufzeit
ZUE	Übergangszeit

# 1 Einleitung

Das erste Kapitel dieser Masterarbeit betrachtet die Motivation und die Vorgehensweise zur Erstellung dieser Arbeit. Diese gliedern sich in die Beschreibung der Ausgangssituation und die damit verbundene Problemstellung, in die daraus ableitbare Zielsetzung, in die methodische Vorgehensweise und in die Gliederung der gesamten Arbeit. Das folgende Kapitel hat das Ziel, einen Überblick über die Inhalte und Erwartungen an diese Arbeit zu geben.

## 1.1 Ausgangssituation und Problemstellung

Anstoß zu dieser Masterarbeit ist die Problemstellung des Unternehmen EPCOS OHG. Die EPCOS OHG entwickelt, fertigt und vertreibt elektronische Bauelemente und gehört zur TDK Corporation, einem Elektronikonzern mit Sitz in Japan. Der Standort in Deutschlandsberg ist der größte Standort in Europa und das Kompetenzzentrum für keramische Bauelemente der Division Piezo and Protection Device Business Group. Von hier aus werden mehrere Standorte weltweit gesteuert und unterstützt.<sup>1</sup> Dazu gehört auch der Standort in Šumperk, Tschechien, der gemeinsam mit dem Standort in Deutschlandsberg an der Produktion von PTC-Widerständen (Kaltleiter, engl. Positive Temperature Coefficient), im Weiteren mit PTC-Bauteile bezeichnet, beteiligt ist. PTC-Bauteile sind elektronische Bauelemente, deren Widerstand mit zunehmender Temperatur ansteigt. Aus diesem Grund werden sie sowohl als Heizelemente als auch als Temperaturfühler eingesetzt.<sup>2</sup>

Das Erreichen des spezifizierten Sollwert für den gesinterten Widerstand ( $R_{25}$ ) ist derzeit in der Produktion nicht ohne mehrere Versuche möglich. Entlang des im Abbildung 1 dargestellten mehrstufigen Produktionsprozesses der Herstellung von PTC-Bauteilen gibt es daher mehrere Freigaben, die den  $R_{25}$  überprüfen, um ggf. eine Anpassung vornehmen zu können.

---

<sup>1</sup> Vgl. EPCOS HR (2016), S. 3

<sup>2</sup> Vgl. Reisch, M. (2007), S. 298 f.



**Abbildung 1: Produktionsprozess der PTC-Bauteile<sup>3</sup>**

Die Granulatproduktion in Deutschlandsberg (in Abbildung 1 dunkelblau hinterlegt) beinhaltet das Mischen der Komponenten, die in mehreren Schritten zu einem Granulat (Grundmasse) verarbeitet werden, das abschließend mit einer Granulatfreigabe überprüft wird. Dieses Granulat wird anschließend nach Šumperk transportiert. Die in Šumperk stattfindenden Prozessschritte sind in Abbildung 1 hellblau hinterlegt. Hier wird entweder aus mehreren Grundmassen eine neue Mischung homogenisiert oder die Grundmasse direkt weiterverarbeitet. Sollte eine Mischung entstehen, erfolgt eine Mischmassenfreigabe, wenn nicht, werden sowohl Mischung als auch Grundmasse bei der Sinterfreigabe erneut überprüft. Hierbei werden die für das Entbinden und Sintern verwendeten Durchstoßöfen unmittelbar vor der tatsächlichen Produktion bestückt. Allerdings wird dabei zunächst nur jeder fünfzigste Platz ausgelastet, um eine Fehlproduktion zu vermeiden. Anschließend wird der Ofen abhängig vom Resultat eingestellt und es gibt laufende Kontrollen der Sinterqualität. Für Granulat- und Mischmassenfreigabe werden aus den Massen standardisierte Scheiben gepresst, die sich in ihrem Verhalten von dem fertigen Bauteil aufgrund ihrer Geometrie unterscheiden können. Die Sinterfreigabe erfolgt mit der Geometrie des Endproduktes. Alle drei Freigaben unterscheiden sich in der Bearbeitung und Messung der PTC-Bauteile von der Endfreigabe (D9-Freigabe), da die Messteile nicht den gesamten Prozess durchlaufen und unter anderen Voraussetzungen behandelt werden. Sie dauern deutlich länger, da die Sinterung bei der D9-Freigabe bereits abgeschlossen ist.

Die Freigabemessungen, wobei die Sinterfreigabe mehrmals wiederholt werden muss, um die Einstellungen entsprechend anzupassen, sind zeitaufwendig. Durch diese Wiederholungen werden die Fertigungskapazitäten schlecht ausgelastet und die Durchlaufzeit des Fertigungsprozesses wird in die Länge gezogen.

<sup>3</sup> Quelle: in Anlehnung an EPCOS PPD PTC PD (2016), S. 10

## 1.2 Zielsetzung und Forschungsfrage

Ziel dieser Arbeit ist es, die im Prozess erhobenen Werte für einzelne Produktionschargen, wie Komponenten und Prozessparameter, in Bezug zu den jeweiligen bei den Freigaben gemessenen Zielwerten  $R_{25}$  zu setzen. Dies soll dazu führen, Vorhersagen für künftige Lose treffen zu können und den Prozess des Sinterns dementsprechend abzuändern. Zusätzlich soll dies einen positiven Effekt auf die Kapazitätsauslastung haben, da die Sinterfreigabe nicht wiederholt werden müsste bzw. ganz auf sie verzichtet werden und somit direkt mit der Sinterung begonnen werden könnte.

Aus diesen Zielen und der in Kapitel 1.1 erläuterten Ausgangssituation können folgende Forschungsfragen abgeleitet werden:

- Ist es möglich, den Widerstandswert  $R_{25}$  der PTC-Bauteile für zukünftige Lose auf Basis der bereits gefertigten vorherzusagen?
- Wie wirkt sich eine exakte Vorhersage auf die Kapazitätsauslastung aus?

## 1.3 Vorgehensweise und Aufbau der Arbeit

Um die Forschungsfragen aus Kapitel 1.2 bestmöglich zu beantworten, wurden verschiedene Schritte unternommen. Nach einer Besichtigung der Produktionsstandorte Deutschlandsberg und Šumperk wurden die gesammelten Daten gesichtet und eine Literaturrecherche hinsichtlich möglicher bereits anderorts verwendeter Ansätze durchgeführt. Aus dieser entstand die Grundlage für die Vorgehensweise zur Lösung des Problems. Dieses erste Konzept ist in Abbildung 2 ersichtlich. Durch das Schaffen einer einheitlichen Datenbasis, einer anschließenden Dimensionalitätsreduktion und durch Anwenden einer oder mehrerer Methoden aus dem Data Mining werden Modelle zur Vorhersage der Zielwerte der Granulat- und Sinterfreigabe entwickelt.

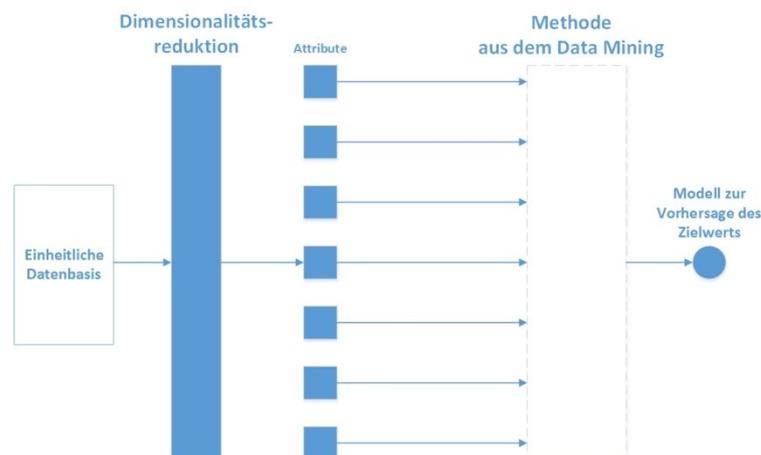
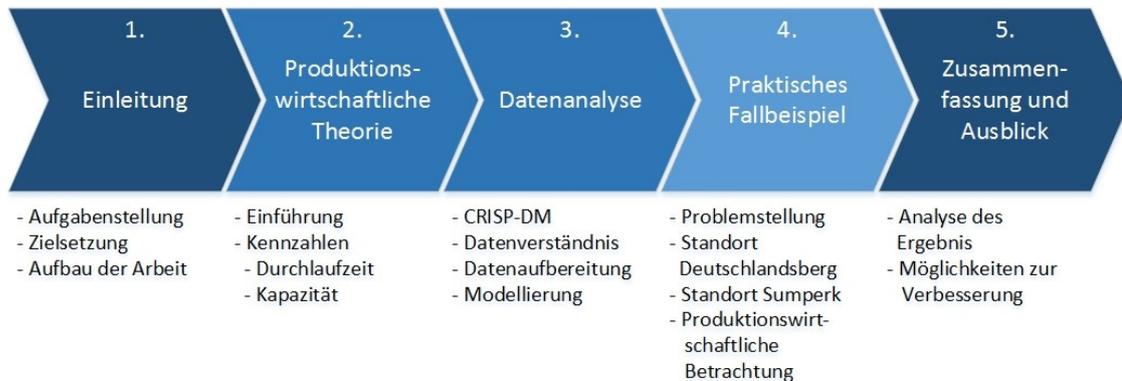


Abbildung 2: Konzept zur Lösung der Forschungsfragen<sup>4</sup>

<sup>4</sup> Quelle: in Anlehnung an Yuan, B. et al. (2000), S. 679

An diesem Konzept orientiert sich der Aufbau der Arbeit, der in Abbildung 3 dargestellt wird. Dunkelblaue Elemente kennzeichnen die grundlegenden Elemente der wissenschaftlichen Arbeit, mittelblaue für die zur Beantwortung der Aufgabenstellung notwendigen Theorie und das hellblaue für die praktische Umsetzung dessen.



**Abbildung 3: Aufbau der Arbeit<sup>5</sup>**

Kapitel 1 dient zur Einführung in die Aufgabenstellung. Dabei soll das Problem verstanden, das Ziel definiert und der Aufbau der Arbeit vorgestellt werden.

Um das Ergebnis messen zu können, werden in Kapitel 2 die produktionswirtschaftliche Theorie und die Kennzahlen definiert.

Innerhalb von Kapitel 3 werden mit dem CRISP-DM als Referenzmodell Methoden und Ansätze zum Verständnis und der Aufbereitung der Daten und zur Modellierung auf Basis derer aufgezeigt.

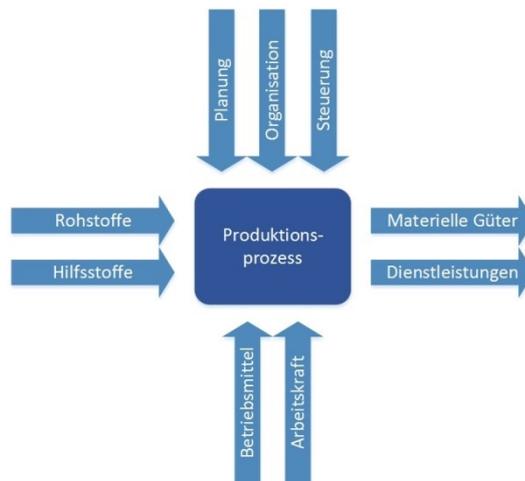
Die praktische Umsetzung erfolgt in Kapitel 4. Hierbei werden die in den beiden vorherigen Kapitel Methoden und Ansätze angewendet und Ergebnisse präsentiert.

Abschließend werden die Arbeit und deren Ergebnisse in Kapitel 5 zusammengefasst. Zusätzlich werden der in Kapitel 4 präsentierte Lösungsweg kritisch hinterfragt und andere Ansätze in Betracht gezogen.

<sup>5</sup> Quelle: eigene Darstellung

## 2 Produktionswirtschaftliche Theorie

In der Aufgabenstellung in Kapitel 1.2 werden lange Durchlaufzeiten und schlechte Auslastung der Fertigungskapazitäten in der Produktion als Grund für diese Arbeit angegeben. Die daraus abgeleitete Forschungsfrage, wie sich eine exakte Vorhersage auf die Kapazitätsauslastung auswirkt, muss in einen messbaren Kontext gebracht werden. Ziel dieses Kapitels ist es, die betriebswissenschaftliche Theorie und Kennzahlen zu diesen Vorgaben zu betrachten.

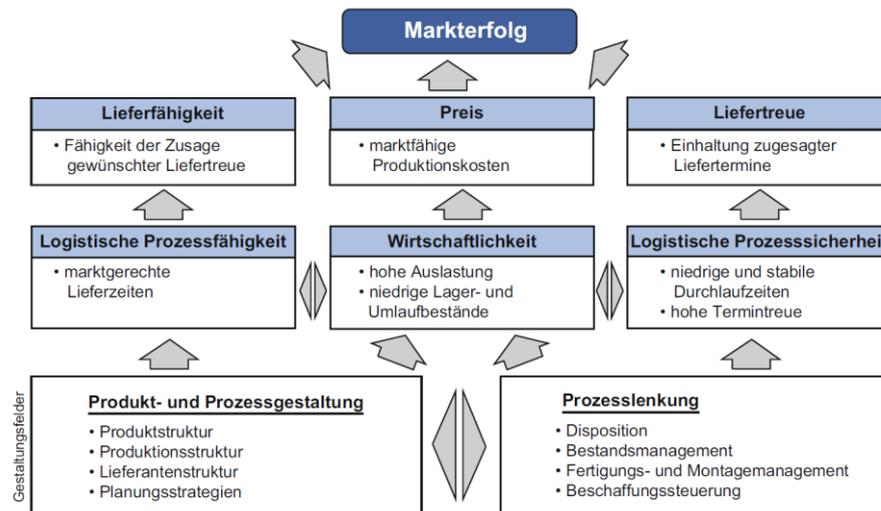


**Abbildung 4: Produktion als Kombinations- und Transformationsprozess<sup>6</sup>**

Die Grundlage für diese Betrachtung ist die Produktion als Einheit. Eine Produktion als solche ist ein Kombinations- und Transformationsprozess. In diesem werden, wie in Abbildung 4 dargestellt, verschiedene Produktionsfaktoren (Rohstoffe, Hilfsstoffe, Betriebsmittel, u. a.) zu Produkten kombiniert bzw. transformiert. Das Produktionssystem kann in Subsysteme aufgeteilt werden, die kleinsten werden als Leistungsstelle, die aus Maschinen oder technischen Anlagen und Maschinenbedienern bestehen.<sup>7</sup>

<sup>6</sup> Quelle: Zsifkovits, H. E. (2013), S. 106 (leicht modifiziert)

<sup>7</sup> Vgl. ebda., S. 106 ff.



**Abbildung 5: Logistische Erfolgsfaktoren von Produktionsunternehmen<sup>8</sup>**

Die logistischen Erfolgsfaktoren, mit denen ein Produktionsunternehmen beschrieben werden kann, sind in Abbildung 5 dargestellt. Für das oberste Ziel, den Markterfolg, sind die Lieferfähigkeit, die Liefertreue und der Preis von Bedeutung. Die Lieferfähigkeit wird über die Produkt-, Produktions- und Lieferantenstrukturen beeinflusst. Bei optimaler Gestaltung dieser wird eine logistische Prozessfähigkeit und in weiterer Folge die Prozesssicherheit des Unternehmens gewährleistet. Die dadurch erreichten niedrigen Durchlaufzeiten sollen auf stabilem Niveau und im laufenden Prozess eingehalten werden können. Dies führt zu einer hohen Liefertreue. Die Wechselwirkungen zwischen Leistungs- und Kostenzielen beeinflussen die Wirtschaftlichkeit der Produktion. Marktfähige Produktionskosten können durch eine hohe Auslastung der bereitgestellten Kapazitäten und durch Reduktion der Kapitalbindungskosten, die durch Lager- und Umlaufbeständen anfallen, realisiert werden.<sup>9</sup>

## 2.1 Das Dilemma der Ablaufplanung

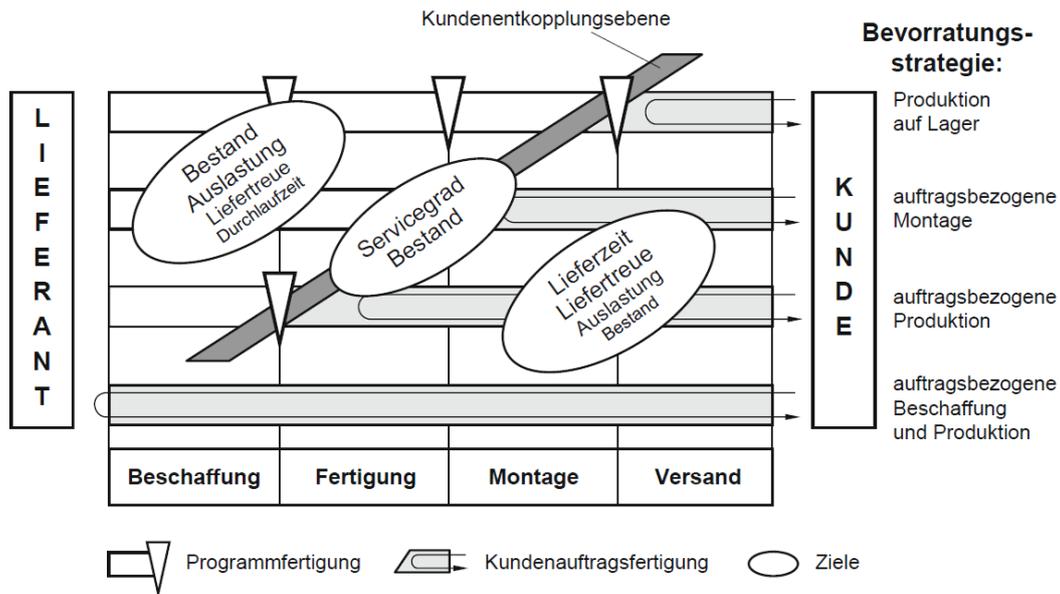
Das Bestreben einzelne logistische Erfolgsfaktoren zu optimieren, führt zu einem Zielkonflikt. Die Zielsetzungen und Anforderungen sind oftmals gegensätzlich und schwanken lokal und temporär. So sind, um eine hohe Auslastung zu erzielen, hohe Bestände notwendig. Dies führt zu langen und schwankenden Durchlaufzeiten, wodurch die Terminalsicherheit nicht gewährleistet werden kann. Diese Konstellation der Zielsetzungen wird als Dilemma der Ablaufplanung bezeichnet.<sup>10</sup> Kein einzelner Wert wird optimiert, sondern es muss ein Gleichgewicht zwischen den Teilzielen gefunden werden. Dazu müssen alle berücksichtigt und gewichtet werden.<sup>11</sup>

<sup>8</sup> Quelle: Gläßner, J. (1995) zitiert nach Nyhuis, P.; Wiendahl, H.-P. (2012), S. 3

<sup>9</sup> Vgl. Nyhuis, P.; Wiendahl, H.-P. (2012), S. 3

<sup>10</sup> Vgl. Gutenberg, E. (1971), S. 216

<sup>11</sup> Vgl. Nyhuis, P.; Wiendahl, H.-P. (2012), S. 4



**Abbildung 6: Gewichtung von logistischen Zielgrößen<sup>12</sup>**

Bei unterschiedlichen Bevorratungsstrategien und unterschiedlichen Kundenentkopplungspunkten innerhalb der Produktion unterscheidet sich die Gewichtung der Zielgrößen. Bei einer Produktion ohne den konkreten Kundenauftrag stehen geringe Bestände und eine hohe Auslastung als Ziele im Vordergrund, die direkt die Wirtschaftlichkeit der Produktion beeinflussen. Die Ziele der Liefertreue und Durchlaufzeit werden in einem zweiten Schritt betrachtet, da diese die Lagerhaltung betreffen. Das Ziel einer geringen Termintreue und großen Durchlaufzeit bedingt einen großen Lagerbestand. Bei einer kundenbezogenen Produktion verändern sich die Zielgewichtungen, sodass die Lieferzeit und die Liefertreue in den Vordergrund treten. Bei nicht Einhaltung zugesagter Lieferzeiten und -termine ist der Kunde direkt betroffen. Fraglich ist jedoch, ob kurze Lieferzeiten zugunsten von geringer Auslastung der Produktionssysteme sinnvoll sind. Die dadurch notwendigen Kapazitätsausweitungen (Betriebsmittel und/oder Personal) führen zu einem Anstieg der Stückkosten, die an den Kunden weiterverrechnet werden müssen.<sup>13</sup>

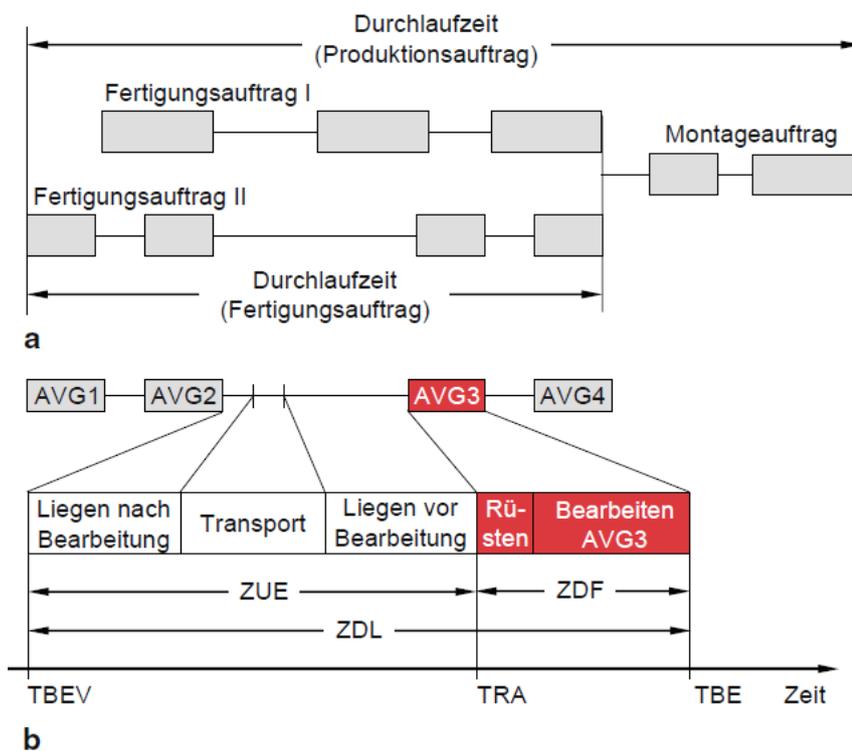
Im Hinblick auf das Ziel dieser Arbeit (Kapitel 1.2) wird genauer auf den Faktor Betriebsmittel und dessen Potenzial eingegangen. Die Ziele der Auslastung und der Durchlaufzeit, die in dieser Arbeit im Vordergrund stehen, bedingen eine Auseinandersetzung mit der Kapazität und ihrer Zusammensetzung und der Durchlaufzeit als Kenngröße der Produktionslogistik.

<sup>12</sup> Quelle: in Anlehnung an Eidenmüller, B. (1995) zitiert nach Nyhuis, P.; Wiendahl, H.-P. (2012), S. 4

<sup>13</sup> Vgl. Nyhuis, P.; Wiendahl, H.-P. (2012), S. 4 f.

## 2.2 Durchlaufzeit

Nyhuis und Wiendahl (2012)<sup>14</sup> unterscheiden zwischen der Durchlaufzeit eines Produktionsauftrags, der Durchlaufzeit eines Fertigungsauftrags, im Weiteren als Auftrag bezeichnet, und dem Durchlaufelement. Diese werden in Abbildung 7 dargestellt. Der obere Teil zeigt den Durchlauf eines Produktionsauftrages, der aus zwei Fertigungsaufträgen und einem Montageauftrag besteht. Handelt es sich um eine losweise Fertigung, folgt nach der Beendigung eines Arbeitsvorgangs auf dem Arbeitssystem mit einer möglichen Liegezeit vor Ort der Transport zum Folgearbeitssystem. Hier kann der Auftrag bearbeitet werden, wenn die Kapazitäten zur Bearbeitung frei sind und die notwendigen Rüstungsmaßnahmen durchgeführt worden sind.<sup>15</sup>



TBEV : Bearbeitungsende Vorgänger      ZDL = TBE - TBEV : Durchlaufzeit  
 TRA : Rüstanfang      ZUE = TRA - TBEV : Übergangszeit  
 TBE : Bearbeitungsende      ZDF = TBE - TRA : Durchführungszeit

**Abbildung 7: Durchlaufzeitanteile und Durchlaufelement<sup>16</sup>**

Im unteren Teil der Abbildung 7 sind die Ablaufschritte für einen Arbeitsvorgang, das arbeitsvorgangsbezogenes Durchlaufelement, dargestellt.<sup>17</sup> Als die Durchlaufzeit  $ZDL$  eines Arbeitsvorgangs gilt die Zeitspanne, „die ein Auftrag von der Beendigung des vorhergehenden Arbeitsvorgangs bzw. vom Einstoßzeitpunkt des Auftrages (beim

<sup>14</sup> Vgl. Nyhuis, P.; Wiendahl, H.-P. (2012), S. 21 f.

<sup>15</sup> Vgl. ebda., S. 21

<sup>16</sup> Quelle: ebda., S. 22

<sup>17</sup> Vgl. ebda., S. 21

ersten Arbeitsvorgang)  $TBEV$  bis zum Bearbeitungsende des betrachteten Arbeitsvorgangs  $TBE$  selbst benötigt“<sup>18</sup> (2.1).

$$ZDL = TBE - TBEV \quad (2.1)$$

Nach Nyhuis und Wiendahls (2012)<sup>19</sup> Definition werden Liegen nach Bearbeiten, Transportzeit und das Liegen vor Bearbeitung dem jeweiligen untersuchten Arbeitsvorgang zugeordnet. Diese werden zur Übergangszeit  $ZUE$  zusammengefasst werden, wobei diese durch die Differenz der Durchlaufzeit  $ZDL$  und Durchführungszeit  $ZDF$  berechnet wird (2.2).<sup>20</sup>

$$ZUE = ZDL - ZDF \quad (2.2)$$

Ein relatives Maß für die Durchlaufzeit ist der mittlere Flussgrad  $FG_m$ . Dieser stellt das Verhältnis der mittleren Durchlaufzeit  $ZDL_m$  zur mittleren Durchführungszeit  $ZDF_m$  dar (2.3):<sup>21</sup>

$$FG_m = \frac{ZDL_m}{ZDF_m} \quad (2.3)$$

Gudehus (2010)<sup>22</sup> definiert die Aufteilung der Zeitanteile anders. Er unterscheidet in Wartezeit, Rüstzeit, Leistungszeit und Verfahrenszeit (2.4).<sup>23</sup>

$$\text{Auftragsdurchlaufzeit} = \text{Wartezeit} + \text{Rüstzeit} + \text{Leistungszeit} + \text{Verfahrenszeit} \quad (2.4)$$

So gelten jene Zeiten als Rüstzeiten, die vor und nach der Leistungszeit anfallen. Dazu zählen u. a. folgende:<sup>24</sup>

- Auftragsannahmezeiten
- Vor- und Nachbereitungszeiten
- Materialbereitstellungszeiten
- Umschalt- und Räumzeiten
- Be- und Entladezeit
- Datenerfassungszeiten

Die Leistungszeit hingegen ist die Zeit, die eine einsatzbereite Leistungsstelle für die Erbringung der im Auftrag geforderten Leistung benötigt und ist in der Regel von der Auftragsgröße abhängig. Typische Beispiele für diese sind:<sup>25</sup>

- Fertigungszeiten
- Montagezeiten
- Bearbeitungszeiten
- Reparaturzeiten

<sup>18</sup> Nyhuis, P.; Wiendahl, H.-P. (2012), S. 21

<sup>19</sup> Vgl. ebda., S. 21 f.

<sup>20</sup> Vgl. ebda., S. 23

<sup>21</sup> Vgl. ebda., S. 23

<sup>22</sup> Vgl. Gudehus, T. (2010), S. 216

<sup>23</sup> Vgl. ebda., S. 216

<sup>24</sup> Vgl. ebda., S. 216

<sup>25</sup> Vgl. ebda., S. 216

Die Verfahrenszeit ist verfahrenstechnisch bedingt und muss verstreichen, bevor am Auftragsgegenstand der nächste Bearbeitungsschritt durchgeführt werden darf. Beispiele hierfür sind:<sup>26</sup>

- Trocknungszeiten
- Aushärtungszeiten
- Ablagerungszeiten
- Reifezeiten
- Gärungszeiten

Aus den Minima der Rüst-, Leistungs- und Verfahrenszeit kann die minimale Durchlaufzeit berechnet werden. Die Differenz zwischen der tatsächlichen und minimalen Durchlaufzeit wird als Wartezeit bezeichnet. Diese ergibt sich aus:<sup>27</sup>

- Ausfallzeiten durch Betriebsunterbrechung, Störung oder fehlende Personalbesetzung
- Totzeiten wegen fehlender Daten, Informationen, Entscheidungen oder Anweisungen
- Nachbearbeitungszeiten zur Beseitigung von Fehlern und Mängeln am Auftragsgegenstand
- Stauzeiten infolge der Belegung einer Leistungsstelle durch vorherige Aufträge
- Blockierzeiten durch einen Rückstau in einer nachfolgenden Leistungsstelle
- Materialbeschaffungszeiten, die durch fehlendes Material verursacht werden
- Unterbrechungszeiten infolge eines Ausfalls oder Nichtverfügbarkeit einer vorangehenden Leistungsstelle
- Pufferzeiten, die disponiert werden, um Aufträge anzusammeln, um die optimale Auslastung zu erreichen

Auffällig ist bei Gudehus (2010)<sup>28</sup> Einteilung, dass Transportzeiten nicht als solche deklariert werden. Sie könnten der Materialbereitstellungszeit zugeordnet werden, und somit in die Rüstzeit fallen. Nyhuis und Wiendahl (2012)<sup>29</sup> unterscheiden hingegen nicht, ob eine Liegezeit verfahrensbedingt ist oder nicht. Bei einer Betrachtung der Durchlaufzeit ist daher zu überprüfen, wie genau einzelne Zeitanteile definiert werden müssen und welche Definition als Referenz dienen soll.

---

<sup>26</sup> Vgl. Gudehus, T. (2010), S. 217

<sup>27</sup> Vgl. ebda., S. 217 f.

<sup>28</sup> Vgl. ebda., S. 216

<sup>29</sup> Vgl. Nyhuis, P.; Wiendahl, H.-P. (2012), S. 21 f.

**Tabelle 1: Kenngrößen der Produktionslogistik<sup>30</sup>**

<p style="text-align: center;"><b>Zeit</b></p> <ul style="list-style-type: none"> <li>• Durchlaufzeit/Lieferzeit</li> <li>• Warte-/Liegezeiten vs. produktive Zeiten</li> <li>• Wiederbeschaffungszeit</li> <li>• Kapitalbindung</li> </ul>	<p style="text-align: center;"><b>Qualität</b></p> <ul style="list-style-type: none"> <li>• Kundenzufriedenheit</li> <li>• Kundennutzen</li> <li>• Produkt-/Dienstleistungsqualität</li> <li>• Lieferqualität</li> </ul>
<p style="text-align: center;"><b>Kosten/Wirtschaftlichkeit</b></p> <ul style="list-style-type: none"> <li>• Input: Material, Personal, Kapital, Overhead</li> <li>• Output: Umsätze</li> <li>• Produktivität, Rentabilität, Effizienzkennzahlen</li> </ul>	<p style="text-align: center;"><b>Flexibilität</b></p> <ul style="list-style-type: none"> <li>• Anpassungsumfang</li> <li>• Kosten/Ertrag aus Anpassungen</li> <li>• Zeitbedarf für Anpassungen</li> </ul>

Die Wechselwirkungen der Durchlaufzeit mit den Kenngrößen der Produktionslogistik, wie sie in Tabelle 1 dargestellt sind, sind eine Ausprägung des in Kapitel 2.1 vorgestellten Dilemmas der Ablaufplanung. Auch hier wird ein Zielkonflikt beim Optimierungswunsch mehrere Zielgrößen verursacht. So kann eine Veränderung der Durchlaufzeit die Veränderung anderer Größen, wie Bestände, Kundenzufriedenheit, Produktionskosten und Flexibilität, beeinflussen.<sup>31</sup>

Eine Beeinträchtigung in negative Richtung liegt dann vor, wenn die Durchlaufzeit ansteigt. Dadurch nimmt die Planungsqualität ab und Risiken und Kosten steigen. Unvorhergesehene Entwicklungen und Störungen können eintreten und lange Lieferzeiten können Kundenunzufriedenheit und Auftragsänderungen mit sich bringen.<sup>32</sup>

Eine Verkürzung der Durchlaufzeit hat verschiedene Auswirkung. So wird die Logistikleistung verbessert. Dies führt zu einer Verkürzung von Lieferzeiten und einer Erhöhung der Liefertreue, -fähigkeit, Flexibilität, Lieferqualität und Informationsbereitschaft. Ebenso erfolgt eine Reduktion der Logistikkosten, die durch eine Verringerung der Lenkungs-, Handlings- und Lagerkosten und der Verringerung der Kapitalbindung und Wagnissen erreicht wird. Um dies zu erreichen sollte eine Kernzeitanalyse vorgenommen werden.<sup>33</sup>

Eine Erhöhung der Flexibilität kann nur in dem Maß erfolgen, in dem Produktionskapazität bereit steht. Daher stehen im nächsten Kapitel die Kapazität und die Möglichkeiten der Anpassung im Fokus.

<sup>30</sup> Quelle: Zsifkovits, H. E. (2013), S. 112

<sup>31</sup> Vgl. ebda., S. 112

<sup>32</sup> Vgl. ebda., S. 113

<sup>33</sup> Vgl. ebda., S. 117

## 2.3 Kapazität

Der Produktionsfaktor Betriebsmittel kann durch verschiedene Größen beschrieben werden. Eine davon ist die Kapazität. In diesem Kapitel soll erklärt werden, wie sie definiert wird und wie Kapazität und Durchlaufzeit in Zusammenhang stehen. Zusätzlich wird die Möglichkeit der Anpassung mit ihren Zielen, Herausforderungen, Voraussetzungen und verschiedenen Formen erläutert.

### 2.3.1 Einführung

Die Kapazität ist das Leistungsvermögen einer technischen Einheit in einem definierten Zeitabschnitt. Die Ausnutzung in Prozent dieses Leistungsvermögens kann als *Beschäftigungsgrad (BG) (Kapazitätsnutzungsgrad)* (2.5) angegeben werden.<sup>34</sup>

$$\text{Beschäftigungsgrad} = \frac{\text{tatsächlich in Anspruch genommene Leistung}}{\text{maximal mögliche Leistung}} * 100\% \quad (2.5)$$

Die Kapazität kann in betriebswirtschaftliche und technische Komponenten zerlegt werden.<sup>35</sup>

$$K = C \times I \times D \quad (2.6)$$

Die Multiplikation von Kapazität  $C$ , Produktionsgeschwindigkeit  $I$ , und Produktionsdauer  $D$  bildet die Kapazität  $K$  (2.6).<sup>36</sup> Die Kapazität  $C$  bildet die technische Komponente und beschreibt Anzahl und Eigenschaften der Betriebsmittel. Als technisch-wirtschaftliche Komponente wird die Produktionsgeschwindigkeit  $I$  (Intensität der Betriebsmittelnutzung) angesehen, die innerhalb bestimmter Grenzen ( $I_{min} \leq I \leq I_{max}$ ) veränderlich ist und die dadurch die Höhe der variablen Kosten und die Durchlaufzeit beeinflussen kann. Die Produktionsdauer  $D$  gibt die Länge des Zeitraums, der zur Normierung verwendet wird, an.<sup>37</sup>

Die Ausprägungen des Kapazitätsbegriffs sind in Tabelle 2 zusammengefasst. Soll die Kapazität als betriebliche Kennzahl erfasst werden, muss zuerst der Kapazitätsbegriff als solches definiert werden. Die Kapazität selbst und ihre Inanspruchnahme sind wichtige Einflussfaktoren für Planung, Steuerung und Überwachung der Produktion.<sup>38</sup>

<sup>34</sup> Vgl. Blohm, H. et al. (2008), S. 168

<sup>35</sup> Vgl. Männel, W. (1979) zitiert nach Blohm, H. et al. (2008), S. 168

<sup>36</sup> Vgl. ebda., S. 168

<sup>37</sup> Vgl. Blohm, H. et al. (2008), S. 168

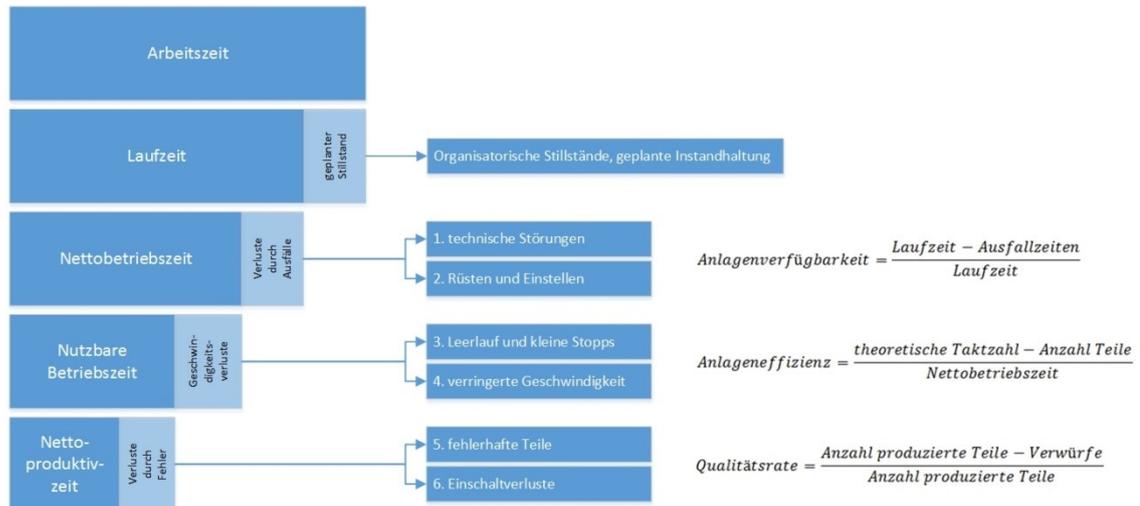
<sup>38</sup> Vgl. ebda., S. 168 f.

Tabelle 2: Ausprägungen des Kapazitätsbegriffs<sup>39</sup>

Unterscheidungsmerkmal	Kapazitätsbegriff	Definition/Charakterisierung
(1) Teilbarkeit	a) Ganzheitliche Kapazität	Die Kapazität des Betriebes wird als unteilbares Ganzes gesehen
	b) Zusammengesetzte Kapazität	Die Kapazität des Betriebes setzt sich aus den Kapazitäten der einzelnen Aggregate zusammen. Bei Aggregaten, die unabhängig voneinander eingesetzt werden: Addition der Teilkapazitäten. Bei verbundenen Aggregaten: Gesamtkapazität = Kapazität des Engpasses
(2) Menge und Art der Leistung	a) Quantitative Kapazität	Mengenmäßige Leistungsfähigkeit
	b) Qualitative Kapazität	Leistungsfähigkeit in Bezug auf die Eigenschaften der abgegebenen Leistungen.
(3) Technische Gesichtspunkte	a) Technische maximal mögliche Leistungsfähigkeit	Theoretisch mögliche Leistungsfähigkeit unter den günstigsten Bedingungen bei voller Auslastung sämtlicher Produktionssysteme; Vorzug: feste Vergleichsgröße für die erbrachte Mengenleistung
	b) Durchschnittliche Dauerleistungsfähigkeit	Gegenüber (3a) sind hier die normalen Ausfälle infolge Wartung und Reparaturen sowie Abwesenheit von Belegschaftsangehörigen berücksichtigt. Nachteil bei der Verwendung als Bezugsgröße: zeitweise Überschreitung durch tatsächlich erbrachte Leistung möglich (BG > 100%)
	c) Momentane (jeweilige) Leistungsfähigkeit	Leistungsfähigkeit in einem bestimmten, in der Regel kurzen Zeitraum (konkreter Tag, konkrete Woche). Wegen ihrer Schwankungen als Bezugsgröße ungeeignet, aber als Vergleichsmaß für die Beurteilung des auf die technische Höchstleistungsfähigkeit bezogenen BG wertvoll.
(4) Bezugszeitraum	a) Totalkapazität	Leistungsfähigkeit während der gesamten technischen oder wirtschaftlichen Nutzungsdauer
	b) Periodenkapazität	Leistungsfähigkeit während einer Rechnungs- oder Planungsperiode

<sup>39</sup> Quelle: Blohm, H. et al. (2008), S. 169 (leicht modifiziert)

Wird in einem Betrieb, einem definierten Betriebsteil oder auf einer Anlage ein einheitliches Produkt erzeugt, kann die Kapazität einfach gemessen werden. Dies passiert dann über die erzeugten Mengeneinheiten, wie Stück, Liter, Meter, etc. pro Zeiteinheit. Bei einer Sortenfertigung können verschiedene Sorten unter Zuhilfenahme von Äquivalenzzahlen miteinander verglichen werden.<sup>40</sup>



**Abbildung 8: Berechnung des OEE<sup>41</sup>**

Die Kapazität kann durch ungeplante und geplante Stillstände, die durch Wartung, Urlaub, Krankheit oder technische Ausfälle bedingt sind, geschmälert werden. Eine Möglichkeit, um die verfügbare Kapazität zu bestimmen, ist die Berechnung der *Gesamtanlageneffektivität OEE* (Overall Equipment Efficiency). Diese setzt sich aus den drei Größen Anlagenverfügbarkeit, Anlageneffizienz und der Qualitätsrate zusammen (2.7). Ihre Berechnung wird in Abbildung 8 zusammengefasst.

$$\text{OEE} = \text{Anlagenverfügbarkeit} * \text{Anlageneffizienz} * \text{Qualitätsrate} \quad (2.7)$$

Die Planung von betrieblichen Teilkapazitäten kann Engpässe aufzeigen. Ein Engpass ist jenes Betriebsmittel bzw. eine Betriebsmittelgruppe, deren Kapazität zur maximalen Gesamtleistung des Betriebs wird, wenn gleichzeitig in anderen Bereichen ungenutzte Kapazitäten vorliegen.<sup>42</sup> Dies und eine Differenz zwischen Kapazitätsbedarf und verfügbaren Kapazitäten führen zur Notwendigkeit der Anpassung der Kapazität.<sup>43</sup> Hierzu sollten verschiedene Möglichkeiten und die damit einhergehenden Herausforderungen betrachtet werden, auf die im nächsten Kapitel eingegangen wird.

<sup>40</sup> Vgl. Blohm, H. et al. (2008), S. 171

<sup>41</sup> Quelle: Fraunhofer-Institut Produktionstechnik und Automatisierung IPA zitiert nach Zsifkovits, H. E. (2013), S. 151

<sup>42</sup> Vgl. Blohm, H. et al. (2008), S. 171

<sup>43</sup> Vgl. Zsifkovits, H. E. (2013), S. 152

### 2.3.2 Anpassungen der Kapazität

Das Ziel der Anpassung der Kapazität liegt „in der Angleichung des Zeitfonds des Kapazitätsangebots und des Kapazitätsbedarfs je Kapazitätseinheit, um eine bestmögliche Bedarfsbefriedigung abzusichern.“<sup>44</sup> Zusätzlich sollten die dabei anfallenden Kosten betrachtet werden. Auszugehen ist von den *fixen Kosten*  $Q$ , die bei Aufrechterhaltung der Betriebsbereitschaft anfallen, und die von der Inanspruchnahme der betrieblichen Anlagen und von der Art der Anpassung an die Beschäftigungsschwankungen unabhängig sind. Die fixen Kosten setzen sich aus Nutzkosten  $K_n$  und Leerkosten  $K_l$  zusammen. Nutzkosten fallen für die tatsächlich genutzte Kapazität der Anlage an, wohingegen Leerkosten für die nicht genutzte Kapazität stehen.<sup>45</sup>

$$Q = K_l + K_n \quad (2.8)$$

Die Nutz- und Leerkosten können unter Zuhilfenahme der effektiv erzeugten Produktmenge  $x$  und dem Maximum der herstellbaren Produktmenge  $m$  folgendermaßen ausgedrückt werden:<sup>46</sup>

$$K_l(x) = (m - x) \frac{Q}{m} \quad (2.9)$$

$$K_n(x) = x \frac{Q}{m} \quad (2.10)$$

Der Zusammenhang der Nutz- und Leerkosten wird in Abbildung 9 dargestellt. Bei voller Ausnutzung der verfügbaren Kapazität fallen keine Leerkosten an und die Fixkosten entsprechen den Nutzkosten.<sup>47</sup>

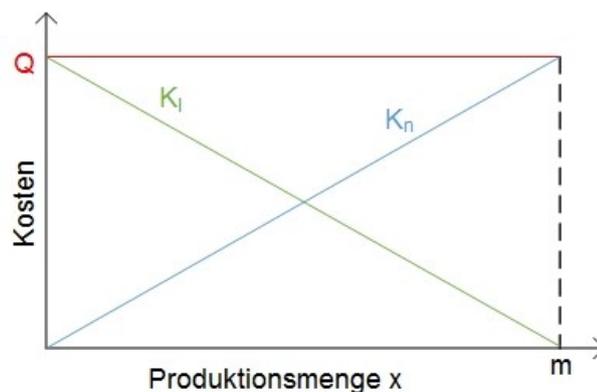


Abbildung 9: Verlauf von Nutz- und Leerkosten<sup>48</sup>

Das Ziel Leerkosten abzubauen kann nicht mit sämtlichen Anpassungsformen erreicht werden. Dies ist der Fall bei der später vorgestellten intensitätsmäßigen Anpassung, da sich hier nur die Inanspruchnahme verändert, jedoch nicht der Faktorbestand.<sup>49</sup>

<sup>44</sup> Vgl. Nebl, T. (2011), S. 224

<sup>45</sup> Vgl. Gutenberg, E. (1971), S. 348 f.

<sup>46</sup> Vgl. ebda., S. 349

<sup>47</sup> Vgl. ebda., S. 349

<sup>48</sup> Vgl. ebda., S. 349

<sup>49</sup> Vgl. ebda., S. 349

### Voraussetzungen und Herausforderungen

Die Voraussetzungen und Herausforderungen bei der Anpassung sind vielfältig. Dazu zählen das Anpassungsintervall, die Kostenremanenz, die Anpassungskosten und die Elastizität als Anpassungsfähigkeit, die im Weiteren näher erklärt werden.

Bei dem *Anpassungsintervall* handelt es sich um den Zeitraum zwischen der Anpassungsnotwendigkeit und der erfolgten Ausführung der Anpassungshandlung. Er beinhaltet sowohl Leerzeiten, wie die Zeitspanne des Erkennens der Anpassungsnotwendigkeit oder die Dauer des informationssammelnden, -bearbeitenden und ausführenden Prozess. Die Länge ist abhängig von der Organisation des Entscheidungsprozesses und auch vom Auslöser der Anpassungsnotwendigkeit. Zusätzlich sind Mechanismen zur Kontrolle, ob eine Anpassung notwendig ist, einzusetzen.<sup>50</sup>

Unter *Kostenremanenz* wird der Umstand bezeichnet, dass sich Kosten nicht proportional zum Beschäftigungsgrad verhalten. Bei rückläufiger Beschäftigung laufen die Gesamtkosten nicht auf der gleichen Kostenkurve  $K_a$  zurück, die die Zunahme der Kosten bei steigender Beschäftigung beschreibt, sondern sie fallen nach der Kostenkurve  $K_r$ , die darüber liegt.<sup>51</sup>

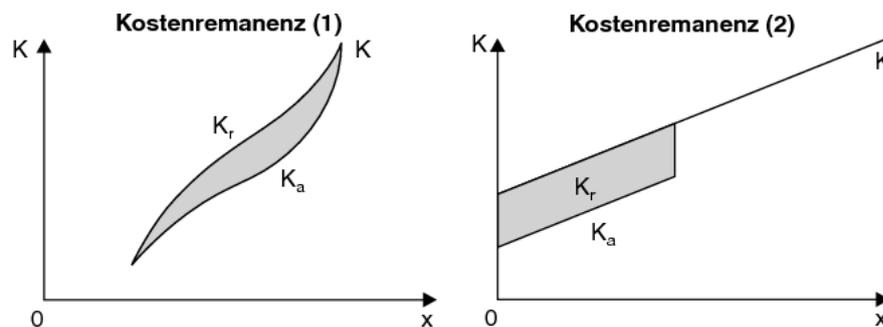


Abbildung 10: Kostenremanenz: (1) Hysteresis-Schleife, (2) remanenter Sprung<sup>52</sup>

Die Kostenremanenz kann sich als Hysteresis-Schleife oder als remanenter Kostensprung manifestieren und ist in Abbildung 10 zu sehen.<sup>53</sup> Dies hat eine Steigerung der Kosten bei gleichem Beschäftigungsgrad aber zu unterschiedlichen Zeitpunkten zur Folge.

Die Anpassungsentscheidung und -ausführung verursacht Kosten. Diese können in drei Gruppen eingeteilt werden: die Elastizitätskosten, die anfallen, um zukünftige Anpassungsmaßnahmen zu erleichtern, jene Kosten, die bei der Wirksamkeit der Anpassungsmaßnahmen entstehen (z. B. Betriebskosten der neuen Anlage), und die *Anpassungskosten*, die gemeinsam mit den Elastizitätskosten durch die Anpassungsentscheidung und -ausführung entstehen. Sie beginnen mit der Erfassung der Anpassungsnotwendigkeit und erreichen ihr Maximum bei der vollen Wirksamkeit der Maßnahme. Die Anpassungskosten gelten als „wertmäßige Trägheitsschwelle“ und

<sup>50</sup> Vgl. Swoboda, P. (1964), S. 56

<sup>51</sup> Vgl. Weber, J.; Steven, M., <http://wirtschaftslexikon.gabler.de/Archiv/7702/kostenremanenz-v6.html> (Zugriff: 25.04.2017)

<sup>52</sup> Quelle: ebda.

<sup>53</sup> Vgl. ebda.

können ausschlaggebend für die Verzögerung oder Verhinderung einer sonst rentablen Anpassungsmaßnahme sein.<sup>54</sup>

Die *Elastizität* wird als Reagibilität eines produktiven Wirtschaftsgebildes gegenüber veränderten Produktionsaufgaben bezeichnet. Sie kann in der Anpassung der Mittel und Kräfte resultieren und wird durch Verfahren und Ziele des Wirtschaftsgebildes gegenüber Änderungen in der Umwelt und im eigenen Betrieb definiert.<sup>55</sup> Die Elastizität der Anpassungsfähigkeit kann in quantitative, qualitative, zeitliche und räumliche Elastizität unterteilt werden.<sup>56</sup> Die Elastizität schafft die Voraussetzungen für eine Anpassung, wobei jedoch die Wirtschaftlichkeit dieser ausschlaggebend ist.<sup>57</sup>

### **Arten der Anpassung**

Nach der Betrachtung der Voraussetzungen und Herausforderungen besteht die Frage nach den Möglichkeiten zur Anpassung. Gutenberg (1971)<sup>58</sup> hat mit der intensitätsmäßigen, quantitativen und zeitlichen Anpassung die Änderung des Beschäftigungsgrades abgedeckt. Swoboda (1964)<sup>59</sup> ergänzte diese drei Grundtypen mit der qualitativen Anpassung.

Eine *intensitätsmäßige Anpassung* liegt vor, wenn der gesamte fertigungstechnische Apparat bei gleichbleibender Betriebsdauer unterschiedlich stark beschäftigt wird. Es wird die Intensität der Betriebsmittelnutzung angepasst und so die Kapazität verändert (siehe (2.6)). Diese Art der Anpassung kommt bei Anlagen und Anlagenkomplexen zum Einsatz, die nicht in mehrere selbstständige Teileinheiten getrennt werden können.<sup>60</sup> Hierbei muss beachtet werden, dass bei einer rückläufigen Beschäftigung die Leerkosten grundsätzlich konstant bleiben, da nicht der Faktoreinsatz sondern die Inanspruchnahme variiert.<sup>61</sup>

---

<sup>54</sup> Vgl. Swoboda, P. (1964), S. 59 f.

<sup>55</sup> Vgl. Moxter, A. et al. (1954), S. 87 zitiert nach Swoboda, P. (1964), S. 61

<sup>56</sup> Vgl. Moxter, A. et al. (1954), S. 105 ff. zitiert nach Swoboda, P. (1964), S. 61

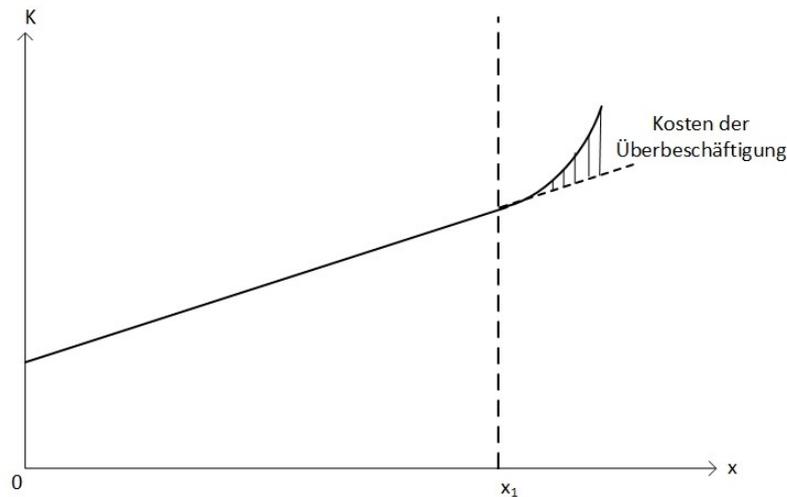
<sup>57</sup> Vgl. Swoboda, P. (1964), S. 61

<sup>58</sup> Vgl. Gutenberg, E. (1971), S. 355 f.

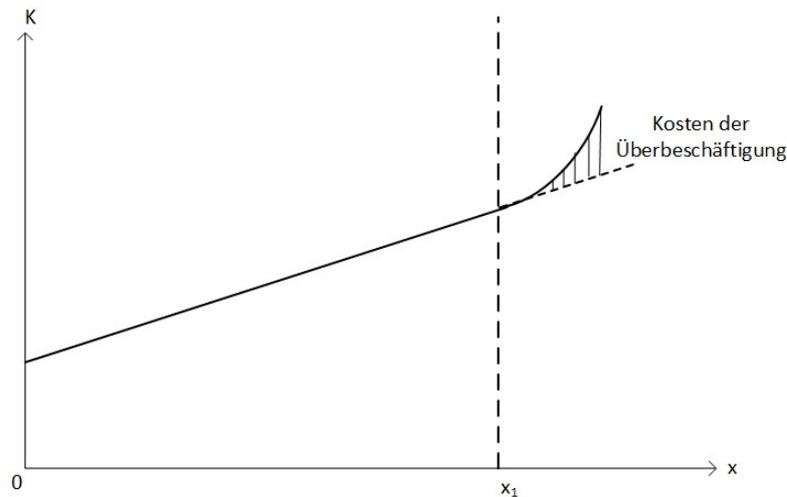
<sup>59</sup> Vgl. Swoboda, P. (1964), S. 49

<sup>60</sup> Vgl. Gutenberg, E. (1971), S. 355

<sup>61</sup> Vgl. ebda., S. 349



**Abbildung 11: Gesamtkostenverlauf beim Übergang<sup>62</sup>**

Ein möglicher Kostenverlauf wie in  liegt vor, wenn bis zur Ausbringungsmenge  $x_1$  innerhalb der maximalen möglichen Betriebszeit und bei optimaler Leistung zeitlich angepasst wird und Mengen, die darüber hinausgehen, durch intensitätsmäßige Anpassungen erzeugt werden.<sup>63</sup>

Werden fertigungstechnische Anpassungen im Sinne einer Stilllegung oder Wiederinbetriebnahme von selbstständigen Teileinheiten an die sich ändernde Beschäftigungslage durchgeführt, wird diese als *quantitative* Anpassung bezeichnet. Es wird die Kapazität angepasst (siehe (2.6)). Die Voraussetzungen hierfür ist die Selbstständigkeit der Teileinheiten.<sup>64</sup> Zusätzlich kann die quantitative Anpassung in zwei Formen, die *rein quantitative* und die *selektive* Anpassung unterschieden werden. Bei der rein quantitativen Anpassung hat ein Betrieb über gleichartige Aggregate und ein beliebiges stillgesetzt bzw. in Betrieb genommen. Bei der selektiven Anpassung wird jenes Aggregat ausgewählt, wenn aus einem heterogenen Maschinenbestand hinsichtlich Kriterien, wie Leistungsvermögen und Kostenstrukturen, jene ausgewählt werden, die in Betrieb genommen oder stillgesetzt werden.<sup>65</sup> Diese Anpassungen sind mit Kosten, entweder solche für die Investition oder die für die kurzfristige Bereitstellung von Kapazitäten, wie z. B. durch Anmietung, verbunden.<sup>66</sup>

Bei einer *zeitlichen Anpassung* wird die Betriebszeit je nach Beschäftigungslage verkürzt oder verlängert. Die Betriebsanlagen werden in dieser Zeit voll ausgelastet. In diesem Fall kommt es zu Veränderungen der Produktionsdauer (siehe (2.6)).<sup>67</sup> Dies geschieht durch Überstunden und/oder durch die Anpassung des Schichtmodells, der Pausenregelung, der Fertigungszeiten. Zusätzlich verbessert die Abstimmung der Taktung der einzelnen Aggregate die Betriebszeit.<sup>68</sup> Durch Überstundenzuschläge

<sup>62</sup> Quelle: Blohm, H. et al. (2008) (leicht modifiziert), S. 81

<sup>63</sup> Vgl. Blohm, H. et al. (2008), S. 81

<sup>64</sup> Vgl. Gutenberg, E. (1971), S. 355 f.

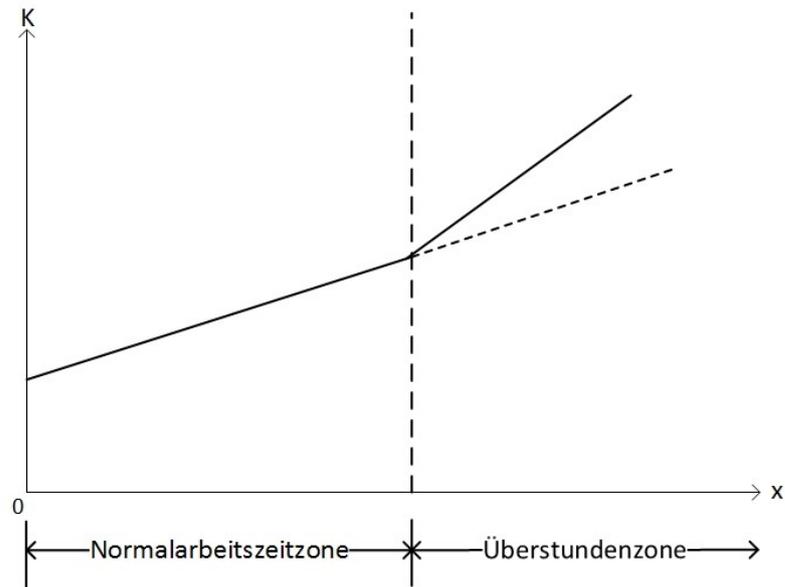
<sup>65</sup> Vgl. Blohm, H. et al. (2008), S. 83

<sup>66</sup> Vgl. Zsifkovits, H. E. (2013), S. 152

<sup>67</sup> Vgl. Gutenberg, E. (1971), S. 356

<sup>68</sup> Vgl. Zsifkovits, H. E. (2013), S. 152

steigen die Lohnkosten außerhalb der Normalarbeitszone stärker, jedoch weiterhin linear an.<sup>69</sup>



**Abbildung 12: Lohnkostenverlauf bei zeitlicher Anpassung<sup>70</sup>**

Die *qualitative Anpassung* umfasst u. a. die Änderung der Produktions- und Faktorenqualität (Schulung der Arbeitskräfte, Erneuerung von Maschinen), die Einführung anderer Verfahren und die Substitution von Faktoren, wie z. B. Mensch durch Betriebsmittel – Automatisierung von Abläufen, im Erzeugungsprozess. Je nach Betrachtungspunkt können Anpassungen unterschiedlich kategorisiert werden. So kann die intensitätsmäßige Anpassung einer Maschine zu einer qualitativen Anpassung eines Produkts führen.<sup>71</sup>

Diese Einführung in die produktionswirtschaftliche Theorie soll den Hintergrund für die Datenanalyse und die Motivation zur Durchführung einer solchen beleuchten. Erhöht sich allein die Intensität der Auslastung eines Betriebsmittels, das einen Engpass darstellt, können so die Gesamtdurchlaufzeit und Kosten reduziert werden. Im nächsten Kapitel soll nun die Methodik, mit der diese Ziele erreicht werden sollen, vorgestellt werden.

<sup>69</sup> Vgl. Blohm, H. et al. (2008), S. 82

<sup>70</sup> Quelle: ebda. (leicht modifiziert), S. 82

<sup>71</sup> Vgl. Swoboda, P. (1964), S. 49

## 3 Datenanalyse

Das in Abbildung 2 vorgestellte Konzept beruht auf den Grundlagen und Methoden der Datenanalyse. Daher ist das Ziel dieses Kapitels, ein grundlegendes Verständnis und die theoretische Basis zu schaffen, die zur Lösung der Aufgabenstellung aus Kapitel 1 und zur Beantwortung der Forschungsfrage benötigt wird.

Zu diesem Zweck werden eine allgemeine Einführung gegeben und einige Begriffsabgrenzungen vorgenommen. Es wird eine Einführung in den Begriff Data Science gegeben, der als Überbegriff für die Methoden und Ansätze im Umgang mit Daten und daraus zu gewinnenden Erkenntnisse im Kontext von Informatik, Mathematik und den einzelnen Anwendungsgebieten steht. Anhand des CRISP-DM soll aufgezeigt werden, wie ein solches Projekt strukturell aufgebaut werden kann. Dessen Unterteilung in Phasen ist verantwortlich für die Struktur dieses Kapitels. Es werden die Grundlagen aus der Statistik behandelt, bevor das Schaffen einer einheitlichen Datenbasis in der Theorie behandelt wird. Mit dem abschließenden Vorstellen der Methoden, die verwendet werden können, endet dieses Kapitel.

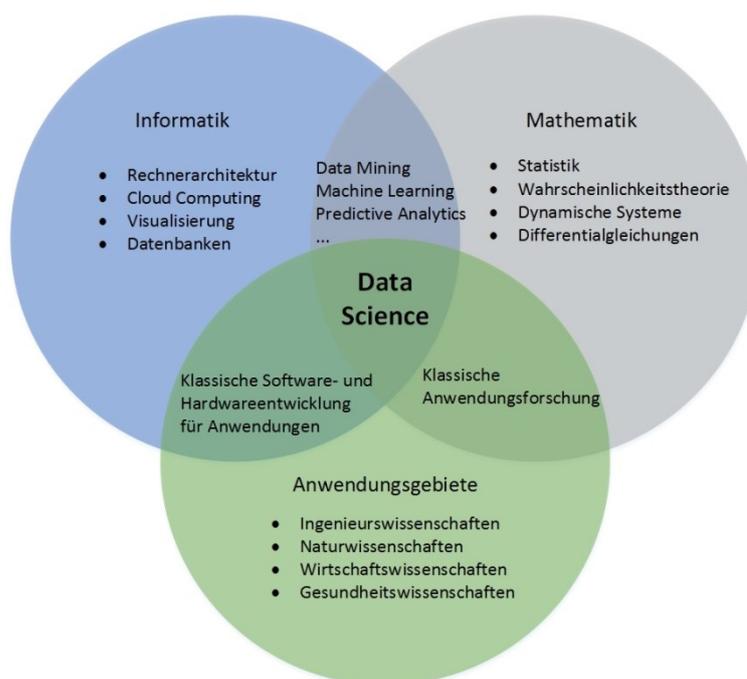
### 3.1 Einführung

Die Verknüpfung datengestützter Produktion und Logistik sind in der 4. industriellen Revolution auf eine hohe Vernetzung und zielorientierte Auswertung von Daten zurückzuführen.<sup>72</sup> Dies soll zu einer effektiveren Datenverarbeitung und -analyse führen, welche wiederum neue Optimierungspotenziale aufzeigen. Durch diese Fähigkeit erhält die Datenanalyse einen immer höheren Stellenwert in den Bereichen Produktion und Logistik.<sup>73</sup>

---

<sup>72</sup> Vgl. Bauernhansl, T. (2014), S. 14

<sup>73</sup> Vgl. Freitag, M. et al. (2015), S. 22 f.



**Abbildung 13: Komponenten von Data Science<sup>74</sup>**

Der seit 40 Jahren existierende Begriff Data Science gibt der Kombination von Mathematik, Informatik und den jeweiligen Anwendungsgebieten einen Überbegriff.<sup>75,76</sup> Dies wird in Abbildung 13 dargestellt. Durch eben die Anwendungsgebiete ergeben sich neue Potenziale für die Datenanalyse, bestehende Ansätze zu verbessern. Es werden Ansätze verfolgt, wie Daten die Transparenz von Prozessen erhöhen und Entscheidungen zu deren Planung und Steuerung auf Basis dieser getroffen werden können.<sup>77</sup>

Es gibt eine Fülle an Theorien und Ansätzen in den Bereichen Mathematik und Informatik sowie in deren Zusammenspiel. Dazu gehören u. a. Data Mining, Statistik, explorative Datenanalyse, Predictive Analytics und Big Data. Unterschiedliches Vorgehen und Ziele kennzeichnen diese, dennoch sind die Übergänge oft fließend. So kann eine Methode unterschiedlichen Bereichen zugeordnet werden und ein anderes Ziel in diesem Kontext verfolgen. Dennoch sind diese Überbegriffe oftmals anders definiert.<sup>78</sup> So bezeichnen die Begriffe Data Mining<sup>79,80</sup>, explorative Datenanalyse<sup>81,82</sup> und Predictive Analytics<sup>83</sup> ein Konzept, ohne jedoch die in diesem Zusammenhang verwendeten Methoden zur Erreichung eines Ziels zu spezifizieren. Maschinelle Lernverfahren hingegen werden als eine Klasse von Verfahren beschrieben, die in Modell (Algorithmus), Verlustfunktion und Vorgehensweise zur Parameteroptimierung

<sup>74</sup> Quelle: in Anlehnung an Freitag, M. et al. (2015), S. 23 und Schutt, R.; O'Neil, C. (2013), S. 7

<sup>75</sup> Vgl. Freitag, M. et al. (2015), S. 22

<sup>76</sup> Vgl. Schutt, R.; O'Neil, C. (2013), S. 7

<sup>77</sup> Vgl. Freitag, M. et al. (2015), S. 23

<sup>78</sup> Vgl. ebda., S. 23

<sup>79</sup> Vgl. Witten, I. H. et al. (2011), S.4

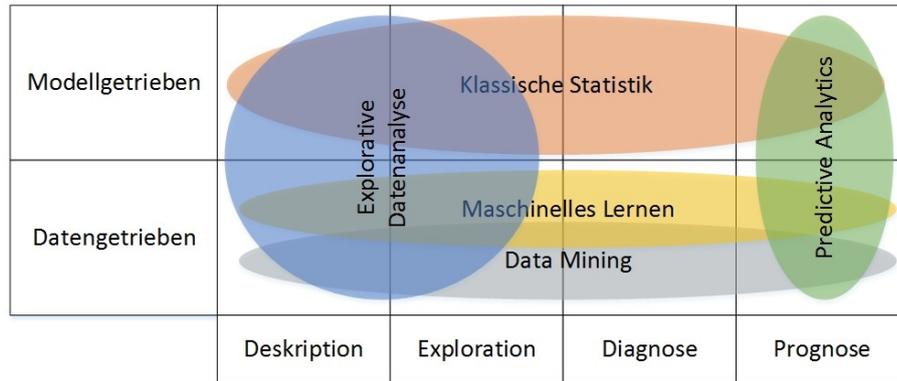
<sup>80</sup> Vgl. Fayyad, U. et al. (1996), S. 37

<sup>81</sup> Vgl. Fahrmeir, L. et al. (2004), S. 12

<sup>82</sup> Vgl. Shmueli, G. (2010), S. 297

<sup>83</sup> Vgl. Abbott, D. (2014), S. 3

unterschieden werden. Eine Aneinanderreihung und Anwendung dieser Verfahren mit unterschiedlichen Ergebnissen wird als Data Mining Prozess bezeichnet.<sup>84</sup>



**Abbildung 14: Qualitative Unterscheidung verschiedener Ansätze zur Datenanalyse<sup>85</sup>**

Es gibt verschiedene Ansätze zur Klassifizierung der Methoden zur Datenanalyse. So definiert IBM über die Eigenschaften der Daten<sup>86</sup>, wohingegen Abbildung 14 die Herangehensweise an die Problemstellung und das Ziel der Anwendung zur Unterscheidung heranzieht. Bei der Herangehensweise kann in modellgetrieben, Einsatz der klassischen Statistik, und datengetrieben, Data Mining und maschinelles Lernen, differenziert werden. Beim modellgetriebenen Ansatz liegt die Annahme einem stochastischen Modell zugrunde, wohingegen beim datengetriebenen Ansatz der datengenerierende Prozess als unbekannt betrachtet und keine Annahme über ein stochastisches Modell getroffen wird.<sup>87</sup>

Ebenfalls in Abbildung 14 ist die Differenzierung durch das Ziel der Anwendung ersichtlich. Hierbei wird nach Datenanalyse zur Beschreibung (Deskription), Erforschung (Exploration), Erklärung (Diagnose) und Prognose der vorhandenen Daten unterschieden. Die Deskription beschreibt anhand vorliegender Daten historische Systemzustände und Ereignisse. Die Deskription ist eng verknüpft mit der Exploration, so dass diese häufig zusammen ausgeführt werden. Die Diagnose baut auf der Frage auf, warum etwas passiert ist, und die Prognose darauf, was passieren wird.<sup>88</sup>

Sämtliche in Abbildung 14 dargestellten Analyseziele können mit unterschiedlichen Methoden der klassischen Statistik, des Data Minings und des maschinellen Lernens erreicht werden. Daneben befasst sich die explorative Datenanalyse mit der Erforschung der Daten, um ein Verständnis für diese zu schaffen.<sup>89</sup> Predictive Analytics beschäftigt sich mit der Prognose der Daten, wobei dieser unterstellt wird, dass die modell- und datengetriebenen Ansätze zur Prognose ungesehener Daten, eingesetzt werden kann. Die Annahme, dass diese geeignete Ergebnisse liefert, ist oftmals falsch.

<sup>84</sup> Vgl. Alpaydin, E. (2010), S. 2 f.

<sup>85</sup> Quelle: Freitag, M. et al. (2015), S. 24 (leicht modifiziert)

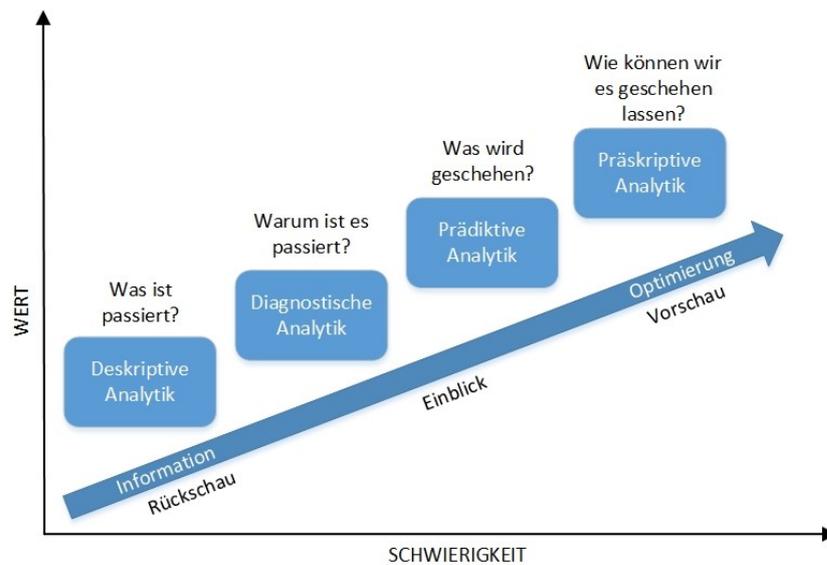
<sup>86</sup> IBM Big Data & Analytic Hub, <http://www.ibmbigdatahub.com/infographic/four-vs-big-data> (Zugriff: 09.02.2017)

<sup>87</sup> Vgl. Breiman, L. (2001), S. 1

<sup>88</sup> Vgl. Fahrmeir, L. et al. (2004), S. 11

<sup>89</sup> Vgl. Freitag, M. et al. (2015), S. 24

Hierbei muss der erwartete Prognosefehler, der sich in die Komponenten zerlegen lässt, betrachtet werden.<sup>90</sup>



**Abbildung 15: Die evolutionstechnischen Stufen der Analytik<sup>91</sup>**

Neben den bereits genannten Ansätzen ergibt sich aus der Betrachtung der in Abbildung 15 dargestellten evolutionstechnischen Stufen der Analytik ein weiterer Ansatz. Bei der präskriptiven Analytik sind das Ziel sowohl eine Prognose als auch eine Handlungsanleitung, um Prognosen Realität werden zu lassen bzw. diese ggf. zu verändern.<sup>92</sup>

All diese Ansätze können auf die Daten, die in Produktions- und Logistiksystemen anfallen, angewendet werden. Durch die Wahl der geeigneten Methoden kann ein klares Verständnis des Systems und des Optimierungspotenzials aufgezeigt werden. Dies soll mit den Beispielen in Tabelle 3 veranschaulicht werden. Die aufgezählten Aufgaben, Ausreißerererkennung, Assoziation, Clustering (Segmentierung), Klassifikation, Prognose und Regression, sind typische Aufgaben der Datenanalyse. Jeweils einer Aufgabe werden verschiedenen Methoden zugewiesen, um das Ziel, im Anwendungsbeispiel dargestellt, zu erreichen.<sup>93</sup>

<sup>90</sup> Vgl. Hastie, T. et al. (2009), S. 38

<sup>91</sup> Quelle: in Anlehnung an Gartner, <http://www.gartner.com/it-glossary/predictive-analytics> (Zugriff: 17.05.2015) zitiert nach Freitag, M. et al. (2015)

<sup>92</sup> Vgl. Freitag, M. et al. (2015), S. 25

<sup>93</sup> Vgl. ebda., S. 26

Tabelle 3: Beispiele für die Aufgaben der Datenanalyse<sup>94</sup>

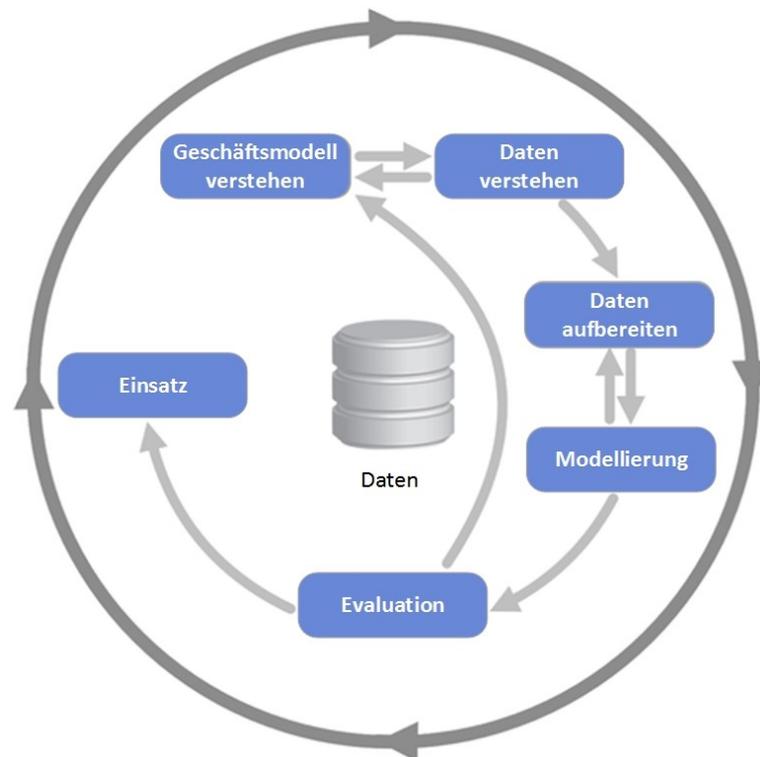
Aufgabe	Beschreibung	Methodenbeispiele	Anwendungsbeispiele
Ausreißererkennung	Entdeckung auffälliger Objekte oder Werte	<ul style="list-style-type: none"> <li>• Ausreißertests</li> <li>• Lineare Regression</li> </ul>	<ul style="list-style-type: none"> <li>• Filterung von Sensordaten</li> <li>• Identifikation von Kreditkartenbetrügern</li> </ul>
Assoziation	Untersuchung der Zusammenhänge und durch Abhängigkeiten durch Wenn-dann-Regeln	<ul style="list-style-type: none"> <li>• Assoziationsregeln</li> <li>• Bayessche Netze</li> </ul>	<ul style="list-style-type: none"> <li>• Warenkorbanalyse</li> <li>• Analyse von Kaufverhalten und Kundenbedürfnissen</li> </ul>
Clustering/ Segmentierung	Bildung von a priori unbekanntem Klassen aufgrund von Ähnlichkeiten	<ul style="list-style-type: none"> <li>• Clusteranalyse</li> <li>• Neuronale Netze</li> <li>• Selbstorganisierende Karten</li> </ul>	<ul style="list-style-type: none"> <li>• Bildung verschiedener Liefergebiete</li> <li>• Bildung von Produktklassen ähnlicher Eigenschaften</li> </ul>
Klassifikation	Zuordnung von Objekten durch Vergleiche von Objekteigenschaften mit den Eigenschaften vorgegebener Klassen	<ul style="list-style-type: none"> <li>• Diskriminanzanalyse</li> <li>• Entscheidungsbäume</li> <li>• Neuronale Netze</li> </ul>	<ul style="list-style-type: none"> <li>• Prozess- und Qualitätsanalyse</li> <li>• Zuordnung von neuen Produkten zu vorgegebenen Produktklassen</li> </ul>
Prognose	Berechnung zukünftiger Werte auf Basis historischer Daten	<ul style="list-style-type: none"> <li>• ARIMA</li> <li>• Exponentielle Glättung</li> <li>• Neuronale Netze</li> <li>• Prognoseverfahren der nichtlinearen Dynamik</li> </ul>	<ul style="list-style-type: none"> <li>• Bedarfsprognosen für die Produktionsplanung</li> <li>• Prognosen der Restlebensdauer von Produktionsanlagen</li> </ul>
Regression	Modellierung der Beziehung zwischen verschiedenen Variablen	<ul style="list-style-type: none"> <li>• Lineare Regression</li> <li>• Neuronale Netze</li> <li>• Support Vector Regression</li> </ul>	<ul style="list-style-type: none"> <li>• Bestimmung des Zusammenhangs zwischen Absatz und Marketing</li> <li>• Bestimmung des Zusammenhangs zwischen Maschineneigenschaften und Lebensdauern</li> </ul>

<sup>94</sup> Quelle: Freitag, M. et al. (2015), S. 26

Diese Einführung in Data Science und dessen Umfeld, in dem sich diese Arbeit bewegt, gibt erste Einblicke in die Möglichkeiten zur Lösung der Aufgabenstellung. Daran schließt das Kapitel 3.5 an und geht auf geeignete Modelle des maschinellen Lernens ein. Zunächst wird das CRISP-DM-Modell vorgestellt. Dieses gibt dem Prozess der Datenanalyse bis hin zur Modellanwendung einen standardisierten Rahmen.

## 3.2 Cross-Industry Standard Process for Data Mining

1996 entstand die Idee für ein nicht eigentumsrechtlich geschütztes, dokumentiertes, branchenunabhängiges und frei verfügbares Data Mining Modell. Dies wird mit dem *CRISP-DM (Cross-Industry Standard Process for Data Mining)* verwirklicht. Das umfassende Modell gliedert den Prozess des Data Minings in die sechs Phasen *Geschäftsmodell verstehen*, *Daten verstehen*, *Daten aufbereiten*, *Modellierung*, *Evaluation* und *Einsatz*. Diese und ihre Beziehungen zueinander sind in Abbildung 16 dargestellt.<sup>95</sup>



**Abbildung 16: Die sechs Phasen des CRISP-DM<sup>96</sup>**

In Abbildung 17 wird aufgelistet, welche Aufgaben bzw. Ergebnisse jede Phase beinhalten. So beschäftigt sich die erste Phase mit dem Umfeld und den generellen Zielen des Prozesses. Die zweite Phase schafft die Grundlage für spätere Phasen, indem die Daten gesammelt, beschrieben und ihre Qualität beurteilt wird. Um in der vierten Phase ein Modell erstellen zu können, müssen in der dritten Phase die Daten

<sup>95</sup> Vgl. Shearer, C. (2000), S. 13

<sup>96</sup> Quelle: in Anlehnung an ebda., S. 14

aufbereitet werden. Hierbei werden Daten gelöscht, Attribute abgeleitet und neuen Datensets erstellt. In der fünften Phase erfolgen die Evaluation der Ergebnisse und das daraus abgeleitete weitere Vorgehen. Das Modell schließt mit der Phase sechs, in die sich mit dem Einsatz der Erkenntnisse bzw. des Modells und der Sicherstellung der weiteren Gültigkeit der Ergebnisse befasst.<sup>97</sup>

Geschäftsmodell verstehen	Daten verstehen	Daten aufbereiten	Modellierung	Evaluation	Einsatz
<b>Bestimmen der Unternehmensziele</b> - Hintergrund - Unternehmensziele - Kriterien für den Unternehmenserfolg  <b>Beurteilung der Situation</b> - Bestandsaufnahme der Ressourcen - Anforderungen, Annahmen und Einschränkungen - Kriterien für den Unternehmenserfolg  <b>Festlegen der Data Mining Ziele</b> - Ziele - Erfolgskriterien  <b>Erstellen eines Projektplans</b> - Projektplan - Erste Einschätzung von zu verwendenden Methoden	<b>Sammeln von Ausgangsdaten</b> - Report über Sammlung der Daten  <b>Beschreibung der Daten</b> - Report zur Beschreibung der Daten  <b>Erforschen der Data</b> - Report über Erforschung der Daten  <b>Datenqualität beurteilen</b> - Report über die Datenqualität	<b>Datenset</b> - Beschreibung des Datensets  <b>Auswahl der Daten</b> - Gründe für Einbeziehung/Ausschluss  <b>Löschen von Daten</b> - Report über das Löschen von Daten  <b>Erschaffen von Daten</b> - Abgeleitete Attribute - Erschaffen neuer Daten  <b>Zusammenführen der Daten</b> - neues Datenset  <b>Formatieren der Daten</b> - formatierte Daten	<b>Auswahl der Methodik</b> - Methodik - Annahmen  <b>Erstellen eines Test Designs</b> - Test Design  <b>Erstellen des Modells</b> - Parametereinstellung - Modelle - Modellbeschreibung  <b>Bewerten des Modells</b> - Bewertung des Modells - Überarbeitung der Parametereinstellungen	<b>Evaluieren der Ergebnisse</b> - Bewertung der Data Mining Ergebnisse in Bezug auf die Unternehmensziele - Bewilligtes Modell  <b>Überprüfen des Prozess</b> - Überprüfung des Prozess  <b>Festlegen nächster Schritte</b> - Liste möglicher Maßnahmen - Entscheidungen	<b>Planen des Einsatzes</b> - Einsatzplan  <b>Planen der Überwachung und Wartung</b> - Überwachungs- und Wartungsplan  <b>Erstellen des finalen Berichts</b> - finaler Bericht - finale Präsentation  <b>Bewerten des Projekts</b> - Dokumentation des Projekts/Erfahrungen

Abbildung 17: Aufgaben und Output des CRISP-DM Modells<sup>98</sup>

Das CRISP-DM Modell basiert auf dem *KDD-Modell*<sup>99</sup> (*Knowledge Discovery in Databases*), das ebenso für das *SEMMA-Modell*<sup>100</sup> (*Sample, Explore, Modify, Model, Assess*) als Grundlage dient. Die beiden Modelle, CRIPS-DM und SEMMA, unterscheiden sich nur marginal. Ein Vergleich der einzelnen Phasen ist in Tabelle 4 zu finden. Alle drei dienen dem Zweck den Anwendern bei ihren Data Mining Projekten eine Struktur zu geben.<sup>101</sup>

Tabelle 4: Vergleich der Schritte von KDD, SEMMA und CRISP-DM<sup>102</sup>

KDD	SEMMA	CRISP-DM
Pre KDD	—	Geschäftsmodell verstehen
Selektion	Datenauswahl	Daten verstehen
Datenaufbereitung	Exploration	
Transformation	Modifizierung	Daten aufbereiten
Data Mining	Modellieren	Modellierung
Interpretation/Evaluation	Bewertung	Evaluation
Post KDD	—	Einsatz

<sup>97</sup> Vgl. Shearer, C. (2000), S. 13

<sup>98</sup> Quelle: in Anlehnung an ebda., S. 14

<sup>99</sup> Vgl. Fayyad, U. et al. (1996), S. 37

<sup>100</sup> Vgl. Fernandez, G. (2003), S. 5 f.

<sup>101</sup> Vgl. Azevedo, A.; Santos, M. F. (2008), S. 185

<sup>102</sup> Quelle: in Anlehnung an ebda., S. 185

CRISP-DM ist vollständig und dokumentiert. Sämtliche Stufen wurden nachvollziehbar strukturiert und definiert. Dies führt zu einem guten Verständnis und erleichtert eine spätere Überarbeitung.<sup>103</sup> Daher wurde entschieden, dass sich die folgenden Kapitel am Aufbau des CRISP-DM orientieren, und dieses als Referenz für das praktische Fallbeispiel einzusetzen. Als erstes werden die Grundlagen der Datenaufbereitung und ihre Methoden aufgezeigt.

### 3.3 Datenverständnis

Das Verständnis des Rohmaterials - der Daten - wird, wie in Kapitel 3.2 erklärt, in der Phase zwei "Daten verstehen" erzielt.

"Daten werden stets an gewissen Objekten beobachtet."<sup>104</sup> Diese Objekte werden als statistische Einheiten bezeichnet. An ihnen werden Merkmale beobachtet, die in weiterer Folge als Attribute bezeichnet werden, und deren konkreter Wert für eine bestimmte statistische Einheit Ausprägung genannt wird.<sup>105</sup> Eine inhaltlich zusammenhängende Menge von Ausprägungen der Attribute einer statistischen Einheit wird sowohl logisch als auch physisch zu Datensätzen zusammengefasst und kann einen Index zur eindeutigen Zuordnung enthalten.<sup>106</sup>

Im Folgenden werden mehrere Ansätze, Methoden und Kriterien vorgestellt, Daten zu analysieren und für die Aufbereitung wichtige Erkenntnisse zu gewinnen. Hierbei liegt der inhaltliche Fokus auf den verschiedenen Attributtypen, in denen Daten vorliegen können, und den Maßen der uni- und multivariaten Deskription. Hierbei handelt es sich um statistische Maße, die für ein Attribut berechnet werden können, sowie um Zusammenhangsmaße, die die Beziehung zwischen den Attributen beschreiben können. Mit Hilfe dieser können Schlüsse über die Verteilung und Zusammenhänge innerhalb der Daten erhalten werden.

#### 3.3.1 Attributtypen

Ausprägungen von Attributen werden durch Messen bzw. Beobachten nach festgelegten Regeln einen Zahlenwert zugewiesen. Grundsätzlich gilt dabei, dass die Messbarkeit der jeweiligen Ausprägung eines Attributs von der Art der Eigenschaft abhängig ist. Das bedeutet, wie gut sie in Zahlen ausgedrückt werden kann. So ist der Durchmesser eines Werkstücks leicht messbar, allerdings die subjektive Wahrnehmung eines Menschen nicht. Um diesen Ausprägungen einen Rahmen zu geben, gibt es Skalen. Diese können durch ihre unterschiedlichen Niveaus, den Skalenniveaus, voneinander unterschieden werden und geben an, in welcher Art und Weise die Eigenschaften eines Objekts gemessen worden sind.<sup>107</sup>

---

<sup>103</sup> Vgl. Azevedo, C. S.; Santos, M. F. (2005) zitiert nach Azevedo, A.; Santos, M. F. (2008), S. 185

<sup>104</sup> Fahrmeir, L. et al. (2004), S. 14

<sup>105</sup> Vgl. ebda., S. 14

<sup>106</sup> Vgl. Mertens, P. et al. (2012), S. 39

<sup>107</sup> Vgl. Backhaus, K. et al. (2011), S. 10

In Tabelle 5 sind die Nominalskala, die Ordinalskala, die Intervallskala und die Ratioskala mit ihren Merkmalen sowie die Anwendbarkeit von Rechenoperationen auf die Daten angegeben. Es ist zu bedenken, dass die Daten immer eine Information darstellen. Somit können Ausprägungen, die unterschiedlich skaliert sind, nicht wahllos miteinander verglichen werden.<sup>108</sup>

**Tabelle 5: Skalenniveau<sup>109</sup>**

Skala		Merkmale	Mögliche rechnerische Handhabung
Nicht-metrische Skalen	<i>Nominalskala</i>	Klassifizierung qualitativer Eigenschaftsausprägungen	Bildung von Häufigkeiten
	<i>Ordinalskala</i>	Rangwert mit Ordinalskalen	Median, Quantile
Metrische Skalen	<i>Intervallskala</i>	Skala mit gleichgroßen Abschnitten ohne natürlichen Nullpunkt	Subtraktion, Mittelwerte
	<i>Ratioskala</i>	Skala mit gleichgroßen Abschnitten und natürlichem Nullpunkt	Summe, Division, Multiplikation

So stellt die *Nominalskala* die primitivste Grundlage des Messens da. Die Klassifizierung von qualitativen Eigenschaften, wie z. B. die Bestimmung der Farbe (rot – blau – grün – ...), kann anschließend in Zahlen übersetzt werden; dennoch sind keine arithmetischen Operationen, wie Addition, Subtraktion, Multiplikation oder Division, erlaubt. Einzig die Häufigkeit des Auftretens der einzelnen Ausprägungen, die Häufigkeit, ist zählbar.<sup>110</sup>

Die zweite nicht-metrische Skala ist die *Ordinalskala*. Bei ihr wird eine Rangordnung mit Hilfe von Rangwerten (d. h. ordinalen Zahlen) ermittelt. Ein Beispiel für eine Ordinalskala ist der Vergleich von Produkt A mit Produkt B, wobei Produkt A besser als Produkt B ist. Dieser Vergleich sagt nichts über die Abstände zwischen den einzelnen Objekten aus. Auch bei dieser Skala sind daher keine arithmetischen Operationen erlaubt. Neben der Häufigkeit sind der Median oder die Quantile erlaubte statistische Größen.<sup>111</sup>

Die *Intervallskala* stellt das nächsthöhere Messniveau dar. Gleichgroße definierte Skalenabschnitte, wie auf z. B. auf der Celsius-Skala zur Temperaturmessung, erlauben erstmals die Anwendung von arithmetischen Operationen (Addition und Subtraktion) sowie statistische Größen wie das arithmetische Mittel und die Standardabweichung. So kann der Temperaturunterschied zwischen zwei Messungen durch dessen Differenz berechnet werden, die wiederum einen eigenen Informationsgehalt hat.<sup>112</sup>

<sup>108</sup> Vgl. Backhaus, K. et al. (2011), S. 10

<sup>109</sup> Quelle: ebda., S. 12

<sup>110</sup> Vgl. ebda., S. 10

<sup>111</sup> Vgl. ebda., S. 11

<sup>112</sup> Vgl. ebda., S. 11

Das höchste Messniveau ist die *Ratio- (oder Verhältnis-)skala*. Der Unterschied zwischen ihr und der Intervallskala liegt in der Existenz eines natürlichen Nullpunkts, der für das betreffende Attribut die Aussage "nicht vorhanden" wiedergibt. Dies ist bei der oben als Beispiel aufgeführten Celsius-Skala nicht der Fall. Die Ratioskala kann daher auf physische Attribute, wie Länge und Gewicht, angewendet werden. Aus diesen so skalierten Daten können sämtliche arithmetische Operationen und die oben genannten statistischen Größen berechnet werden. Zusätzlich können das geometrische Mittel und der Variationskoeffizient berechnet werden.<sup>113</sup>

Grundsätzlich ist zu erkennen, dass je höher das Skalenniveau ist, desto höher ist der Informationsgehalt der betreffenden Daten, und desto mehr Rechenoperationen und statistische Größen können mit den Daten ausgeführt werden. Die Umwandlung von einem höheren auf ein niedrigeres Skalenniveau ist generell möglich. Sinnvoll ist das, wenn die Daten für die Analyse vereinfacht werden müssen. Allerdings geht diese Transformation mit einem Informationsverlust einher.<sup>114</sup>

Zusätzlich zur Einteilung in die Skalen können Attribute als *diskret* und *stetig* bezeichnet werden. Diese geben darüber Auskunft, ob Ausprägungen eines Attributs endlich viele Werte annehmen können (diskret) oder alle Werte eines Intervalls (stetig). Liegen diskrete Attribute in sehr feinen Abstufungen vor, werden diese als *quasistetig* bezeichnet.<sup>115</sup>

Eine Unterscheidung in *qualitativ (oder kategorial)* und *quantitativ* wird dann getroffen, wenn zwischen Qualität und Ausmaß getrennt wird. Bei qualitativen Attributen handelt es sich um Größen mit endlich vielen Ausprägungen, die höchstens ordinalskaliert sind. Quantitative Attribute kennzeichnen sich durch die Angabe der Ausprägung in deren Intensität oder Ausmaß. Ordinalskalierte Attribute haben in dieser Ordnung eine Zwitterstellung; sie besitzen sowohl einen schwach quantitativen als auch einen dominierenden qualitativen Anteil.<sup>116</sup>

Nach dieser Einführung, in welcher Form Merkmale skaliert sein können, beschäftigt sich das nächste Kapitel mit den statistischen Maßen, die von einem Attribut abhängig sind. Ziel ist es, zu vermitteln, wie sie berechnet werden können, welche Bedeutung sie haben und welche Erkenntnisse gewonnen werden können.

### 3.3.2 Univariate Deskription

Wird bei einer Erhebung vom Umfang  $n$  ausgegangen, bei der an  $n$  Einheiten die Werte  $x_1, \dots, x_n$  eines Attributes erfasst wurden, so werden diese als Urliste, Roh- oder Primärdaten bezeichnet. Bereits bei kleinem oder mittlerem Umfang kann die Darstellung der Rohdaten unübersichtlich werden. Um eine übersichtliche und zusammenfassende Darstellung zu ermöglichen, können verschiedene Methoden angewendet bzw. Werte errechnet werden.<sup>117</sup>

---

<sup>113</sup> Vgl. Backhaus, K. et al. (2011), S. 11

<sup>114</sup> Vgl. ebda., S. 11 f.

<sup>115</sup> Vgl. Fahrmeir, L. et al. (2004), S. 16 f.

<sup>116</sup> Vgl. ebda., S. 19

<sup>117</sup> Vgl. ebda., S. 31

In einem ersten Teil werden die Berechnung von Häufigkeit und die Darstellung von Verteilungen betrachtet. Diese kann auch durch Maßzahlen bzw. Parameter ausgedrückt werden.<sup>118</sup> Zu diesem Zweck wird eine kurze Einführung zu Lagemaßen sowie zu den Streuungsmaßen gegeben.

### Häufigkeiten und Graphische Darstellungen

Die Urliste kann nach *Häufigkeiten* der Ausprägungen  $a_1, a_2, \dots, a_k, k \leq n$  untersucht werden. Die Anwendung auf die Nominalskala bringt keine inhaltliche Bedeutung. Bei kategorialen Attributen ist  $k$  die Anzahl der Kategorien und zudem meist maßgeblich kleiner als  $n$ . Liegen metrischen Werte in der Urliste vor, ist  $k$  meist fast genauso groß wie  $n$ . Die relative Häufigkeit  $f_j$  ist das Verhältnis zwischen der absoluten Häufigkeit  $h_j$  zu  $n$ . Häufigkeiten und Ausprägungen zusammen werden als Häufigkeitsdaten bezeichnet.<sup>119</sup>

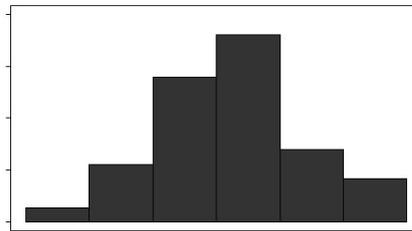


Abbildung 18: Beispiel eines Histogramms<sup>120</sup>

Häufigkeitsdaten können durch ein *Histogramm* visualisiert werden.<sup>121</sup> Zusätzlich können sie in Intervalle zusammengefasst werden. Dieser Vorgang wird gruppieren, auch klassieren oder kategorisieren, genannt.<sup>122</sup> Die dadurch entstandenen Klassen werden auf der Abszisse und die relative oder absolute Häufigkeit auf der Ordinate aufgetragen. Die entstandenen Blöcke sind proportional zu der absoluten oder der relativen Häufigkeit. Eine Darstellung, wie sie in Abbildung 18 zu sehen ist, ist insbesondere von der Klassenbreite bzw. von der Anzahl der Klassen abhängig. Wird die Klassenbreite größer, sinkt die Anzahl der Sprünge. Der Extremfall wäre eine maximale Breite, bei dem nur ein einziger Block dargestellt wird. Es gibt mehrere Ansätze die Klassenanzahl zu berechnen, so wie (3.1), (3.2) und (3.3). Zusätzlich soll der subjektive optische Eindruck des Histogramms berücksichtigt werden.<sup>123</sup>

$$k = \sqrt{n} \quad (3.1)$$

$$k = 2\sqrt{n} \quad (3.2)$$

$$k = 10 \log n \quad (3.3)$$

<sup>118</sup> Vgl. Fahrmeir, L. et al. (2004), S. 53

<sup>119</sup> Vgl. ebda., S. 32

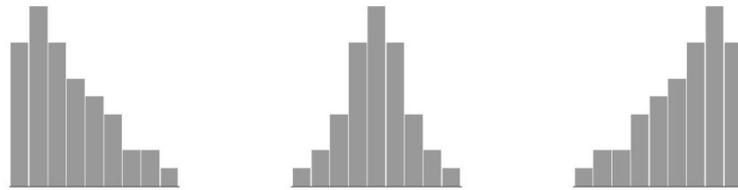
<sup>120</sup> Quelle: ebda., S. 43

<sup>121</sup> Vgl. ebda., S. 41

<sup>122</sup> Vgl. ebda., S. 17

<sup>123</sup> Vgl. ebda., S. 42 f.

Verteilungen können durch diese Darstellung in *unimodal*, *bimodal* und *multimodal* eingeteilt werden. Dies gibt an, ob ein, zwei oder mehrere Gipfel in einer Verteilung vorliegen. Abbildung 18 stellt eine unimodale Verteilung dar.<sup>124</sup>



**Abbildung 19: Eine linkssteile (a), symmetrische (b) und rechtssteile Verteilung (c)**<sup>125</sup>

Sind Verteilungen wie in Abbildung 19 gegeben, wird von *Schief* bzw. *Symmetrie* gesprochen. Wenn die rechte und linke Hälfte einer Verteilung annähernd zueinander gleich sind, heißt diese *symmetrisch*. Bei empirischen Verteilungen ist eine exakte Symmetrie eher selten. Sind diese deutlich unsymmetrisch, werden sie als *schief* bezeichnet. Die Verteilung wird als *linkssteil* (oder *rechtsschief*) bezeichnet, wenn der mehrheitliche Teil der Messwert linksseitig konzentriert ist. Fällt die Verteilung nach rechts steiler und links flacher ab, wird diese *rechtssteil* bzw. *linksschief* genannt.<sup>126</sup>

### Lagemaße

Lagemaße sind Maßzahlen zur Beschreibung des Schwerpunkts einer Verteilung. Sie werden in numerischen Werten angegeben. Es gibt verschiedene Lagemaße und die Auswahl, welche sinnvoll zur Beantwortung einer Fragestellung ist, hängt von Kontext, der Datensituation und dem Skalenniveau des Attributs ab.<sup>127</sup>

Das *arithmetische Mittel*  $\bar{x}$  (3.4) ist das gebräuchlichste Maß für das Festlegen des Zentrums einer Verteilung.<sup>128</sup> Es wird aus der Urliste berechnet.<sup>129</sup>

$$\bar{x}_{arithm} = \frac{1}{n}(x_1 + x_2 + \dots + x_n) = \frac{1}{n} \sum_{i=1}^n x_i \quad (3.4)$$

Das arithmetische Mittel wird von extremen Werten und Ausreißern in den Daten beeinflusst. Solche können entstehen, wenn eine empfindliche Einheit betrachtet wird, oder ein Fehler bei der Datenerhebung oder -aufbereitung unterlaufen ist.<sup>130</sup> Die Erkennung solcher Ausreißer wird im Kapitel 3.4.1 behandelt.

Es gibt *robuste* bzw. *resistente* Lagemaße, wie den *Median* (3.5), die den Einfluss von Extremwerten beschränkt. Der Median wird ermittelt, in dem die Urliste mit den Werten  $x_1, \dots, x_n$  geordnet wird. In der dadurch enthaltenen geordneten Urliste  $x_1 \leq \dots \leq x_i \leq \dots \leq x_n$  wird der Median in der Mitte platziert. Daraus folgt, dass mindestens 50 % der

<sup>124</sup> Vgl. Fahrmeir, L. et al. (2004), S. 46

<sup>125</sup> Vgl. ebda., S. 49

<sup>126</sup> Vgl. ebda., S. 48

<sup>127</sup> Vgl. ebda., S. 53

<sup>128</sup> Vgl. Haseloff, O. W.; Hoffmann, H. J. (1965), S.43

<sup>129</sup> Vgl. Fahrmeir, L. et al. (2004), S. 53

<sup>130</sup> Vgl. ebda., S. 55

Daten kleiner oder gleich und mindestens 50 % größer oder gleich  $x_{med}$  sind. Dies ist ein weiterer Vorteil des Median gegenüber dem arithmetischen Mittel.<sup>131</sup>

$$x_{med} = \begin{cases} x_{(n+\frac{1}{2})} & \text{für } n \text{ ungerade} \\ \frac{1}{2} \left( x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)} \right) & \text{für } n \text{ gerade} \end{cases} \quad (3.5)$$

Ebenfalls zu den Lagemaßen wird der *Modus*  $x_{mod}$  gezählt. Dieser gibt Auskunft darüber, welche Ausprägung am häufigsten vorkommt. Der Modus ist eindeutig, wenn die Häufigkeitsverteilung ein eindeutiges Maximum aufweist.<sup>132</sup>

Handelt es sich um metrisch skalierte Attribute, kann eine Aussage über Symmetrie und Schiefe mithilfe des arithmetischen Mittels, Median und Modus getroffen werden. Dazu dienen die *Lageregeln*.<sup>133</sup>

- Symmetrische Verteilung:  $\bar{x} \approx x_{med} \approx x_{mod}$
- Linkssteile Verteilung:  $\bar{x} > x_{med} > x_{mod}$
- Rechtssteile Verteilung:  $\bar{x} < x_{med} < x_{mod}$

Die größte Bedeutung haben diese Lageregeln für unimodale Verteilungen. Je größer die Unterschiede zwischen  $\bar{x}$ ,  $x_{med}$  und  $x_{mod}$  sind, desto schiefere sind die Verteilungen. Dieses Wissen kann auf den Umgang mit Situationen, in denen ein "Mittelwert" benötigt wird, Einfluss haben.<sup>134</sup>

### Quantile, Standardabweichung und Varianz

Um eine vollständige Beschreibung von Verteilungen zu liefern, sollten die Lagemaße um Angaben zur Streuung ergänzt werden.<sup>135</sup> Hierzu gibt es mehrere Möglichkeiten. Innerhalb dieser Arbeit werden die Quantile, die Standardabweichung und die Varianz betrachtet.

Quantile und deren graphische Darstellung als Boxplots können zur Charakterisierung von Streuungen verwendet werden. Das  $p$ -Quantil teilt die Verteilung von Daten in zwei Teile, sodass circa  $p * 100$  % der Daten unterhalb und circa  $(1 - p) * 100$  % oberhalb liegen. Zur Berechnung der Quantile wird die geordnete Urliste als Basis verwendet. Dies kann verallgemeinert werden: Ein  $p$ -Quantil wird durch den Wert  $x_p$  mit  $0 < p < 1$  definiert.  $p$  stellt den Anteil der Daten dar, die kleiner/gleich  $x_p$  sind (3.6), und der Wert  $1 - p$  jenen, die größer/gleich  $x_p$  sind (3.7).  $[np]$  stellt dabei die nächste zu  $np$  kleinere Zahl dar.<sup>136</sup>

$$\frac{\text{Anzahl } x - \text{Werte} \leq x_p}{n} \geq p \quad (3.6)$$

$$\frac{\text{Anzahl } x - \text{Werte} \geq x_p}{n} \geq 1 - p \quad (3.7)$$

<sup>131</sup> Vgl. Fahrmeir, L. et al. (2004), S. 55 f.

<sup>132</sup> Vgl. ebda., S. 57

<sup>133</sup> Vgl. ebda., S. 60

<sup>134</sup> Vgl. ebda., S. 60

<sup>135</sup> Vgl. ebda., S. 64

<sup>136</sup> Vgl. ebda., S. 64

Daher gilt für das  $p$ -Quantil (3.8), wenn  $np$  nicht ganzzahlig ist, und (3.9), wenn  $np$  ganzzahlig ist.<sup>137</sup>

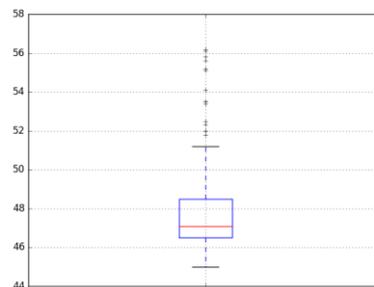
$$x_p = x_{([np]+1)} \quad (3.8)$$

$$x_p \in [x_{(np)}, x_{(np+1)}] \quad (3.9)$$

Das 50 %-Quantil ist der Median. Als *Quartil* werden das *untere Quartil*  $Q_1$  und *obere Quartil*  $Q_3$  bezeichnet. Sie sind das 25 %-Quantil bzw. 75 %-Quantil. Eine weitere Definitionsmöglichkeit wäre, die die Datenhälfte, die über und unter dem Median liegen, wiederum zu halbieren.<sup>138</sup> Gemeinsam mit dem Median sind Quartile Indikatoren für eine schiefe Verteilung. Dies ergibt sich aus der Annahme, dass bei einer symmetrischen Verteilung der Abstand zwischen Median und den beiden Quartilen jeweils gleich groß ist. Ist dies nicht der Fall, weist dies auf eine steilere Verteilung bei der kleineren Differenz hin. Ebenfalls ist an den Quartilen zu erkennen, wie weit eine Verteilung streut. Daraus lässt sich die Maßzahl der Streuung, der *Interquartilsabstand (IQR)* (3.10) ableiten.<sup>139</sup>

$$d_Q = x_{0,75} - x_{0,25} \quad (3.10)$$

Die Quartile werden durch die Lage der Daten, die kleiner als das untere Quartil und größer als das obere Quartil sind, nicht beeinflusst. Daher ist der Interquartilsabstand resistent gegen Ausreißer.<sup>140</sup> Dieser Umstand wird sich in der Ausreißeranalyse im Kapitel 3.4.1 zu Nutze gemacht.



**Abbildung 20: Beispiel eines Boxplots<sup>141</sup>**

Graphisch werden die Quartile mit Boxplots, wie einer in Abbildung 20 beispielhaft zu sehen ist, dargestellt. Die rote Linie steht für den Median, das untere Ende der Box für das untere Quartil, die obere für das obere Quartil. Die Länge der Linien, die auch „whiskers“ genannt werden, können selbst definiert werden. Diese könnten den minimalen und maximalen Wert der Verteilung einschließen oder eine definierte Länge, wie z. B.  $Q_1 - 1,5 \times IQR$ ,  $Q_3 + 1,5 \times IQR$  haben.<sup>142</sup>

Neben dem IQR ist die Standardabweichung bzw. ihr Quadrat, die Varianz die bekannteste Maßzahl für die Streuung einer Verteilung. Es wird dabei die Streuung der

<sup>137</sup> Vgl. Fahrmeir, L. et al. (2004), S. 64

<sup>138</sup> Vgl. ebda., S. 65

<sup>139</sup> Vgl. ebda., S. 66

<sup>140</sup> Vgl. ebda., S. 66 f.

<sup>141</sup> Quelle: eigene Darstellung

<sup>142</sup> Vgl. Fahrmeir, L. et al. (2004), S. 67

Daten um ihr Mittel  $\bar{x}$  gemessen und ist somit nur für metrische Attribute sinnvoll anwendbar. Die *Varianz*  $\tilde{s}^2$  (3.11) ist das mittlere arithmetische Mittel der quadratischen Abweichung.<sup>143</sup>

$$\tilde{s}^2 = \frac{1}{n} [(x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2] = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (3.11)$$

Zusätzlich kann die (*korrigierte*) *Stichprobenvarianz* (3.12) geschätzt werden. Dieser Schätzer gibt darüber Aufschluss, ob eine Stichprobe erwartungstreu, also unverzerrt, ist.<sup>144,145</sup>

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (3.12)$$

Sofern Varianzen nicht miteinander verglichen werden, sind sie auf Grund der Verwendung der Quadrate der Abweichung kein geeignetes Maß zur direkten Beurteilung der Streuung. Zu diesem Zweck wird die *Standardabweichung*  $\tilde{s}$  (3.13) herangezogen.<sup>146</sup>

$$\tilde{s} = +\sqrt{\tilde{s}^2} = \frac{1}{x} \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (3.13)$$

### Maßzahlen für Schiefe und Wölbung

Neben den Merkmalen Lage und Streuung können Verteilungen auch in Bezug auf ihre Symmetrie bzw. Schiefe und die Wölbung unterschieden werden. Durch die Wölbung wird gezeigt, ob sich die Daten im Zentrum der Verteilung häufen, oder ob die Verteilung eher flach ist. Das bedeutet, dass die Daten um das Zentrum streuen.<sup>147</sup>

Neben der Ermittlung der Schiefe oder Symmetrie einer Verteilung durch die Lageregeln, gibt es die Möglichkeit für metrische Attribute den *Momentenkoeffizienten*  $g_m$  (3.14) zur Beurteilung der Schiefe zu berechnen.<sup>148</sup>

$$g_m = \frac{m_3}{\tilde{s}^4} \text{ mit } m_3 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3 \quad (3.14)$$

$g_m$  ist maßstabsunabhängig und durch den Erhalt der positiven bzw. negativen Abweichungen vom Mittelwert kann nach folgenden Regeln die Schiefe, wie sie bereits in Abbildung 19 vorgestellt wurde, beurteilt werden.<sup>149</sup>

- $g_m = 0$  für symmetrische Verteilungen
- $g_m > 0$  für linkssteile Verteilungen
- $g_m < 0$  für rechtssteile Verteilungen

<sup>143</sup> Vgl. Haseloff, O. W.; Hoffmann, H. J. (1965), S. 53

<sup>144</sup> Vgl. Czado, C.; Schmidt, T. (2011), S. 5

<sup>145</sup> Vgl. ebda., S. 104

<sup>146</sup> Vgl. Haseloff, O. W.; Hoffmann, H. J. (1965), S. 54

<sup>147</sup> Vgl. Fahrmeir, L. et al. (2004), S. 74

<sup>148</sup> Vgl. ebda., S. 74 f.

<sup>149</sup> Vgl. ebda., S. 75

Wie stark oder wie schwach der zentrale Bereich bzw. die Randbereiche der Daten besetzt sind, kann durch das *Wölbungsmaß nach Fischer* (3.15) angezeigt werden. Bei spitzen Verteilungen sind die Randbereiche stärker besetzt als bei flachen Verteilungen. Die Normalverteilung dient als Vergleichsverteilung.<sup>150</sup>

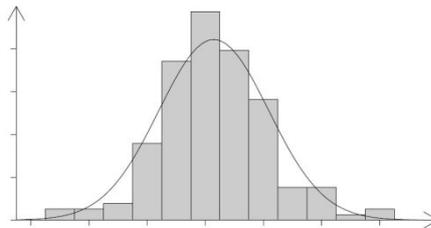
$$\gamma = \frac{m_4}{\bar{s}^4} \text{ mit } m_4 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^4 \quad (3.15)$$

Dabei gilt:<sup>151</sup>

- $\gamma = 0$  bei Normalverteilung
- $\gamma > 0$  bei spitzeren Verteilungen
- $\gamma < 0$  bei flacheren Verteilungen

### Dichtekurven und Normalverteilung

Bei stetigen Merkmalen kommt es vor, dass Histogramme Informationen unterdrücken. Dies geschieht u. a. durch die Wahl der Klassenbreite. Diese empirische Verteilung verhindert ein Erkennen von Besonderheiten in der Verteilung. Daher kann das Histogramm durch eine glatte Kurve, die sogenannte *Dichtekurve*, approximiert werden.<sup>152</sup>



**Abbildung 21: Beispiel für eine Dichtefunktion in einem Histogramm<sup>153</sup>**

In Abbildung 21 ist ersichtlich, dass diese stetige Kurve einige Details glättet. Ihren Zweck, eine kompakte optische Visualisierung zu liefern, erfüllt sie und darf als Dichtekurve bezeichnet werden, wenn sie die Anforderung (3.16) erfüllt. Also  $f(x) \geq 0$  und die Gesamtfläche unter  $f(x)$  gleich 1 ist.<sup>154</sup>

$$\int f(x) dx = 1 \quad (3.16)$$

Bei Dichtekurven kann wieder in symmetrische und schiefe sowie in unimodale und multimodale Dichtefunktionen unterschieden werden. Genauso können Maßzahlen der Häufigkeitsverteilung, wie Median und Quantil, auf die Dichtefunktion übertragen werden. Hierbei steht die Fläche anstatt der Anzahl der jeweiligen Werte im Vordergrund. So definiert das  $p$ -Quantil  $x_p$  für  $0 < p < 1$  jenen Wert, der die Gesamtfläche unter  $f(x)$  in die beiden Flächen  $p * 100\%$  und Flächen  $(1 - p) * 100\%$  teilt. Somit ist der Median  $x_{0,5}$  jener Wert, der die Gesamtfläche in zwei gleich große

<sup>150</sup> Vgl. Fahrmeir, L. et al. (2004), S. 76

<sup>151</sup> Vgl. ebda., S. 76

<sup>152</sup> Vgl. ebda., S. 87

<sup>153</sup> Quelle: ebda., S. 87 (leicht modifiziert)

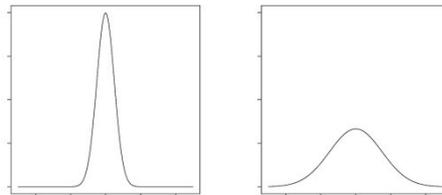
<sup>154</sup> Vgl. ebda., S. 87 f.

Hälften aufteilt. Handelt es sich um eine symmetrische, unimodale Dichtekurve ist der Median die Mitte und repräsentiert gleichzeitig den Modus und den Mittelwert der Dichtefunktion  $f(x)$  bzw. des Attributes  $X$ .<sup>155</sup>

Eine wichtige Klasse der Dichtefunktionen sind die *Normalverteilungen* oder auch *Gauss-Verteilungen*, die symmetrisch, unimodal und glockenförmig sind. Die in Abbildung 21 dargestellte Dichtefunktion gehört zu dieser Klasse. Sie werden durch (3.17) definiert, wobei  $\mu$  der Mittelwert und  $\sigma$  die Standardabweichung ist. Für gegebene Werte dieser beiden Parameter und für die  $\mu \in \mathbb{R}$  und  $\sigma > 0$  gilt, ist  $f(x|\mu, \sigma)$  eindeutig.<sup>156</sup>

$$f(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right) \quad (3.17)$$

Die Exponentialfunktion ist verantwortlich für die Form von  $f(x)$ . Das Maximum der Dichtefunktion ist mit  $x = \mu$  erreicht. Links und rechts davon fällt sie symmetrisch und glockenförmig gegen Null ab. Je kleiner der Wert  $\sigma$  ist, desto steiler fällt diese Kurve ab und umso enger ist die Kurve um den Mittelwert  $\mu$ . Im Gegensatz dazu ist die Kurve umso flacher, je größer  $\sigma$  ist. Dieser Zusammenhang wird in Abbildung 22 veranschaulicht.<sup>157</sup>



**Abbildung 22: Zwei Normalverteilungsdichtekurven mit kleinem und großem  $\sigma$** <sup>158</sup>

Der erste Faktor von (3.17) wird als Normalisierungskonstante bezeichnet. Ihr Zweck ist die Gesamtfläche unter  $f(x)$  auf 1 zu setzen und sie hat keine Auswirkung auf die Form der Kurve.<sup>159</sup>

Der bisher verwendete Mittelwert  $\mu$  und die Standardabweichung  $\sigma$  entsprechen dem arithmetischen Mittel  $\bar{x}$  und der Standardabweichung  $\tilde{s}$  von Beobachtungen  $x_1, \dots, x_n$  eines Attributes  $X$ . Durch die Variation von  $\mu$  und  $\sigma$  können verschiedene Normalverteilungen entstehen. Alle diese können auf die sogenannte *Standardnormalverteilung*, eine Art der Normalverteilung, zurückgeführt werden.<sup>160</sup>

Die Standardnormalverteilung ist als eine Normalverteilung mit Dichtefunktion  $f(x)$ , Mittelwert  $\mu$  und Standardabweichung  $\sigma$  definiert, deren Attribut  $Z$ , wie in (3.18) beschrieben, standardisiert wurde.<sup>161</sup>

$$Z = \frac{X - \mu}{\sigma} \quad (3.18)$$

<sup>155</sup> Vgl. Fahrmeir, L. et al. (2004), S. 88 f.

<sup>156</sup> Vgl. ebda., S. 90

<sup>157</sup> Vgl. ebda., S. 90 f.

<sup>158</sup> Quelle: ebda., S. 92

<sup>159</sup> Vgl. ebda., S. 91

<sup>160</sup> Vgl. ebda., S. 92

<sup>161</sup> Vgl. ebda., S. 92

Durch die Standardisierung entsteht eine Dichtekurve einer Normalverteilung mit  $\mu = 0$  und  $\sigma = 1$ . Das Attribut  $Z$  standardnormalverteilt. Die Dichtekurve wird mit  $\phi(z)$  (3.19) bezeichnet.<sup>162</sup>

$$\phi(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) \quad (3.19)$$

Durch die Standardisierung eines Attributs werden Beobachtungen als Abweichungen vom Mittel  $\mu$  mit der Standardabweichung  $\sigma$  als Maßeinheit angesehen. Dadurch können Berechnungen, wie etwa die der Quantile, auf die Standardnormalverteilung zurückgeführt werden. Diese Werte sind in Tabellen gesammelt und somit entfällt die Berechnung.<sup>163</sup>

### Kolmogorov-Smirnov-Test

Da manche Methoden nur anwendbar sind, wenn Attribute bzw. deren Werte normalverteilt sind, ist es notwendig, diese auf Normalverteilung zu testen.<sup>164</sup> Diese Tests sind ein wichtiges Instrument der quantitativen Analyse.<sup>165</sup> Ziel einer solchen Untersuchung ist es, die Frage zu klären, ob eine tatsächliche Verteilung mit einer vorgegebenen Verteilung übereinstimmt.<sup>166</sup> Diese Tests werden daher auch Anpassungs- oder Goodness-of-fit-Tests genannt.<sup>167</sup> Es kann daraus die folgende Nullhypothese (3.20) abgeleitet werden.<sup>168</sup>

$$H_0: X \sim P_\theta \quad (3.20)$$

$P_\theta$  steht für die gegebene Verteilung, die durch  $\theta$  näher erklärt wird.  $\theta$  kann als Parameter(-vektor) fixiert sein oder nicht. Würde eine Normalverteilung vorliegen, bestünde  $\theta$  aus den Werten  $\mu$  und  $\sigma$ .<sup>169</sup>

Einer dieser Tests ist der *Kolmogorov-Smirnov-Test*. Er dient zur Überprüfung, ob die tatsächliche Verteilung mit der einer stetigen Prüfungsverteilung, wie der Normalverteilung, übereinstimmt. Er basiert auf der empirischen Verteilungsfunktion (3.21).<sup>170</sup>

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x) \quad (3.21)$$

$$F_n(x) = P(X \leq x) \quad (3.22)$$

Mit ihr (3.21) wird jener Teil der Ausprägungen vom Attribut  $X$  angegeben, die kleiner oder gleich  $x_i$  sind. Dieser wird anschließend der theoretischen Verteilungsfunktion, welche die Prüfungsverteilung (3.22) repräsentiert, gegenübergestellt. Letztere ist die Wahrscheinlichkeit, einen Wert kleiner (oder gleich)  $x$  zu erlangen, im Falle, dass die

<sup>162</sup> Vgl. Fahrmeir, L. et al. (2004), S. 92

<sup>163</sup> Vgl. ebda., S. 93

<sup>164</sup> Vgl. Gertheiss, J.; Tutz, G. (2009), S. 448

<sup>165</sup> Vgl. ebda., S. 439

<sup>166</sup> Vgl. ebda., S. 448

<sup>167</sup> Vgl. Fahrmeir, L. et al. (2004), S. 445

<sup>168</sup> Vgl. Gertheiss, J.; Tutz, G. (2009), S. 448

<sup>169</sup> Vgl. ebda., S. 448

<sup>170</sup> Vgl. ebda., S. 450

Zufallsgröße  $X$  der zu überprüfenden Verteilung zugrunde liegt. Hierbei ist die Prüfgröße der maximale Abstand  $D$  (3.23) von theoretischer und empirischer Verteilungsfunktion.<sup>171</sup>

$$D = \sup_x |F_n(x) - F(x)| \quad (3.23)$$

Ist diese Distanz  $D$  hinreichend groß, wird die Nullhypothese  $H_0$  verworfen.<sup>172</sup> Bei der Anwendung der Anpassungstests und auf Hinblick auf die Nullhypothese (3.20) ist zu beachten, dass dieser nur überprüft, ob signifikante Abweichungen von der theoretischen Verteilungsannahme existieren. Das Ergebnis trifft die Aussage, ob etwas dagegen spricht, dass die beobachteten Werte der angenommenen Verteilung zugrunde liegen.<sup>173</sup> Auch ist der Test unabhängig von der jeweiligen zu testenden Verteilungsfamilie.<sup>174</sup>

### 3.3.3 Multivariate Deskription

Neben der Möglichkeit ein Attribut zu analysieren, werden oftmals mehrere Attribute und deren Zusammenhänge beschrieben bzw. aufgedeckt.<sup>175</sup> Daher werden in diesem Abschnitt der Korrelationskoeffizient als Zusammenhangsmaß bei metrischen Attributen und die Möglichkeit der Darstellung in einem Streudiagramm vorgestellt.

#### Streudiagramm

Das *Streudiagramm* stellt die einfachste Darstellung der gemeinsamen Messwerte  $(x_i, y_i)$ ,  $i = 1, \dots, n$  zweier stetiger Attribute dar. Dabei werden die Messwerte in ein  $x - y$ -Koordinatensystem eingetragen.<sup>176</sup>

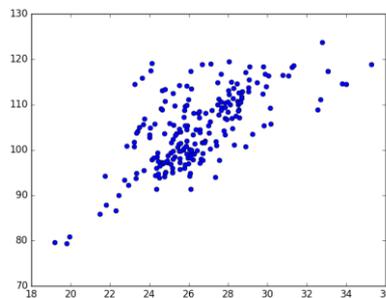


Abbildung 23: Beispiel eines Streudiagramms<sup>177</sup>

Aus einem wie in Abbildung 23 dargestellten Streudiagramm kann ein erster Eindruck über die Zusammenhänge von zwei Attributen gewonnen werden. So ist aus diesem ersichtlich, dass sobald die Werte des einen Attributes ansteigen, die des anderen ebenso ansteigen.

<sup>171</sup> Vgl. Gertheiss, J.; Tutz, G. (2009), S. 450

<sup>172</sup> Vgl. Wilcox, R. (2005), S. 1

<sup>173</sup> Vgl. Gertheiss, J.; Tutz, G. (2009), S. 449

<sup>174</sup> Vgl. ebda., S. 450

<sup>175</sup> Vgl. Fahrmeir, L. et al. (2004), S. 109

<sup>176</sup> Vgl. ebda., S. 128

<sup>177</sup> Quelle: eigene Darstellung

### Zusammenhangsmaße bei metrischen Attributen

Das Streudiagramm ist eine Hilfe zur Visualisierung der Anordnung von Beobachtungspunkten. Die aus Abbildung 23 a Tendenz zwischen den Werten des jeweiligen Attributs, lässt vermuten, dass ein Zusammenhang besteht. Diese Stärke kann durch einen *empirischen Korrelationskoeffizient* ausgedrückt werden. Hierbei handelt es sich um ein Maß für die Stärke des Zusammenhangs zwischen zwei Attributen. Es gibt verschiedene Methoden, um diesen zu berechnen.<sup>178</sup>

Einer dieser Korrelationskoeffizienten ist der *Bravais-Pearson-Korrelationskoeffizient*  $r$ . Er misst die Stärke des linearen Zusammenhangs und wird aus den Daten  $(x_i, y_i)$ ,  $i = 1, \dots, n$  nach (3.24) berechnet.<sup>179</sup>

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \quad (3.24)$$

Als Beispiele für unterschiedliche Zusammenhänge sind verschiedene Attributpaare als Punktwolken mit unterschiedlichen Korrelationskoeffizienten in Abbildung 24 visualisiert. Der Korrelationskoeffizient kann Werte im Wertebereich von  $-1 \leq r \leq 1$  annehmen, welche folgende Bedeutung haben:<sup>180</sup>

- $r > 0$ : positive Korrelation, gleichsinniger linearer Zusammenhang, Tendenz: Werte  $(x_i, y_i)$  um eine Gerade positiver Steigung liegend
- $r < 0$ : negativer Korrelation, gegensinniger linearer Zusammenhang, Tendenz: Werte  $(x_i, y_i)$  um eine Gerade negativer Steigung liegend
- $r = 0$ : keine Korrelation, unkorreliert, kein linearer Zusammenhang

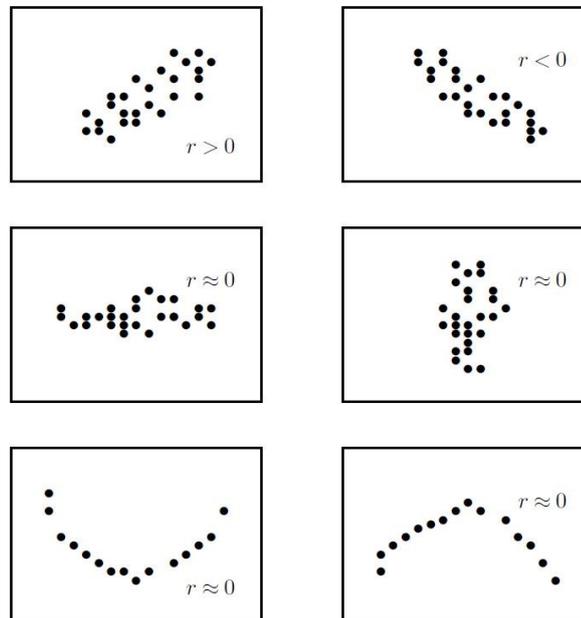


Abbildung 24: Punktekongfiguration und Korrelationskoeffizienten<sup>181</sup>

<sup>178</sup> Vgl. Fahrmeir, L. et al. (2004), S. 135 ff.

<sup>179</sup> Vgl. ebda., S. 139

<sup>180</sup> Vgl. ebda., S. 139

<sup>181</sup> Quelle: ebda., S. 138

Der Bravais-Pearson-Korrelationskoeffizient kann nur bei metrischen Attributen angewendet werden. Sofern nichtmetrische Attribute in nur zwei Ausprägungen, sogenannte dichotome oder binäre Attribute, vorliegen, können diese durch  $X, Y \in [0,1]$  ausgedrückt werden.<sup>182</sup>

Um den Korrelationskoeffizienten von metrisch und ordinalskalierten Attributen zu berechnen, kann der *Spearman's Korrelationskoeffizient* herangezogen werden.<sup>183</sup> Bei der Berechnung dieses wird nicht auf die bisher verwendeten Wertepaare  $(x_i, y_i)$  zurückgegriffen sondern auf deren Ränge. Dabei werden jeweils unabhängig voneinander jedem  $x$ -Wert aus  $x_1, \dots, x_n$  und jedem  $y$ -Wert aus  $y_1, \dots, y_n$  der Rang als Platzzahl zugeordnet, die der jeweilige Wert bei einer Ordnung nach der Größe erhält. Damit ergeben sich aus den ursprünglichen Messpaaren die neuen Rangdaten  $rg(x_i), rg(y_i), i = 1, \dots, n$ . Treten innerhalb der Werte identische Werte auf, kann eine Rangvergabe nicht eindeutig sein. Daher werden Durchschnittsränge verwendet. Das bedeutet, dass den identischen Messwerten, die als Bindungen oder Ties bezeichnet werden, der Durchschnitt der in Frage kommenden Ränge als Rang zugewiesen wird. Mit diesen Grundlagen wird Spearman's Korrelationskoeffizient nach (3.25) berechnet.<sup>184</sup>

$$r = \frac{\sum_{i=1}^n (rg(x_i) - \bar{rg}_X)(rg(y_i) - \bar{rg}_Y)}{\sqrt{\sum_{i=1}^n (rg(x_i) - \bar{rg}_X)^2 \sum_{i=1}^n (rg(y_i) - \bar{rg}_Y)^2}} \quad (3.25)$$

Der Korrelationskoeffizient kann wieder den Wertebereich  $-1 \leq r_{SP} \leq 1$  annehmen. Die Bedeutung ist im Folgenden zusammengefasst, wobei zu beachten ist, dass der Spearman's Korrelationskoeffizient im Gegensatz zum Bravais-Pearson-Korrelationskoeffizient den monotonen Zusammenhang betrachtet. Dies ist in der Abbildung 25 dargestellt, wobei hier die Extremfälle mit  $r_{SP} = 1$  (oben) und mit  $r_{SP} = -1$  (unten) zu sehen sind.<sup>185</sup>

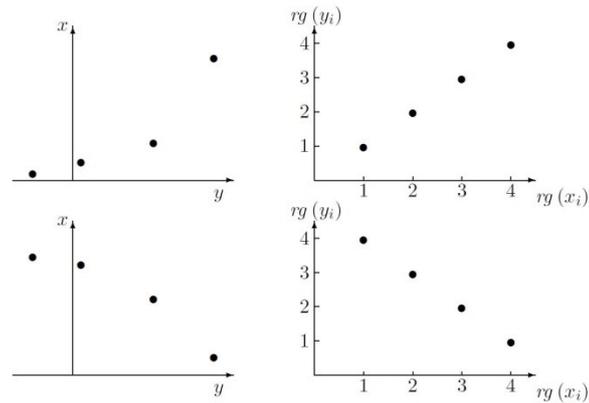
- $r_{SP} > 0$ : gleichsinniger monotoner Zusammenhang, Tendenz:  $x$  groß  $\Leftrightarrow y$  groß,  $x$  klein  $\Leftrightarrow y$  klein
- $r_{SP} < 0$ : gegensinniger monotoner Zusammenhang, Tendenz:  $x$  groß  $\Leftrightarrow y$  klein,  $x$  klein  $\Leftrightarrow y$  groß
- $r_{SP} \approx 0$ : kein monotoner Zusammenhang

<sup>182</sup> Vgl. Fahrmeir, L. et al. (2004), S. 140

<sup>183</sup> Vgl. ebda., S. 144

<sup>184</sup> Vgl. ebda., S. 142

<sup>185</sup> Vgl. ebda., S. 143 f.



**Abbildung 25: Extremfälle für Spearmans Korrelationskoeffizienten**<sup>186</sup>

Durch die beschriebenen Korrelationskoeffizienten wird die Richtung der Wirkung, sofern vorhanden, nicht ausgedrückt. Die *Korrelation* ist nur ein Maß für die Stärke des Zusammenhangs zwischen zwei Attributen  $X, Y$ .

Auf *Kausalzusammenhänge* kann daher niemals allein durch statistische Analyse bzw. hohe Werte ein Zusammenhangsmaß geschlossen werden. Zu diesem Zweck sind sachlogische Überlegungen einzubeziehen.<sup>187</sup> Dabei kann es passieren, dass wesentliche Merkmale übersehen werden. Das kann zum Auftreten von Scheinkorrelationen aber auch zu verdeckten Korrelationen führen. Eine Scheinkorrelation ist eine solche, die eine hohe Korrelation zwischen zwei Attributen angibt, aber inhaltlich nicht vertretbar ist. Bei einer verdeckten Korrelation kann durch das Auslassen eines Merkmals eine tatsächlich vorhandene Korrelation verschleiert oder hinsichtlich des Vorzeichens umgekehrt werden.<sup>188</sup>

Mit dem Wissen aus diesem Kapitel werden im nachfolgenden die Daten aufbereitet. Die vorgestellten statistischen Maße bilden die Grundlage für einige der Methoden.

<sup>186</sup> Quelle: Fahrmeir, L. et al. (2004), S. 144

<sup>187</sup> Vgl. ebda., S. 148

<sup>188</sup> Vgl. ebda., S. 149 ff.

## 3.4 Datenaufbereitung

Die Aufbereitung stellt die dritte Phase des CRISP-DM dar. Die Grundlage sind oftmals Daten, die keinerlei Qualitätskontrolle unterliegen, und die aus diesem Grund vor der Verwendung einer Analyse unterzogen und hinsichtlich des Vorhabens aufbereitet werden sollten. Dazu werden im Folgenden die Identifikation von Ausreißern, die Analyse von fehlenden Daten, die Gruppierung und die Dimensionalitätsreduktion thematisiert.

### 3.4.1 Ausreißeranalyse

Ausreißer sind eine einzigartige Kombination von Charakteristiken, die sich von anderen Beobachtungen klar abgrenzen. Die Auswirkungen, die diese auf die Analyse haben könnten, können von verschiedenen Standpunkten aus betrachtet werden. So stellt sich die Frage, ob die Grundgesamtheit der Daten als solche die Situation vollständig repräsentiert, oder ob es sich dabei schon um einen Ausreißer als solchen handelt. Ebenso ist zu hinterfragen, ob ein einzelner Wert z. B. die Lagemaße so beeinflussen kann, dass sie den Großteil der Daten nicht korrekt wiedergibt.<sup>189</sup>

Ausreißer können somit nicht als „gut“ oder „schlecht“ kategorisiert werden. Vielmehr sollten sie im Kontext der Analyse betrachtet werden und welche Art von Informationen sie repräsentieren. So können Ausreißer, obwohl sie sich vom Großteil der Grundgesamtheit unterscheiden, vorteilhaft für die Analyse sein. Sie können bestimmte Zusammenhänge repräsentieren, die sonst nicht erkannt worden wären. Im Gegensatz dazu stellen sie nicht die Grundgesamtheit dar, können aber statistische Tests verfälschen und das Erreichen des Analyseziels verhindern.<sup>190</sup>

Das Auftreten von Ausreißern kann in vier Klassen eingeteilt werden. Als Merkmal zur Einteilung dieser dient der Grund für ihr Auftreten. Hierbei kann in Verfahrensfehler, außergewöhnliche Ereignisse, außergewöhnliche Beobachtungen und Einzigartigkeit in ihrer Kombination unterschieden werden. Verfahrensfehler treten auf, wenn Daten falsch erfasst wurden. Solche Fehler sollten bereits bei der Aufbereitung der Daten erkannt und eliminiert bzw. korrigiert werden. Außergewöhnliche Ereignisse sind einzigartig und können, falls es dem Ziel der Analyse dient, in der Gesamtmenge belassen werden. Wenn dies nicht der Fall ist, sollte der Datensatz gelöscht werden. Bei außergewöhnlichen Beobachtungen liegt keine Erklärung vor, und es muss individuell entschieden werden, ob diese in der Gesamtmenge belassen werden oder nicht. Datensätze, die sich durch ihre einzigartige Kombination von Werten als Ausreißer repräsentieren, sollten in der Gesamtmenge behalten werden.<sup>191</sup>

Zum Entdecken von Ausreißern gibt es verschiedene Methoden, wie den Ausreißertest nach Grubbs<sup>192</sup> oder den David-Hartley-Pearson-Test<sup>193</sup>, die die Normalverteilung der Daten voraussetzen. Eine weitere Möglichkeit ist die graphische Auswertung mit

---

<sup>189</sup> Vgl. Hair, J. F. et al. (2014), S. 62 f.

<sup>190</sup> Vgl. ebda., S. 63

<sup>191</sup> Vgl. ebda., S. 63

<sup>192</sup> Vgl. Grubbs, F. E. (1950), S. 27

<sup>193</sup> Vgl. David, H. A. et al. (1954), S. 482

Boxplots.<sup>194</sup> Wie bereits in Kapitel 3.3.2 besprochen, kann die Länge der Whisker selbst gewählt werden. Tukey (1977)<sup>195</sup> schlägt die Whiskerlänge von  $Q_1 - 1,5 \times IQR$ ,  $Q_3 + 1,5 \times IQR$  vor, wobei diese dann bei dem letzten innerhalb dieser Grenze gelegenen Wert enden. Dies kann zu unterschiedlich langen Whiskers führen. Alle außerhalb dieser gelegenen Werte werden als „extreme“ Ausreißer identifiziert werden. Werte die außerhalb der Quantile - der Box - liegen, werden als „milde“ Ausreißer bezeichnet.<sup>196</sup>

Ausreißer können nur dann erkannt werden, wenn genügend Daten vorhanden sind. Das Fehlen von Daten und die Möglichkeiten zur Ergänzung werden im nachfolgenden Kapitel betrachtet.

### 3.4.2 Analyse der Fehlenden Daten

Bei der Exploration der Daten mit denen in Kapitel 3.3 vorgestellten Methoden, waren fehlende Daten kein Grund für die Nichtanwendung. Oftmals ist es aber erforderlich, dass die Daten vollständig ohne Lücken vorliegen. Würden Datensätze aufgrund von fehlenden Daten ausgeschlossen werden, könnte dies die Anzahl der Datensätze so verringern, dass keine Analyse mit diesen Daten als Grundlage durchgeführt werden könnte. Zusätzlich besteht das Problem der Verzerrung der Daten bzw. Maßzahlen der Attribute, sollten die Daten absichtlich oder unabsichtlich verfälscht worden sein. Es sollten daher die Fragen beantwortet werden, ob Daten zufällig oder nach einem bestimmten Muster fehlen.<sup>197</sup>

Fehlende Daten können bei Analysen und vor allem bei multivariaten Analysen eine signifikante Auswirkung auf die Ergebnisse haben. Die zu betrachtenden Beziehungen zwischen den Daten werden komplexer, und die Wahrscheinlichkeit fehlende Daten zu übersehen, wird größer. Daher ist es von Vorteil fehlende Daten vorab zu finden und zu überprüfen. Zu diesem Zweck hat Hair et al. (2014)<sup>198</sup> einen vier Schritte umfassenden Prozess zum Identifizieren und Bereinigen von fehlenden Daten entworfen und wird im Folgenden vorgestellt. Der englische Originaltitel „A Four-Steps Process for Identifying Missing Data and Applying Remedies“ lässt Spielraum bei der Übersetzung des Wort „data“ zu. Im Deutschen wurde daher das Wort „Daten“ gewählt, wobei im späteren praktischen Teil der Prozess für fehlende Werte angewendet wird.

Die folgenden Ausführungen dieses Prozesses wurden, wenn nicht anders angegeben, von Hair et al. (2014)<sup>199</sup> übernommen und in Abbildung 26 zusammengefasst.

---

<sup>194</sup> Vgl. Fahrmeir, L. et al. (2004), S. 67

<sup>195</sup> Vgl. Tukey, J. W. (1977), S. 44 ff.

<sup>196</sup> Vgl. ebda., S. 44 ff.

<sup>197</sup> Vgl. Hair, J. F. et al. (2014), S. 40

<sup>198</sup> Vgl. ebda., S. 42

<sup>199</sup> Vgl. ebda., S. 42 ff.

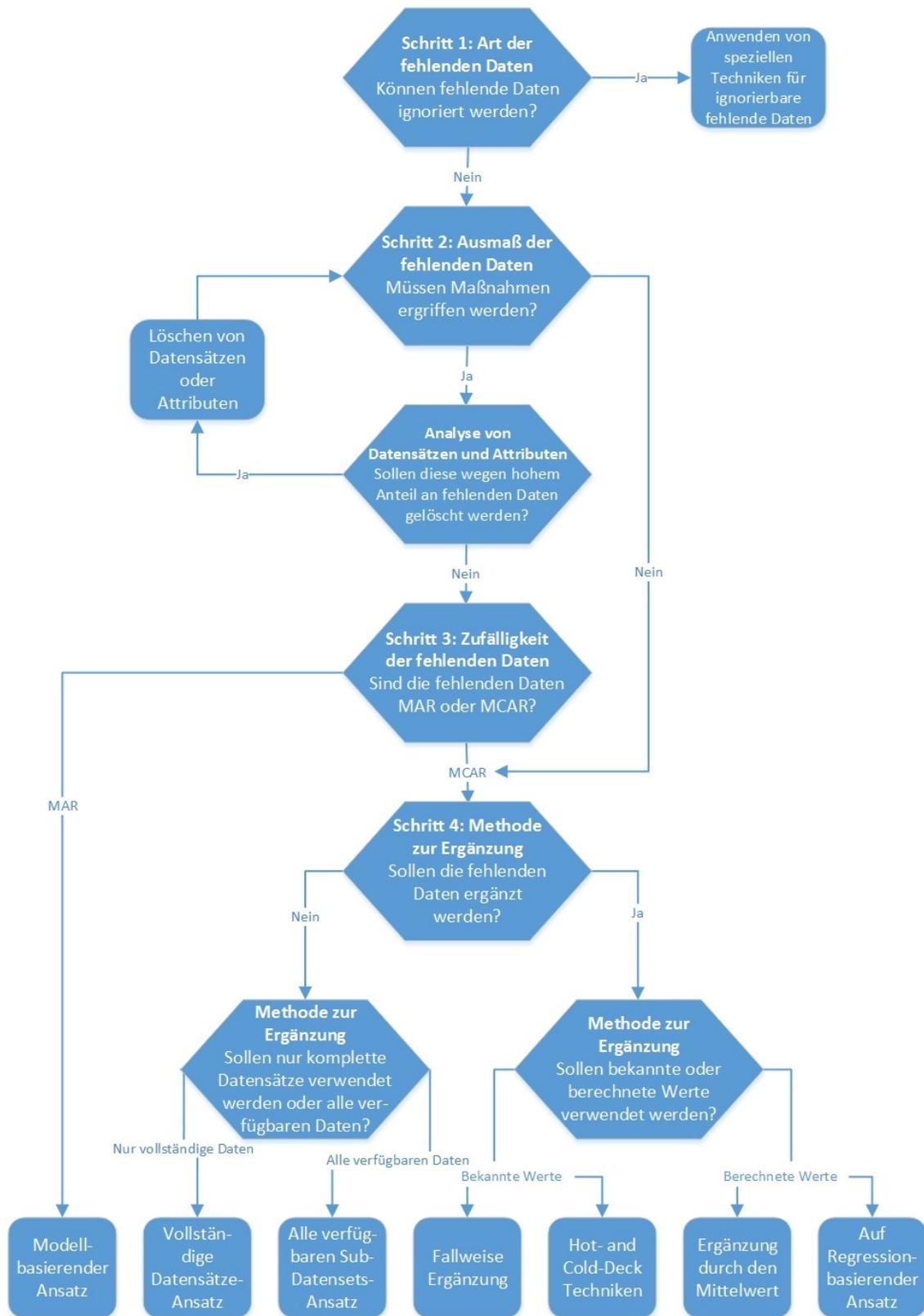


Abbildung 26: Prozess zum Identifizieren und Bereinigen fehlender Daten<sup>200</sup>

<sup>200</sup> Quelle: in Anlehnung an Hair, J. F. et al. (2014), S. 43

### Schritt 1: Feststellen der Art der fehlenden Daten

In diesem ersten Schritt soll festgestellt werden, ob Daten, die fehlen, ignorierbar sind oder nicht. Ignorierbare fehlende Daten können folgende sein:

- Daten, die nicht in die Stichprobe einfließen
- Daten, die durch die Struktur der Erhebung entstanden sind
- Zensierte Daten: Diese können aufgrund ihrer Stellung im Erhebungsprozess nicht vollständig sein.

Nicht ignorierbare fehlende Daten können als solche durch mehrere Gründe oder Situationen bezeichnet werden. Im Allgemeinen lassen sie sich in die Kategorien „bekannte“ und „unbekannte“ Ursachen, die während des Datenerhebungsprozess auftreten, einteilen.

Werden die fehlenden Daten als nicht ignorierbar eingeschätzt, wird zu Schritt 2 übergegangen. Ansonsten ist die Analyse beendet.

### Schritt 2: Feststellen des Ausmaßes der fehlenden Daten

Die Hauptaufgabe in diesem Schritt ist, festzustellen, ob die Menge der fehlenden Daten gering genug ist, um das Ergebnis nicht zu beeinflussen. Wenn der Anteil gering ist, können die fehlenden Daten durch Ergänzen nach bestimmten Regeln eliminiert werden. Hierbei wird die Frage aufgeworfen, wie das Ausmaß von „gering“ definiert ist. Diese Grenze wird individuell festgelegt. Durch Löschen von Datensätzen oder Attributen kann die Menge von fehlenden Daten verkleinert werden und somit Daten hervorbringen, bei denen das Auffüllen der fehlenden Daten keine Verzerrung der Resultate bewirkt. Hierzu gibt es keine definierten Regeln.

Hair et al. (2014)<sup>201</sup> haben folgende Faustregeln zusammengefasst:

- Grenzen für fehlende Daten
  - Der Anteil fehlender Daten von unter 10 % können generell ignoriert werden. Es sei denn Daten fehlen nach einem bestimmten Muster.<sup>202</sup>
  - Die Anzahl der Datensätze ohne fehlende Daten ist ausreichend für die gewählte Analyseverfahren, wenn die fehlenden Daten nicht ergänzt wurden.
- Löschen von Attributen und Datensätzen
  - Hair et al. (2014)<sup>203</sup> beziehen sich auf Hertel (1976)<sup>204</sup> mit der Aussage, dass Attribute mit einem Anteil von weniger als 15 % fehlender Daten Kandidaten zum Löschen sind. Bei Anteilen von 20 - 30 % können die Daten oftmals ergänzt werden. Hertel (1976)<sup>205</sup> hingegen verweist auf Sharp und Feldt (1959)<sup>206</sup>, dass bei einem Komplettheitsgrad von mindestens 85 % die Daten gelöscht, aber auch andere Möglichkeiten in Betracht gezogen werden sollten. Bei einem geringeren Komplettheitsgrad sollte das Löschen in Betracht gezogen werden.

<sup>201</sup> Vgl. Hair, J. F. et al. (2014), S. 45

<sup>202</sup> Vgl. Malhotra, N. K. (1987), S. 75

<sup>203</sup> Vgl. Hair, J. F. et al. (2014), S. 45

<sup>204</sup> Vgl. Hertel, B. R. (1976), S. 461

<sup>205</sup> Vgl. ebda., S. 461

<sup>206</sup> Vgl. Sharp, H.; Feldt, A. (1959), S. 652

- Es sollte sichergestellt werden, ob die erreichbare Verbesserung des Anteils der fehlenden Daten an der Gesamtmenge durch das Löschen eines Datensatzes oder Attribut groß genug ist.
- Datensätze mit fehlenden Daten in abhängigen Attributen werden üblicherweise gelöscht, um künstliche Verstärkungen in den Beziehungen zwischen unabhängigen Attributen zu vermeiden.
- Wenn ein Attribut gelöscht wird, sollte sichergestellt sein, dass alternative Attribute, die stark miteinander korrelieren und somit die gelöschten Attribute repräsentieren können, zurückbleiben.
- Es sollte in Betracht gezogen werden, die Analyse mit und ohne die fehlenden Daten durchzuführen, um den Unterschied zu erkennen.

Datensätze oder Attribute mit mehr als 50 % fehlender Daten sollten gelöscht werden. Sobald der Prozentsatz darunter ist, muss individuell nach Ziel der Analyse und dem Motto „trial and error“ entschieden werden.

Ist das Ausmaß der fehlenden Daten akzeptabel niedrig und liegt fehlenden Daten kein Muster zugrunde, kann direkt zum Schritt 4 übergegangen werden. Ansonsten wird mit Schritt 3 fortgefahren.

### **Schritt 3: Feststellen der Zufälligkeit der fehlenden Daten**

Sollte beschlossen werden, dass diese fehlenden Daten ergänzt werden sollen, muss der Grad Zufälligkeit der fehlenden Daten bestimmt werden. Hierfür geht man von Daten, die mindestens zwei Attribute repräsentieren, aus. Enthält nun eines fehlende Daten, andere nicht, so kann das Attribut bzw. dessen fehlende Daten als „Missing At Random“ (MAR) oder als „Missing Completely At Random“ (MCAR) definiert werden.<sup>207</sup> *MAR* bedeutet in diesem Fall, dass die fehlenden Werte aus dem einen Attribut mit dem anderen Attribut in Zusammenhang stehen und nicht mit den anderen Werten des eigenen Attributs. *MCAR* hingegen sind solche fehlenden Daten, deren Datensätze sich nicht von anderen Datensätzen unterscheiden lassen. Attribute mit fehlenden Daten können durch Begutachtung bzw. Kenntnisse oder durch vorgefertigte Tests der jeweiligen Kategorie zugeteilt werden.

### **Schritt 4: Auswahl der Methode zur Ergänzung**

Das Ergänzen als solches wird als Prozess zum Schätzen von Daten auf Basis der vorhandenen Daten angesehen. Als Ziel gilt es, Beziehungen zwischen den Daten zu erkennen, die zum Ergänzen verwendet werden können. In diesem Schritt wird die Methode ausgewählt, mit der die fehlenden Daten ergänzt werden. Hierbei wird die in Schritt 3 getroffene Entscheidung, ob es sich um MAR oder MCAR handelt, als Basis verwendet. Unabhängig von dem Grad der Zufälligkeit gibt es verschiedene Methoden zum Ergänzen der Daten.<sup>208,209,210</sup> Wie in Abbildung 26 zuerkennen ist, wird als erstes die Entscheidung getroffen, ob die Daten ergänzt werden sollen. Ist dies nicht der Fall können entweder nur solche Daten verwendet werden, die ohne fehlende Daten sind, oder es wird anstelle der fehlenden Daten ein Lage-, Streuungs- oder

<sup>207</sup> Vgl. Little, R. J. A.; Rubin, D. B. (2002), S. 12

<sup>208</sup> Vgl. ebda., S. 19

<sup>209</sup> Vgl. Heitjan, D. F. (1997), S. 548 ff.

<sup>210</sup> Vgl. Roth, P. L. (1994), S. 539 ff.

Zusammenhangsmaß eingesetzt. Sollte die Daten ersetzt werden, gibt es die Möglichkeit bekannte oder berechnete Werte einzusetzen.

Abschließend ist festzuhalten, dass nur die Ergänzung von metrischen Daten nach diesen Regeln sinnvoll ist. Nicht-metrische Daten sollten anders behandelt werden. Zusätzlich ist die Frage, ob nicht-metrische Attribute für die gewählte Analyseverfahren überhaupt zulässig sind.

Nach der Analyse der Ausreißer innerhalb der Daten werden im nächsten Kapitel das Zielattribut und seine Eigenschaften betrachtet.

### 3.4.3 Gruppierung

Die in Kapitel 3.5 verwendeten Methoden sollen klassieren. Daher ist es sinnvoll, das Zielattribut zu gruppieren. Dies bedeutet, die vorhandenen Werte in Klassen mit definierten Grenzen zuzuordnen.<sup>211</sup> Bei der Gruppierung werden jeweils mehrere Ausprägungen des ursprünglichen Attributes zusammengefasst und in eine Klasse bzw. Ausprägung des transformierten Merkmals überführt. Dieser Vorgang geht mit einer Skalentransformation einher. Bei der Klassenbildung wird das Gesamtintervall in Teilintervalle zerlegt, die eindeutig durch ihre Klassenmitte und ihre obere und untere Klassengrenze definiert werden. Bei der Festlegung dieser Klassen gibt es zwei Ansätze: Durch Definition der Klassenbreite, wie schon bei den Histogrammen in Kapitel 3.3.2 erläutert, oder durch die Bestimmung der Klassenanzahl. Das Festlegen der Klassenbreite kann zu Verzerrungen führen, da homogene Intervalle innerhalb der Beobachtungswerte zerschnitten werden könnten. Bei der Bestimmung der Klassenanzahl ist darauf zu achten, dass bei zu vielen Klassen nicht die gewünschte Informationsverdichtung erreicht wird, wohingegen zu wenige die ursprüngliche Verteilung der Attributausprägungen verschleiern.<sup>212</sup>

Zum Festlegen der Klassenanzahl gibt es verschiedene Ansätze. Wird eine Einteilung nach Anzahl der Klassen getroffen, so gibt die DIN 55302<sup>213</sup> an, dass bei 100 Beobachtungswerten mindestens 10 Klassen, bei 1000 mindestens 13 und bei 10.000 mindestens 16 Klassen gebildet werden sollten. Daneben sei zu beachten, dass die Anzahl der Klassen bei weniger als 100 Beobachtungswerten nicht größer sein sollte als die Quadratwurzel aus der Zahl der Beobachtungswerte, siehe (3.1). Eine weitere Methode ist die von Sturges aufgestellte Formel (3.26), wobei  $k$  die Anzahl der Klassen und  $n$  die Anzahl der Beobachtungswerte darstellt. Bei der Bildung der Klassen sollte auf jeden Fall darauf geachtet werden, die Daten nicht zu sehr zu verändern, da so Informationen verloren gehen könnten.<sup>214</sup>

$$k = 1 + 3.3 \log n \quad (3.26)$$

Das Gruppieren dient der Verdichtung des Zielattributs. Die anderen Attribute können ebenso zusammengefasst werden. Diese Möglichkeit wird im nachfolgenden Kapitel behandelt.

<sup>211</sup> Vgl. Fahrmeir, L. et al. (2004), S. 17

<sup>212</sup> Vgl. Degen, H.; Lorscheid, P. (2012), S. 19 ff.

<sup>213</sup> Vgl. DIN 55302 (Blatt 1) zitiert nach Degen, H.; Lorscheid, P. (2012), S. 21

<sup>214</sup> Vgl. Degen, H.; Lorscheid, P. (2012), S. 21

### 3.4.4 Dimensionalitätsreduktion

In Hinblick auf das Ziel und die Methode der Datenanalyse kann es sein, dass das Verhältnis von Attributen zu Datensätzen nicht optimal ist oder dass die Anzahl von Attributen einen höheren Rechenaufwand nach sich ziehen würde. Um ein optimales Verhältnis zu erhalten, gibt es mehrere Möglichkeiten. Eine davon ist das unwiderrufliche Löschen von Attributen, das zu einem Informationsverlust führt. Die andere ist die Transformation der ursprünglichen Attribute in eine geringere Anzahl neuer Attribute, die annähernd gleiche Informationen enthalten, aber ein besseres Verhältnis von Attributen zu Datensätzen bzw. eine einfachere Struktur liefert.

Es gibt verschiedene Methoden zur Dimensionalitätsreduzierung. Dazu gehören die Erstellung eines Sub-Datensets, die Faktorenanalyse oder die multidimensionale Skalierung.<sup>215</sup> Eine der bekanntesten Methoden ist die *Hauptkomponentenanalyse*, die im Folgenden als *PCA (Principal Component Analysis)* bezeichnet wird. Die zentrale Idee der PCA ist das Verringern der Attributanzahl, die in einem Datenset enthalten sind. Zu diesem Zweck werden aus ursprünglichen Attributen neue, die sogenannten *Hauptkomponenten (PCs – Principal Components)*, die untereinander nicht korrelieren, gebildet. Dabei soll der Informationsverlust möglichst gering gehalten werden, um die Varianz zu maximieren. Die PCs sind so geordnet, dass die ersten PCs die meisten Variation enthalten, die in allen ursprünglichen Attributen vorhanden ist.<sup>216</sup>

Es gab eine Vielzahl von Mathematikern, die sich seit Ende des 19. Jahrhunderts<sup>217,218</sup> mit der der PCA zugrunde liegenden Singulärwertzerlegung (*SVD – Singular Value Decomposition*) in unterschiedlichsten Kontexten beschäftigt haben. Die PCA als solche wurde als erstes von Pearson (1901)<sup>219</sup> und Hotelling (1933)<sup>220</sup> beschrieben. Dennoch haben sie unterschiedliche Ansätze. Pearsons Ansatz beruht auf der geometrischen Betrachtung der Lage der Datensätze in einem  $p$ -dimensionalen Raum und wie diese durch Linien und Ebenen dargestellt werden können.<sup>221</sup> Hotellings Ansatz begründet sich auf der Faktorenanalyse, einem Verfahren der multivariaten Statistik, die die PCA verwendet.<sup>222</sup> Sein Ergebnis unterscheidet sich stark von diesem heute noch genutzten Verfahrens. Seine Idee war, dass es ein kleineres Set von Attributen geben muss, das die Werte der ursprünglichen  $p$  Attribute repräsentiert.<sup>223</sup> Hotellings Methodik ähnelt der in diesem Kapitel vorgestellten allgemeinen Methodik für die PCA. Seit den Anfängen der Forschung im Bereich der PCA, die immer noch anhält, bietet diese ein breites Anwendungsspektrum in den unterschiedlichsten Branchen.<sup>224</sup>

<sup>215</sup> Vgl. Alpaydin, E. (2010), S. 110 ff.

<sup>216</sup> Vgl. Jolliffe, I. T. (2002), S. 1

<sup>217</sup> Vgl. Beltrami, E. (1873)

<sup>218</sup> Vgl. Fisher, R. A.; Mackenzie, W. A. (1923)

<sup>219</sup> Vgl. Pearson, K. (1901)

<sup>220</sup> Vgl. Hotelling, H. (1933)

<sup>221</sup> Vgl. Pearson, K. (1901), S. 559

<sup>222</sup> Vgl. Backhaus, K. et al. (2011), S. 356

<sup>223</sup> Vgl. Hotelling, H. (1933), S. 417

<sup>224</sup> Vgl. Jolliffe, I. T. (2002), S. 7 ff.

Für die Anwendung der PCA wird für das Verhältnis von Attributen zu Beobachtungen eine absolute Untergrenze von eins zu fünf gefordert. Ein wünschenswertes Verhältnis wäre ein zu zehn.<sup>225</sup>

### Methodik der PCA

Die folgende Beschreibung der Methodik der PCA basiert auf der Eigenvektorzerlegung und zur Gänze auf Shlens (2014)<sup>226</sup>. Es gibt andere mathematische Ansätze, die z. B. auf der SVD beruhen. Der Grundgedanke ist, die beste Basis zu finden, um die Attribute in geringerem Ausmaß ausdrücken zu können. Dabei sollen Rauschen herausgefiltert und neue Strukturen entdeckt werden. Hierzu wird ein neues Achsensystem eingeführt, und überprüft, ob eine neue lineare Kombination der Attribute das ursprüngliche Datenset repräsentieren kann. Dazu wird angenommen, dass die Daten linear vorliegen, eine hohe Varianz in den Daten neuen Strukturen verbirgt und die daraus errechneten Hauptkomponenten orthogonal zueinander sind.

Dieser Grundgedanke kann mathematisch formuliert werden: Das Ziel ist es eine orthogonale Matrix  $\mathbf{P}$  in  $\mathbf{Y} = \mathbf{P}\mathbf{X}$ , sodass  $\mathbf{C}_Y = \frac{1}{n}\mathbf{Y}\mathbf{Y}^T$  eine Diagonalmatrix ist. Die Zeilen von  $\mathbf{P}$  sind die Hauptkomponenten von  $\mathbf{X}$ . Die  $m \times n$ -Matrix  $\mathbf{X}$  stellt in diesem Fall den Datensatz dar, wobei  $m$  die Anzahl der Attribute und  $n$  die Anzahl der Datensätze ist.  $\mathbf{C}_Y$  ist in diesem Fall die Kovarianz (3.27), die eine symmetrische  $m \times m$ -Matrix ist. Ihre Diagonalelemente sind die Varianz der einzelnen Attribute und die Elemente abseits der Diagonalen sind die Kovarianzen zwischen den einzelnen Attributen. Sie spiegeln das Rauschen und die Redundanz der Daten wieder. Große Werte in den Elementen abseits der Diagonalen weisen auf große Redundanz hin.

Das daraus abzuleitende Ziel ist, die Varianz zu maximieren und die Redundanz zu minimieren. Dies geschieht, indem die Zielgleichung nach der unbekanntem Kovarianz umgeformt wird:

$$\begin{aligned}\mathbf{C}_Y &= \frac{1}{n}\mathbf{Y}\mathbf{Y}^T \\ &= \frac{1}{n}(\mathbf{P}\mathbf{X})(\mathbf{P}\mathbf{X})^T \\ &= \frac{1}{n}\mathbf{P}\mathbf{X}\mathbf{X}^T\mathbf{P}^T \\ &= \mathbf{P}\left(\frac{1}{n}\mathbf{X}\mathbf{X}^T\right)\mathbf{P}^T \\ \mathbf{C}_Y &= \mathbf{P}\mathbf{C}_X\mathbf{P}^T\end{aligned}\tag{3.27}$$

Jede symmetrische Matrix  $\mathbf{A}$  wird durch eine orthogonale Matrix ihrer Eigenvektoren diagonalisiert. Für diese gilt:  $\mathbf{A} = \mathbf{E}\mathbf{D}\mathbf{E}^T$ , wobei  $\mathbf{D}$  die Diagonalmatrix ist. Die Matrix  $\mathbf{E}$  besteht aus Eigenvektoren der Matrix  $\mathbf{A}$ .

<sup>225</sup> Vgl. Rinkenburger, R. (2009), S. 463

<sup>226</sup> Vgl. Shlens, J. (2014), S. 2 ff.

Die Matrix  $\mathbf{P}$  wird so gewählt, dass jede Reihe  $\mathbf{p}_i$  ein Eigenvektor von  $\frac{1}{n}\mathbf{X}\mathbf{X}^T$  ist. Durch diese Wahl gilt  $\mathbf{P} \equiv \mathbf{E}^T$ . Dies führt zu (3.28):

$$\begin{aligned} \mathbf{C}_Y &= \mathbf{P}\mathbf{C}_X\mathbf{P}^T \\ &= \mathbf{P}(\mathbf{E}^T\mathbf{D}\mathbf{E})\mathbf{P}^T \\ &= \mathbf{P}(\mathbf{P}^T\mathbf{D}\mathbf{P})\mathbf{P}^T \\ &= (\mathbf{P}\mathbf{P}^T)\mathbf{D}(\mathbf{P}\mathbf{P}^T) \\ &= (\mathbf{P}\mathbf{P}^{-1})\mathbf{D}(\mathbf{P}\mathbf{P}^{-1}) \\ \mathbf{C}_Y &= \mathbf{D} \end{aligned} \tag{3.28}$$

Die Ergebnisse der PCA können in den Matrizen  $\mathbf{P}$  und  $\mathbf{C}_Y$  zusammengefasst werden:

- Die Hauptkomponenten von  $\mathbf{X}$  sind die Eigenvektoren von  $\mathbf{C}_X = \frac{1}{n}\mathbf{X}\mathbf{X}^T$ .
- Das  $i$ -te-Diagonalelement von  $\mathbf{C}_Y$  ist die Varianz von  $\mathbf{X}$  entlang  $\mathbf{p}_i$ .

Die tatsächliche Anzahl, wie viele Hauptkomponenten dieser Lösung verwendet werden, kann über Kriterien und Tests angegeben werden, wie z. B. den Anteil der erklärten Varianz (Faustregel für Naturwissenschaften 0,95), den Scree-Test oder das Kaiser-Kriterium.<sup>227</sup>

### Eignung zur der Daten PCA

Zur Überprüfung der Eignung der vorliegenden Daten für die PCA können verschiedene Kriterien herangezogen werden. In diesem Fall wird die Korrelationsmatrix, die Inverse der Korrelationsmatrix und der Kaiser-Meyer-Olkin-Test vorgestellt.

Die *Korrelationsmatrix*  $R$ , die sich aus den in Kapitel 3.3.3 berechneten Korrelationskoeffizienten zusammensetzt, sollte wenige kleine Werte enthalten. Dies würde auf eine heterogene Datenstruktur hindeuten; in diesem Fall wäre die Anwendung in Frage zu stellen.<sup>228</sup>

Die *Inverse der Korrelationsmatrix*  $R^{-1}$  sollte eine Diagonalmatrix darstellen. Das bedeutet, dass die Elemente außerhalb der Diagonalen möglichst nahe bei Null liegen. In wie weit diese Werte allerdings von Null abweichen dürfen und wie häufig dies vorkommen darf, ist nicht allgemein geklärt und liegt im eigenen Ermessen.<sup>229</sup>

Guttman geht davon aus, dass sich die Varianz eines Attributes in zwei Teile zerlegen lässt: das *Image* und das *Anti-Image*. Das *Image* stellt jenen Teil der Varianz dar, der durch die restlichen Attribute beschrieben werden kann. Dagegen ist das *Anti-Image* jener Teil, der von den verbleibenden Attributen unabhängig ist.<sup>230</sup> Beim *Kaiser-Meyer-Olkin-Test* wird die *Anti-Image-Korrelationsmatrix* berechnet. Anhand dieser kann dieser über eine Eignung sämtlicher Attribute und einzelner Attribute für die Faktorenanalyse entschieden werden. Die Gesamtheit der Attribute wird mit dem *MSA-Kriterium (measure of sampling adequacy)*, (3.29) bewertet, das auch *KMO-Kriterium*

<sup>227</sup> Vgl. Rinckenburger, R. (2009), S. 469

<sup>228</sup> Vgl. Backhaus, K. et al. (2011), S. 339

<sup>229</sup> Vgl. ebda., S. 340

<sup>230</sup> Vgl. Guttman, L. (1953), S. 277

(Kaiser-Meyer-Olkin) genannt wird und auch für jedes einzelne Attribut (3.30) berechnet werden kann.<sup>231</sup>

$$MSA = \frac{\sum_i \sum_{j \neq i} r_{ij}^2}{\sum_i \sum_{j \neq i} r_{ij}^2 + \sum_i \sum_{j \neq i} q_{ij}^2} \quad (3.29)$$

$$MSA(i) = \frac{\sum_{j \neq i} r_{ij}^2}{\sum_{j \neq i} r_{ij}^2 + \sum_{j \neq i} q_{ij}^2} \quad (3.30)$$

Bei der Berechnung ist darauf zu achten, dass die partielle Korrelationsmatrix  $Q$  sich mit  $Q = SR^{-1}S$  berechnen lässt, wobei  $S = (diag R^{-1})^{-1/2}$  ist. Die so berechneten Werte können die Werte im Intervall  $[0; 1]$  annehmen. Wünschenswert wäre ein Wert von mindestens 0,8.<sup>232</sup> Und ein  $MSA < 0,5$  gibt die Untauglichkeit der Daten bzw. des jeweiligen Attributs an.<sup>233</sup> Zudem kann die Eignung der Gesamtheit bzw. einzelner Attribute, wie in Tabelle 6 angegeben, eingestuft werden.<sup>234</sup>

**Tabelle 6: Einstufung des MSA-Kriterium hinsichtlich der Eignung für die PCA**

$MSA \geq 0,9$	marvelous	„erstaunlich“
$MSA \geq 0,8$	meritorious	„verdienstvoll“
$MSA \geq 0,7$	middling	„ziemlich gut“
$MSA \geq 0,6$	mediocre	„mittelmäßig“
$MSA \geq 0,5$	miserable	„kläglich“
$MSA < 0,5$	unacceptable	„untragbar“

Der KMO-Test gilt als einer der besten Tests zur Bestimmung der Eignung von Daten für die PCA.<sup>235</sup>

Mit der Dimensionalitätsreduktion ist die Vorbereitung der Daten abgeschlossen, und diese können, wie in Abbildung 2 dargestellt, als Basis für die Modellerstellung verwendet werden. Welche Methoden innerhalb des maschinellen Lernens zur Verfügung stehen, wird im nachfolgenden Kapitel vorgestellt.

### 3.5 Modellierung

Dieses Kapitel soll die in Kapitel 4 angewendeten Methoden des maschinellen Lernens und deren Theorie vorstellen, die im Rahmen der CRISP-DM Phase „Modellierung“ zum Einsatz kommen. Im Rahmen dieser Arbeit liegt der Fokus auf dem Klassieren. Wie bereits aus Tabelle 3 für diese Methode zu entnehmen war, ist das Ziel, einen Datensatz zu klassieren bzw. die Klasse eines Attributes zu ermitteln.

<sup>231</sup> Vgl. Kaiser, H. F. (1970), S. 405

<sup>232</sup> Vgl. ebda., S. 405

<sup>233</sup> Vgl. Edward E. Cureton, R. B. D. (1993), S. 391

<sup>234</sup> Vgl. Backhaus, K. et al. (2011), S. 343

<sup>235</sup> Vgl. Stewart, D. W. (1981), S. 57

### 3.5.1 Einführung

Um Probleme mit dem Computer zu lösen, werden Algorithmen benötigt. Da die Arbeitsweise nicht mit der eines Menschen vergleichbar ist, ist es ihnen oftmals nicht möglich, den ganzen Prozess zu erfassen. Aber es können Muster und Regelmäßigkeiten identifiziert werden. Das ist der Ansatzpunkt des maschinellen Lernens. Dadurch können Erkenntnisse aus den Daten die Grundlagen für Vorhersagen über die nähere Zukunft sein. Oder sie können die Daten beschreiben und so Wissen über den Prozess generieren.<sup>236</sup>

Die Ausgangsbasis stellen Daten dar. Diese, wie in den vorherigen Kapitel beschriebenen, aufbereiteten Daten werden in zwei Gruppen geteilt: die Trainings- und die Testdaten. *Trainingsdaten* sind jene, von denen gelernt und deren Wissen extrahiert wird. Sie sind die Grundlage für das Modell bzw. dessen Parametrierung. Mit den *Testdaten* wird die Genauigkeit bzw. die Leistung des gelernten Modells überprüft, indem der klassierte Wert mit dem bereits eingetretenen Wert verglichen wird.<sup>237</sup>

Das Lernen von den Trainingsdaten ist mit den im folgenden Kapitel 0 beschriebenen Methoden ein *überwachtes Lernen*. Es wird mit Datensätzen trainiert, die den Zielwert enthalten. Dies beeinflusst den Lernprozess so, dass, wenn etwas falsch gelernt wurde, darauf aufmerksam gemacht wird, und dies bei zukünftigen Klassierungen berücksichtigt wird.<sup>238</sup> Das Problem bei dieser Art des Lernens ist, die Generalisierung des Modells. Sind die Trainingsdaten keine gute Repräsentation der Gesamtmenge der Daten, kann es passieren, dass das Modell nicht generalisierbar ist.<sup>239</sup>

Eine weitere Gefahr ist das die *Überanpassung* (engl. *overfitting*). Dies kann auftreten, wenn das Verhältnis von Attributen zu Datensätzen zu klein ist.<sup>240</sup> Um ein optimales Verhältnis erreichen, kann auf die in Kapitel 3.4.4 Dimensionalitätsreduktion zurückgegriffen werden. Des Weiteren kann zur zusätzlichen Verbesserung des Modells die *Kreuzvalidierung* angewendet werden. Bei der auch *K-fache Kreuzvalidierung* genannten Methode, werden die  $X$  Datensätze in  $K$  gleich große Teile,  $X_i, i = 1, \dots, K$  geteilt. Dabei wird ein Teil als Testdaten zurückgehalten und die restlichen  $K - 1$  Teile sind die Trainingsdaten. Dies wird  $K$ -mal wiederholt und so wird immer ein Teil ausgelassen.<sup>241</sup>

Die später im praktischen Fallbeispiel im Kapitel 4 verwendeten Methoden und deren Ansätze werden im folgenden Kapitel vorgestellt.

---

<sup>236</sup> Vgl. Alpaydin, E. (2010), S. 1 ff.

<sup>237</sup> Vgl. Ertel, W. (2013), S. 181

<sup>238</sup> Vgl. Backhaus, K. et al. (2015), S. 299

<sup>239</sup> Vgl. Alpaydin, E. (2010), S. 24

<sup>240</sup> Vgl. ebda., S. 154

<sup>241</sup> Vgl. ebda., S. 486 f.

### 3.5.2 Methoden

Alle Methoden haben das Ziel, eine Klassifikation durchzuführen, und somit eine Prognose über den Zielwert durch die Zuordnung eines Datensatzes zu einer Klasse abzugeben. Die Klassifikation beruht auf den Zusammenhängen zwischen den einzelnen Attributen und der Bedeutung ihrer Kombination in Bezug auf den Zielwert.<sup>242</sup>

**Tabelle 7: Anwendbarkeit auf Attribute unterschiedlicher Skalenniveaus<sup>243</sup>**

Methode	Skala		
	metrisch	nicht-metrisch	gemischt
Entscheidungsbaum	Ja	Ja	Ja
Diskriminanzanalyse	Ja	Nein	Nein
SVM	Ja	Ja	Ja
KNN	nur euklidische Distanz	Nur Hamming-Distanz	Nein
Ensemblemethoden	Ja	Ja, außer Subspace Diskriminanzanalyse	Ja, außer Subspace Methoden

Nicht jede der nachfolgend vorgestellten Methoden kann mit Attributen unterschiedlichsten Skalenniveaus (siehe Kapitel 3.3.1) arbeiten. Tabelle 7 enthält daher eine Zusammenstellung über die Anwendbarkeit von Methoden auf Attribute in metrischen, nicht-metrischen und gemischten Skalenniveaus.<sup>244</sup>

#### Entscheidungsbaum

*Entscheidungsbäume* (engl. *decision trees*) sind äußerst einfache aber effiziente Verfahren, um aus Daten Wissen zu generieren. Der Vorteil ist, dass diese Bäume graphisch dargestellt werden können. Dies erleichtert das Verständnis, die Interpretierbarkeit und Kontrolle.<sup>245</sup> Ein Entscheidungsbaum besteht aus Knoten und Blättern, die die Zielwerte darstellen. Jeder Entscheidungsknoten  $m$  hat eine Testfunktion  $f_m(x)$ , deren jeweiliges Ergebnis in einen der zwei abzweigenden Äste mündet. Beispielhaft ist dies in Abbildung 27 dargestellt, wobei die Knoten mit den Funktionen beschriftet wurden und die Blätter durch die Rechtecke mit den jeweiligen Zielwerten  $C_1, C_2$  dargestellt werden.<sup>246</sup>

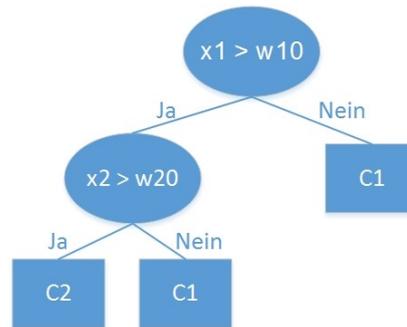
<sup>242</sup> Vgl. Alpaydin, E. (2010), S. 5

<sup>243</sup> Quelle: MathWorks, <https://de.mathworks.com/help/stats/choose-a-classifier.html> (Zugriff: 21.02.2017) (leicht modifiziert)

<sup>244</sup> Vgl. ebda.

<sup>245</sup> Vgl. Ertel, W. (2013), S. 202 f.

<sup>246</sup> Vgl. Alpaydin, E. (2010), S. 185



**Abbildung 27: Beispiel eines Entscheidungsbaums<sup>247</sup>**

Aus einem Entscheidungsbaum können nun Entscheidungsregeln abgelesen werden. Diese sind für das Beispiel in Abbildung 15:

- WENN  $(x_1 > w_{10})$  UND  $(x_2 > w_{20})$  DANN Zielwert = C2
- WENN  $(x_1 > w_{10})$  UND  $(x_2 < w_{20})$  DANN Zielwert = C1
- WENN  $(x_1 < w_{10})$  DANN Zielwert = C1

Dies macht einen Entscheidungsbaum leicht interpretierbar. Zusätzlich kann dieser durch Experten verifiziert werden, ob dieses Modell bzw. Entscheidungsregeln sinnvoll ist/sind oder nicht.<sup>248</sup> Zusätzlich gibt es die Möglichkeit mehrere Bäume miteinander zu verbinden. Dies fällt unter die Ensemblemethoden.<sup>249</sup>

### Diskriminanzanalyse

Bei der *Diskriminanzanalyse* wird angenommen, dass die Datensätze einer Klasse linear von den Datensätzen anderer Klassen getrennt werden können. Dabei wird auf eine Betrachtung der Dichte der Trainingsdaten innerhalb der Klassen verzichtet. Die Methode schätzt die Parameter dieser Trennfunktion durch die Analyse der Trainingsdaten und konzentriert sich einzig auf die Lage dieser Trennfunktion.<sup>250</sup>

### Support Vector Machine

Die *Support Vector Machine (SVM)* gehört zur Klasse der Kernel Maschinen und basiert auf der Diskrimanten-Methode. Auch hier wird versucht, Datensätze, die als Objekte im Raum vorliegen, durch Trennebene in Klassen zu teilen.<sup>251</sup>

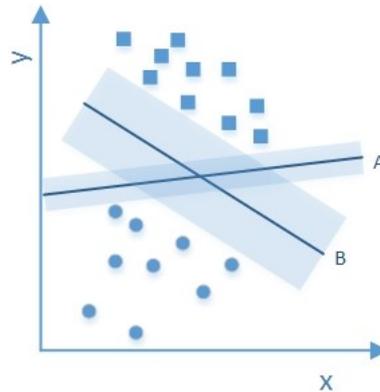
<sup>247</sup> Quelle: in Anlehnung an Alpaydin, E. (2010), S. 186

<sup>248</sup> Vgl. ebda., S. 198

<sup>249</sup> Vgl. Witten, I. H. et al. (2011), S. 356

<sup>250</sup> Vgl. Alpaydin, E. (2010), S. 209 f.

<sup>251</sup> Vgl. ebda., S. 319



**Abbildung 28: Mögliche Platzierungen der Trennebene<sup>252</sup>**

Die Möglichkeiten diese zu platzieren sind vielfältig, wie in Abbildung 28 dargestellt wird. Dieses Optimierungsproblem, einen möglichst großen Abstand zwischen den Objekten am Rand der Klassen und der Ebene zu schaffen, wird nicht durch Schätzung sondern durch Berechnung mit Hilfe von Lagrange-Multiplikatoren für die Testdaten gelöst und die daraus abgeleiteten Supportvektoren anschließend zur Verwendung an den Testdaten im Modell gespeichert. Die Abstände zwischen Ebene und Randobjekten soll dazu führen, dass Testdaten richtig klassiert werden.<sup>253</sup> Sind die Daten nicht linear trennbar, kann mit dem Kernel-Trick das lineare Modell in ein neues Achsensystem überführt werden. Durch die Lösung in diesem System und anschließende Rücktransformation wird ein nicht-lineares Problem in einem linearen Modell gelöst.<sup>254</sup>

### **k-Nearest-Neighbor**

Der k-Nearest-Neighbor-Algorithmus (KNN, Nächste-Nachbarn) gehört zur Klasse der parameterfreien Klassifikationsverfahren. Während des Trainings erfolgt nur eine Speicherung der Daten und keine Klassifikationsvorschrift. Erst bei der Anwendung des Modells auf Testdaten, werden diese Datensätze ausgewertet und mit dem zu klassierenden Datensatz abgeglichen. Dabei wird nach dem ähnlichsten gesucht und dem Testdatensatz dessen Klasse zugeordnet. Diese Methode wird auch als „faules Lernen“ (engl. Lazy Learning) bezeichnet.<sup>255</sup>

### **Ensemblemethoden**

Die Kombination von mehreren Klassifikationsmethoden werden Ensemblemethoden genannt. Es gibt verschiedene Ansätze, dazu gehören Bagging und Boosting.<sup>256</sup> *Bagging* kombiniert mehrere Klassenzuordnungen aus verschiedenen Modellen, wobei die einzelnen Klassenzuordnungen gleich gewichtet und am Ende der Durchschnitt bestimmt wird.<sup>257</sup> Beim *Boosting* entsteht aus vielen schwachen Klassifikatoren ein

<sup>252</sup> Quelle: in Anlehnung an Alpaydin, E. (2010), S. 314

<sup>253</sup> Vgl. ebda., S. 309 ff.

<sup>254</sup> Vgl. ebda., S. 319

<sup>255</sup> Vgl. Shakhnarovich, G. et al. (2005), S. 1

<sup>256</sup> Vgl. Witten, I. H. et al. (2011), S. 351 f.

<sup>257</sup> Vgl. Breiman, L. (1996), S. 123

starker Klassifikator. Dieses allgemeine Verfahren hat sehr viele verschiedene Implementierungen, die bekannteste ist AdaBoost.<sup>258</sup>

Im praktischen Fallbeispiel im folgenden Kapitel wird die Entwicklung eines Modells mit allen im Kapitel 3.2 bzw. Kapitel 3 vorgestellten Schritten durchlaufen. Ziel ist es, das Produktionslos durch korrekte Zuordnung der Klasse die Sinterung ohne weitere Einstellung der Öfen durchläuft, und dadurch die Intensität der Betriebsmittelnutzung zu erhöhen. Schlussendlich ist neben der Tauglichkeit des Modells die Kapazität bzw. deren Veränderung die Kennzahl, mit der der Einsatz des Modells begründet wird.

---

<sup>258</sup> Vgl. Witten, I. H. et al. (2011), S. 358

## 4 Praktische Fallstudie

Ziel dieses Kapitels ist es, die Vorgehensweise und Ergebnisse von der praktischen Bearbeitung der in Kapitel 1.1 vorgestellten Aufgabenstellung bzw. des Problems zu erläutern. Als erstes wird die Problemstellung wiederholt und ein Grundverständnis für die Art und Aufbau der Datenstruktur geschaffen. Anschließend werden die im Laufe der praktischen Bearbeitung dieser verwendeten Softwareprogramme vorgestellt und die Datenaufzeichnung und deren Zusammenhänge betrachtet.

Danach erfolgt eine Zweiteilung der Problemstellung in die Standorte Deutschlandsberg und Šumperk, wobei in die Bearbeitung des Standort Šumperk die Daten und Erkenntnisse aus Deutschlandsbergs einfließen. Zweck dieser Betrachtung ist die praktische Umsetzung der theoretischen Hintergründe und die Beantwortung der in Kapitel 1.2 formulierten Forschungsfrage, ob und wie der Widerstandswert  $R_{25}$  vorausgesagt werden könnte. Dabei orientiert sich die Vorgehensweise am CRISP-DM. Danach wird auf die Forschungsfrage, wie sich eine exakte Vorhersage auf die Kapazitätsauslastung auswirken könnte, eingegangen.

### 4.1 Einführung

Wie bereits in Kapitel 1.1 erläutert, gibt es mehrere Freigaben entlang des Prozesses zur Herstellung von PTC-Bauteilen. Dieser ist zur besseren Übersicht nochmals in Abbildung 29 dargestellt wird. Sie wird dem Umstand angepasst, dass vorerst nur für einen Produkttypen je ein Modell für die beiden Standorte Deutschlandsberg und Šumperk mit zwei unterschiedlichen Zielwerten erstellt werden soll. Dieser Typ, der im Betrachtungszeitraum eines Jahres am öftesten produziert wurde, besteht aus einer in Deutschlandsberg gefertigten Grundmasse, die als reine Masse weiterverarbeitet wird. Daher entfallen das Homogenisieren und die danach notwendige Mischmassenfreigabe.



Abbildung 29: Produktionsprozess der PTC-Bauteile<sup>259</sup>

Entlang des in Abbildung 29 dargestellten Prozesses fallen eine Vielzahl von Messwerten an. Generell stehen die Daten der Vor- und Nachmahlung der Granulatproduktion, sowie der Granulatfreigabe, und die Daten der Sinterfreigabe für diese Arbeit zur Verfügung. Zusätzlich ist es notwendig, eine Verknüpfung zwischen den Chargen- und Produktionslosnummern der Granulatproduktion und der Frontend Production herzustellen. Die Zusammenhänge und der Umstand, dass mehrere Chargen aus den vorhergehenden Produktionsabschnitten in die nachfolgenden einfließen können, werden in Abbildung 30 dargestellt.



Abbildung 30: Anzahl einfließender Chargen je Produktionsabschnitts<sup>260</sup>

Die Auftragseinplanung erfolgt in Chargen, die beim Start des Loses in der Vormahlung zwischen 1020 und 1200 kg umfasst. Durch Schwankungen im Produktionsprozess werden Endmengen von 490 bis 1010 kg erreicht. Diese werden in den Nachmahlungsprozess eingebracht, wobei sich die Chargen- und Losnummer ändert, da bis zu zwei Chargen aus der Vormahlung in die Nachmahlung einfließen können. In diesem Fall gibt es keine klare Startmenge, wobei sich diese in der Größenordnung von der Liefermenge der Vormahlung bewegt. Die Nachmahlung wird als solches in Šumperk wiederum mit einer neuen Chargen- und Losnummer versehen, bei denen bis zu vier Chargen der Nachmahlung gemischt werden. Hier beträgt die Startmenge ca. 2000 kg. Die Sinterfreigabemessungen erfolgen fortlaufend während des Sinterns. Für die Beantwortung der Problemstellung ist nur die erste Messung relevant, da diese

<sup>259</sup> Quelle: in Anlehnung an EPCOS PPD PTC PD (2016), S. 10

<sup>260</sup> Quelle: eigene Darstellung

das Zielattribut für Šumperk darstellt. Da bei den nachfolgenden Messungen die Sinterereinstellungen an die Ergebnisse angepasst werden, liegt ab diesem Zeitpunkt kein standardisierter Prozess mehr vor.

Vormahlung	Nachmahlung	Materialverwendung	Sinterfreigabe
Stammdaten Komponenten Rückmeldungen Prozesse Messwerte	Stammdaten Komponenten Rückmeldungen Prozesse Messwerte Granulatfreigabe RFA-Analyse	Stammdaten Prozesse	Stammdaten Prozessgrenzen Messwerte
↔ = Chargennummer	↔ = Chargennummer	↔ = Losnummer	↔ = Lieferschein
$\Sigma = 162$ Attribute	$\Sigma = 116$ Attribute	$\Sigma = 43$ Attribute	$\Sigma = 36$ Attribute

Abbildung 31: Aufbau der Ausgangsdaten<sup>261</sup>

Die Ausgangsdaten sind auf 4 einzelne Datensets verteilt. Diese sind in Abbildung 31 übersichtsmäßig dargestellt, da eine vollständige Auflistung aller enthaltenen Attribute zu diesem Zeitpunkt nicht zielführend wäre. Daher werden die Kategorien, in die die Attribute zusammengefasst werden können, das Schlüsselattribut und am unteren Bildrand jeweils Gesamtanzahl der vorliegenden Attribute pro Datenset angegeben. Die Daten der Vor- und Nachmahlung sind bereits aggregierte Daten, die aus unterschiedlichen Quellen für unternehmensinterne Auswertungen zusammengestellt werden. Nach der Begutachtung der Rohdaten für diese Auswertung wird entschieden, mit den aggregierten Daten zu arbeiten, da die Aufbereitung der Rohdaten eben jene aggregierten Daten zur Folge hätte.

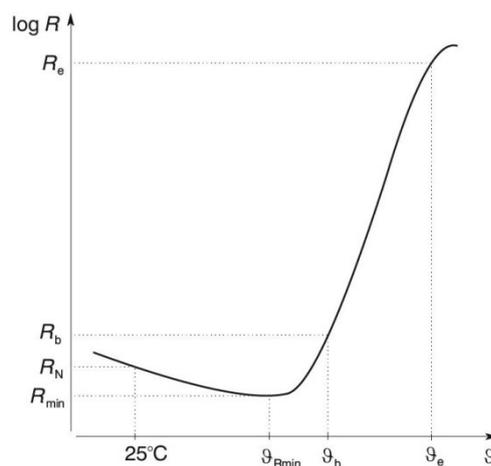


Abbildung 32:  $R(\theta)$ -Kennlinie eines Kaltleiters (schematisch)<sup>262</sup>

Die Zielattribute sind in den Datensets der Nachmahlung und der Sinterfreigabe zu finden. Bei beiden handelt es sich um den Nennwiderstand  $R_N$ , der bei 25°C

<sup>261</sup> Quelle: eigene Darstellung

<sup>262</sup> Quelle: Reisch, M. (2007), S. 294 (leicht modifiziert)

gemessen wird und daher als  $R_{25}$  bezeichnet wird. In Abbildung 32 ist der schematische Verlauf zum Verhalten von Widerstand zu Temperatur eines PTC-Bauteile schematisch dargestellt. Bei Temperaturen  $\vartheta$  kleiner als  $\vartheta_{min}$  ist der Temperaturkoeffizient negativ, zudem besitzt das PTC-Bauteil hier den geringsten Widerstand  $R_{min}$  und ab diesem Punkt erhöht sich der Widerstandswert mit ansteigender  $\vartheta$ . Die Bezugstemperatur  $\vartheta_b$  markiert den Beginn des steilen Anstiegs der Widerstandskurve, die über  $R_b = R(\vartheta_b) = 2 R_{min}$  definiert ist. Die Bezugstemperaturen können in einem Wertebereich von  $-30^\circ\text{C}$  bis  $340^\circ\text{C}$  liegen. Das Ende des steilen Anstiegs liegt bei  $\vartheta_e$  mit dem Endwiderstand  $R_e$ . Bei der Herstellung der PTC-Bauteile kann über die chemische Zusammensetzung die Bezugstemperatur eingestellt werden. Eine Anpassung des  $R_{25}$  erfolgt über die im Herstellungsprozess eingesetzten Sintertemperaturen.<sup>263</sup> Durch eine korrekte Vorhersage des zu erwartenden  $R_{25}$  können die Sintereinstellungen so angepasst werden, dass der gewünschte  $R_{25}$  erreicht wird und keine Nachjustierung der Einstellungen nötig ist.

Die im Folgenden vorgestellten (Zwischen-)Ergebnisse und Berechnungen wird unter Zuhilfenahme der Softwareprogramme Microsoft Excel 2010, PyCharm Community Edition und Matlab 2015b erstellt. Microsoft Excel 2010 als Tabellenkalkulationsprogramm mit der Möglichkeit zur Darstellung der Daten und Erstellen von Diagrammen<sup>264</sup>, sowie der Verknüpfung mit Visual Basic for Applications (VBA), der Microsofteigenen Skriptsprache<sup>265</sup>, bilden das Grundgerüst für diese Arbeit. PyCharm Community Edition bildet mit der Programmiersprache Python bilden einen kostenlos nutzbare Plattform zur Verarbeitung großer Anzahl von Datensätzen.<sup>266</sup> Dies gilt ebenfalls für Matlab 2015b, wobei die kommerziellen Lizenzen mehrere 1000 Euro kosten können.<sup>267</sup> Daher wurde versucht, Matlab weitestgehend nicht zu verwenden.

Mit diesen Werkzeugen und den Ausgangsdaten wird nun im zweiten Kapitel der Weg zum Erlangen des Verständnisses für die Daten, die Aufbereitung und die Modellierung für die Daten für Deutschlandsberg erläutert.

## 4.2 Standort Deutschlandsberg

In Deutschlandsberg (DL) wird die Granulatfreigabe betrachtet. Der Input für diese sind Probeteile aus dem freizugebenen Granulat und standardisierte Kontrollteile, die für alle Lose diesen Typs gleich sind. Diese werden gemeinsam gesintert und anschließend deren Widerstände gemessen. Aus dem Verhältnis zueinander wird das Zielattribut „R25 Vergl.%“ berechnet. Ziel ist es, diesen Wert mit einem Modell vorherzusagen.

<sup>263</sup> Vgl. Reisch, M. (2007), S. 293 ff.

<sup>264</sup> Vgl. Microsoft, <https://products.office.com/en-us/excel> (Zugriff: 13.02.2017)

<sup>265</sup> Vgl. Microsoft, <https://msdn.microsoft.com/en-us/library/office/gg264383.aspx> (Zugriff: 13.02.2017)

<sup>266</sup> Vgl. JetBrains s.r.o., <https://www.jetbrains.com/pycharm/> (Zugriff: 13.02.2017)

<sup>267</sup> Vgl. MathWorks, <https://de.mathworks.com/pricing-licensing.html?prodcode=ML> (Zugriff: 13.02.2017)

### 4.2.1 Analyse der Fehlenden Daten

Die Daten vom Standort Deutschlandsberg umfassen die in Abbildung 31 dargestellten Datensets zur Vor- und Nachmahlung im Zeitraum von September 2015 bis August 2016. Wie aus Abbildung 31 hervorgeht, beträgt die Gesamtanzahl der Attribute 278. Sämtliche Attribute, für das Ziel ein Modell wie in Kapitel 0 zu entwickeln, heranzuziehen, ist nicht ratsam. Gründe hierfür können sein, dass diese entweder leer, und somit überflüssig sind, oder zu viele Ausreißer enthalten. Ebenfalls kann das Skalenniveau (siehe Kapitel 3.3.1) unpassend für die gewählten Modelle sein. Genauso müssen die Datensätze überprüft werden, ob Chargen zu Forschungs- und Entwicklungszwecken produziert worden sind. Diese müssen ebenfalls aussortiert werden.

Die Aufbereitung der Daten ist im CRISP-DM (siehe Kapitel 3.2) als dritter Schritt angesiedelt. In diesem Fall macht es Sinn ihn vorzuziehen, um gerade leere Attribute zu entfernen und somit die Übersichtlichkeit zu erhöhen. Außerdem werden solche Datensätze entfernt, die zur weiteren Analyse nicht geeignet sind. Diese Bereinigung wurde mit der in Kapitel 3.4.2 vorgestellten Analyse der fehlenden Daten kombiniert.

#### Vormahlung

Bei den Daten der Vormahlung (VM) werden sämtliche Forschungs- und Entwicklungslose entfernt sowieso solche Datensätze, bei denen kein Ist-Ende angegeben ist, da diese offensichtlich noch nicht fertig gestellt sind, und somit noch nicht in eine Nachmahlung eingeflossen sind. Im anschließenden Schritt 1 werden solche Attribute gelöscht, die Kommentare enthalten oder zu 100 % leer sind. Dazu zählen z. B. Komponenten, die nicht für diesen Produkttyp verwendet werden. Im nächsten Schritt werden Datensätze mit mehr als 69 % fehlender Daten aussortiert. Die 69 %-Schwelle ist als solches nicht in der Literatur zu finden und wird nach Betrachtung der Prozentsätze aller Datensätze gewählt. Abschließend werden solche Attribute eliminiert, die durch das Entfernen der Datensätze keine Werte mehr enthielten.

**Tabelle 8: Veränderung der Anzahl und Prozente der fehlenden Werte (VM)<sup>268</sup>**

	Anzahl der fehlenden Werte	Prozent der fehlenden Werte
<b>Start</b>	20973	54%
<b>Ende Schritt 1</b>	4478	20%
<b>Ende Schritt 2</b>	3428	15%
<b>Ende Schritt 3</b>	2976	14%

Das Verhältnis von fehlenden zu vorhandenen Daten verbessert sich durch diesen Vorgang von 54 % auf 14 %. Dies ist in Tabelle 8 ersichtlich.

Die verbleibenden Attribute werden in die Kategorien MAR und MCAR eingeteilt. Diese Einteilung befindet sich aus Gründen der besseren Darstellung im Anhang Tabelle 27. Die Einteilung wird nach subjektiver Einschätzung vorgenommen. So werden solche Attribute als MAR kategorisiert, die Komponenten repräsentieren. Oftmals gibt es für

<sup>268</sup> Quelle: eigene Darstellung

eine Komponente mehrere Attribute. So ist es möglich, dass bei der Zugabe von z. B. Blei zwei unterschiedliche Rohstoffe verwendet werden, die jeweils mit Menge und Chargennummer vermerkt werden. Dies hat zur Folge, dass der Komponente Blei vier Attribute zugeordnet sind. Sollte allerdings der Bleibedarf mit einer Rohstoffcharge abgedeckt werden können, bleiben die anderen beiden Felder frei. Diese werden in weiterer Folge durch eine Null ergänzt.

Attribute, die als MCAR bezeichnet werden, sind solche, die eigentlich gemessen hätten werden sollen, aber nicht gemessen worden sind. Dies trifft z. B. auf eine fehlenden Vermerk der Bleimenge oder auf eine nicht gemessene oder nicht eingetragene Messung während eines Prozesses. Hier wird mit dem arithmetischen Mittel (siehe Kapitel 3.3.2) des jeweiligen Attributs ergänzt.

Das Ergebnis dieses Vorgangs ist ein Datenset bestehend aus 226 Datensätzen mit 90 Attributen.

### Nachmahlung

Die Analyse der fehlenden Werte in den Daten der Nachmahlung (NM) verläuft ähnlich wie bei der Vormahlung. Bevor die Analyse der fehlenden Werte begonnen wird, werden die Forschungs- und Entwicklungslose und solche Lose, die keinen Zielwert in „R25 Vergl.“ und keinen Eintrag in „Umsatz1“ besitzen, gelöscht. Letzteres ist die Vormahlungscharge, die in der Nachmahlung als Grundlage eingesetzt wird, und ohne die eine Verknüpfung von Vor- und Nachmahlung nicht möglich ist. Im ersten Schritt werden die Attribute mit 100 % fehlenden Werten und die Kommentare gelöscht. In einem zweiten Schritt wird der Prozentsatz der fehlenden Werte der einzelnen Datensätze betrachtet, wobei hier solche Attribute ausgelassen werden, die nach dem „skip lot“-Verfahren gemessen werden. Dabei wird das Attribut nur zu bestimmten Zeiten gemessen und nicht für jedes Los. Daher entstehen hier bewusst Lücken. Diese werden beim Auffüllen anders behandelt. Bei einem Anteil von mehr als 42 % fehlender Werte wurden Datensätze gelöscht.

**Tabelle 9: Veränderung der Anzahl und Prozente der fehlenden Werte (NM)<sup>269</sup>**

	Anzahl der fehlenden Werte	Prozent der fehlenden Werte
<b>Start</b>	8472	33%
<b>Ende Schritt 1</b>	1904	13%
<b>Ende Schritt 2</b>	1860	13%

Die Entwicklung des Verhältnisses von fehlenden zu vorhandenen Werten ist in Tabelle 9 ersichtlich und verbessert sich um 20 %.

Wieder werden die verbleibenden Attribute in die MAR und MCAR aufgeteilt. Die genaue Aufschlüsselung ist im Anhang in Tabelle 28 zu finden. Die Attribute, die als MAR bezeichnet werden, sind die bereits erwähnten nach „skip lot“-Verfahren gemessenen Attribute. Diese werden so aufgefüllt, dass die Datensätze ohne Werte mit dem als letztes gemessenen Wert annehmen. Bei den mit MCAR gekennzeichneten Werten wird, wie bei der Vormahlung, das arithmetische Mittel gemessen über sämtliche Werte des Attributs ergänzt.

<sup>269</sup> Quelle: eigene Darstellung

Das Ergebnis dieses Vorgangs ist ein Datenset bestehend aus 224 Datensätzen mit 63 Attributen.

#### 4.2.2 Datenzusammenführung

Da für die Modellierung die Daten aus Vor- und Nachmahlung verknüpft werden müssen, wird erneut die Abbildung 30 betrachtet. Es können ein oder zwei Chargen der Vormahlung in die Nachmahlung einfließen. Außerdem fließen diese in veränderlichen Gewichtsanteilen ein. Die Vormahlung hat, wie in Kapitel 4.1 erwähnt, unterschiedliche Startgewichte (Ist-Start). Aus diesem Grund müssen vor der Zusammenführung alle Attribute, die sich auf diese Startmenge beziehen, auf diese genormt werden. Gleiches gilt für Attribute der Nachmahlung, die sich auf deren Startmenge beziehen. Bei zwei einfließenden Chargen liegen jeweils zwei Messwerte für ein Attribut vor. Dies wird unterschiedlich gehandhabt:

- Bei Attributen, die gemessene Werte enthalten, wird von beiden der Mittelwert, der mit den Gewichtsanteilen der jeweiligen Chargen gewichtet wird, verwendet.
- Bei Chargennummern der Rohstoffe werden beide beibehalten. (Bei nur einer einfließenden Charge wird der entsprechend zweite Wert freigelassen.)
- Bei Aufzeichnungen über Maschinen (Fryma, Sprühturm, UD (Ofen)) wird jene der ersten einfließenden Charge gewählt.

Zusätzlich werden Rohstoffanteile, wie z. B. die bereits erwähnte Komponente Blei, die durch diesen Prozess der Zusammenführung pro Datensatz bis zu vier unterschiedliche Ausprägungen haben kann, zu einem Wert zusammengefasst. Genauso werden die Ringtemperatur 1 und Ringtemperatur 2 in einem Attribut gemittelt, da diese nicht klar unterschieden werden können.

Hinsichtlich der Modellierung wird entschieden, Chargennummern der Rohstoffe aus der Analyse auszuschließen, da sich diese in Zukunft ändern können und nicht das passende Skalenniveau haben. Genauso werden die Maschinen, die für die Produktion verwendet wurden, vorerst ausgeschlossen und nicht zielführende Attribute wie Materialnummer, Ist-Start, oder Umsatz Menge entfernt.

Durch diese Zusammenführung ergibt sich ein Datenset mit 99 Attributen, der Chargennummer als Index zur Identifizierung der Datensätze und 218 Datensätzen. Sämtliche Attribute sind im Anhang in Tabelle 29 aufgelistet, wobei die hellblauen hinterlegten solche Attribute sind, die aus dem Datenset der Nachmahlung stammen, die dunkelblau hinterlegten aus der Vormahlung (mit dem Suffix „\_VM“) und der Zielwert „R25 Vergl.“ ist rot hinterlegt und stammt aus dem Datenset der Nachmahlung.

Dieses neue Datenset dient als Grundlage für die Trainingsdaten für das Modell für den Standort Deutschlandsberg. Bevor ein solches trainiert werden kann, müssen für die Daten, wie in Kapitel 3.3 beschrieben, verstanden werden.

### 4.2.3 Auswertung

Aus diesem neu entstandenen Datenset werden die in Kapitel 3.3.2 und Kapitel 3.3.3 vorgestellten statistische Maße mit Hilfe von Python berechnet. Hierbei ist das Ziel, ein Verständnis für die Daten zu entwickeln. Dafür werden der arithmetische Mittelwert  $\bar{x}$ , der Median  $x_{med}$ , das Minimum, das Maximum, die Varianz  $\tilde{s}^2$ , die Standardabweichung  $\tilde{s}$ , die Wölbung  $\gamma$ , die Schiefe  $g_m$  und der Bravais-Pearson-Korrelationskoeffizient zwischen dem jeweiligen Attribut und dem Zielwert ermittelt. Die Werte für Varianz, Standardabweichung, Schiefe und Korrelationskoeffizient sind aus der Tabelle 30 zu entnehmen.

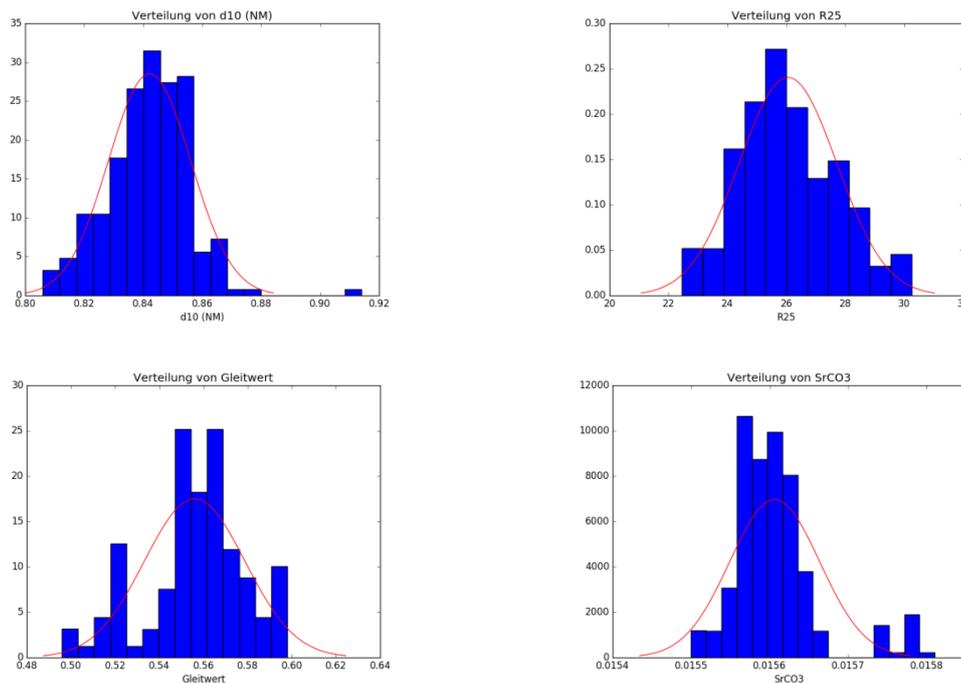
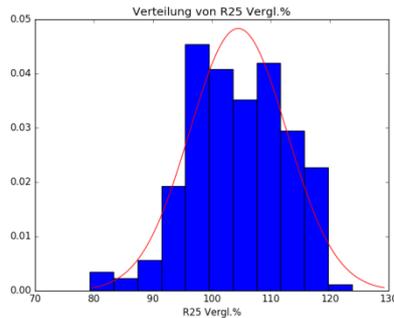


Abbildung 33: Beispiele für Histogramme für DL mit Ausreißer<sup>270</sup>

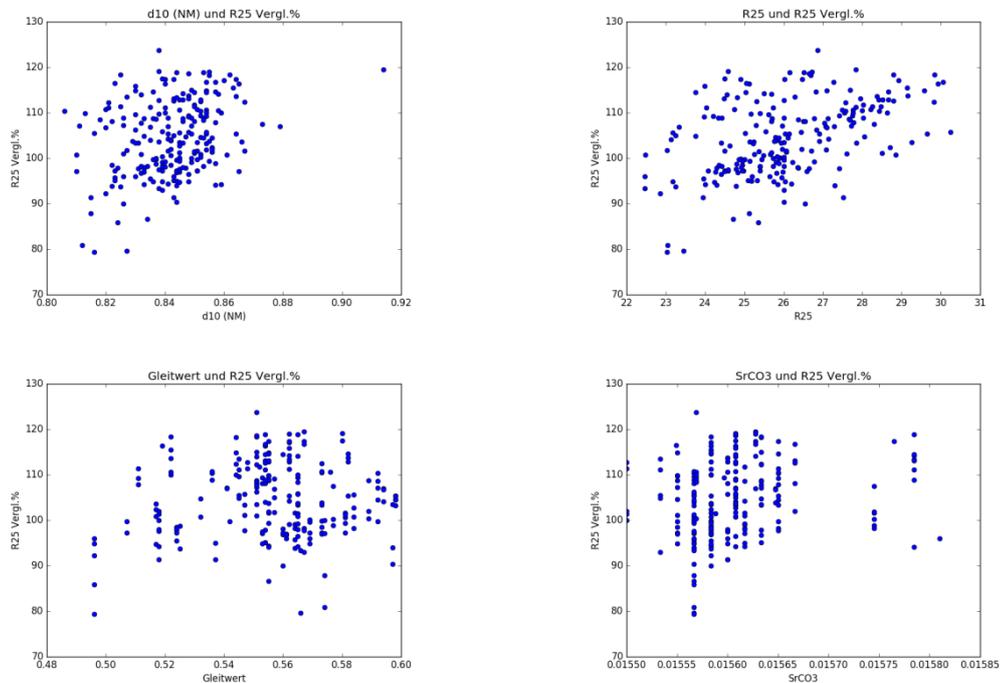
Zusätzlich werden Histogramme für sämtliche Attribute erstellt, in die die Kurven der Normalverteilung eingefügt wird. In Abbildung 33 sind diese für die Attribute d10 (NM) (links oben), R25 (rechts oben), Gleitwert (links unten) und SrCO3 (rechts unten) zu sehen. Sie weisen unterschiedliche Verteilungen auf, die nicht der Normalverteilung entsprechen. Dies ist ebenfalls in der Tabelle 30 durch die Werte  $g_m$  und  $\gamma$  ersichtlich. Diese Erkenntnis gilt auch für den Zielwert R25 Vergl.%, zu sehen in Abbildung 34.

<sup>270</sup> Quelle: eigene Darstellung



**Abbildung 34: Histogramm für R25 Vergl.% mit Ausreißern<sup>271</sup>**

In Abbildung 35 sind die Streudiagramme der Wertepaare (Attribut, Zielwert) zu sehen. Als Beispiel werden wieder die Attribute aus Abbildung 33 herangezogen. Für die Attribute d10 (NM), R25 und SrCO3 kann eine positive Korrelation zwischen dem jeweiligen Attribut und dem Zielwert festgestellt werden. Beim Attribut R25 war dies zu erwarten, da dieser als Grundlage für die Berechnung von R25 Vergl.% verwendet wird. Bei dem Attribut Gleitwert gibt es nur eine sehr schwache Korrelation.

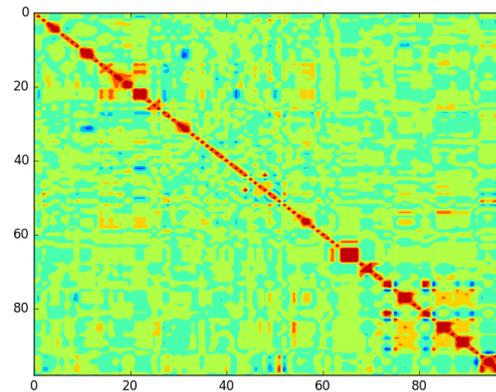


**Abbildung 35: Beispiele für Streudiagramme für DL mit Ausreißer<sup>272</sup>**

Die Zusammenhänge zwischen den einzelnen Attributen kann ebenso in der Korrelationsmatrix in Abbildung 36 dargestellt werden. Die dunkelroten und dunkelblauen Areale sind solche mit hoher Korrelation. Die helleren Farben sind solche mit geringer Korrelation. Die Achsenbeschriftung erfolgt in Nummern, die die Attribute bezeichnen, in der Reihenfolge wie sie in Tabelle 30 zeilenweise von links nach rechts aufgezählt wurden.

<sup>271</sup> Quelle: eigene Darstellung

<sup>272</sup> Quelle: ebda.



**Abbildung 36: Korrelationsmatrix (graphisch) für DL mit Ausreißern<sup>273</sup>**

Die hohe Korrelation zwischen den Attributen um den Punkt (10,10) in Abbildung 36 entsteht durch die verschiedenen Durchmesser der Körnung der Rohstoffe und die dadurch resultierende Schlickerdichte für die Nachmahlung. Die Attribute um Punkt (20, 20) werden bei der Granulatfreigabe gemessen. Ebenso korrelieren die gemeinsam bei der Röntgenfluoreszenzanalyse gemessenen Attribute um Punkt (50,50). Die Korrelation um Punkt (60,60) ist mit der um Punkt (10,10) zu vergleichen, wobei hier die Attribute der Vormahlung aufgetragen sind. Um Punkt (67,67) liegt die Ursache der Korrelation im Zusammenhang zwischen den einzelnen Heizgruppen. Die hohe Korrelation in der rechten unteren Ecke der Korrelationsmatrix folgt aus dem Zusammenhang zwischen den Attributen, die von jeweils einem Attribut das Minimum, Maximum und Standardabweichung angeben.

Zusätzlich wird der Kolmogorov-Smirnov-Test für jedes Attribut durchgeführt, um festzustellen, ob bei diesen eine Normalverteilung vorliegt. Dies ist nicht der Fall. Dies wird auch durch die in Tabelle 30 berechneten Werte bestätigt.

Wie bereits in den Histogrammen in Abbildung 33 ersichtlich, gibt es Ausreißer. Diese befinden sich außerhalb der Normalverteilungskurve bzw. stellen die zweite Spitze der bimodalen Verteilung dar (siehe Attribut Gleitwert). Diese gilt es im nächsten Kapitel zu identifizieren. Zusätzlich werden die Veränderungen, die das Eliminieren von Ausreißern in der Datenstruktur mit sich bringt, anhand der fünf in diesem Kapitel als Beispiele besprochenen Attribute aufgezeigt.

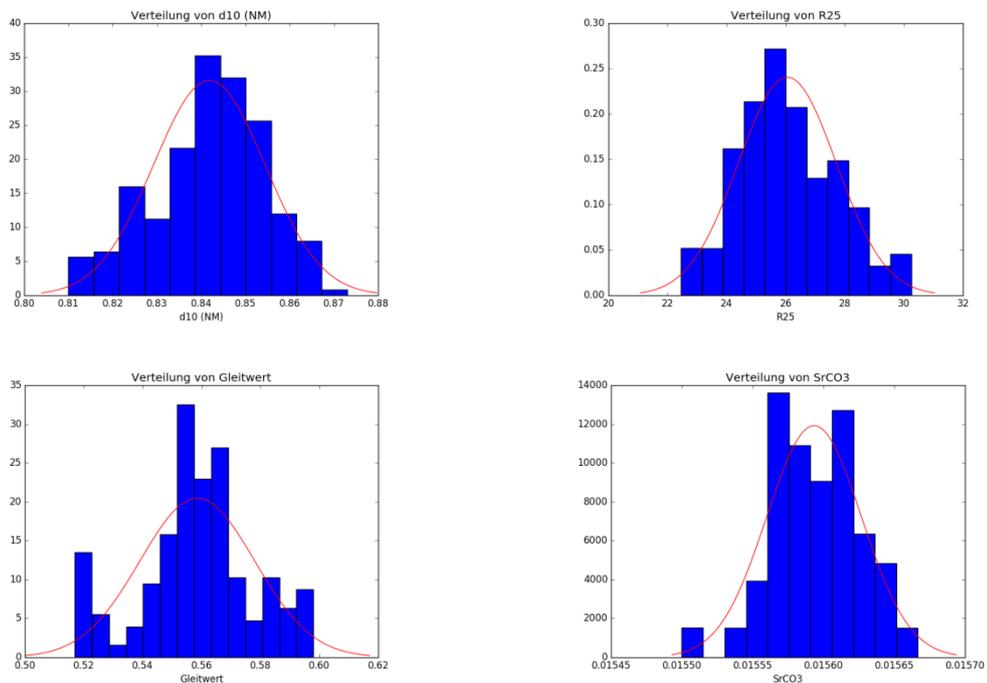
#### 4.2.4 Ausreißeranalyse

Die Identifikation der Ausreißer erfolgte nach der logischen Überlegung, die hinter der Erstellung von Boxplots steht. Dies ist Inhalt des in Kapitel 3.4.1. Die Berechnung des  $IQR$  und die Identifikation von  $Q_1$  und  $Q_3$  für jedes Attribut ist erforderlich. Anschließend werden jene Messwerte innerhalb der einzelnen Attribute markiert, die kleiner als  $Q_1 - 1,5 \times IQR$  oder größer als  $Q_3 + 1,5 \times IQR$  sind. Lediglich 21 von 99 Attributen und 33 von 218 Datensätzen enthalten keine Ausreißer. Eine Aufschlüsselung hinsichtlich Anzahl bzw. Prozentsatz von Ausreißern je Attribut ist in Tabelle 31 im Anhang zu

<sup>273</sup> Quelle: eigene Darstellung

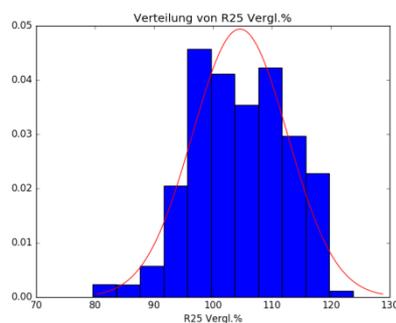
finden. Die Datensätze mit Ausreißern gänzlich auszuschließen, ist keine Option. Somit wird ab diesem Punkt mit zwei Datensets für Deutschlandsberg weitergearbeitet. Das eine ist das ursprüngliche inklusive der Ausreißer. Bei dem anderen werden die Ausreißer mit dem ohne Ausreißer berechneten arithmetischen Mittelwert der jeweiligen Attribute ersetzt.

In Abbildung 37 sind die Histogramme der gleichen Attribute wie in Abbildung 33 dargestellt. Das Attribut R25 enthielt keine Ausreißer, somit gibt es keine Veränderung im Histogramm. Bei den anderen drei Attributen sind die Veränderungen in den Randbereichen und somit die Erscheinungsbilder klar erkennbar. Die andere Häufigkeitsverteilung von SrCO3 ist die wohl die Folge der automatischen Berechnung der Klassenbreiten, die sich aufgrund der neuen Gesamtheit nun auch ändert.



**Abbildung 37: Beispiele für Histogramme für DL ohne Ausreißer<sup>274</sup>**

Auch der Zielwert R25 Vergl.% enthält Ausreißer. Durch die die Beseitigung dieser wandelt sich das Histogramme in Abbildung 34 und zu jenem in Abbildung 38.



**Abbildung 38: Histogramm für R25 Vergl.% ohne Ausreißer<sup>275</sup>**

<sup>274</sup> Quelle: eigene Darstellung

Werden die Streudiagramme in Abbildung 35 und Abbildung 39 miteinander verglichen, fällt auf, dass die Randbereiche weniger stark besetzt sind und sich die Wertepaare in der Mitte konzentrieren.

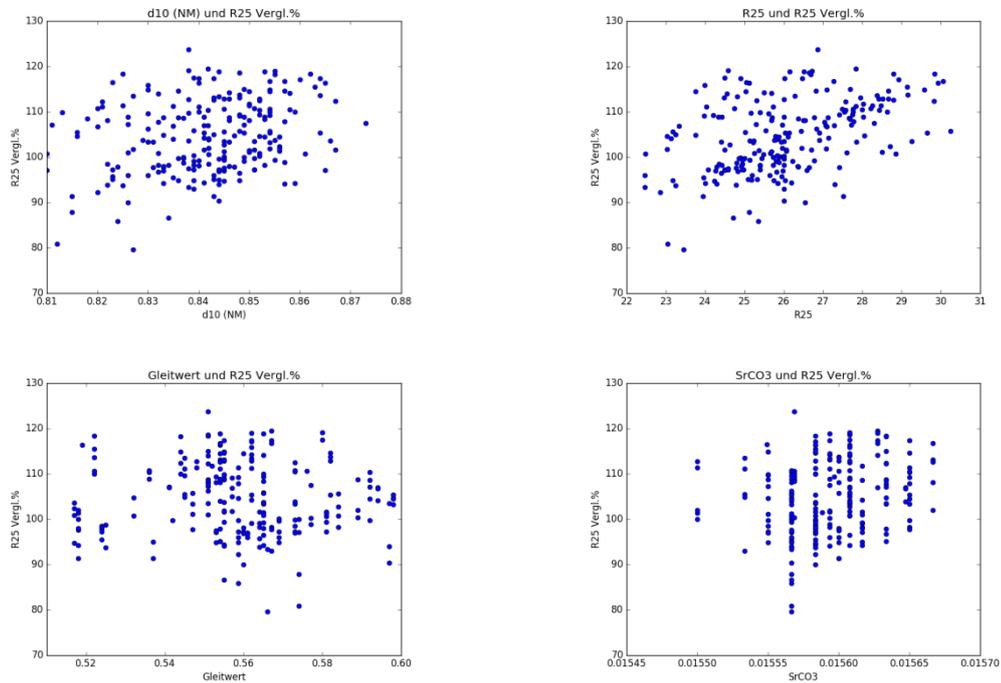


Abbildung 39: Beispiele für Streudiagramme für DL ohne Ausreißer<sup>276</sup>

Auch auf die Korrelationsmatrix hat die Eliminierung der Ausreißer Auswirkungen. Der Vergleich der Korrelationsmatrix mit Ausreißern, in Abbildung 36 und links in der Abbildung 40 zusehen, mit der ohne Ausreißern, rechts in Abbildung 40, zeigt, dass die Korrelation in der rechten unteren Ecke abnimmt. Dies ist auf den Umstand zurückzuführen, dass gerade bei diesen Attributen eine hohe Anzahl von Ausreißern auftritt (siehe Tabelle 31).

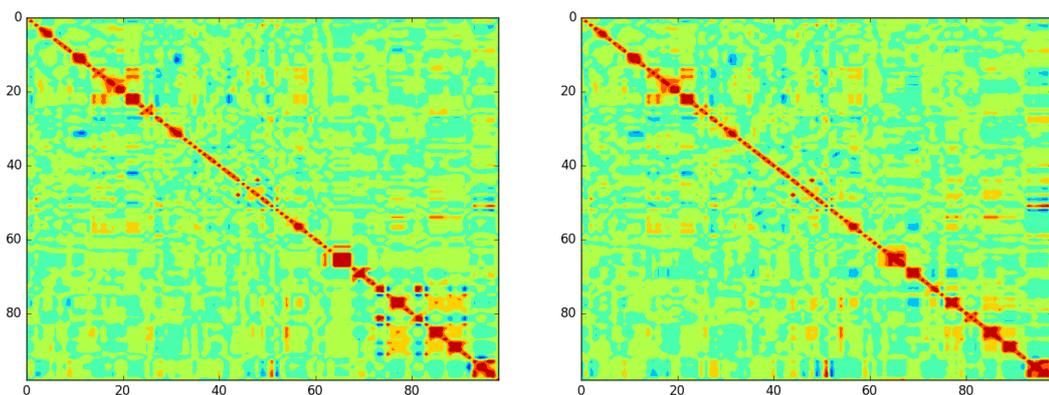


Abbildung 40: Korrelationsmatrizen (graphisch) für DL mit und ohne Ausreißer<sup>277</sup>

<sup>275</sup> Quelle: eigene Darstellung

<sup>276</sup> Quelle: ebda.

<sup>277</sup> Quelle: ebda.

Am Ende dieses Schritts wird entschieden, dass obwohl im Zielattribut ein Ausreißer vorliegt, diesen nicht zu ändern und die anschließende Gruppierung des Zielwerts nur auf Basis der unveränderten Ausprägungen des Attributs vorzunehmen. Dies wird im nächsten Kapitel erläutert.

#### 4.2.5 Gruppierung

In Kapitel 3.4.3 wurden mehrere Ansätze zur Gruppierung erläutert. Für dieses Fallbeispiel wurde die Berechnung der Klassenanzahl nach (3.1) für den Zielwert R25 Vergl.% durchgeführt:

$$k_{DL} = \sqrt{n_{DL}} = \sqrt{218} \approx 15$$

Durch diese Berechnung entstehen 15 Klassen, die durch eine Bezeichnung (arithmetisches Mittel aller in der Klasse vertretenen Werte), die Klassenmitte, die Unter- und Obergrenze definiert werden, siehe Tabelle 10. Zusätzlich wird die Breite angegeben, die in der Skala des Zielattributs, in diesem Fall Prozentpunkte, vorliegt.

**Tabelle 10: Einteilung von R25 Vergl.% in Klassen<sup>278</sup>**

Bezeichnung	Klassenmitte	Untergrenze	Obergrenze	Anzahl von Datensätzen pro Klasse	Breite
88.62466667	86.57	79.32	93.82	15	14.5
94.87333333	95.02	94.02	96.02	15	2
96.936	96.66	96.02	97.3	15	1.28
98.03066667	98.03	97.53	98.53	15	1
99.516	99.48	98.8	100.17	15	1.37
101.05266667	100.96	100.27	101.66	15	1.39
102.92466667	102.71	101.7	103.72	15	2.02
104.642	104.58	103.73	105.43	15	1.7
106.6053333	106.59	105.44	107.74	15	2.3
108.566	108.51	107.92	109.1	15	1.18
110.0333333	109.99	109.2	110.78	15	1.58
111.99666667	111.86	110.8	112.93	15	2.13
114.066	114.03	113.13	114.93	15	1.8
116.8253333	116.64	114.94	118.34	15	3.4
119.52375	121.095	118.41	123.78	8	5.37

Dieses Kapitel widmet sich einzig dem Zielwert und die Grundlage für die Klassifikation künftiger Datensätze geschaffen. Durch die vorgenommene Einteilung gibt es 15 bzw. 8 Datensätze je Klasse, die für das Training des Modells verwendet werden können. Hierbei muss beachtet werden, dass die Einteilung in Klassen an dem Datenset mit Ausreißern vorgenommen und für das bereinigte übernommen wird. Die nachfolgend durchgeführte Überprüfung auf Eignung für die PCA bezieht sich wieder auf alle Attribute und ist der letzte Schritt vor der Erstellung eines Modells.

<sup>278</sup> Quelle: eigene Darstellung

#### 4.2.6 Eignung für die PCA

Die in Kapitel 3.4.4 vorgestellten Vortests, die über die Eignung eines Datensets für die PCA entscheiden, werden für beide Datensets durchgeführt. Es werden die Korrelationsmatrix und ihre Inverse betrachtet. In der Literatur beinhaltet keine aussagekräftigen Grenzwerte. Die Korrelationsmatrix sollte wenige kleine Werte enthalten, wobei hier die Frage ist, in wie weit „klein“ definiert ist und ob klein große negative Zahlen einschließt oder nicht. Ähnliche Überlegungen werden für die Überprüfung der Inversen der Korrelationsmatrix auf die Eigenschaft der Diagonalmatrix angestellt: Was bedeutet klein und in wie weit betrifft dies größere negative Zahlen? Daher wird im Folgenden nur auf die Ergebnisse des KMO-Tests eingegangen, der iterativ durchgeführt wurde, bis zumindest jedes Attribut einen  $MSA \geq 0,5$  hatte. Um dies zu erreichen, werden jene Attribute gelöscht, die  $MSA < 0,5$  hatten, und der Test wiederholt. Beide Datensets haben zu Beginn dieser Überprüfung dieselben Attribute, die in Tabelle 29 ersichtlich sind, wobei das Zielattribut R25 Vergl.% ausgeschlossen wird, da dieses nicht in der nachfolgenden PCA eingerechnet werden darf. Zusätzlich werden die Datenwerte nach (3.18) standardisiert.

##### Datenset mit Ausreißern

Für das Datenset mit Ausreißern wird der KMO-Test viermal durchgeführt, bis für sämtliche Attribute  $MSA \geq 0,5$  gilt. Im Anhang in Tabelle 32 werden diese genau aufgelistet. Attribute mit  $MSA < 0,5$  sind rot markiert und gelöscht. Zusätzlich sind sie durchgestrichen. Die Anzahl der Attribute nach dem Löschen ändert sich nach jedem Durchgang, wobei mit 98 gestartet wurde. Die Entwicklung ist in Tabelle 11 zu sehen.

**Tabelle 11: Entwicklung der Anzahl der Attribute (VM+NM mit Ausreißern)<sup>279</sup>**

	Durchgang			
	1	2	3	4
<b>Anzahl zu löschende Attribute</b>	33	2	1	0
<b>Anzahl Attribute nach Löschen</b>	65	63	62	62

Das MSA-Kriterium für das gesamte Datenset hat sich während der vier Durchgänge wie in Tabelle 12 dargestellt verhalten.

**Tabelle 12: MSA für Datenset (VM+NM mit Ausreißern)<sup>280</sup>**

	Durchgang			
	1	2	3	4
<b>MSA</b>	0,5957	0,7301	0,7356	0,7382

Das MSA-Kriterium für das Datenset verbessert sich mit jedem Durchgang. Die hohe Anzahl von 36 zu löschenden Attributen nach dem ersten Durchgang hat einen positiven Einfluss auf das das MSA-Kriterium. Dieser Wert ist nach der in Tabelle 6 angegebenen Skala als für die PCA geeignet einzustufen.

Das Ergebnis dieses Schritts ist ein Datenset mit 63 Attributen (Tabelle 11 plus den Zielwert) und 218 Datensätzen.

<sup>279</sup> Quelle: eigene Darstellung

<sup>280</sup> Quelle: ebda.

### Datenset ohne Ausreißer

Für das Datenset ohne Ausreißer wird genauso wie bei dem mit Ausreißern vorgegangen. Es gibt vier Durchgänge, wobei beim ersten die Inverse der Korrelationsmatrix nicht berechnet werden kann. Daher kann in weiterer Folge das MSA-Kriterium nicht ermittelt werden. Um dies zu ermöglichen, werden nach dem ersten fehlgeschlagenen Versuch jene fünf Attribute gelöscht, die bei der Berechnung Probleme bereitet haben und der Test wiederholt. Die Auflistung aller Attribute ist im Anhang in Tabelle 33 zu finden. In Tabelle 13 ist die Entwicklung der Anzahl der Attribute zu sehen.

**Tabelle 13: Entwicklung der Anzahl der Attribute (VM+NM ohne Ausreißer)<sup>281</sup>**

	Durchgang			
	1	2	3	4
<b>Anzahl zu löschende Attribute</b>	5	9	1	0
<b>Anzahl Attribute nach Löschen</b>	93	84	83	83

Der MSA für den gesamten Datensatz verhält sich während der vier Durchgänge wie in Tabelle 14 dargestellt. Auch hier ist eine Veränderung hin zu einem höheren MSA-Kriterium zu sehen.

**Tabelle 14: MSA für Datenset (VM+NM ohne Ausreißer)<sup>282</sup>**

	Durchgang			
	1	2	3	4
<b>MSA</b>	nan	0,7077	0,7411	0,7441

Aus dieser Abfolge von Durchgängen ergibt sich als Endresultat ein Datenset mit 84 Attributen (Tabelle 13 plus dem Zielwert) und 218 Datensätzen.

Nach der Durchführung des KMO-Tests folgt die Anwendung der PCA. Die Ausführung dieser erfolgte im gleichen Schritt wie das Trainieren der Modelle. Beides wird im nächsten Kapitel besprochen und die Ergebnisse präsentiert.

### 4.2.7 Modellierung

Das Erstellen der Modelle, die PCA und die Anwendung der Modelle erfolgen in Matlab mittels der App Classification Learner (CL). In dieser App stehen verschiedene Modelle, die mit Hilfe vom überwachten Lernen trainiert werden können, zur Verfügung. Dazu zählen die Kategorie der Entscheidungsbäume, der Diskriminanzanalyse, der SVM, der KNN und der Ensemblemethoden. Jede Kategorie enthält mehrere Variationen der jeweiligen Methode.<sup>283</sup> Zusätzlich können die Einstellungen für die Methode variiert werden, worauf aber in diesem Kontext verzichtet wird und die Standardeinstellungen verwendet werden. In Tabelle 15 und Tabelle 16 wird eine Übersicht über die in Matlab zur Verfügung stehenden Methoden zur Klassifikation gegeben.

<sup>281</sup> Quelle: eigene Darstellung

<sup>282</sup> Quelle: ebda.

<sup>283</sup> Vgl. MathWorks, <https://de.mathworks.com/help/stats/choose-a-classifier.html> (Zugriff: 21.02.2017)

Tabelle 15: Übersicht über Methoden in Matlab (1)<sup>284</sup>

Methoden	Vorhersage- geschwindigkeit	Benötigter Speicherplatz	Interpretier- barkeit	Flexibilität des Modells
Complex Tree	Schnell	Klein	Leicht	Hoch; viele Blätter für eine sehr feine Unterscheidungen zwischen den Klassen (maximale Anzahl der Splits = 100)
Medium Tree	Schnell	Klein	Leicht	Medium; mittlere Anzahl Blätter zur feineren Unterscheidung zwischen Klassen (maximale Anzahl der Splits = 20)
Simple Tree	Schnell	Klein	Leicht	Gering; nur wenige Blätter, um grobe Unterscheidungen zwischen den Klassen zu machen (maximale Anzahl der Splits = 4)
Linear Discriminant	Schnell	Klein	Leicht	Gering; erschafft lineare Grenzen zwischen Klassen
Quadratic Discriminant	Schnell	Groß	Leicht	Gering; erschafft nicht-lineare Grenzen zwischen Klassen (elliptisch, parabolisch oder hyperbolisch)
Linear SVM	binär: schnell; multiclass: medium	Medium	Leicht	Gering; einfache lineare Trennung zwischen den Klassen
Quadratic SVM	binär: schnell; multiclass: langsam	binär: Medium; multiclass: groß	Schwer	Medium
Cubic SVM	binär: schnell; multiclass: langsam	binär: Medium; multiclass: groß	Schwer	Medium
Fine Gaussian SVM	binär: schnell; multiclass: langsam	binär: Medium; multiclass: groß	Schwer	Hoch; Kerneleinstellung = $\sqrt{P}/4$ (P...Anzahl Prädiktoren)
Medium Gaussian SVM	binär: schnell; multiclass: langsam	binär: Medium; multiclass: groß	Schwer	Medium; Kerneleinstellung = $\sqrt{P}$ (P...Anzahl Prädiktoren)
Coarse Gaussian SVM	binär: schnell; multiclass: langsam	binär: Medium; multiclass: groß	Schwer	Gering; Kerneleinstellung = $\sqrt{P} \cdot 4$ (P...Anzahl Prädiktoren)

<sup>284</sup> Quelle: MathWorks, <https://de.mathworks.com/pricing-licensing.html?prodcode=ML> (Zugriff: 13.02.2017) (leicht modifiziert)

Für das gegebene Fallbeispiel sind die Geschwindigkeit der Klassifikation und der benötigte Speicherplatz von fast keiner Bedeutung, da die zu analysierenden Datenmenge mit 218 Datensätzen, einem zu klassierenden Attribut mit 15 Klassen und mehr als 60 Prädiktoren eher gering gegenüber der Testmenge ( $\leq 7.000$  Datensätze, 80 Prädiktoren, 50 Klassen) ist, auf der diese Angaben beruhen. Zur Vollständigkeit werden sie dennoch angegeben. Bei der Geschwindigkeit wird in Schnell (0,01 Sekunden), Medium (1 Sekunde) und Langsam (100 Sekunden) unterschieden. Der benötigte Speicherplatz wird in Klein (1 MB), Medium (4 MB) und Groß (100 MB) eingeteilt.<sup>285</sup>

**Tabelle 16: Übersicht über Methoden in Matlab (2)<sup>286</sup>**

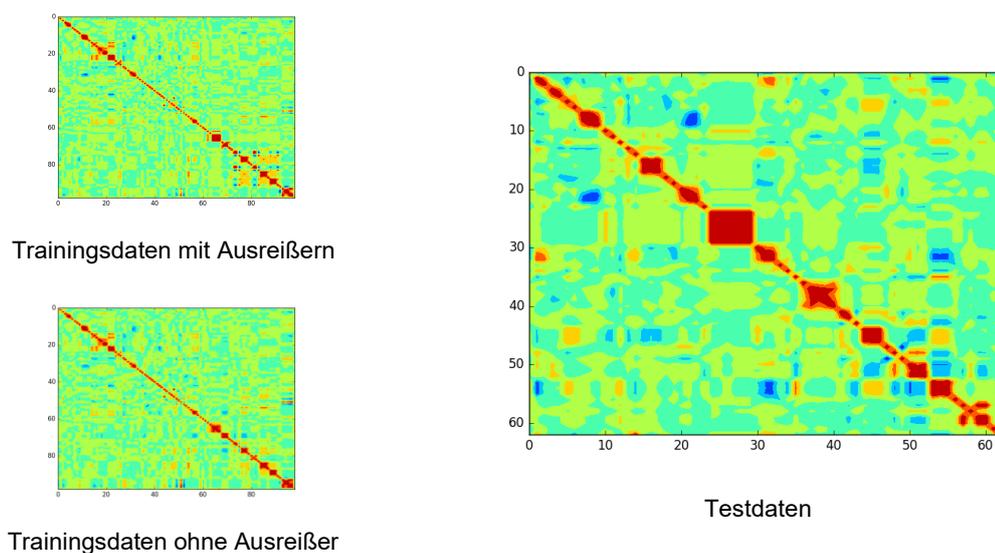
Methode	Vorhersagegeschwindigkeit	Benötigter Speicherplatz	Interpretierbarkeit	Flexibilität des Modells
Fine KNN	Medium	Medium	Schwer	Feine detaillierte Unterscheidung zwischen Klassen. Anzahl der Nachbarn = 1
Medium KNN	Medium	Medium	Schwer	Mittelfeine Unterscheidung zwischen Klassen. Anzahl der Nachbarn = 10
Coarse KNN	Medium	Medium	Schwer	Grobe Unterscheidung zwischen Klassen. Anzahl der Nachbarn = 100
Cosine KNN	Medium	Medium	Schwer	Mittelfeine Unterscheidung zwischen Klassen unter Verwendung einer Cosinus-Distanz-Metrik Anzahl der Nachbarn = 10
Cubic KNN	Langsam	Medium	Schwer	Mittelfeine Unterscheidung zwischen Klassen unter Verwendung einer Kubische-Distanz-Metrik Anzahl der Nachbarn = 10
Weighted KNN	Medium	Medium	Schwer	Mittelfeine Unterscheidung zwischen Klassen unter Verwendung der gewichteten Distanz Anzahl der Nachbarn = 10
Boosted Trees	Schnell	Klein	Schwer	Medium bis Hoch - Methode: AdaBoost mit Decision Tree
Bagged Trees	Medium	Groß	Schwer	Hoch
Subspace Discriminant	Medium	Klein	Schwer	Medium
Subspace KNN	Medium	Medium	Schwer	Medium
RUSBoosted Trees	Schnell	Klein	Schwer	Medium - gut für schiefe Daten

<sup>285</sup> Vgl. MathWorks, <https://de.mathworks.com/help/stats/choose-a-classifier.html> (Zugriff: 21.02.2017)

<sup>286</sup> Quelle: MathWorks, <https://de.mathworks.com/pricing-licensing.html?prodcode=ML> (Zugriff: 13.02.2017) (leicht modifiziert)

Die Überlegung bei der Anwendung der App ist, eine erste Abschätzung über die Sinnhaftigkeit des Einsatzes der einzelnen Methoden ohne aufwendige Implementierung treffen zu können. Zu diesem Zweck werden verschiedene Datensets in einem ersten Schritt getestet, und die Ergebnisse miteinander verglichen. In einem zweiten Schritt werden zwei Datensets ausgewählt und mit denen alle zur Verfügung stehenden Modelle trainiert werden. Diese Modelle werden anschließend auf die Testdaten angewendet.

Bei den Testdaten handelt es sich um die 78 Datensätze, die aus Zeitraum von September 2016 bis Anfang Jänner 2017 stammen und nicht in den Trainingsdaten enthalten sind. Die Aufbereitung der Testdaten wird auf das Zusammenführen der Daten (siehe Kapitel 4.2.2) und das Standardisieren beschränkt. Zusätzlich wird das Zielattribut R25 Vergl.% entfernt.



**Abbildung 41: Korrelationsmatrix der VM+NM<sup>287</sup>**

Da im praktischen Einsatz das Modell nur auf einen einzigen Datensatz angewendet wird, wird auf eine umfassende Auswertung verzichtet. Einzig die Korrelationsmatrix für die Trainingsdaten wird berechnet. In Abbildung 41 sind sowohl diese als auch die bereits bekannten Korrelationsmatrizen auf Basis der Trainingsdaten aus Abbildung 36 bzw. Abbildung 40 zu sehen. Aus dem direkten Vergleich wird ersichtlich, dass die Testdaten anders korrelieren als in die Trainingsdaten, wobei dies für die Daten mit und ohne Ausreißer gilt. Diese Erkenntnisse bilden die Grundlage für den Zweifel an der Generalisierbarkeit der auf den Trainingsdaten basierenden Modelle.

<sup>287</sup> Quelle: eigene Darstellung

**Tabelle 17: Auswertung der Testdaten (Anzahl je Klasse)<sup>288</sup>**

Bezeichnung	Anzahl Datensets je Klasse
88.62	0
94.87	1
96.94	4
98.03	3
98.78	2
99.52	4
101.05	9
102.92	15
104.64	18
106.61	11
107.8	2
108.57	3
110.03	2
112.00	4
114.07	0
116.83	0
119.52	0

Zur späteren Überprüfung der Genauigkeit der Modelle werden die Ausprägungen des Zielattributs in die in Tabelle 10 angeführten Klassen eingeordnet. Dies ist eine Voraussetzung für einen Vergleich der vorhergesagten und tatsächlichen bzw. Ist-Werten. In Tabelle 17 ist eine Aufschlüsselung zu sehen, wie sich die 78 Datensätze auf die 15 Klassen verteilen. Da es zwischen den Klassen Intervalle gibt, die keiner Klasse zugeordnet sind, ist es möglich, dass die in den Testdaten vorhandenen Werte zwischen zwei Klassen liegen. Sie sind in Tabelle 17 rot hinterlegt und müssen als Grenzfälle in der Auswertung gesondert angeführt werden.

### Vergleich verschiedener Modelle und Datensets

Für die Daten der Vor- und Nachmahlung liegen bis zu diesem Punkt drei Datensets vor: beide durch den KMO-Test verkleinerten Datensets: jener mit Ausreißern, jener ohne Ausreißer und das ursprüngliche Datenset. Zusätzlich wird die Überlegung seitens des Unternehmens eingebracht, dass das Attribut UD, der Ofen, in dem die PTC-Bauteile gesintert werden, Einfluss auf das Ergebnis hat. Ebenso gibt es die Möglichkeit die Anzahl der Anteile, die bei der Kreuzvalidierung (siehe Kapitel 3.5.1) gebildet werden, zu beeinflussen. Die PCA kann als Option, die bei der Modellerstellung verwendet wird, gewählt werden. Die Varianz wurde hierbei auf 0,95 gesetzt, siehe Kapitel 3.4.4. Aus diesen Kombinationsmöglichkeiten ergaben sich die Datensets in Tabelle 18, die getestet wurden.

<sup>288</sup> Quelle: eigene Darstellung

**Tabelle 18: Kombinationen für VM+NM<sup>289</sup>**

	<b>Ausreißer</b>	<b>Kreuz- validierung</b>	<b>PCA</b>	<b>UD</b>	<b>Anzahl Prädiktoren</b>
5-fold PCA inkl. Ausreißer (Prädiktoren: 98)	mit	5	ja	nein	98
5-fold PCA inkl. Ausreißer + UD (Prädiktoren: 99)	mit	5	ja	ja	99
5-fold PCA inkl. Ausreißer (Prädiktoren: 62)	mit	5	ja	nein	62
10-fold PCA inkl. Ausreißer (Prädiktoren: 62)	mit	10	ja	nein	62
5-fold PCA exkl. Ausreißer (Prädiktoren: 83)	ohne	5	ja	nein	83
10-fold PCA exkl. Ausreißer (Prädiktoren: 83)	ohne	10	ja	nein	83
5-fold inkl. Ausreißer (Prädiktoren: 62)	mit	5	nein	nein	62
5-fold inkl. Ausreißer + UD (Prädiktoren: 63)	mit	5	nein	ja	63

Die in Tabelle 18 angegebenen Datensets wurden als Trainingsdaten für die Modellerstellung verwendet. In Abbildung 42 sind die Ergebnisse zu sehen. Die angegebenen Prozent auf der Abszisse geben die Genauigkeit des Modells an (siehe Kapitel 3.5.1). Auf der Ordinate sind die in Tabelle 15 und Tabelle 16 vorgestellten Methoden aufgetragen. Die Ergebnisse schwanken in einem Bereich von 7 % bis 23 %. Die Genauigkeit von 0 % steht für die Nicht-Anwendung der Methode. Ursache hierfür ist das Skalenniveau der Daten, das sich durch das Einbeziehen des Attributs UD ändert (siehe Tabelle 7).

<sup>289</sup> Quelle: eigene Darstellung

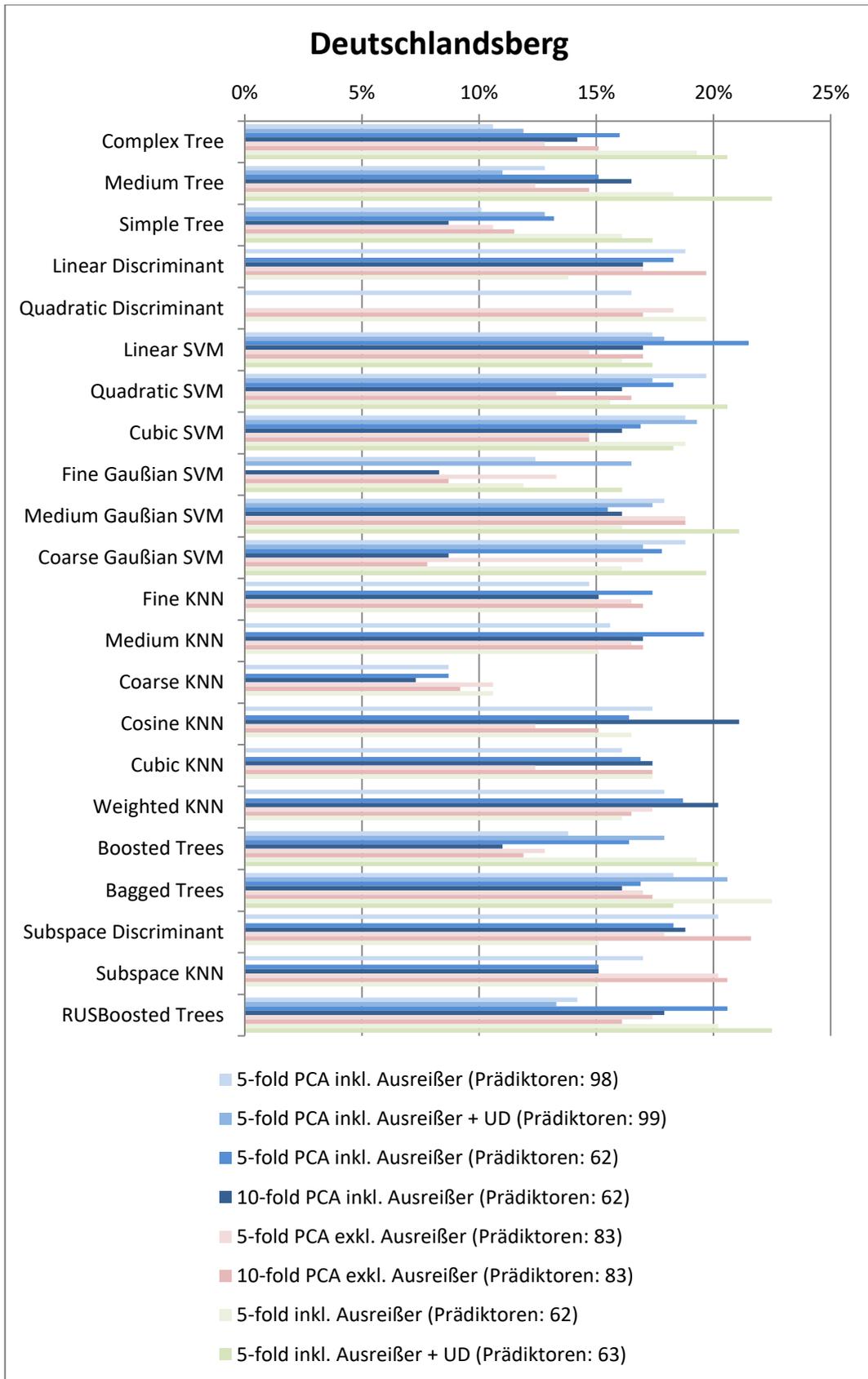


Abbildung 42: Vergleich verschiedener Kombinationsmöglichkeiten<sup>290</sup>

<sup>290</sup> Quelle: eigene Darstellung

Durch diese erste Abschätzung werden folgende Behauptungen aufgestellt:

- Die PCA trägt nicht unbedingt zu einer höheren Genauigkeit bei. Auch die Rechengeschwindigkeit wird nicht wesentlich verbessert, da es sich um so wenige Datensätze handelt.
- Das Attribut UD verändert die Genauigkeit nicht wesentlich; allerdings können nicht alle Methoden verwendet werden.
- Der Unterschied zwischen dem Datenset mit und ohne Ausreißer ist nicht signifikant.
- Die Anzahl der Kreuzvalidierungsteilmengen ist hier nicht ausschlaggebend.

Zusätzlich stellte sich die Frage, ob die beiden Datensets mit und ohne Ausreißer tatsächlich miteinander verglichen werden können, da eine unterschiedliche Anzahl von Attributen vorliegt. Ebenso ist zu bedenken, dass ein neuerliches Durchführen der gleichen Methode zu Modellen mit unterschiedlichen Genauigkeiten führen kann.

Die Abschätzung ist somit nur bedingt maßgebend für die Auswahl der endgültigen Trainingsdaten. Es wird die Entscheidung getroffen, auf die PCA zu verzichten. Dies liegt einerseits daran, dass ähnliche Ergebnisse bei Anwendung bzw. Nicht-Anwendung erzielt wurden, andererseits daran, dass die in Kapitel 3.4.4 empfohlene Untergrenze unterschritten wird. Eine 5-fache Kreuzvalidierung für die Trainingsdaten wird verwendet, um ein mögliches Overfitting zu vermeiden, und es wird mit den 62 Attributen aus dem Datenset mit Ausreißern und dem Attribut UD als Prädiktoren weitergearbeitet. Prädiktoren sind jene Attribute, auf denen die Modelle basieren. Zusätzlich wird das Datenset ohne Ausreißer auf diese 63 Prädiktoren beschränkt, um eine bessere Vergleichbarkeit der Ergebnisse zu gewährleisten. Mit diesen beiden Datensets werden Modelle mit den zur Verfügung stehenden Methoden erstellt und mit den Testdaten getestet. Diese Ergebnisse werden im Folgenden vorgestellt.

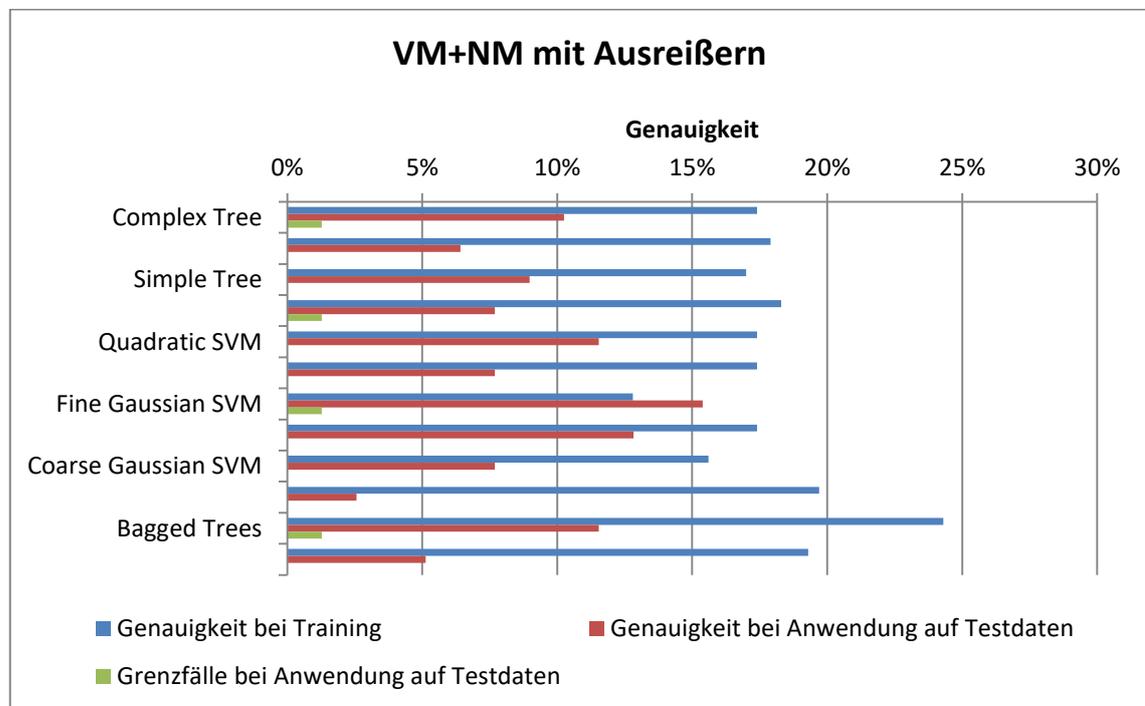
### **Ergebnisse des Datenset mit Ausreißern**

Als erster Schritt werden erneut sämtliche Modelle trainiert, deren Methoden Daten in gemischten Skalenniveaus verarbeiten können. Dies führt zu den verschiedenen Methoden der Entscheidungsbäume (Complex Tree, Medium Tree, Simple Tree), den SVM (Linear SVM, Quadratic SVM, Cubic SVM, Fine Gaussian SVM, Medium Gaussian SVM, Coarse Gaussian SVM) und den Ensemblemethoden (Boosted Trees, Bagged Trees, RUSBoosted Trees), die auf den Entscheidungsbäumen basieren. Die Entscheidungsbäume unterscheiden sich im Detaillierungsgrad voneinander, die SVM in der Art der Kernels, die verwendet werden und die Ensemblemethoden in der Art des Kombinierens der einzelnen Bäume. Durch die Anwendung all dieser Methoden wird eine erste Aussage über die Genauigkeit der Modelle getroffen. Diese sind in Abbildung 43 durch die blauen Balken dargestellt.

Bei der Anwendung der Modelle auf die Testdaten wird die Klasse, zu denen ein Datensatz gehört, vorhergesagt. Diese vorhergesagte Klasse wird mit der Ist-Klasse verglichen. Sind diese dieselbe, wird die Vorhersage als Treffer gewertet. Durch die Anzahl der richtig vorhergesagten Klassen je Modell, die in Bezug zur Gesamtanzahl der Testdatensätze gesetzt wird, ergibt sich die Genauigkeit in Prozent je Modell. Diese werden durch die roten Balken in Abbildung 43 dargestellt.

Die grünen Balken in Abbildung 43 symbolisieren die Grenzfälle. Diese entstehen, wenn die tatsächliche Klasse nicht eindeutig zugeordnet werden kann, da der Messwert von  $R_{25}$  Vergl.% zwischen den zwei Klassen liegt. Bei der vorhergesagten Klasse handelt es sich um jene, die unmittelbar an diesen Wert angrenzen. Es kann daher nicht definitiv gesagt werden, ob richtig oder falsch zugeordnet wird. Grenzfälle, die einen größeren Abstand zwischen der vorhergesagten und tatsächlichen Klasse haben, werden nicht berücksichtigt. Daher resultieren die Schwankungen zwischen 0 und 1,3 %. Dieser Prozentsatz entspricht einem Grenzfall.

In Abbildung 43 ist zu sehen, dass die Genauigkeit bei der Anwendung auf die Testdaten im Durchschnitt 9 % niedriger als bei der Abschätzung nach dem Lernen. Ein Grund dafür könnte ein Overfitting der Methode bei der Kreuzvalidierung sein. Die Bagged Trees erreichen im Training die höchste Genauigkeit (24,3 %), und die Fine Gaussian SVM die niedrigste. Bei der Anwendung auf die Testdaten erzielen diese aber die höchste Genauigkeit. Die dem Diagramm zugrunde liegenden Daten sind in Tabelle 34 im Anhang aufgeführt.



**Abbildung 43: Ergebnisse des CL für VM+NM mit Ausreißern<sup>291</sup>**

Eine genaue Aufschlüsselung über die vorhergesagte und Ist-Klasse je Datensatz bzw. Charge und Methode ist in Abbildung 62 bis Abbildung 73 im Anhang zu finden. Im Durchschnitt ist die Abweichung von vorhergesagter zur Ist-Klasse 2,47 bis 3,09 Klassen (siehe Tabelle 34), wobei das Maximum von der quadratischen SVM erreicht wird. Dies ist interessant, da die Vorhersagegenauigkeit mit 11,5 % im oberen Drittel liegt. Ebenfalls ist zu bemerken, dass die Fine Gaussian SVM bei der Anwendung auf die Testdaten bessere Ergebnisse als für die Trainingsdaten liefert.

Die Interpretierbarkeit der Modelle wurde bereits in Tabelle 15 und Tabelle 16 angesprochen. Lediglich die Entscheidungsbäume und die Lineare SVM seien leicht zu

<sup>291</sup> Quelle: eigene Darstellung

interpretieren. Die Lineare SVM bezieht sich auf eine lineare Trennung der Datensätze im Raum. Da für die Platzierung dieser 63 Attribute betrachtet werden, ist es dennoch nicht möglich, diese graphisch darzustellen. Die Entscheidungsbäume können wie in Kapitel 0 als Graph oder Code dargestellt werden. Daher wird der Simple Tree als Beispiel für das Ergebnis einer solchen Methode verwendet und für alle Datensets graphisch dargestellt. Zusätzlich werden die Regeln, die der Medium Tree und Complex Tree aufgestellt hat, betrachtet.

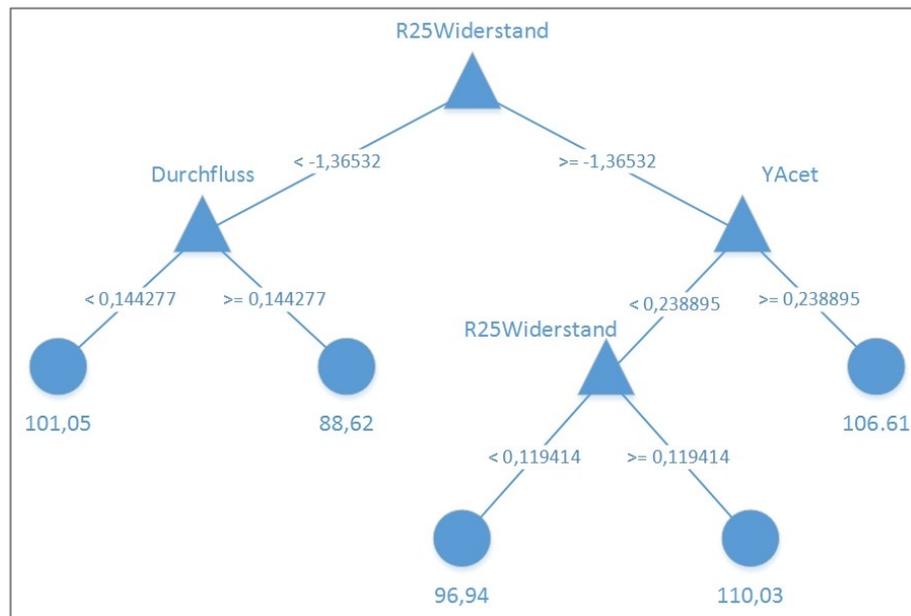


Abbildung 44: Simple Tree für Datenset VM+NM mit Ausreißern<sup>292</sup>

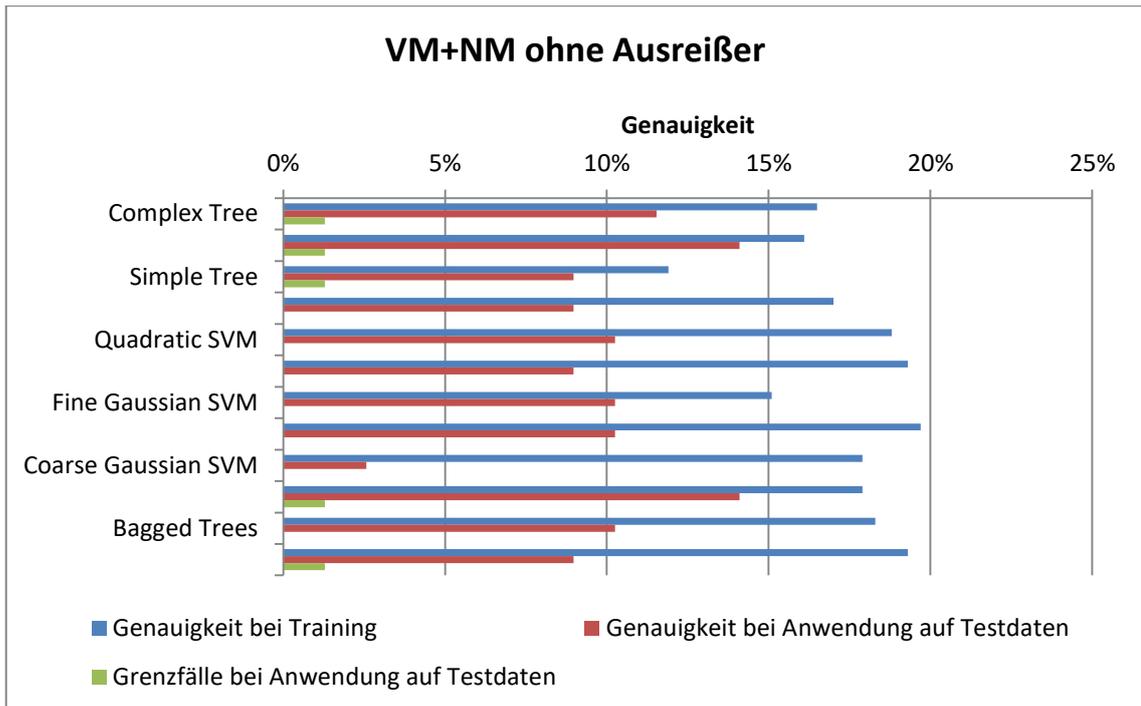
In Abbildung 44 ist der Simple Tree graphisch dargestellt. Es ist ersichtlich, dass die Anzahl der Blätter, also der Endpunkte, nur fünf beträgt. Dies bedeutet, dass nur ein Drittel der Klassen abgebildet werden und somit korrekt klassiert werden können. Grund dafür sind die Standardeinstellungen von Matlab, in denen die Anzahl der Teilungen angegeben werden kann. Das Modell ist somit trotz seiner Genauigkeit von 17 % in der Lernphase nicht für die Testdaten geeignet, da zwei der fünf Endknoten in den Daten nicht enthalten sind. Dieser Umstand wird auch in den Abbildung 64 deutlich, da der Linie der Vorhersage wenig differenziert. Zusätzlich ist die Verwendung für zukünftige Datensätze nicht sinnvoll, da das Modell nicht alle Klassen abbildet.

Ab Seite w im Anhang sind sowohl die Regeln für den Medium Tree als auch die für den Complex Tree angegeben. Diese sind wesentlich detaillierter verzweigt, dennoch liefern sie nur bedingt bessere Ergebnisse.

### Ergebnisse des Datensets ohne Ausreißer

Wie bereits bei den Ergebnissen des Datensets mit Ausreißern besprochen, werden sämtliche Methoden auf die Testdaten angewendet. Die Ergebnisse sind nach dem gleichen Prinzip aufbereitet worden und in Abbildung 45 dargestellt.

<sup>292</sup> Quelle: eigene Darstellung



**Abbildung 45: Ergebnisse des CL für VM+NM ohne Ausreißer<sup>293</sup>**

In Abbildung 45 sind die Ergebnisse des Datensets ohne Ausreißer zu sehen. Sie sind schlechter als für das Datenset mit Ausreißern. Dennoch weicht die Genauigkeit bei der Anwendung mit einem Durchschnitt von 7,5 % weniger vom Ergebnis des Trainings ab. Der durchschnittliche Unterschied zwischen vorhergesagter und Ist-Klasse ist mit 2,06 bis 3,37 Klassen sehr viel größer. Der Trend vom vorher betrachteten Datenset, dass die SVM bessere Ergebnisse liefern, setzt sich bei diesem Test nicht fort und in diesem Fall erzielen die Bäume bzw. die Ensemblemethoden bessere Ergebnisse. Eine genaue Aufschlüsselung ist der Tabelle 35 zu entnehmen. Zusätzlich werden die einzelnen Ergebnisse in den Abbildung 74 bis Abbildung 85 im Anhang über die Chargin als Diagramm dargestellt.

Die andere Datenstruktur und Korrelationsmatrix der Datensets ohne Ausreißer hat ebenso einen anderen Entscheidungsbaum in Abbildung 46 als Ergebnis. Es werden andere Attribute für den Baum gewählt oder andere Grenzwerte für die gleichen Attribute eingesetzt. Ebenso werden andere Klassen als Blätter ausgewählt, wovon eine nicht in den Testdaten vorkommt. Die schlechten Ergebnisse bei der Anwendung des Simple Trees auf die Testdaten sind wieder durch den Umstand der nicht vollständigen Abbildung aller Klassen zu erklären. Die Regeln für den Medium und Complex Tree, die bei diesem Datenset deutlich besser abschnitten und vor allem der Medium Tree die maximale Genauigkeit dieses Datensets liefert, sind im Anhang ab Seite gg zu finden.

<sup>293</sup> Quelle: eigene Darstellung

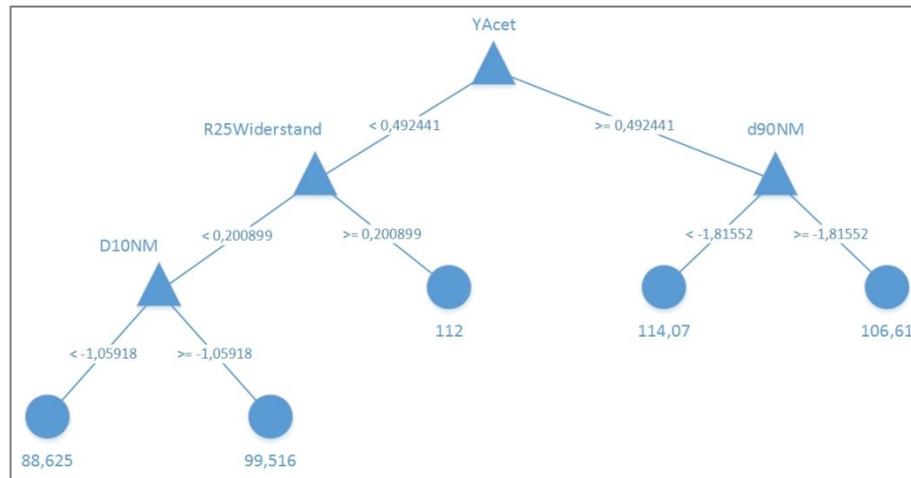


Abbildung 46: Simple Tree für Datenset VM+NM ohne Ausreißer<sup>294</sup>

Die Ergebnisse des Deutschlandsberger Datensets sind auf Grund der geringen Genauigkeit nicht zufriedenstellend. In einem weiteren Schritt werden die in diesem Kapitel entstandenen Datensets als Grundlage für die Erstellung von Modellen für Šumperk bzw. einer Datenbasis für diese verwendet. Die Vorgehensweise wird im nächsten Kapitel erläutert.

### 4.3 Standort Šumperk

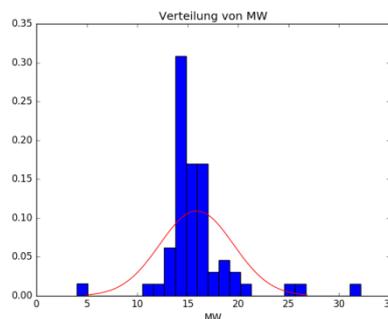
In Šumperk wird die Sinterfreigabe betrachtet. Diese findet unmittelbar vor und während des Sinterprozesses im für die Produktion eingesetzten Durchstoßofen statt und bildet die Grundlage für eine eventuelle Änderung der Prozesseinstellungen, um die Produktspezifikationen u. a. hinsichtlich  $R_{25}$  zu treffen. Die PTC-Bauteile haben, anders als bei der Granulatfreigabe, die für diesen Produkttypen vorgesehene Geometrie. Durch diese anderen Maße bzw. andere Oberfläche gibt es Änderungen zwischen dem Verhalten bei der Sinterung für die Granulatfreigabe und der Sinterung für Sinterfreigabe. Daher kann es zu anderen Korrelationskoeffizienten zwischen den Attributen und dem neuen Zielattribut, dem gesinterten Widerstand, als zwischen den Attributen und  $R_{25}$  Vergl.% kommen. Das Ziel ist es, den gesinterten Widerstand vorherzusagen, und somit die Sinterfreigabe nicht mehr durchführen zu müssen. Um ein neues Modell, das auf diesem Ziel und diesen Zusammenhängen beruht, erstellen zu können, wird im Folgenden die Aufbereitung des Zielattributs erklärt. Anschließend wird das Schaffen einer neuen Datenbasis für die Modellerstellung beschrieben und die entstandenen Modelle ausgewertet.

<sup>294</sup> Quelle: eigene Darstellung

### 4.3.1 Aufbereitung des Zielattributs

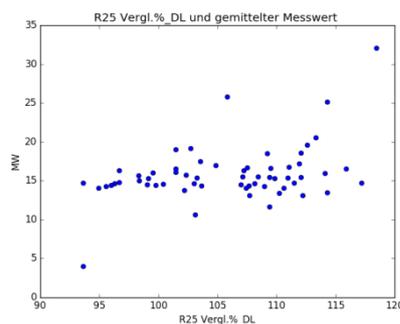
Das Zielattribut, der gesinterte Widerstand, wird für die Sinterfreigabe aufgezeichnet. Dies ist eine laufende Überprüfung während des Sinterprozesses. Je Freigabemessung werden 15 Bauteile gemessen, die je Charge aufgezeichnet und mit einer laufenden Nummer versehen werden. Für die Erstellung des Modells und das Ziel die Sinterfreigabe obsolet werden zu lassen, ist nur die erste Messung je Charge zu betrachten. Aus diesem Grund liegen je Charge bzw. Datensatz 15 Messwerte für den gesinterten Widerstand vor. Dieser gesinterte Widerstand hat den Sollwert von 15,01.

Als erstes werden diese 15 Messwerte von Ausreißern bereinigt. Es werden jene Messwerte als Ausreißer bezeichnet, die den Wert „-99“ (Messung ungültig) oder einen Wert größer als 50 haben. Da nur ein Wert vorhergesagt werden soll, wird aus diesen 15 Messwerten je Datensatz das arithmetische Mittel gebildet. Aus diesem Grund wird das Zielattribut, der gemittelte gesinterte Widerstand, im Weiteren als MW bezeichnet. Die Verteilung von MW ist in Abbildung 47 zusehen.



**Abbildung 47: Histogramm für gemittelten gesinterten Widerstand (MW)<sup>295</sup>**

Die Korrelation von MW mit dem Zielattribut R25 Vergl.% des Kapitels 4.2 ist in Abbildung 48 zu sehen und hat einen Korrelationskoeffizienten von 0,39. Es liegt eine positive Korrelation vor (siehe Kapitel 3.3.3), die im mittleren Ausprägungsbereich liegt.



**Abbildung 48: Streudiagramm für R25 Vergl.% und MW<sup>296</sup>**

Der MW wird keiner weiteren Ausreißeranalyse unterzogen, obwohl die Notwendigkeit in Abbildung 47 deutlich wird. Hierbei wird an der Überlegung aus Kapitel 4.2.4 und

<sup>295</sup> Quelle: eigene Darstellung

<sup>296</sup> Quelle: ebda.

Kapitel 4.2.5 festgehalten, dass auf Basis der unveränderten Ausprägungen des Zielattributs die Gruppierung vorgenommen wird.

**Tabelle 19: Einteilung von MW in Klassen<sup>297</sup>**

Bezeichnung	Klassenmitte	Untergrenze	Obergrenze	Anzahl von Datensätzen pro Klasse	Breite
11.66375	8.9	4.03	13.77	8	9.74
14.24125	14.25	14.08	14.42	8	0.34
14.58875	14.58	14.45	14.71	8	0.26
15.08875	15.07	14.74	15.4	8	0.66
15.68	15.75	15.45	16.05	8	0.6
16.5	16.44	16.08	16.8	8	0.72
18.33625	18.31	16.97	19.65	8	2.68
25.9075	26.335	20.56	32.11	4	11.55

Die in Tabelle 19 angeführten Klassen beruhen wie die Einteilung in Kapitel 4.2.5 auf (3.1) und der Tatsache, dass 60 Datensätze vorhanden sind:

$$k_{\text{Šumperk}} = \sqrt{n_{\text{Šumperk}}} = \sqrt{60} \approx 8$$

Nach dieser Gruppierung und Einteilung werden das Zielattribut MW und die Deutschlandsberger Daten der Vor- und Nachmahlung zusammengeführt.

### 4.3.2 Datenzusammenführung

Die Deutschlandsberger Daten dienen als Grundlage für die Erstellung eines Datensets für Šumperk. In Abbildung 31 wird der Aufbau der Ausgangsdaten beschrieben. Durch die Zusammenführung der Vor- und Nachmahlung stellt sich dieser Aufbau für Šumperk wie in Abbildung 49 dar. Einerseits ist es möglich, nur jene 62 Attribute (ohne UD und R25 Vergl.%) zu verwenden, die in Kapitel 4.2.7 die Basis für die Modellierung sind. Andererseits kann mit sämtlichen Attributen, die nach der Zusammenführung der Daten in Kapitel 4.2.2 (angeführt in Tabelle 29), ein erneuter KMO-Test durchgeführt werden.

<sup>297</sup> Quelle: eigene Darstellung

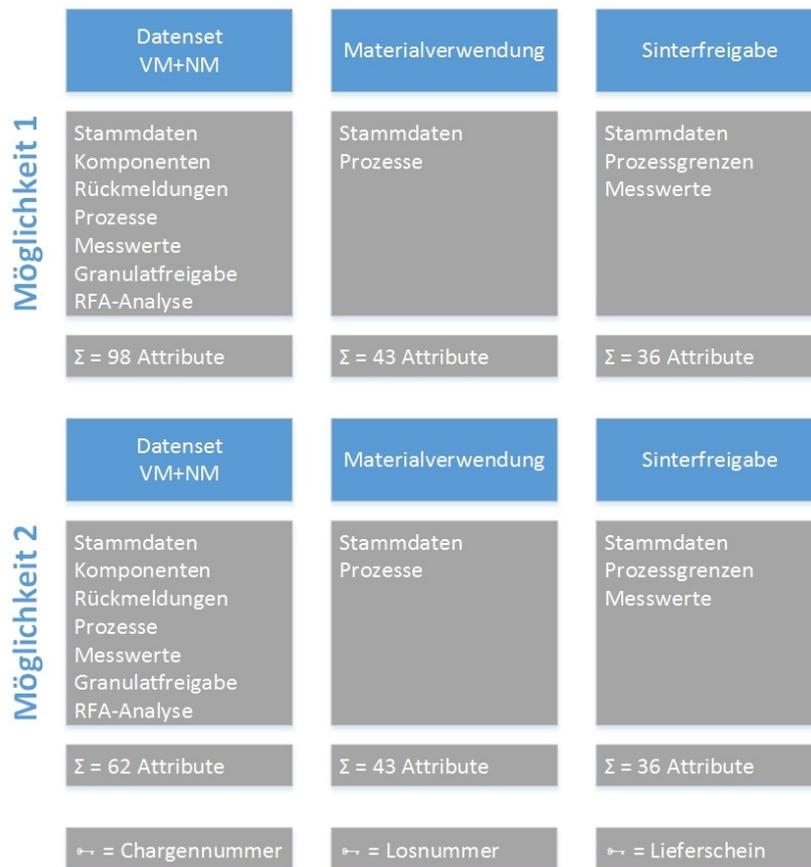


Abbildung 49: Aufbau der Ausgangsdaten für Šumperk<sup>298</sup>

Ebenso muss beachtet werden, ob die Datensets mit oder ohne Ausreißer verwendet werden. Zu diesem Zeitpunkt wird entschieden, dass die Datensets mit Ausreißern verwendet werden, da bis zu vier Chargen aus der Nachmahlung in eine Charge für die Sinterfreigabe einfließen und die Ausreißer dadurch ausgeglichen werden könnten. Die Chargen- und Losnummern werden im Datenset Materialverwendung miteinander verknüpft. Diese Beziehungen werden verwendet, um das arithmetische Mittel aller 98 bzw. 62 Attribute der einfließenden Chargen, das über die Gewichtsanteile gewichtet wird, zu berechnen. Zusätzlich wird der Ofen, auf dem die Freigabe stattgefunden hat, als Attribut aufgenommen.

Die entstandenen Datensets bestehen aus einer veränderlichen Anzahl von Attributen (64 bzw. 100) und 60 Datensätzen. Sämtliche Attribute, die in diesem Datenset das Suffix „\_DL“ in ihrer Bezeichnung (siehe Tabelle 36) enthalten, sind jenen im Datenset der Vor- und Nachmahlung ohne Suffix gleichzusetzen. Im nächsten Schritt werden die neuentstandenen Datensets ausgewertet.

<sup>298</sup> Quelle: eigene Darstellung

### 4.3.3 Auswertung

Wie bereits in Kapitel 4.2.3 wird das Datenset hinsichtlich der statistischen Maße aus Kapitel 3.3.2 und 3.3.3 ausgewertet. Dies wurde zuerst am Datenset mit Ausreißern durchgeführt. Der Vollständigkeit halber wird hier bereits die Auswertung des Datensets ohne Ausreißer, das durch das Zusammenführen der Daten der Vor- und Nachmahlung ohne Ausreißer entstanden ist, angeführt. Allerdings wird erst zu einem späteren Zeitpunkt in dieser Arbeit entschieden, dieses zu verwenden. Auf eine explizite Auswertung des Datensets mit 64 Attributen und 60 Datensätzen wird verzichtet, da dieses ein Subset des 100 Attribute/60 Datensätze-Datensets mit Ausreißern darstellt. Allerdings wird der Ofen nicht mit einbezogen, da dieser ein nicht-metrisches Attribut darstellt und die Maße daher nicht berechnet werden können (siehe Kapitel 3.3.2 und Kapitel 3.3.3).

#### Datenset mit Ausreißern

In diese Auswertung fließen 99 Attribute und 60 Datensätze ein, deren Grundlage das Deutschlandsberger Datenset mit Ausreißern ist.

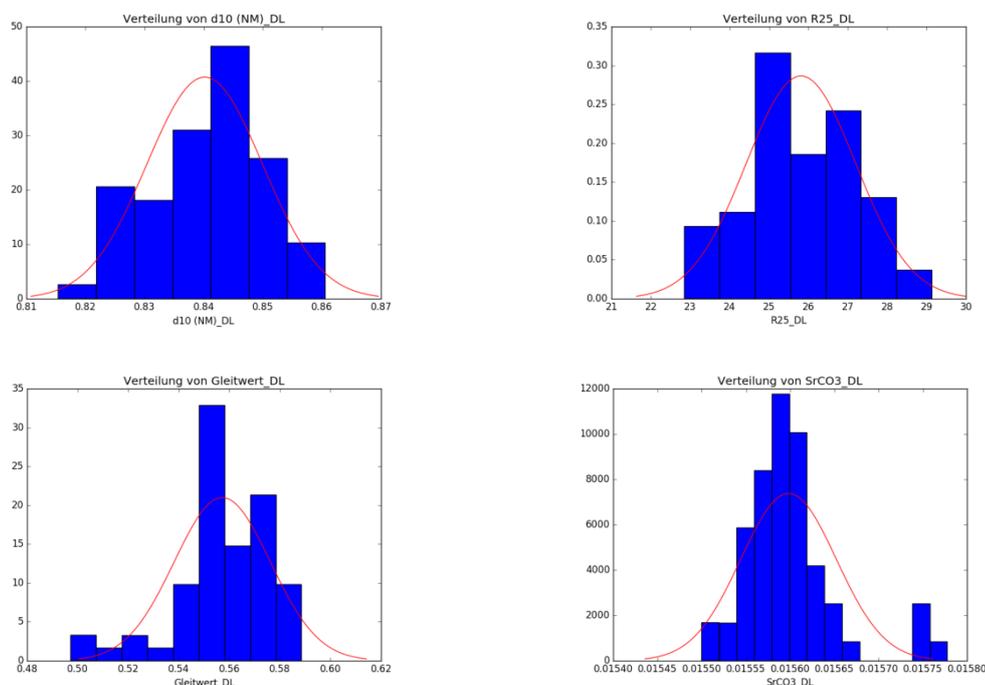
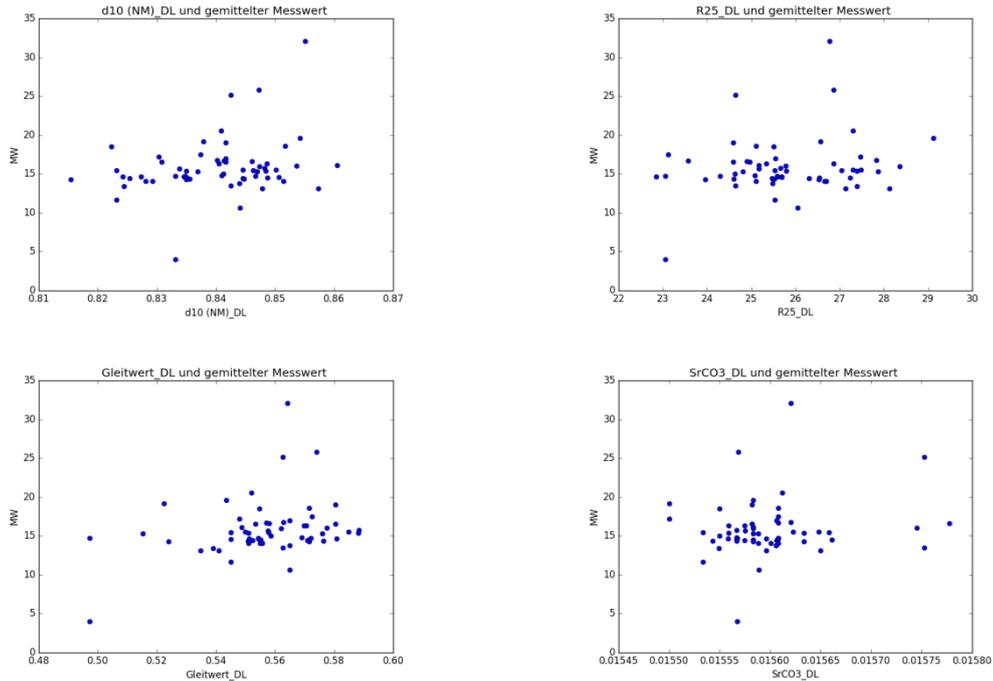


Abbildung 50: Beispiele für Histogramme für Šumperk mit Ausreißern<sup>299</sup>

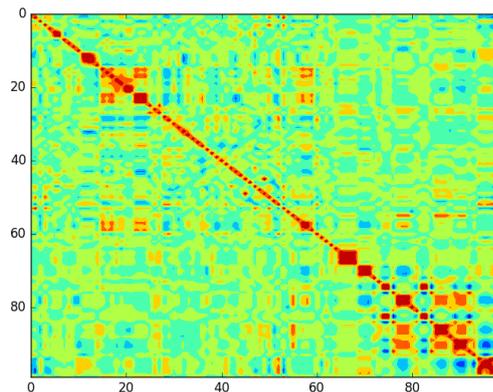
Aufgrund der geringeren Anzahl von Datensätzen verändert sich die Verteilung der Attribute. In Abbildung 50 werden wieder die vier gleichen Attribute wie in Abbildung 33 und Abbildung 37 dargestellt. Die Ausreißer sind trotz der Mittelwertbildung vor allem im Gleitwert\_DL und SrCO3\_DL wiederzufinden. Ob es sich um normalverteilte Daten handelt, wird nicht explizit mittels Kolmogorov-Smirnov-Test überprüft. Es werden lediglich die im Anhang in Tabelle 36 angeführte Wölbung und Schiefe betrachtet. Diese lassen den Schluss zu, dass die Daten nicht normalverteilt sind.

<sup>299</sup> Quelle: eigene Darstellung



**Abbildung 51: Beispiele für Streudiagramme für Šumperk mit Ausreißern<sup>300</sup>**

Die Zusammenhänge zwischen MW und den ausgewählten Attributen sind in den Streudiagrammen in Abbildung 51 zu sehen. Die Beziehungen zwischen den MW und den jeweiligen Attributen ist nicht so klar ersichtlich wie in Abbildung 35 und Abbildung 39. Dies liegt u. a. an der geringeren Anzahl von Datensätzen und der nicht stark ausgeprägten Korrelation von R25 Vergl.% und MW. Bei der Betrachtung der Korrelationskoeffizienten in Tabelle 36 wird klar, dass zwischen Gleitwert\_DL und MW die am stärksten ausgeprägte Korrelation mit 0,29 innerhalb dieser vier Attribute herrscht.



**Abbildung 52: Korrelationsmatrix (graphisch) für Šumperk mit Ausreißern<sup>301</sup>**

Bei der Betrachtung der Korrelationsmatrix sind Parallelen zu Abbildung 36 zu erkennen. Die Reihenfolge der Attribute ist zu Tabelle 29 bzw. Abbildung 36

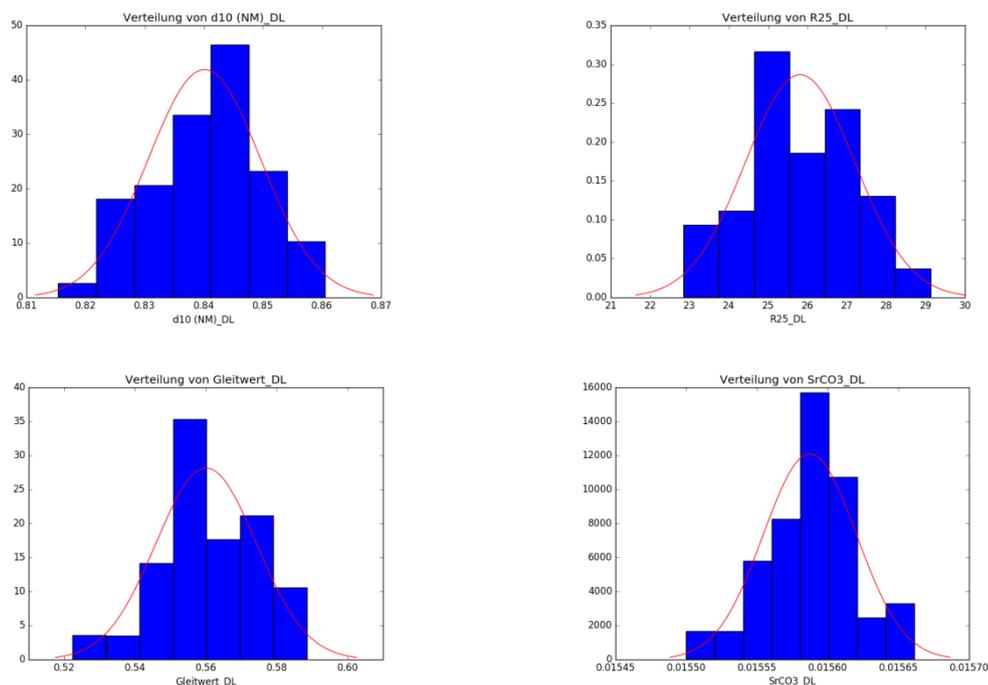
<sup>300</sup> Quelle: eigene Darstellung

<sup>301</sup> Quelle: ebda.

unverändert, einzig das Attribut MW steht nun an erste Stelle (siehe Tabelle 36). Hier bestätigt sich, dass die Trends aus dem Datenset der Vor- und Nachmahlung weiter geführt werden (siehe Kapitel 4.2.3). Zusätzlich gibt es eine höhere Anzahl negativer Korrelationskoeffizienten (blau dargestellt). Die Korrelation um Punkt (20,20) sind jene Attribute, die bei der Granulataufgabe erhoben werden. Diese korrelieren ebenso um den Punkt (60,20) mit den Attributen d50 (VM)\_VM\_DL, d90 (VM)\_VM\_DL und pH Wert Schlicker\_VM\_DL, wobei die ersteren beiden auch hoch miteinander korrelieren. Um den Punkt (68,68) sind wieder die Zusammenhänge zwischen den Heizgruppen ersichtlich. Die diagonal anschließende dunkelrote Fläche steht für die Korrelation der statistischen Maße, die für den Differenzdruck aufgezeichnet wurden. Dieser Trend setzt sich in den nachfolgenden diagonalen auffälligen Flächen fort, da hier jeweils die Maße Minimum, Maximum und Standardabweichung eines Attributs angegeben wurde. Die Korrelation dieser Attribute untereinander ist sehr viel negativer ausgeprägt als bei der Korrelationsmatrix des Datensets der Vor- und Nachmahlung in Abbildung 36.

### Datenset ohne Ausreißer

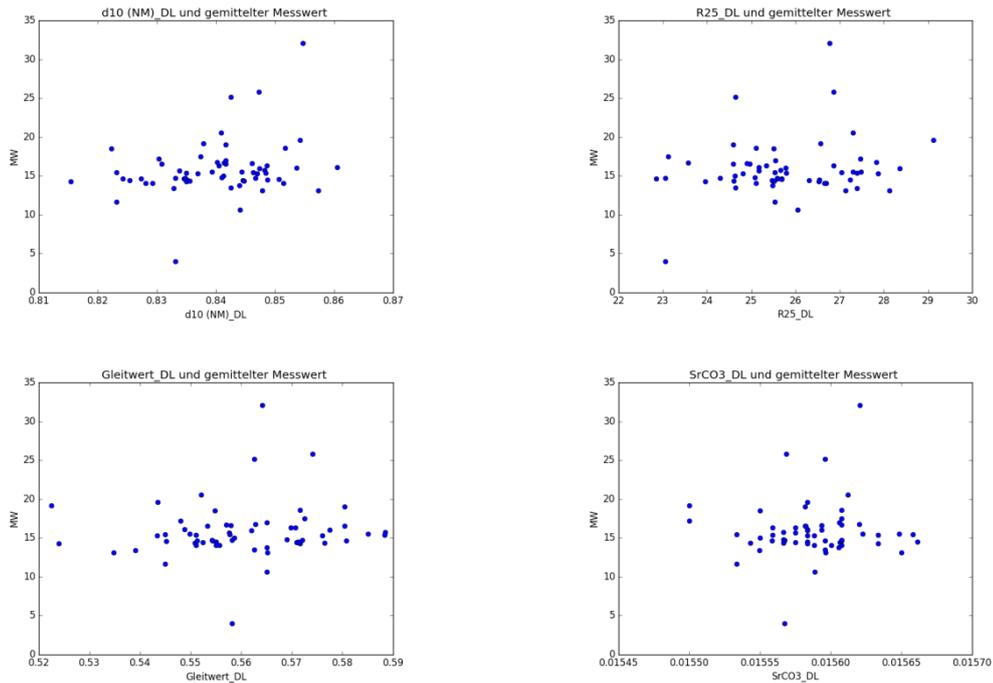
In die folgende Betrachtung fließen 100 Attribute und 60 Datensätze ein, deren Grundlage das Deutschlandsberger Datenset ohne Ausreißer ist.



**Abbildung 53: Beispiele für Histogramme für Šumperk ohne Ausreißer<sup>302</sup>**

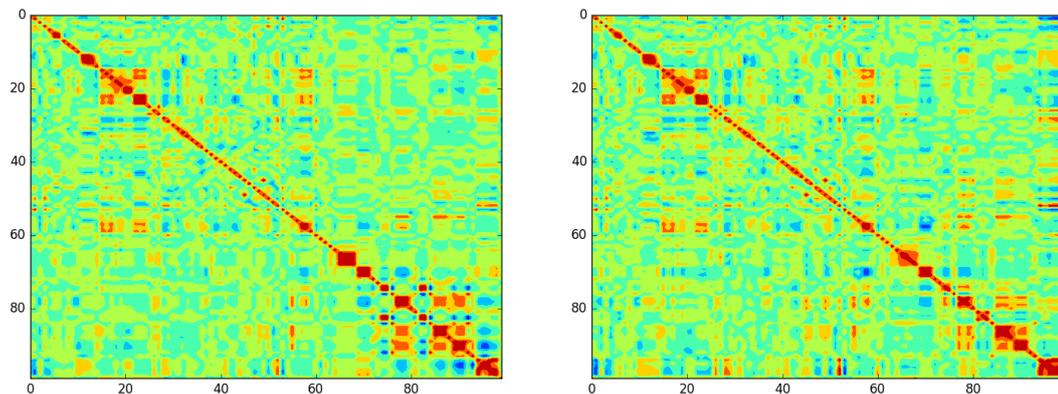
In Abbildung 53 sind die Histogramme des neu entstandenen Datensets zu sehen. Im Vergleich mit Abbildung 50 ist eine Veränderung für die Attribute Gleitwert\_DL und SrCO3\_DL ersichtlich. Die Ursache hierfür ist die Eliminierung der Ausreißer in Kapitel 4.2.4. Gleiches gilt für die in Abbildung 54 zu sehenden Streudiagramme.

<sup>302</sup> Quelle: eigene Darstellung



**Abbildung 54: Beispiele für Streudiagramme für Šumperk ohne Ausreißer<sup>303</sup>**

Die Änderung der Streudiagramme in Abbildung 54 ist ebenso in der Korrelationsmatrix in Abbildung 55 wiederzufinden. Hier ist die links die Korrelationsmatrix mit Ausreißern aus Abbildung 52 und rechts die Korrelationsmatrix ohne Ausreißer zu dargestellt. Die Veränderung der Korrelationsmatrizen in Abbildung 55 folgen demselben Muster wie jene in Abbildung 40. Die ausgeprägten Zusammenhänge in der rechten unteren Ecke werden aufgebrochen.



**Abbildung 55: Korrelationsmatrizen (graphisch) für Šumperk mit und ohne Ausreißer<sup>304</sup>**

Nach dieser kurzen Betrachtung der beiden Datensets wird die Eignung des Datensets mit Ausreißern für die PCA überprüft.

<sup>303</sup> Quelle: eigene Darstellung

<sup>304</sup> Quelle: ebda.

#### 4.3.4 Eignung für PCA

Die Vortests für die PCA wurden in Kapitel 3.4.4 theoretisch betrachtet und in Kapitel 4.2.6 am Beispiel der Deutschlandsberger Daten durchgeführt. Die Vorgehensweise stimmt mit jener in Kapitel 4.2.6 überein. Die Basis für den KMO-Test bildete das Datenset mit 100 Attributen/60 Datensätzen mit Ausreißern, wobei das Attribut Ofen ausgenommen wird, und wird nicht mit dem Datenset ohne Ausreißer durchgeführt. Der Grund hierfür liegt in der in Kapitel 4.2.7 gewonnenen Erkenntnis, dass die Ergebnisse noch schlechter vergleichbar sind, wenn unterschiedliche Attribute für die Modellerstellung herangezogen werden. Für den Zweck der KMO werden der Ofen und das Zielattribut MW ausgeschlossen.

Für das vorliegende Datenset wurde der KMO-Test viermal durchgeführt. Attribute mit  $MSA < 0,5$  wurden ausgeschlossen. Eine genaue Aufschlüsselung über die Werte sind der Tabelle 37 im Anhang zu entnehmen. Die Entwicklung über die Anzahl der Attribute ist in Tabelle 20 zusammengefasst.

**Tabelle 20: Entwicklung der Anzahl der Attribute (Šumperk mit Ausreißern)<sup>305</sup>**

	Durchgang			
	1	2	3	4
<b>Anzahl zu löschende Attribute</b>	61	28	1	0
<b>Anzahl Attribute nach Löschen</b>	38	10	9	9

Das MSA-Kriterium für das gesamte Datenset hat sich, wie in Tabelle 21 dargestellt, verändert. Wie bereits beim Datenset der Vor- und Nachmahlung ohne Ausreißer in Kapitel 4.2.6, ist es im ersten Durchgang nicht möglich, das MSA-Kriterium für das gesamte Datenset zu berechnen. Dies wird auf gleiche Art gelöst: Durch Löschen jener Attribute, die die Berechnung der Inversen der Korrelationsmatrix nicht zulassen.

**Tabelle 21: MSA für Datenset (Šumperk mit Ausreißern)<sup>306</sup>**

	Durchgang			
	1	2	3	4
<b>MSA</b>	nan	0,3730	0,7009	0,7325

Es ist eine deutliche Verbesserung des MSA-Kriteriums zu sehen und ist nur unwesentlich schlechter als jener der Vor- und Nachmahlung ohne Ausreißer in Tabelle 14. Das Ergebnis dieses Tests ist ein Datenset mit 9 Attributen und 60 Datensätzen, dem die Attribute MW und Ofen wieder hinzugefügt werden. Dadurch entsteht für die anschließende Modellierung ein Datenset mit 10 Prädiktoren, einem vorherzusagenden Zielattribut mit 60 Datensätzen.

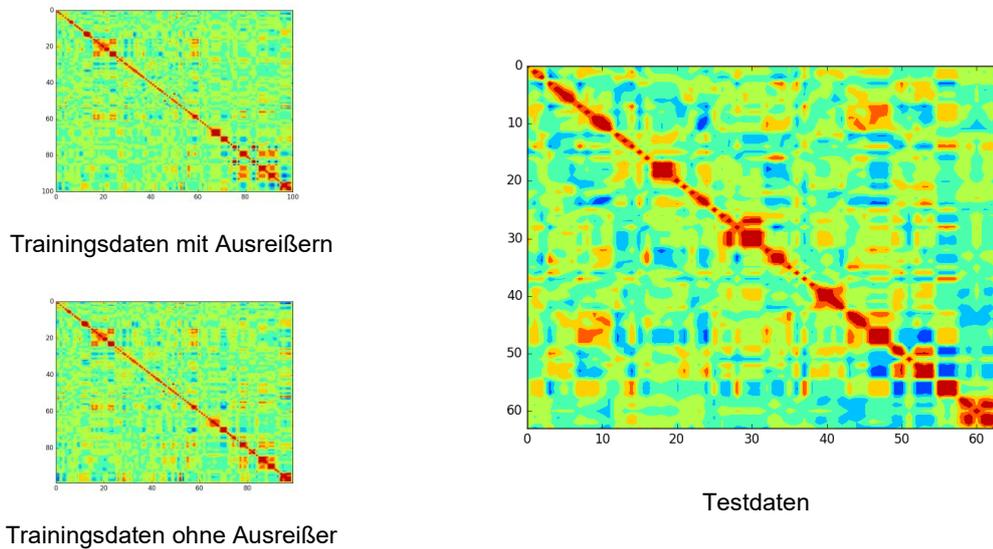
#### 4.3.5 Modellierung

Die Modellierung erfolgte wie in Kapitel 4.2.7 wieder in zwei Schritten und wurde in Matlab in der Classification Learner App durchgeführt. Es stehen dieselben Methoden zur Auswahl (siehe Tabelle 15 und Tabelle 16) und in einem ersten Schritt wird eine generelle Abschätzung getroffen, welches Datenset verwendet werden soll.

<sup>305</sup> Quelle: eigene Darstellung

<sup>306</sup> Quelle: ebda.

Anschließend werden auf dieses und das passende Gegenstück sämtliche Methoden angewendet, um Modelle zu erstellen. Diese Modelle werden zudem mit den Testdaten verifiziert. Bei den Testdaten handelt es sich um 25 Datensätze, die wie bereits die Testdaten in Kapitel 4.2.7 aus dem Zeitraum von September 2015 bis Anfang Jänner 2017 stammen. Sie werden wie in Kapitel 4.3.2 beschrieben, aufbereitet und standardisiert, wobei auch hier keine Ausreißeranalyse vorgenommen wird.



**Abbildung 56: Korrelationsmatrizen für Šumperk<sup>307</sup>**

Abbildung 56 zeigt die bereits bekannten Korrelationsmatrizen aus Abbildung 52 bzw. Abbildung 55. Wie in Kapitel 4.2.7 fällt die anders strukturierte Korrelationsmatrix der Testdaten gegenüber jenen der Trainingsdaten auf. Auch hier kann die als Indiz dafür gewertet werden, dass die Generalisierbarkeit des Modells, das auf den Trainingsdaten beruht, nicht gegeben ist.

**Tabelle 22: Auswertung der Testdaten (Anzahl je Klasse)<sup>308</sup>**

Bezeichnung	Anzahl Datensets je Klasse
11.66	2
14.24	1
14.59	3
14.73	1
15.09	7
15.68	1
16.50	5
18.34	4
25.91	1

Ebenso wird das Zielattribut in Klassen eingeteilt. Eine Aufschlüsselung über die Anzahl der Datensets je Klasse sind in Tabelle 22 zu sehen. Wieder gibt es einen Grenzfall (rot hinterlegt), der nicht zugeordnet werden kann.

<sup>307</sup> Quelle: eigene Darstellung

<sup>308</sup> Quelle: ebda.

### Vergleich verschiedener Modelle und Datensets

Für die Daten für Šumperk liegen zu diesem Zeitpunkt 2 Datensets vor: jenes, das durch die Ausführung des KMO Tests entsteht, und jenes, dessen Attribute von der Vor- und Nachmahlung übernommen werden. Die Kombinationsmöglichkeiten, mit denen eine erste Abschätzung stattfindet, sind in Tabelle 23 zu sehen. Auf die Durchführung der PCA wird bei dieser Abschätzung verzichtet, da bereits in Kapitel 4.2.7 auf dies verzichtet wurde.

**Tabelle 23: Kombinationen für Šumperk<sup>309</sup>**

	Datenset	Ausreißer	Kreuz-validierung	PCA	Ofen	Anzahl Prädiktoren
CZ 5-fold inkl. Ausreißer + Ofen (Prädiktoren: 10)	KMO	mit	5	nein	ja	10
CZ 10-fold inkl. Ausreißer + Ofen (Prädiktoren: 10)	KMO	mit	10	nein	ja	10
CZ 5-fold inkl. Ausreißer + Ofen (Prädiktoren: 63)	VM+NM	mit	5	nein	ja	63
CZ 10-fold inkl. Ausreißer + Ofen (Prädiktoren: 63)	VM+NM	mit	10	nein	ja	63

Die in Tabelle 23 angegebenen Datensets werden als Trainingsdaten für die Modellerstellung verwendet. In Abbildung 57 sind die Ergebnisse sehen. Auf der Abszisse sind ist die Genauigkeit der Modelle aufgetragen und auf der Ordinate die verwendeten Modelle. Hierbei wird sich auf die Modelle, die bereits in Kapitel 4.2.7 bei der endgültigen Modellerstellung verwendet werden, konzentriert. Dazu gehören die verschiedenen Bäume, die Varianten der SVM und die Ensemblemethoden.

<sup>309</sup> Quelle: eigene Darstellung

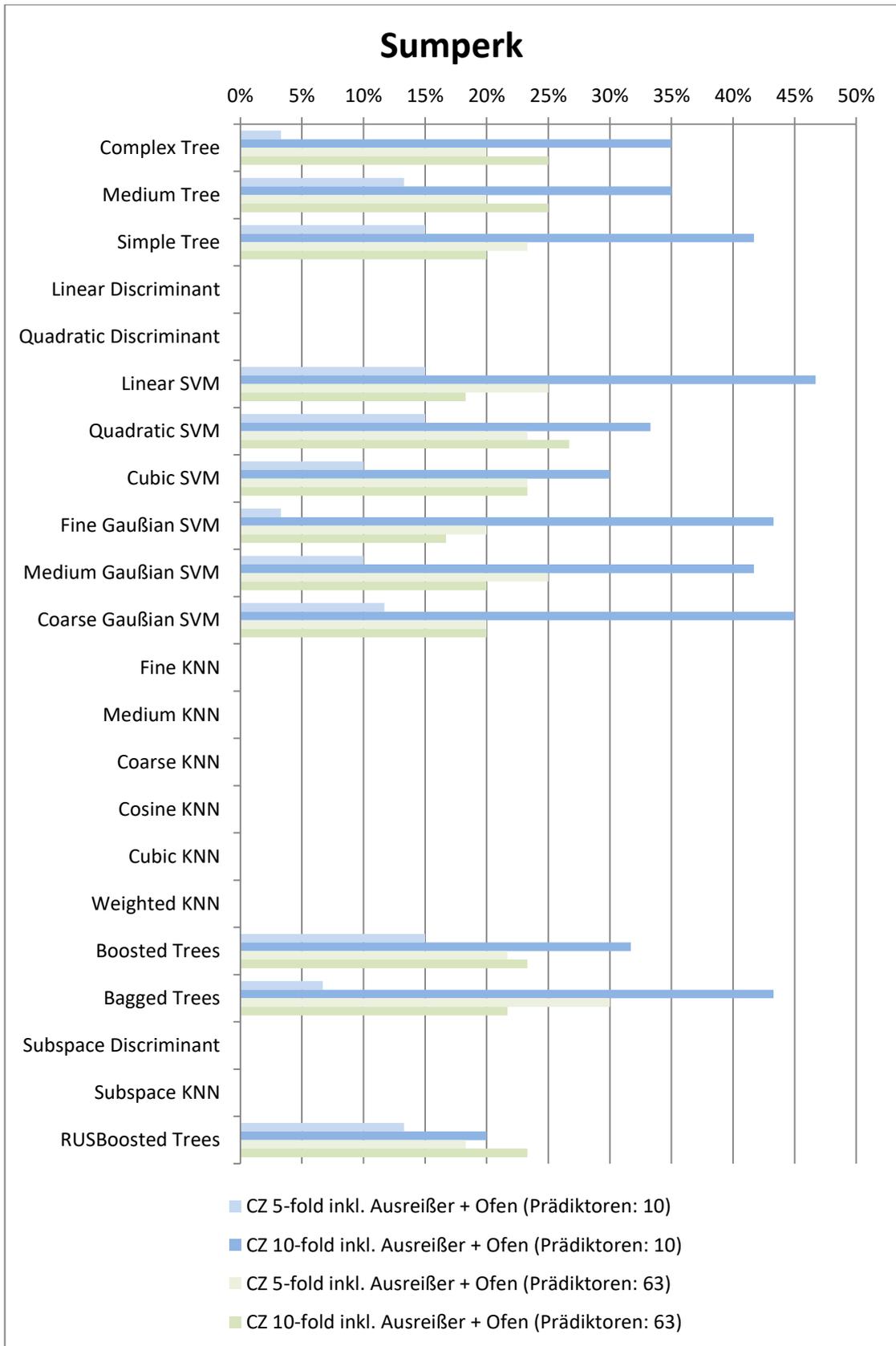


Abbildung 57: Vergleich verschiedener Kombinationsmöglichkeiten für Šumperk<sup>310</sup>

<sup>310</sup> Quelle: eigene Darstellung

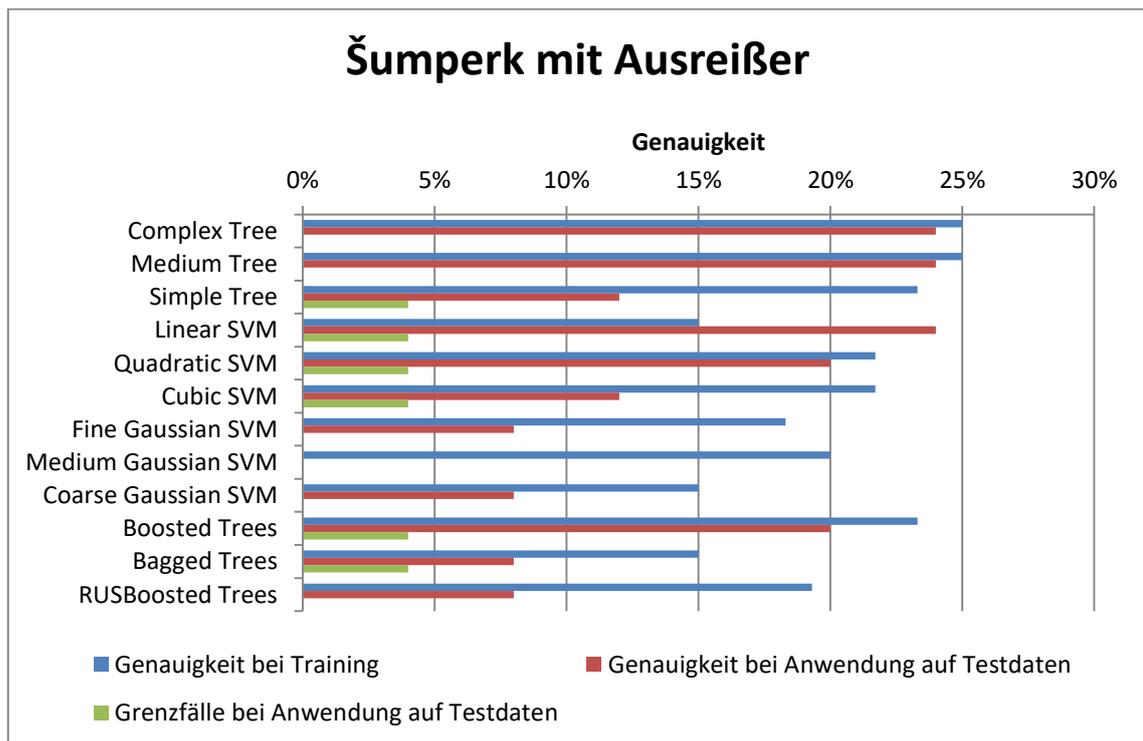
Nach dieser ersten Abschätzung können folgende Überlegungen angestellt werden:

- Die Verbesserung zwischen den Datensets, die 10 Prädiktoren aber unterschiedliche Anzahl an Kreuzvalidierungsmengen haben, beruht höchstwahrscheinlich auf Overfitting.
- Die Anzahl der Kreuzvalidierungsmengen hat auf das Datenset mit 63 Prädiktoren keine nennenswerte Auswirkung.
- Im Vergleich beider Datensets mit der 5-fachen Kreuzvalidierung schneidet jenes mit den 63 Prädiktoren besser ab.

Auf Basis dieser Überlegungen und vorrangig auf der letzten wird entschieden, dass das Datenset mit den 63 Prädiktoren für die Modellerstellung und die anschließende Anwendung auf die Testdaten verwendet wird. In Anlehnung an die Vorgehensweise in Kapitel 4.2.7 werden wieder die Ergebnisse für die Datensets mit und ohne Ausreißer gegenüber gestellt. Daher werden die bereits in Kapitel 4.3.3 aufbereiteten Daten in Folge verwendet.

### Ergebnisse des Datensets mit Ausreißern

Die Ergebnisse für das Datenset mit Ausreißern für Šumperk sind nach dem gleichen Prinzip wie die Ergebnisse aus Kapitel 4.2.7 in Abbildung 43 und Abbildung 45 aufbereitet.

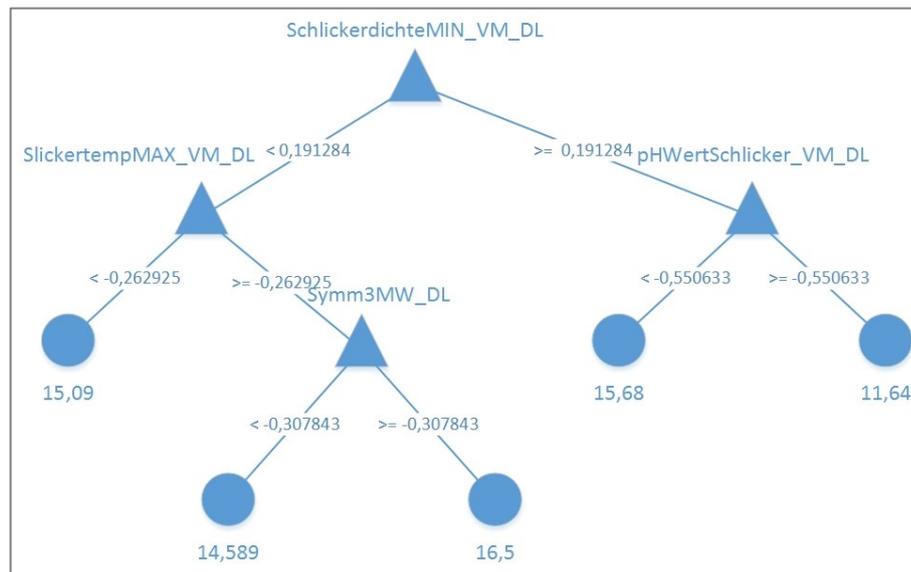


**Abbildung 58: Ergebnisse des CL für Šumperk mit Ausreißern<sup>311</sup>**

In Abbildung 46 sind die Genauigkeiten beim Training (blau) und bei der Anwendung auf die Testdaten (rot) und deren Grenzfälle (grün) dargestellt. Die genauen Daten hinter dem Diagramm sind der Tabelle 38 zu finden. Complex Tree (25 %), Medium Tree (25 %) und Boosted Trees (23,3 %) liefern die besten Ergebnisse beim Training.

<sup>311</sup> Quelle: eigene Darstellung

Dieser Trend setzt sich bei der Anwendung auf die Testdaten fort (24 %, 24 % und 20 %). Zusätzlich hat die Linear SVM dasselbe Ergebnis (24 %). Die Medium Gaussian SVM versagt trotz 20 % bei den Trainingsdaten bei der Anwendung auf die Testdaten. Die durchschnittliche absolute Abweichung zwischen vorhergesagter und Ist-Klasse liegt bei 1,94 bis 3,42 im ähnlichen Bereich wie die der bereits betrachteten Ergebnisse aus Kapitel 4.2.7. Eine genauere Betrachtung der Ergebnisse je Modell und Charge ist den Abbildung 86 bis Abbildung 97 zu entnehmen.



**Abbildung 59: Simple Tree für Šumperk mit Ausreißern<sup>312</sup>**

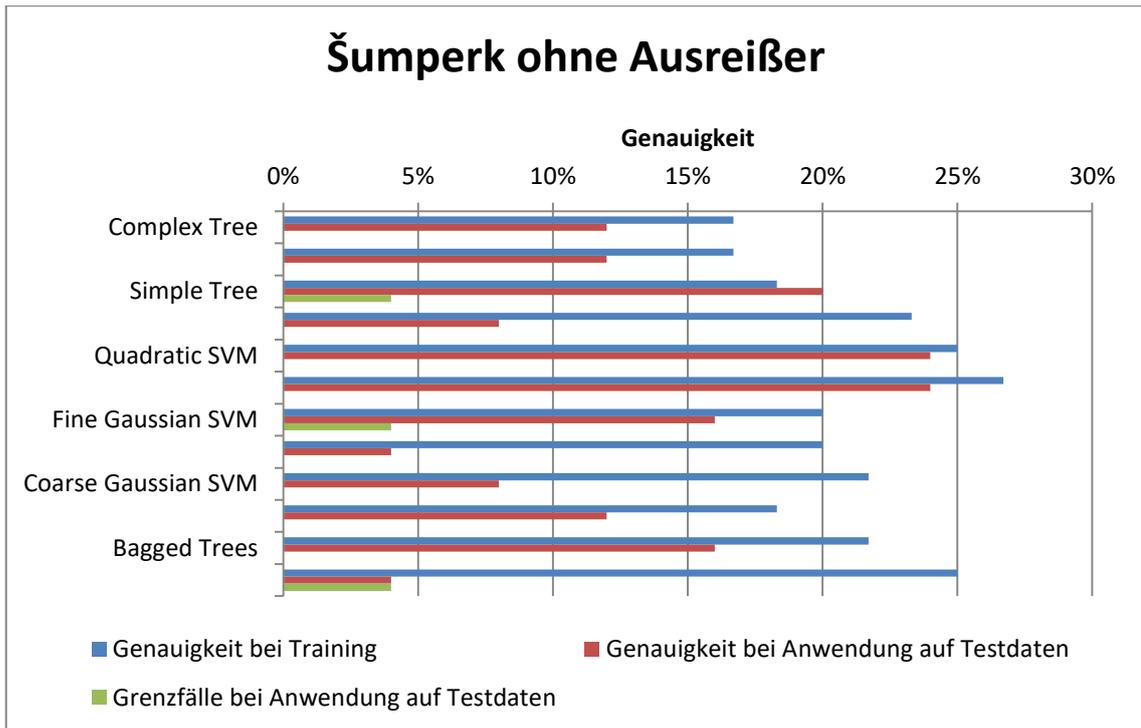
Bei der Betrachtung des Simple Tree in Abbildung 59 wird wieder ersichtlich, dass dieser nur 5 Blätter hat, die 5 der 8 im Datenset enthaltenen Klassen repräsentieren. Diese prozentual höhere Abdeckung der Klassen kann zu einem besseren Ergebnis führen als bei den Datensets der Vor- und Nachmahlung. Dennoch wird diese Methode wegen der Unvollständigkeit verworfen.

Bei der Betrachtung der Regeln für Medium Tree und Complex Tree ab Seite ww im Anhang wird ersichtlich, dass diese weniger verzweigt sind als jene für die Vor- und Nachmahlung. Beide enthalten alle Klassen als Blätter.

### Ergebnisse des Datensets ohne Ausreißer

Das Datenset ohne Ausreißer wird wie die anderen drei Datensets zuvor aufbereitet. Die Aufbereitung der Ergebnisse des Training, und der Anwendung auf die Testdaten wird in Abbildung 60 zusammengefasst.

<sup>312</sup> Quelle: eigene Darstellung



**Abbildung 60: Ergebnisse des CL für Šumperk ohne Ausreißer<sup>313</sup>**

Die Genauigkeiten beim Training (blau) und bei der Anwendung auf die Testdaten (rot) und deren Grenzfälle (grün) sind in Abbildung 60 zu sehen. Die genauen Daten hinter dem Diagramm sind der Tabelle 39 zu finden. Die Ergebnisse beim Training waren relativ konstant, aus denen sich die Linear SVM (23,3 %), Quadratic SVM (25 %) und Cubic SVM (26,7 %) sowie die RUSBoostedTrees (25 %) hervorheben. Dieser Trend wird bei der Anwendung auf die Testdaten nur bei der Quadratic SVM und der Cubic SVM (beide 24 %) beobachtet. Die durchschnittliche absolute Abweichung zwischen vorhergesagter und Ist-Klasse lag bei 1,82 (Cubic SVM) bis 3,62 (Coarse Gaussian SVM) im ähnlichen Bereich wie die der bereits betrachteten Ergebnisse aus Kapitel 4.2.7 und bei den Ergebnissen mit Ausreißer. Eine genauere Betrachtung der Ergebnisse je Modell und Charge ist im Anhang in Abbildung 98 bis Abbildung 109.

<sup>313</sup> Quelle: eigene Darstellung

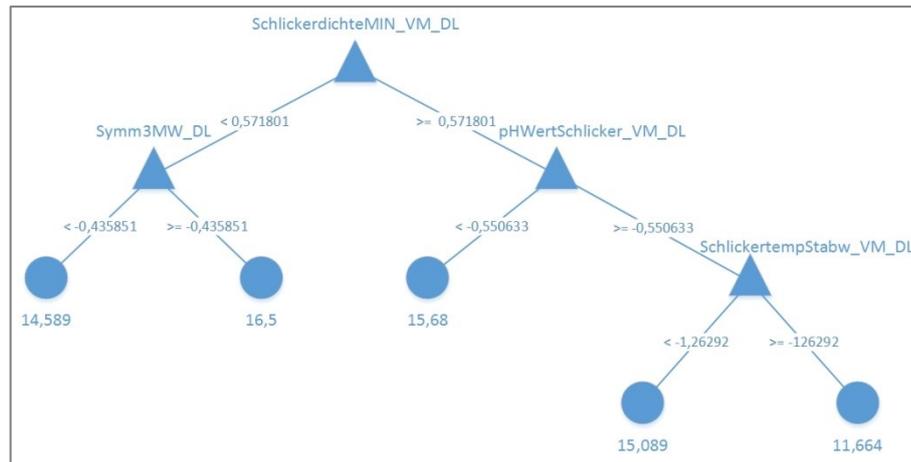


Abbildung 61: Simple Tree für Šumperk ohne Ausreißer<sup>314</sup>

Abbildung 61 zeigt den entstandenen Simple Tree, der wieder nur fünf Blätter hat. Somit werden drei Klassen nicht abgebildet. Dies ist ein Ausschlussgrund für das Modell. Ab Seite fff sind die Regeln für den Medium Tree und Complex Tree zu finden. Wie beim Šumperker Datenset mit Ausreißern ist zu erkennen, dass diese weniger verzweigt sind als jene für die Vor- und Nachmahlung. Beide enthalten alle Klassen als Blätter.

Nach der Aufbereitung der Daten und dem Erstellen der Modelle für Deutschlandsberg und Šumperk wird klar, dass diese nicht die gewünschten Ergebnisse erbringen. Daher ist von einem Einsatz der Modelle abzuraten. Im nächsten Kapitel wird die zusätzliche Fertigungskapazität berechnet, die ein genaues Modell freierwerden lässt.

#### 4.4 Produktionswirtschaftliche Betrachtung

Innerhalb dieses Kapitels soll geklärt werden, welche Auswirkungen ein korrekt vorhergesagter Zielwert auf die Fertigungskapazitäten hätte. Die folgenden Berechnungen konnten nicht in der Praxis überprüft werden.

Im Fall der Granulatfreigabe in Deutschlandsberg würden keine Fertigungskapazitäten freiwerden, da die Sinterung für die Freigabe nicht auf in den Fertigungsprozess integrierten Öfen durchgeführt wird. Einzig könnte es sich auf die Durchlaufzeit auswirken. Allerdings nur unter der Voraussetzung, dass die Granulate direkt nach der Produktion nach Šumperk transportiert und hier ohne Lagerzeit verarbeitet werden.

Im Fall der Sinterfreigabe sind die Voraussetzungen anders. Die Sinterung für die Sinterfreigabe wird auf den Durchstoßöfen durchgeführt, die ebenso in der Produktion eingesetzt werden. Hierbei werden Kapazitäten blockiert, da bis zur korrekten Einstellung der Öfen nur jeder 50. Platz im Ofen besetzt wird. Bei einer korrekten Vorhersage wäre die Sinterfreigabe nicht mehr notwendig und die schlecht ausgelasteten Freigabesinterungen sind hinfällig. Dies wird im Folgenden an Hand der Sinterdaten des bisher betrachteten Produkts auf Basis der Produktionslose, auf denen

<sup>314</sup> Quelle: eigene Darstellung

die Testdaten beruhen, erläutert. Dabei ist der Ofen das einzige betrachtete Aggregat, der gleichzeitig den Engpass darstellt (siehe Punkt 1b) in Tabelle 2).

**Tabelle 24: Auswertung der Bestückung des Ofens<sup>315</sup>**

	Bestückung des Ofens		
	1	10	50
<b>Anzahl der Messungen</b>	1967	107	32

In Tabelle 24 werden die Sinterdaten hinsichtlich der Bestückung der Öfen und die Anzahl der Messungen ausgewertet. Die Bestückungen 1, 10 und 50 bedeuten, dass jeder, jeder zehnte oder jeder fünfzigste Platz belegt wird. Bei der Belegung jedes Platzes wird nur jeder zehnte Platz gemessen. Dies bedeutet für die Summe der benötigten Plätze folgendes:

$$\text{Summe der benötigten Plätze} = 1967 * 10 + 107 + 32 = 19809 \text{ Plätze}$$

Wird vom Idealfall ausgegangen, dass ein Modell die korrekte Sinterung vorhersagt, die Sinterfreigabemessung entfällt und somit jeder Platz belegt werden kann, würden folgende Kapazitäten freiwerden:

$$\text{leere Plätze} = (107 * 9) + (32 * 49) = 2531 \text{ Plätze}$$

Diese leeren Plätze haben einen Anteil von 12,78 % an der benötigten Kapazität. Wird davon ausgegangen, dass ein Los im Mittel 792 Plätze (ohne leere Plätze) im Ofen benötigt, hätten 3 Lose mehr im betrachteten Zeitraum (September 2016 bis Anfang Jänner 2017) produziert werden können.

Zusätzlich wirkt sich der Wegfall der leeren Plätze innerhalb der Produktion auch auf die Durchlaufzeit aus. Dies würde zu einer durchschnittlichen Beschleunigung um 101 Plätze je Los führen.

Würde eine Anpassung der Produktion durch ein funktionierendes Modell vorgenommen werden, würde es sich um eine qualitative Anpassung handeln, die wiederum in einer intensitätsmäßigen Kapazitätsanpassung resultieren würde (siehe Kapitel 0).

Nach dieser produktionswirtschaftlichen Betrachtung des Hintergrundes der Arbeit werden im nächsten Kapitel die gewonnenen Erkenntnisse zusammengefasst. Zusätzlich soll ein Ausblick zu weiteren Möglichkeiten gegeben werden.

<sup>315</sup> Quelle: eigene Darstellung

## 5 Zusammenfassung und Ausblick

Innerhalb dieser Masterarbeit wurde die Thematik der Datenanalyse mit dem produktionswirtschaftlichen Ansatz in Verbindung gesetzt. Anreiz dafür war die Aufgabenstellung der EPCOS OHG, ein Verfahren zur Vorhersage der Sinterergebnisse aus Daten, die entlang der Prozesskette bzw. bei Freigabemessungen gesammelt wurden, zu entwickeln. Hierbei wurden die Produktionsstandorte Deutschlandsberg und Šumperk, die beide am Prozess beteiligt sind, betrachtet. Auf Basis dieser Aufgabenstellung wurden die Forschungsfragen in Kapitel 1.2 abgeleitet, die im Folgenden wiederholt werden:

- Ist es möglich, den Widerstandswert  $R_{25}$  der PTC-Bauteile für zukünftige Lose auf Basis der bereits gefertigten vorherzusagen?
- Wie wirkt sich eine exakte Vorhersage auf die Kapazitätsauslastung aus?

Dieses Kapitel dient dazu, eine Zusammenfassung über die theoretischen Hintergründe und praktischen Ergebnisse, die zur Beantwortung der Forschungsfragen führen, zu geben. Eine klare Beantwortung der ersten Frage ist nur beschränkt möglich. Die praktische Fallstudie lieferte kein zufriedenstellendes Ergebnis, so ist es denkbar, dass der theoretisch betrachtete Lösungsweg für eine andere und/oder größere Datenbasis funktioniert. Zusätzlich wird am Ende des Kapitels ein Überblick über die Einflussfaktoren und andere Ansätze geben. Die Antwort der zweiten Forschungsfrage kann in der Praxis nicht überprüft werden, da sie ein funktionierendes Modell voraussetzt.

Die dazu notwendige produktionswirtschaftliche Theorie wurde in Kapitel 2 vorgestellt. Sie umfasste die Grundlagen eines Betriebes und die Kennzahlen Durchlaufzeit und Kapazität. Eine Verbesserung durch ein funktionierendes Modell konnte so quantifiziert werden. Zusätzlich wurden die Anpassungsformen der Kapazität betrachtet.

Eine Abfolge von Methoden, die es möglich machen sollten, das in der ersten Forschungsfrage geforderte Verfahren zu entwickeln, wurden in Kapitel 3 präsentiert. Hierbei wurden zunächst die verschiedenen Ausprägungen der Datenanalyse betrachtet, bevor mit CRISP-DM ein Referenzmodell zur Durchführung einer solchen vorgestellt wurde. Anhand dieses wurde das Kapitel in weiterer Folge aufgebaut. Nach der Theorie zur Untersuchung der Daten in Bezug auf ihre Skala und statistischen Eigenschaften und anschließenden Aufbereitung hinsichtlich Ausreißern, Gruppierung und Dimensionalitätsreduktion, wurde das maschinelle Lernen vorgestellt. Hierbei wurden zuerst das Konzept und die Grundlagen erläutert, bevor die Methoden zum Klassieren der Daten vorgestellt wurden. Diese beruhen auf unterschiedlichen Ansätzen und sind oftmals nicht für dieselbe Datenbasis geeignet.

Das praktische Fallbeispiel übertrug die theoretischen Überlegungen in die Praxis. Hierfür wurden, orientiert am CRISP-DM, eine Datenaufbereitung und -analyse, an Daten die vom Unternehmen in aggregierter Form bereitgestellt wurden, vorgenommen und anschließend Modelle auf Basis der Daten aus Deutschlandsberg und Šumperk erstellt. Für beide Standorte erbrachten die Modelle in der Evaluierungsphase kein

zufriedenstellendes Ergebnis. Somit bleibt die erste Forschungsfrage in der Praxis unbeantwortet. Die veränderte Kapazitätsauslastung, die in der zweiten Forschungsfrage thematisiert wird, kann auf Basis der Daten beantwortet, aber nicht in der Praxis überprüft werden. Allerdings hätten dies im betrachteten Zeitraum freie Kapazitäten von fast 13 % zur Folge gehabt.

Dass das Modell als solches nicht verwendbar ist, kann als Resultat mehrerer Einflussfaktoren gesehen werden. Der vordringlichste Einflussfaktor besteht in einem nicht zufriedenstellenden kleinen Verhältnis der Datensätze zu den Attributen.

**Tabelle 25: Verhältnisse Datensätze zu Attributen für VM+NM**

	Start	Analyse der fehlende Daten	Zusammenführung	KMO	Modellierung
VM - Datensätze	266	226	-	-	-
VM - Attribute	162	90	-	-	-
NM - Datensätze	265	224	-	-	-
NM - Attribute	115	62	-	-	-
VM+NM - Datensätze	-	-	218	218	218
VM+NM - Attribute	-	-	98	62	63
Verhältnis VM	1,64	2,51	-	-	-
Verhältnis NM	2,30	3,61	-	-	-
Verhältnis VM+NM	-	-	2,22	3,52	3,46

Für das Deutschlandsberger Datenset wurden die Verhältnisse von Datensätzen zu Attributen für die Vor- und Nachmahlung in Tabelle 25 zusammengefasst. Durch die Aufbereitung konnten diese verbessert werden, allerdings nicht im ausreichenden Maße.

**Tabelle 26: Verhältnisse Datensätze zu Attribute für Šumperk**

	Start	Modellieren
VM+NM - Datensätze	218	-
VM+NM - Attribute	63	-
Sinterfreigabe - Datensätze	141	-
Sinterfreigabe - Attribute	2	-
Šumperk - Datensätze	-	60
Šumperk - Attribute	-	63
Verhältnis VM	3,46	-
Verhältnis NM	70,50	-
Verhältnis VM+NM	-	0,95

Dasselbe gilt für das Šumperker Datenset. Modelle benötigen große Datenmengen bzw. ein besseres Verhältnis, um ausreichend gut lernen zu können und Generalisierbarkeit gewährleisten zu können. Dies hat sich auf den hier vorgestellten Lösungsweg ausgewirkt und wird sich auch auf zukünftige Ansätze auswirken, die hinsichtlich der Aufbereitung und der Modellwahl vom vorgestellten Ansatz abweichen. Bei der Aufbereitung können andere Ansätze für die Ausreißerelimination gewählt

werden. Dazu gehören die multivariate Ansätze<sup>316</sup>, die den gesamten Datensatz betrachten. In deren Folge andere Werte als Ausreißer identifiziert oder mit anderen Werten ersetzt werden. Genauso gibt es andere Ansätze für die Gruppierung, wie z. B. die Gruppenbildung über die Entropie.<sup>317</sup> Andere Verfahren, die zu anderen Klassen führen, können das Ergebnis beeinflussen. Der Einsatz der PCA, der wegen zu geringer Datenmenge verworfen wurde, lieferte den Grund für die Anwendung der KMO und der daraus resultierenden Löschung der Attribute. Daher ist die Frage, ob diese Löschung dennoch gerechtfertigt ist. Der in dieser Arbeit gewählte Ansatz des maschinellen Lernens ist jener der Klassifikation. Es muss in Betracht gezogen werden, ob andere Ansätze, wie z. B. das Prognostizieren mittels Neuronaler Netze oder eine andere Kombination von Bäumen, das angestrebte Resultat erbringt. Da dieses Verfahren fast ausschließlich datengetrieben ist, besteht die Frage, ob ein anderes Ergebnis erzielt werden würde, wenn ein modellgetriebener Ansatz verfolgt wird.

Alle genannten Ansätze können jedoch nur in Betracht gezogen werden, wenn die Datenbasis verbessert wird. Ein besseres Verhältnis von Datensätzen zu Attributen kann entweder durch eine Reduktion der Attribute oder durch eine Erweiterung der Anzahl der Datensätze erzielt werden. Bei einer Reduktion der Attribute muss ein Experte jene Attribute löschen, die von ihm als nicht relevant eingestuft werden. Damit verändert sich der datengetriebene Ansatz zu einem modellgetriebenen Ansatz. Die andere Möglichkeit ist eine Erhöhung der Anzahl der Datensätze. Um mehr Daten als Grundlage zu haben, muss der Betrachtungszeitraum in die Vergangenheit ausgeweitet werden. Eine Vergrößerung des Betrachtungszeitraums und somit Erhöhung der Anzahl der Datensätze wurde mit dem Unternehmen diskutiert, jedoch verworfen, da es ungefähr jährlich zu Technologieumbrüchen kommt und somit die Daten nicht vergleichbar sind. Zusätzlich muss die Möglichkeit einer Ausreißeranalyse vor der Aggregation der Daten in Erwägung gezogen werden, um eventuelle Ausreißer vor ihrem Einfließen in die Datenbasis zu erkennen und zu beseitigen.

---

<sup>316</sup> Vgl. Hair, J. F. et al. (2014), S. 64

<sup>317</sup> Vgl. Witten, I. H. et al. (2011), S. 316

## Literaturverzeichnis

- Abbott, D. (2014): Applied Predictive Analytics: Principles and Techniques for the Professional Data Analyst. Indianapolis, IN: Wiley. ISBN 978-1-118-72796-6.
- Alpaydin, E. (2010): Introduction to Machine Learning. 2nd ed, Cambridge, Massachusetts: MIT Press. (Adaptive Computation and Machine Learning). ISBN 978-0-262-01243-0.
- Azevedo, A.; Santos, M. F. (2008): KDD, SEMMA and CRISP-DM: A Parallel Overview. Proceedings of IADIS European Conference on Data Mining 2008, Amsterdam, The Netherlands, July 24-26, 2008. Proceedings, januari 1 2008.
- Azevedo, C. S.; Santos, M. F. (2005): Data Mining - Descoberta de Conhecimento em Bases de Dados. Capa mole. ISBN 978-972-722-509-5.
- Backhaus, K.; Erichson, B.; Plinke, W.; Weiber, R. (2011): Multivariate Analysemethoden: Eine Anwendungsorientierte Einführung. 13., überarbeitete Auflage, Berlin Dordrecht London New York: Springer. (Springer-Lehrbuch). ISBN 978-3-642-16490-3.
- Backhaus, K.; Erichson, B.; Weiber, R. (2015): Fortgeschrittene Multivariate Analysemethoden: Eine anwendungsorientierte Einführung. 3., überarbeitete und aktualisierte Auflage, Berlin Heidelberg: Springer Gabler. (Lehrbuch). ISBN 978-3-662-46086-3.
- Bauernhansl, T. (2014): Die Vierte Industrielle Revolution – Der Weg in ein wertschaffendes Produktionsparadigma. In: Bauernhansl, T.; ten Hompel, M.; Vogel-Heuser, B. (Hrsg.): Industrie 4.0 in Produktion, Automatisierung und Logistik. Wiesbaden: Springer Fachmedien Wiesbaden. ISBN 978-3-658-04681-1, S. 5–35.
- Beltrami, E. (1873): Sulle Funzioni Bilineari. In: Giornale di matematiche, Jg. 11, S. 98–106.
- Blohm, H.; Beer, T.; Seidenberg, U.; Silber, H. (2008): Produktionswirtschaft. 4., vollst. überarb. Aufl, Herne: Verl. Neue Wirtschafts-Briefe. (NWB-Studium-Betriebswirtschaft Lehrbuch). ISBN 978-3-482-63024-8.
- Breiman, L. (1996): Bagging Predictors. In: Machine Learning, Jg. 24, Nr. 2, S. 123–140.
- Breiman, L. (2001): Statistical Modeling: The Two Cultures (with Comments and a Rejoinder by the Author). In: Statistical science, Jg. 16, Nr. 3, S. 199–231.
- Czado, C.; Schmidt, T. (2011): Mathematische Statistik. Berlin: Springer. (Statistik und ihre Anwendungen). ISBN 978-3-642-17260-1.
- David, H. A.; Hartley, H. O.; Pearson, E. S. (1954): The Distribution of the Ratio, in a Single Normal Sample, of Range to Standard Deviation. In: Biometrika, Jg. 41, Nr. 3/4, S. 482–493.
- Degen, H.; Lorscheid, P. (2012): Statistik-Lehrbuch: Methoden der Statistik im wirtschaftswissenschaftlichen Bachelor-Studium. 4., korrigierte Aufl, München: Oldenbourg. ISBN 978-3-486-71588-0.
- Edward E. Cureton, R. B. D. (1993): Factor Analysis: An Applied Approach. Psychology Press. ISBN 0-8058-1546-5.
- Eidenmüller, B. (1995): Die Produktion als Wettbewerbsfaktor: Das Potential der Mitarbeiter nutzen - Herausforderung an das Produktionsmanagement. 3. aktualis. u. erw. Aufl, Köln: TÜV Rheinland. (Leitfaden für Unternehmer und Führungskräfte). ISBN 978-3-8249-0181-4.

- EPCOS HR (2016): EPCOS OHG: A TDK Company: Herzlich Willkommen, unveröffentlichtes Dokument zur Information neuer Mitarbeiter, Deutschlandsberg, 2016.
- EPCOS PPD PTC PD (2016): 2nd Monozukuri Strategy Meeting: PPD Zero Defect Approach: Reduction of resistance variation at PTC, 2016.
- Ertel, W. (2013): Grundkurs Künstliche Intelligenz: Eine praxisorientierte Einführung. 3. Aufl, Wiesbaden: Springer Vieweg. (Lehrbuch). ISBN 978-3-8348-2157-7.
- Fahrmeir, L.; Künstler, R.; Pigeot, I.; Tutz, G. (2004): Statistik: Der Weg zur Datenanalyse. 5., Aufl, Berlin: Springer. (Springer-Lehrbuch). ISBN 978-3-540-21232-4.
- Fayyad, U.; Piatetsky-Shapiro, G.; Smyth, P. (1996): From Data Mining to Knowledge Discovery in Databases. In: AI Magazine, Jg. 17, Nr. 3, S. 37–54.
- Fernandez, G. (2003): Data Mining Using SAS Applications: a Case Study Approach. Cary, N.C.: SAS Pub.
- Fisher, R. A.; Mackenzie, W. A. (1923): Studies in Crop Variation. II. The Manurial Response of Different Potato Varieties. In: The Journal of Agricultural Science, Jg. 13, Nr. 03, S. 311.
- Freitag, M.; Kück, M.; Ait Alla, A.; Lütjen, M. (2015): Potentiale von Data Science in Produktion und Logistik: Teil 1 - Eine Einführung in aktuelle Ansätze der Data Science. In: Industrie 4.0 Management, Jg. 31, Nr. 5, S. 22–26.
- Gartner (2015): Gartner: Predictive Analytics. URL: <http://www.gartner.com/it-glossary/predictive-analytics> (Zugriff: 17.05.2015).
- Gertheiss, J.; Tutz, G. (2009): Statistische Tests. In: Schwaiger, M.; Meyer, A. (Hrsg.): Theorien und Methoden der Betriebswirtschaft: Handbuch für Wissenschaftler und Studierende. München: Vahlen. ISBN 978-3-8006-3613-6, S. 439–454.
- Gläßner, J. (1995): Modellgestütztes Controlling der beschaffungslogistischen Prozeßkette. Als Ms. gedr, Düsseldorf: VDI-Verl. (Fortschritt-Berichte VDI Reihe 2, Fertigungstechnik; 337). ISBN 978-3-18-333702-6.
- Grubbs, F. E. (1950): Sample Criteria for Testing Outlying Observations. In: The Annals of Mathematical Statistics, S. 27–58.
- Gudehus, T. (2010): Logistik: Grundlagen - Strategien - Anwendungen. 4., aktualisierte Aufl, Berlin: Springer. ISBN 978-3-540-89388-2.
- Gutenberg, E. (1971): Grundlagen der Betriebswirtschaftslehre Die Produktion. Berlin, Heidelberg: Springer Berlin Heidelberg. ISBN 978-3-642-61989-2.
- Guttman, L. (1953): Image theory for the structure of quantitative variates. In: Psychometrika, Jg. 18, Nr. 4, S. 277–296.
- Hair, J. F.; Black, W. C.; Babin, B. J.; Anderson, R. E. (2014): Examining Your Data. In: Hair, J. F. (Hrsg.): Multivariate data analysis. 7. ed., Pearson new internat. ed, Harlow: Pearson. ISBN 978-1-292-02190-4 (Pearson custom library), S. 31–88.
- Haseloff, O. W.; Hoffmann, H. J. (1965): Kleines Lehrbuch der Statistik: Für Naturwissenschaftler, Mediziner, Psychologen, Sozialwissenschaftler und Pädagogen. 2., Berlin: Walter De Gruyter & Co.
- Hastie, T.; Tibshirani, R.; Friedman, J. H. (2009): The Elements of Statistical Learning: Data Mining, Inference, and Prediction. 2nd ed, New York, NY: Springer. (Springer Series in Statistics). ISBN 978-0-387-84857-0.
- Heitjan, D. F. (1997): Annotation: What Can Be Done About Missing Data? Approaches to Imputation. In: American Journal of Public Health, Jg. 87, Nr. 4, S. 548–550.

- Hertel, B. R. (1976): Minimizing Error Variance Introduced By Missing Data Routines in Survey Analysis. In: *Sociological Methods & Research*, Jg. 4, Nr. 4, S. 459–474.
- Hotelling, H. (1933): Analysis of a Complex of Statistical Variables with Principal Components. In: *Journal of Educational Psychology*, Jg. 24, S. 417–441.
- IBM Big Data & Analytic Hub Infographic: The Four V's of Big Data | IBM Big Data & Analytics Hub. URL: <http://www.ibmbigdatahub.com/infographic/four-vs-big-data> (Zugriff: 09.02.2017).
- JetBrains s.r.o. (2017): PyCharm. JetBrains. URL: <https://www.jetbrains.com/pycharm/> (Zugriff: 13.02.2017).
- Jolliffe, I. T. (2002): *Principal Component Analysis*. 2nd ed, New York: Springer. (Springer Series in Statistics). ISBN 978-0-387-95442-4.
- Kaiser, H. F. (1970): A Second Generation Little Jiffy. In: *Psychometrika*, Jg. 35, Nr. 4, S. 401–415.
- Little, R. J. A.; Rubin, D. B. (2002): *Statistical Analysis with Missing Data*. 2nd ed, Hoboken, N.J: Wiley. (Wiley Series in Probability and Statistics). ISBN 978-0-471-18386-0.
- Malhotra, N. K. (1987): Analyzing Marketing Research Data with Incomplete Information on the Dependent Variable. In: *Journal of Marketing Research*, S. 74–84.
- MathWorks (2017): Choose Classifier Options - MATLAB & Simulink - MathWorks Deutschland. URL: <https://de.mathworks.com/help/stats/choose-a-classifier.html> (Zugriff: 21.02.2017).
- MathWorks (2017): Pricing and Licensing - MATLAB & Simulink. URL: <https://de.mathworks.com/pricing-licensing.html?prodcode=ML> (Zugriff: 13.02.2017).
- Mertens, P.; Bodendorf, F.; König, W.; Picot, A.; Schumann, M.; Hess, T. (2012): *Grundzüge der Wirtschaftsinformatik*. 11. Aufl, Berlin: Springer Gabler. (Springer-Lehrbuch). ISBN 978-3-642-30514-6.
- Microsoft (2017): Office VBA language reference. URL: <https://msdn.microsoft.com/en-us/library/office/gg264383.aspx> (Zugriff: 13.02.2017).
- Microsoft (2017): Spreadsheet Software Programs | Excel Free Trial. URL: <https://products.office.com/en-us/excel> (Zugriff: 13.02.2017).
- Moxter, A.; Riebel, P.; Wittmann, W.; Muscheid, W. (1954): *Die Elastizität des Betriebes: eine produktions- und marktwirtschaftliche Untersuchung*. Westdeutscher Verlag.
- Männel, W. (1979): Eignung von Produktionsanlagen. In: *HWProd* (Hrsg) Schäffer-Poeschel, Stuttgart, S. 1465–1481.
- Nebl, T. (2011): *Produktionswirtschaft*. 7., vollständig überarbeitete und erweiterte Auflage, München: Oldenbourg Verlag. (Lehr- und Handbücher der Betriebswirtschaftslehre). ISBN 978-3-486-59669-4.
- Nyhuis, P.; Wiendahl, H.-P. (2012): *Logistische Kennlinien: Grundlagen, Werkzeuge und Anwendungen*. 3. Auflage, Berlin Heidelberg Dordrecht London New York: Springer Vieweg. (VDI-Buch). ISBN 978-3-540-92838-6.
- Pearson, K. (1901): LIII. *On Lines and Planes of Closest Fit to Systems of Points in Space*. In: *Philosophical Magazine Series 6*, Jg. 2, Nr. 11, S. 559–572.
- Reisch, M. (2007): *Elektronische Bauelemente: Funktion, Grundsaltungen, Modellierung mit SPICE*. 2., vollst. neu bearb. Aufl, Berlin: Springer. ISBN 978-3-540-34014-0.

- Rinkenburger, R. (2009): Einführung in die Explorative Faktorenanalyse. In: Schwaiger, M.; Meyer, A. (Hrsg.): Theorien und Methoden der Betriebswirtschaft: Handbuch für Wissenschaftler und Studierende. München: Vahlen. ISBN 978-3-8006-3613-6, S. 455–476.
- Roth, P. L. (1994): Missing Data: A Conceptual Review for Applied Psychologists. In: Personnel Psychology, Jg. 47, Nr. 3, S. 537–560.
- Schutt, R.; O’Neil, C. (2013): Doing Data Science. First edition, Beijing ; Sebastopol: O’Reilly Media. ISBN 978-1-4493-5865-5.
- Shakhnarovich, G.; Indyk, P.; Darrell, T. (2005): Introduction. In: Shakhnarovich, G.; Darrell, T.; Indyk, P. (Hrsg.): Nearest-Neighbor Methods in Learning and Vision: Theory and Practice. Cambridge, Mass: MIT Press. ISBN 978-0-262-19547-8 (Neural information processing series), S. 1–12.
- Sharp, H.; Feldt, A. (1959): Some Factors in a Probability Sample Survey of a Metropolitan Community. In: American Sociological Review, Jg. 24, Nr. 5, S. 650.
- Shearer, C. (2000): The CRISP-DM Model: The New Blueprint for Data Mining. In: Journal of Data Warehousing, Jg. 5, Nr. 4, S. 13–22.
- Shlens, J. (2014): A Tutorial on Principal Component Analysis. In: arXiv preprint arXiv:1404.1100.
- Shmueli, G. (2010): To Explain or to Predict? In: Statistical Science, Jg. 25, Nr. 3, S. 289–310.
- Stewart, D. W. (1981): The Application and Misapplication of Factor Analysis in Marketing Research. In: Journal of Marketing Research, Jg. 18, Nr. 1, S. 51.
- Swoboda, P. (1964): Die betriebliche Anpassung als Problem des betrieblichen Rechnungswesens. Wiesbaden: Gabler Verlag. ISBN 978-3-663-13555-5.
- Tukey, J. W. (1977): Exploratory Data Analysis. Reading, Mass: Addison-Wesley Pub. Co. (Addison-Wesley Series in Behavioral Science). ISBN 978-0-201-07616-5.
- Weber, J.; Steven, M. (2017): Definition >>Kostenremanenz<< | Gabler Wirtschaftslexikon. URL: <http://wirtschaftslexikon.gabler.de/Archiv/7702/kostenremanenz-v6.html> (Zugriff: 25.04.2017).
- Wilcox, R. (2005): Kolmogorov-Smirnov Test. In: Encyclopedia of Biostatistics. Chichester, UK: John Wiley & Sons, Ltd. ISBN 978-0-470-84907-1.
- Witten, I. H.; Frank, E.; Hall, M. A. (2011): Data Mining: Practical Machine Learning Tools and Techniques. 3rd ed, Burlington, MA: Morgan Kaufmann. (Morgan Kaufmann Series in Data Management Systems). ISBN 978-0-12-374856-0.
- Yuan, B.; Wang, X. .; Morris, T. (2000): Software Analyser Design Using Data Mining Technology for Toxicity Prediction of Aqueous Effluents. In: Waste Management, Jg. 20, Nr. 8, S. 677–686.
- Zsifkovits, H. E. (2013): Logistik. Konstanz: UVK-Verl.-Ges. (UTB Betriebswirtschaftslehre, Technik, Ingenieurwesen; 3673). ISBN 978-3-8252-3673-1.

## Anhang

Tabelle 27: VM: Einteilung der Attribute in MAR und MCAR (fehlenden Daten in %) <sup>318</sup>

OK	fehlende Werte					
	MAR (nonrandom)		MCAR (random)			
Matnummer	Pb_2	98,2%	Pb_1	0,4%	Durchfluss MW	5,3%
Mat,- Text	Pb_Menge2	98,2%	BaCO3_1	0,4%	Durchfluss Stabw	5,3%
Losnummer	BaCO3_2	89,8%	SrCO3_1	0,4%	Leistung MAX	5,3%
ChargeSAP	BaCO3_Menge2	89,8%	CaCO3_1	0,4%	Leistung MIN	5,3%
Chargennummer	SrCO3_2	98,7%	Pb3O4_1	0,4%	Leistung MW	5,3%
IstStart	SrCO3_Menge2	98,7%	Y-Acetat_1	0,4%	Leistung Stabw	5,3%
IstEnde	CaCO3_2	96,9%	TiO2 HT05SS_1	0,4%	Schlickerdichte MAX	5,3%
Mstart	CaCO3_Menge2	96,9%	Sprühturm	2,2%	Schlickerdichte MIN	5,3%
Mlief	Pb3O4_2	98,2%	pH Wert Schlicker	0,4%	Schlickerdichte MW	5,3%
Pb_Menge1	Pb3O4_Menge2	98,2%	Schüttdichte Granulat	2,2%	Schlickerdichte Stabw	5,3%
BaCO3_Menge1	TiO2 HT05SS_2	83,2%	Granulatfeuchte	2,2%	Schlickerdruck MAX	5,3%
SrCO3_Menge1	TiO2 HT05SS_Menge2	83,2%	Differenzdruck Filter	2,2%	Schlickerdruck MIN	5,3%
CaCO3_Menge1			Ringtemperatur 1	0,0%	Schlickerdruck MW	5,3%
Pb3O4_Menge1			Ringtemperatur 2	0,0%	Schlickerdruck Stabw	5,3%
Y-Acetat_Menge1			Ablufttemperatur	0,4%	Schlickertemp, MAX	5,3%
TiO2 HT05SS_Menge1			Gewölbeunterdruck	3,5%	Schlickertemp, MIN	5,3%
A40_Menge1			Heizgruppe 7	2,2%	Schlickertemp, MW	5,3%
Optapix_Menge1			Heizgruppe 9	2,2%	Schlickertemp, Stabw	5,3%
Fryma			Heizgruppe 11	2,2%	Anzahl Segmente	5,3%
UD			Heizgruppe 13	2,2%	Segmentzeit MAX	5,3%
Kugelmenge			Differenzdruck MAX	5,3%	Segmentzeit MIN	5,3%
Durchfluss			Differenzdruck MIN	5,3%	MW Dauer je Segment	5,3%
Siebrückstand feucht			Differenzdruck MW	5,3%	Stabw Segmentzeit	5,3%
d10 (VM)			Differenzdruck Stabw	5,3%	Gesamtzeit Segmente	5,3%
d50 (VM)			Durchfluss MAX	5,3%	Stabw Segmentzeit	5,3%
d90 (VM)			Durchfluss MIN	5,3%	Gesamtzeit Segmente	5,3%

<sup>318</sup> Quelle: eigene Darstellung

Tabelle 28: NM: Einteilung der Attribute in MAR und MCAR (fehlenden Daten in %) <sup>319</sup>

OK	fehlende Werte					
	MAR (nonrandom)		MCAR (random)			
Matnummer	Umsatz_2	99,1%	Fryma	0,4%	TiO2	0,4%
Mat,- Text	Umsatz_Menge2	99,1%	Datum	0,4%	Y2O3	0,4%
Losnummer	Methocell A4C_1	99,6%	Kugelmenge	0,4%	SPHT3 MW	0,4%
ChargeSAP	Methocell A4C_Menge2	99,6%	Durchfluss	0,4%	Symm3 MW	0,4%
Chargennummer	Gleitwert	67,9%	Schlickertemperatur	0,4%	b/I3 MW	0,4%
IstStart	Grünteilfestigkeit	67,9%	d10 (NM)	0,4%		
IstEnde	Kohäsion	67,9%	d50 (NM)	0,4%		
Mstart	Pressdruck	67,9%	d90 (NM)	0,4%		
Mlief	Radiale Auffederung	67,9%	Schlickerdichte	0,4%		
Umsatz_1	Schüttwinkel	67,9%	pH Wert Schlicker	1,3%		
Umsatz_Menge1			Glühverlust	0,9%		
Methocell A4C_Menge1			Sinterdichte	0,9%		
Contraspum K1012_Menge1			R25	2,7%		
Sprühturm			Rmin	2,7%		
Datum			T-Rmin	2,7%		
Auslaufzeit			Tb	2,7%		
Schüttdichte Granulat			Tcorr	2,7%		
Granulatfeuchte			BaO	0,4%		
Differenzdruck Filter			CaO	0,4%		
d10 Granulat (NM)			MnO	0,4%		
d50 Granulat (NM)			Na2O	0,4%		
d90 Granulat (NM)			PbO	0,4%		
R25 Widerstand			SiO2	0,4%		
R25 Vergl,%			SrO	0,4%		

<sup>319</sup> Quelle: eigene Darstellung

Tabelle 29: Attribute nach der Zusammenführung (VM+NM)<sup>320</sup>

Kugelmenge	Durchfluss	Schlickertemperatur	d10 (NM)	d50 (NM)
d90 (NM)	Auslaufzeit	Schuettdichte Granulat	Granulatfeuchte	Differenzdruck Filter
d10 Granulat (NM)	d50 Granulat (NM)	d90 Granulat (NM)	Schlickerdichte	pH Wert Schlicker
Gluehverlust	Sinterdichte	R25 Widerstand	R25 Vergl.%	R25
Rmin	T-Rmin	Tb	Tcorr	Gleitwert
Gruenteilfestigkeit	Kohaesion	Pressdruck	Radiale Auffederung	Schuettwinkel
SPHT3 MW	Symm3 MW	b I3 MW	Contraspum K1012_Menge1_NM_Gehalt	Methocell A4C_NM_Gehalt
BaO_NM_Gehalt	CaO_NM_Gehalt	MnO_NM_Gehalt	Na2O_NM_Gehalt	PbO_NM_Gehalt
SiO2_NM_Gehalt	SrO_NM_Gehalt	TiO2_NM_Gehalt	Y2O3_NM_Gehalt	Pb_Me
BaCO3	SrCO3	CaCO3	Pb3O4	Y-Acet
TiO2	A40_Me	Optapi	Kugelmenge_VM	Siebrueckstand feucht_VM
d10 (VM)_VM	d50 (VM)_VM	d90 (VM)_VM	pH Wert Schlicker_VM	Schuettdichte Granulat_VM
Granulatfeuchte_VM	Differenzdruck Filter_VM	Ablufttemperatur_VM	Gewoelbeunterdruck_VM	Heizgruppe 7_VM
Heizgruppe 9_VM	Heizgruppe 11_VM	Heizgruppe 13_VM	Differenzdruck MAX_VM	Differenzdruck MIN_VM
Differenzdruck MW_VM	Differenzdruck Stabw_VM	Durchfluss MAX_VM	Durchfluss MIN_VM	Durchfluss MW_VM
Durchfluss Stabw_VM	Leistung MAX_VM	Leistung MIN_VM	Leistung MW_VM	Leistung Stabw_VM
Schlickerdichte MAX_VM	Schlickerdichte MIN_VM	Schlickerdichte MW_VM	Schlickerdichte Stabw_VM	Schlickerdruck MAX_VM
Schlickerdruck MIN_VM	Schlickerdruck MW_VM	Schlickerdruck Stabw_VM	Schlickertemp. MAX_VM	Schlickertemp. MIN_VM
Schlickertemp. MW_VM	Schlickertemp. Stabw_VM	Anzahl Segmente_VM	Segmentzeit MAX_VM	Segmentzeit MIN_VM
MW Dauer je Segment_VM	Stabw Segmentzeit_VM	Gesamtzeit Segmente_VM	Ringtemperatur_VM	

<sup>320</sup> Quelle: eigene Darstellung

Tabelle 30: Statistische Maße der VM+NM (mit Ausreißern)<sup>321</sup>

Bezeichnung	$\tilde{s}^2$	$\tilde{s}$	$\gamma$	$g_m$	Korrelationskoeffizient mit R25 Vergl. %
Kugelmenge	0,22	0,47	50	-5,3	0,20
Durchfluss	0,87	0,93	-3	-0,2	-0,33
Schlickertemperatur	7,39	2,72	-3	-0,4	0,21
d10 (NM)	0,00	0,01	0	0,3	0,31
d50 (NM)	0,00	0,03	11	2,0	0,00
d90 (NM)	0,02	0,15	0	1,0	-0,18
Auslaufzeit	4,85	2,20	0	1,7	-0,11
Schuettdichte Granulat	0,00	0,02	0	-0,9	0,08
Granulatfeuchte	0,00	0,01	-3	0,6	-0,23
Differenzdruck Filter	0,00	0,03	-3	-0,2	0,05
d10 Granulat (NM)	10,71	3,27	-3	0,4	-0,15
d50 Granulat (NM)	31,82	5,64	-3	0,1	-0,14
d90 Granulat (NM)	37,95	6,16	-3	0,0	-0,13
Schlickerdichte	0,00	0,02	-1	-1,0	-0,25
pH Wert Schlicker	0,03	0,17	-4	-0,1	0,23
Gluehverlust	0,00	0,04	1	1,1	0,14
Sinterdichte	0,00	0,01	-4	-0,3	0,37
R25 Widerstand	6,19	2,49	-2	0,5	0,68
R25 Vergl. %	68,09	8,25	-3	-0,2	1,00
R25	2,75	1,66	-3	0,3	0,46
Rmin	1,62	1,27	-3	0,3	0,41
T-Rmin	6,21	2,49	-3	-0,2	0,24
Tb	0,62	0,79	-2	0,0	0,07
Tcorr	5,45	2,34	-3	-0,2	0,21
Gleitwert	0,00	0,02	-3	-0,5	0,09
Gruenteilfestigkeit	0,47	0,69	-2	0,6	-0,09
Kohaesion	0,01	0,12	-3	0,1	0,01
Pressdruck	18,36	4,29	-3	0,0	-0,07
Radiale Auffederung	0,00	0,00	-4	0,1	-0,33
Schuettwinkel	1,00	1,00	0	1,2	0,07
SPHT3 MW	0,00	0,00	3	-0,9	-0,07
Symm3 MW	0,00	0,00	-2	-0,7	0,12
b I3 MW	0,00	0,01	-3	-0,5	0,32
Contraspum					
K1012_Menge1_NM_Gehalt	0,00	0,00	36	6,4	0,03
Methocell A4C_NM_Gehalt	0,02	0,13	210	14,7	0,04
BaO_NM_Gehalt	0,00	0,04	-3	-0,2	0,22
CaO_NM_Gehalt	0,00	0,02	-2	0,5	0,08
MnO_NM_Gehalt	0,00	0,00	-3	-1,3	0,03
Na2O_NM_Gehalt	0,00	0,00	-4	-0,4	-0,45
PbO_NM_Gehalt	0,00	0,03	2	1,1	-0,14
SiO2_NM_Gehalt	0,00	0,00	-2	0,0	-0,20
SrO_NM_Gehalt	0,00	0,01	8	2,9	0,01
TiO2_NM_Gehalt	0,00	0,03	-3	0,3	-0,16
Y2O3_NM_Gehalt	0,00	0,00	-3	0,2	0,06

<sup>321</sup> Quelle: eigene Darstellung

Fortsetzung von Tabelle 30: Statistische Zusammenhangsmaße der VM+NM<sup>322</sup>

Bezeichnung	$\tilde{s}^2$	$\tilde{s}$	$\gamma$	$g_m$	Korrelationskoeffizient mit R25 Vergl. %
Pb Me	0,00	0,00	101	9,7	-0,09
BaCO3	0,00	0,00	35	-4,6	-0,06
SrCO3	0,00	0,00	0	1,5	0,21
CaCO3	0,00	0,00	-3	-0,4	0,02
Pb3O4	0,00	0,00	101	9,7	-0,09
Y-Acet	0,00	0,00	-4	0,3	0,44
TiO2	0,00	0,00	4	-1,6	0,11
A40 Me	0,00	0,00	-4	-0,8	-0,37
Optapi	0,00	0,00	-4	0,8	0,37
Kugelmenge_VM	0,01	0,11	65	-8,3	-0,18
Siebrueckstand feucht_VM	3198,41	56,55	-3	0,4	-0,08
d10 (VM)_VM	0,00	0,01	25	3,4	0,14
d50 (VM)_VM	0,00	0,05	-3	0,0	0,22
d90 (VM)_VM	0,02	0,14	-3	0,0	0,28
pH Wert Schlicker_VM	0,02	0,15	12	2,0	0,23
Schuettdichte Granulat_VM	0,00	0,01	0	-0,9	-0,12
Granulatfeuchte_VM	0,00	0,01	2	1,7	-0,02
Differenzdruck Filter_VM	0,00	0,07	-4	-0,4	0,16
Ablufttemperatur_VM	117,64	10,85	20	-3,2	0,12
Gewoelbeunterdruck_VM	0,00	0,05	0	0,7	-0,16
Heizgruppe 7_VM	250,38	15,82	122	-10,0	0,04
Heizgruppe 9_VM	187,01	13,68	131	-10,6	0,04
Heizgruppe 11_VM	152,56	12,35	133	-10,7	0,03
Heizgruppe 13_VM	162,04	12,73	110	-9,3	0,05
Differenzdruck MAX_VM	0,00	0,05	39	1,2	-0,13
Differenzdruck MIN_VM	0,00	0,04	37	-5,6	-0,10
Differenzdruck MW_VM	0,00	0,03	38	-5,2	-0,13
Differenzdruck Stabw_VM	0,00	0,01	161	12,2	-0,03
Durchfluss MAX_VM	0,04	0,20	3	2,0	-0,15
Durchfluss MIN_VM	0,69	0,83	186	-13,4	0,09
Durchfluss MW_VM	0,03	0,17	99	-8,2	0,09
Durchfluss Stabw_VM	0,08	0,27	196	13,9	-0,07
Leistung MAX_VM	495462,94	703,89	-1	1,4	0,08
Leistung MIN_VM	366392,06	605,30	-2	1,3	0,11
Leistung MW_VM	412585,33	642,33	-2	1,2	0,10
Leistung Stabw_VM	4502,54	67,10	41	5,1	0,01
Schlickerdichte MAX_VM	0,00	0,01	100	-8,3	0,09
Schlickerdichte MIN_VM	0,01	0,09	195	-13,9	0,06
Schlickerdichte MW_VM	0,00	0,02	79	-8,6	0,12
Schlickerdichte Stabw_VM	0,00	0,02	197	14,0	-0,05
Schlickerdruck MAX_VM	0,01	0,09	-4	0,1	-0,16
Schlickerdruck MIN_VM	0,00	0,06	-2	0,4	-0,06
Schlickerdruck MW_VM	0,01	0,07	-3	0,2	-0,13
Schlickerdruck Stabw_VM	0,00	0,02	6	2,6	-0,10
Schlickertemp, MAX_VM	2,96	1,72	22	3,7	-0,03
Schlickertemp, MIN_VM	2,04	1,43	1	1,6	-0,06
Schlickertemp, MW_VM	1,95	1,40	1	1,6	-0,04

<sup>322</sup> Quelle: eigene Darstellung

Fortsetzung von Tabelle 30: Statistische Zusammenhangsmaße der VM+NM<sup>323</sup>

Bezeichnung	$\tilde{s}^2$	$\tilde{s}$	$\gamma$	$g_m$	Korrelationskoeffizient mit R25 Vergl. %
Schlickertemp, Stabw_VM	0,05	0,23	93	8,4	0,00
Anzahl Segmente_VM	0,31	0,56	3	2,3	0,04
Segmentzeit MAX_VM	1183399,42	1087,84	2	-2,0	-0,30
Segmentzeit MIN_VM	19803659,69	4450,13	-3	-1,2	-0,03
MW Dauer je Segment_VM	6965942,86	2639,31	-3	-1,2	-0,10
Stabw Segmentzeit_VM	438016*109	20928831	-1	-1,4	-0,32
Gesamtzeit Segmente_VM	845274,61	919,39	4	-2,1	-0,29
Ringtemperatur_VM	8,00	2,83	-4	0,0	0,30

<sup>323</sup> Quelle: eigene Darstellung

Tabelle 31: Ausreißer in Anzahl und Prozent je Attribut<sup>324</sup>

Attribut	Ausreißer	
	Anzahl	Prozent
Kugelmenge	1	0,5%
Durchfluss	0	0,0%
Schlickertemperatur	1	0,5%
d10 (NM)	3	1,4%
d50 (NM)	7	3,2%
d90 (NM)	6	2,8%
Auslaufzeit	15	6,9%
Schuettdichte Granulat	3	1,4%
Granulatfeuchte	0	0,0%
Differenzdruck Filter	0	0,0%
d10 Granulat (NM)	0	0,0%
d50 Granulat (NM)	0	0,0%
d90 Granulat (NM)	3	1,4%
Schlickerdichte	8	3,7%
pH Wert Schlicker	0	0,0%
Gluehverlust	5	2,3%
Sinterdichte	0	0,0%
R25 Widerstand	9	4,1%
R25 Vergl,%	1	0,5%
R25	0	0,0%
Rmin	0	0,0%
T-Rmin	6	2,8%
Tb	11	5,0%
Tcorr	16	7,3%
Gleitwert	10	4,6%
Gruenteilfestigkeit	11	5,0%
Kohaesion	3	1,4%
Pressdruck	3	1,4%
Radiale Auffederung	0	0,0%
Schuettwinkel	0	0,0%
SPHT3 MW	9	4,1%
Symm3 MW	7	3,2%
b_l3 MW	3	1,4%
Contraspum		
K1012_Menge1_NM_Gehalt	5	2,3%
Methocell A4C_NM_Gehalt	2	0,9%
BaO_NM_Gehalt	2	0,9%
CaO_NM_Gehalt	9	4,1%
MnO_NM_Gehalt	52	23,9%
Na2O_NM_Gehalt	0	0,0%
PbO_NM_Gehalt	13	6,0%
SiO2_NM_Gehalt	1	0,5%
SrO_NM_Gehalt	34	15,6%
TiO2_NM_Gehalt	3	1,4%
Y2O3_NM_Gehalt	12	5,5%
Pb_Me	3	1,4%
BaCO3	4	1,8%

<sup>324</sup> Quelle: eigene Darstellung

Fortsetzung von Tabelle 31: Ausreißer in Anzahl und Prozent je Attribut<sup>325</sup>

Attribut	Ausreißer	
	Anzahl	Prozent
SrCO3	16	7,3%
CaCO3	0	0,0%
Pb3O4	3	1,4%
Y-Acet	0	0,0%
TiO2	24	11,0%
A40_Me	0	0,0%
Optapi	0	0,0%
Kugelmenge_VM	3	1,4%
Siebrueckstand feucht_VM	1	0,5%
d10 (VM)_VM	2	0,9%
d50 (VM)_VM	2	0,9%
d90 (VM)_VM	0	0,0%
pH Wert Schlicker_VM	1	0,5%
Schuettdichte Granulat_VM	3	1,4%
Granulatfeuchte_VM	4	1,8%
Differenzdruck Filter_VM	0	0,0%
Ablufttemperatur_VM	1	0,5%
Gewoelbeunterdruck_VM	1	0,5%
Heizgruppe 7_VM	23	10,6%
Heizgruppe 9_VM	15	6,9%
Heizgruppe 11_VM	28	12,8%
Heizgruppe 13_VM	21	9,6%
Differenzdruck MAX_VM	13	6,0%
Differenzdruck MIN_VM	17	7,8%
Differenzdruck MW_VM	13	6,0%
Differenzdruck Stabw_VM	20	9,2%
Durchfluss MAX_VM	12	5,5%
Durchfluss MIN_VM	5	2,3%
Durchfluss MW_VM	10	4,6%
Durchfluss Stabw_VM	11	5,0%
Leistung MAX_VM	1	0,5%
Leistung MIN_VM	21	9,6%
Leistung MW_VM	3	1,4%
Leistung Stabw_VM	12	5,5%
Schlickerdichte MAX_VM	4	1,8%
Schlickerdichte MIN_VM	11	5,0%
Schlickerdichte MW_VM	3	1,4%
Schlickerdichte Stabw_VM	54	24,8%
Schlickerdruck MAX_VM	0	0,0%
Schlickerdruck MIN_VM	6	2,8%
Schlickerdruck MW_VM	0	0,0%
Schlickerdruck Stabw_VM	18	8,3%
Schlickertemp, MAX_VM	19	8,7%
Schlickertemp, MIN_VM	19	8,7%
Schlickertemp, MW_VM	18	8,3%
Schlickertemp, Stabw_VM	16	7,3%

<sup>325</sup> Quelle: eigene Darstellung

Fortsetzung von Tabelle 31: Ausreißer in Anzahl und Prozent je Attribut<sup>326</sup>

Attribut	Ausreißer	
	Anzahl	Prozent
Anzahl Segmente_VM	47	21,6%
Segmentzeit MAX_VM	4	1,8%
Segmentzeit MIN_VM	47	21,6%
MW Dauer je Segment_VM	49	22,5%
Stabw Segmentzeit_VM	3	1,4%
Gesamtzeit Segmente_VM	2	0,9%
Ringtemperatur_VM	0	0,0%

---

<sup>326</sup> Quelle: eigene Darstellung

Tabelle 32: MSA für einzelne Attribute (VM+NM mit Ausreißern)<sup>327</sup>

Attribute	Durchgang			
	1	2	3	4
Kugelmenge	0,49			
Durchfluss	0,76	0,86	0,86	0,86
Schlickertemperatur	0,70	0,80	0,82	0,82
d10 (NM)	0,59	0,61	0,66	0,65
d50 (NM)	0,60	0,53	0,52	0,52
d90 (NM)	0,55	0,62	0,61	0,61
Auslaufzeit	0,49			
Schuettdichte Granulat	0,55	0,50	0,48	
Granulatfeuchte	0,64	0,59	0,59	0,59
Differenzdruck Filter	0,67	0,74	0,75	0,76
d10 Granulat (NM)	0,53	0,57	0,56	0,56
d50 Granulat (NM)	0,55	0,59	0,59	0,59
d90 Granulat (NM)	0,60	0,65	0,64	0,64
Schlickerdichte	0,62	0,64	0,64	0,64
pH Wert Schlicker	0,83	0,82	0,83	0,83
Gluehverlust	0,73	0,73	0,73	0,75
Sinterdichte	0,80	0,91	0,91	0,91
R25 Widerstand	0,74	0,64	0,64	0,64
R25	0,49			
Rmin	0,48			
T-Rmin	0,81	0,77	0,77	0,78
Tb	0,73	0,75	0,75	0,75
Tcorr	0,64	0,72	0,72	0,72
Gleitwert	0,40			
Gruenteilfestigkeit	0,34			
Kohaesion	0,48			
Pressdruck	0,57	0,84	0,84	0,84
Radiale Auffederung	0,62	0,82	0,82	0,81
Schuettwinkel	0,46			
SPHT3 MW	0,54	0,57	0,56	0,56
Symm3 MW	0,64	0,68	0,67	0,67
b I3 MW	0,71	0,69	0,69	0,68
Contraspum				
K1012_Menge1_NM_Gehalt	0,34			
Methocell A4C_NM_Gehalt	0,26			
BaO_NM_Gehalt	0,36			
CaO_NM_Gehalt	0,29			
MnO_NM_Gehalt	0,31			
Na2O_NM_Gehalt	0,66	0,72	0,73	0,73
PbO_NM_Gehalt	0,20			
SiO2_NM_Gehalt	0,23			
SrO_NM_Gehalt	0,12			
TiO2_NM_Gehalt	0,38			
Y2O3_NM_Gehalt	0,27			
Pb_Me	0,45			
BaCO3	0,47			
SrCO3	0,56	0,75	0,75	0,75

<sup>327</sup> Quelle: eigene Darstellung

Fortsetzung von Tabelle 32: MSA für einzelne Attribute (VM+NM mit Ausreißern)<sup>328</sup>

Attribute	Durchgang			
	1	2	3	4
CaCO3	0,58	0,58	0,58	0,58
Pb3O4	0,45			
Y-Acet	0,77	0,87	0,86	0,87
TiO2	0,54	0,76	0,76	0,76
A40_Me	0,87	0,85	0,84	0,84
Optapi	0,87	0,85	0,84	0,84
Kugelmenge_VM	0,50			
Siebueckstand feucht_VM	0,84	0,89	0,89	0,89
d10 (VM)_VM	0,40			
d50 (VM)_VM	0,62	0,69	0,69	0,69
d90 (VM)_VM	0,63	0,70	0,70	0,70
pH Wert Schlicker_VM	0,77	0,78	0,79	0,80
Schuettdichte Granulat_VM	0,64	0,75	0,76	0,76
Granulatfeuchte_VM	0,56	0,65	0,65	0,67
Differenzdruck Filter_VM	0,40			
Ablufttemperatur_VM	0,78	0,79	0,82	0,83
Gewoelbeunterdruck_VM	0,34			
Heizgruppe 7_VM	0,64	0,66	0,66	0,66
Heizgruppe 9_VM	0,59	0,62	0,62	0,61
Heizgruppe 11_VM	0,59	0,63	0,63	0,63
Heizgruppe 13_VM	0,59	0,64	0,64	0,64
Differenzdruck MAX_VM	0,47			
Differenzdruck MIN_VM	0,50			
Differenzdruck MW_VM	0,46			
Differenzdruck Stabw_VM	0,57	0,49		
Durchfluss MAX_VM	0,48			
Durchfluss MIN_VM	0,76	0,77	0,78	0,78
Durchfluss MW_VM	0,74	0,78	0,77	0,77
Durchfluss Stabw_VM	0,78	0,79	0,80	0,80
Leistung MAX_VM	0,76	0,71	0,71	0,71
Leistung MIN_VM	0,70	0,75	0,76	0,76
Leistung MW_VM	0,77	0,67	0,68	0,68
Leistung Stabw_VM	0,49			
Schlickerdichte MAX_VM	0,42			
Schlickerdichte MIN_VM	0,72	0,73	0,72	0,72
Schlickerdichte MW_VM	0,65	0,68	0,69	0,69
Schlickerdichte Stabw_VM	0,80	0,78	0,77	0,78
Schlickerdruck MAX_VM	0,76	0,77	0,83	0,83
Schlickerdruck MIN_VM	0,68	0,72	0,76	0,77
Schlickerdruck MW_VM	0,71	0,83	0,75	0,75
Schlickerdruck Stabw_VM	0,57	0,49		
Schlickertemp, MAX_VM	0,70	0,74	0,74	0,75
Schlickertemp, MIN_VM	0,70	0,73	0,73	0,73
Schlickertemp, MW_VM	0,68	0,71	0,71	0,71
Schlickertemp, Stabw_VM	0,68	0,74	0,74	0,74
Anzahl Segmente_VM	0,46			
Segmentzeit MAX_VM	0,70	0,75	0,75	0,75

<sup>328</sup> Quelle: eigene Darstellung

Fortsetzung von Tabelle 32: MSA für einzelne Attribute (VM+NM mit Ausreißern)<sup>329</sup>

Attribute	Durchgang			
	1	2	3	4
Segmentzeit MIN_VM	0,48			
MW Dauer je Segment_VM	0,54	0,73	0,75	0,74
Stabw Segmentzeit_VM	0,77	0,75	0,75	0,75
Gesamtzeit Segmente_VM	0,68	0,80	0,83	0,83
Ringtemperatur_VM	0,61	0,64	0,64	0,68

---

<sup>329</sup> Quelle: eigene Darstellung

Tabelle 33: MSA für einzelne Attribute (VM+NM ohne Ausreißer)<sup>330</sup>

Attribute	Durchgang			
	1	2	3	4
Kugelmenge	nan	0,83	0,86	0,86
Durchfluss	nan	0,85	0,87	0,87
Schlickertemperatur	nan	0,76	0,84	0,84
d10 (NM)	nan	0,63	0,68	0,67
d50 (NM)	nan	0,47		
d90 (NM)	nan	0,55	0,64	0,65
Auslaufzeit	nan	0,41		
Schuettdichte Granulat	nan	0,55	0,52	0,54
Granulatfeuchte	nan	0,66	0,65	0,65
Differenzdruck Filter	nan	0,65	0,75	0,75
d10 Granulat (NM)	nan	0,51	0,52	0,52
d50 Granulat (NM)	nan	0,56	0,57	0,58
d90 Granulat (NM)	nan	0,67	0,71	0,71
Schlickerdichte	nan	0,70	0,68	0,67
pH Wert Schlicker	nan	0,89	0,89	0,89
Gluehverlust	nan	0,77	0,79	0,79
Sinterdichte	nan	0,76	0,83	0,83
R25 Widerstand	nan	0,74	0,76	0,79
R25	nan	0,55	0,56	0,56
Rmin	nan	0,52	0,53	0,54
T-Rmin	nan	0,75	0,80	0,80
Tb	nan	0,75	0,78	0,79
Tcorr	nan	0,83	0,83	0,82
Gleitwert	nan	0,57	0,57	0,57
Gruenteilfestigkeit	nan	0,63	0,66	0,66
Kohaesion	nan	0,67	0,64	0,64
Pressdruck	nan	0,79	0,81	0,81
Radiale Auffederung	nan	0,74	0,80	0,80
Schuettwinkel	nan	0,47		
SPHT3 MW	nan	0,71	0,74	0,74
Symm3 MW	nan	0,59	0,60	0,59
b I3 MW	nan	0,67	0,71	0,71
Contraspum				
K1012_Menge1_NM_Gehalt	nan	0,40		
Methocell A4C_NM_Gehalt	nan	0,36		
BaO_NM_Gehalt	nan	0,71	0,75	0,75
CaO_NM_Gehalt	nan	0,71	0,74	0,75
MnO_NM_Gehalt	nan			
Na2O_NM_Gehalt	nan	0,71	0,73	0,72
PbO_NM_Gehalt	nan	0,55	0,58	0,58
SiO2_NM_Gehalt	nan	0,66	0,62	0,63
SrO_NM_Gehalt	nan	0,63	0,62	0,62
TiO2_NM_Gehalt	nan	0,79	0,80	0,80
Y2O3_NM_Gehalt	nan	0,63	0,61	0,61
Pb_Me	nan			
BaCO3	nan	0,70	0,68	0,68
SrCO3	nan	0,68	0,69	0,69

<sup>330</sup> Quelle: eigene Darstellung

Fortsetzung von Tabelle 33: MSA für einzelne Attribute (VM+NM ohne Ausreißer)<sup>331</sup>

Attribute	Durchgang			
	1	2	3	4
CaCO3	nan	0,53	0,57	0,58
Pb3O4	nan			
Y-Acet	nan	0,87	0,91	0,91
TiO2	nan	0,68	0,71	0,73
A40 Me	nan	0,88	0,88	0,88
Optapi	nan	0,88	0,88	0,88
Kugelmengung_VM	nan			
Siebrueckstand feucht_VM	nan	0,78	0,84	0,84
d10 (VM)_VM	nan	0,43		
d50 (VM)_VM	nan	0,69	0,70	0,70
d90 (VM)_VM	nan	0,68	0,69	0,69
pH Wert Schlicker_VM	nan	0,78	0,84	0,84
Schuettdichte Granulat_VM	nan	0,83	0,82	0,82
Granulatfeuchte_VM	nan	0,74	0,73	0,73
Differenzdruck Filter_VM	nan	0,45		
Ablufttemperatur_VM	nan	0,62	0,64	0,65
Gewoelbeunterdruck_VM	nan	0,57	0,60	0,58
Heizgruppe 7_VM	nan	0,73	0,75	0,75
Heizgruppe 9_VM	nan	0,72	0,77	0,78
Heizgruppe 11_VM	nan	0,76	0,78	0,79
Heizgruppe 13_VM	nan	0,77	0,77	0,77
Differenzdruck MAX_VM	nan	0,78	0,78	0,80
Differenzdruck MIN_VM	nan	0,73	0,76	0,79
Differenzdruck MW_VM	nan	0,75	0,78	0,78
Differenzdruck Stabw_VM	nan	0,50	0,50	
Durchfluss MAX_VM	nan	0,50	0,51	0,51
Durchfluss MIN_VM	nan	0,67	0,67	0,66
Durchfluss MW_VM	nan	0,58	0,58	0,58
Durchfluss Stabw_VM	nan	0,69	0,68	0,68
Leistung MAX_VM	nan	0,69	0,72	0,71
Leistung MIN_VM	nan	0,83	0,84	0,85
Leistung MW_VM	nan	0,70	0,74	0,75
Leistung Stabw_VM	nan	0,71	0,72	0,72
Schlickerdichte MAX_VM	nan	0,58	0,58	0,58
Schlickerdichte MIN_VM	nan	0,58	0,65	0,65
Schlickerdichte MW_VM	nan	0,64	0,64	0,64
Schlickerdichte Stabw_VM	nan	0,37		
Schlickerdruck MAX_VM	nan	0,77	0,77	0,77
Schlickerdruck MIN_VM	nan	0,69	0,71	0,72
Schlickerdruck MW_VM	nan	0,75	0,75	0,76
Schlickerdruck Stabw_VM	nan	0,75	0,76	0,75
Schlickertemp, MAX_VM	nan	0,74	0,76	0,76
Schlickertemp, MIN_VM	nan	0,66	0,75	0,75
Schlickertemp, MW_VM	nan	0,64	0,68	0,68
Schlickertemp, Stabw_VM	nan	0,33		
Anzahl Segmente_VM	nan	0,57	0,63	0,61
Segmentzeit MAX_VM	nan	0,83	0,83	0,83

<sup>331</sup> Quelle: eigene Darstellung

Fortsetzung von Tabelle 33: MSA für einzelne Attribute (VM+NM ohne Ausreißer)<sup>332</sup>

Attribute	Durchgang			
	1	2	3	4
Segmentzeit MIN_VM	nan	0,66	0,66	0,66
MW Dauer je Segment_VM	nan	0,79	0,79	0,79
Stabw Segmentzeit_VM	nan	0,84	0,84	0,85
Gesamtzeit Segmente_VM	nan			
Ringtemperatur_VM	nan	0,75	0,76	0,77

Tabelle 34: Ergebnis CL und Testdaten: VM+NM mit Ausreißern<sup>333</sup>

	Genauigkeit bei Training	Genauigkeit bei Anwendung auf Testdaten	Grenzfälle bei Anwendung auf Testdaten	Mittelwert der absoluten Differenz zwischen der vorhergesagten und Ist-Klasse
Complex Tree	17,4%	10,3%	1,3%	2,47
Medium Tree	17,9%	6,4%	0,0%	2,92
Simple Tree	17,0%	9,0%	0,0%	2,85
Linear SVM	18,3%	7,7%	1,3%	3,04
Quadratic SVM	17,4%	11,5%	0,0%	3,09
Cubic SVM	17,4%	7,7%	0,0%	3,14
Fine Gaussian SVM	12,8%	15,4%	1,3%	2,42
Medium Gaussian SVM	17,4%	12,8%	0,0%	2,59
Coarse Gaussian SVM	15,6%	7,7%	0,0%	3,26
Boosted Trees	19,7%	2,6%	0,0%	2,97
Bagged Trees	24,3%	11,5%	1,3%	2,81
RUSBoosted Trees	19,3%	5,1%	0,0%	2,82

<sup>332</sup> Quelle: eigene Darstellung<sup>333</sup> Quelle: ebda.

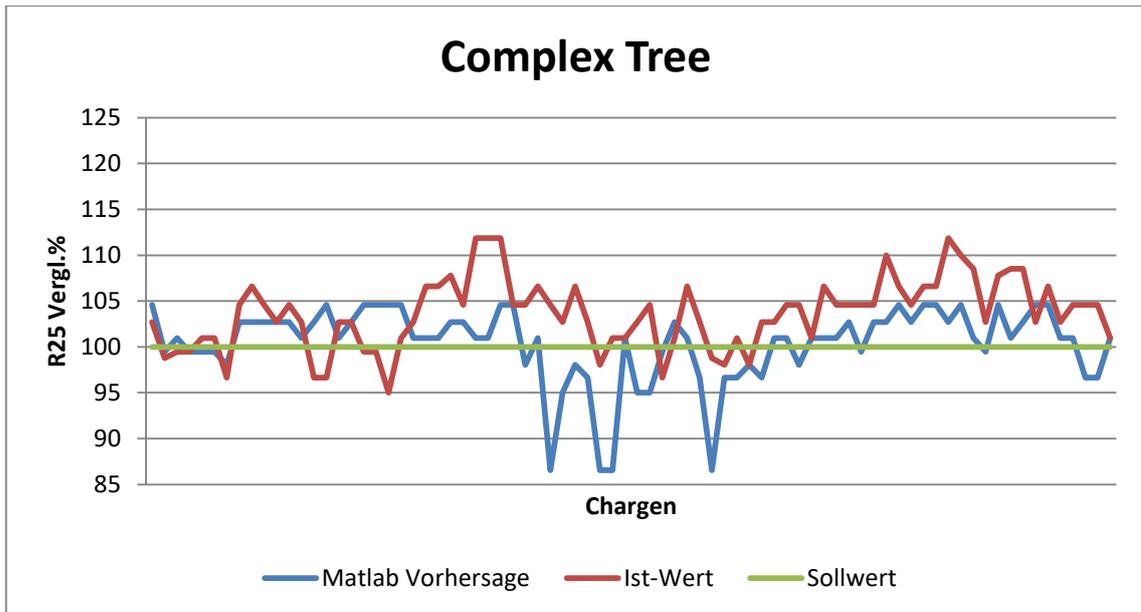


Abbildung 62: CT: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern)<sup>334</sup>

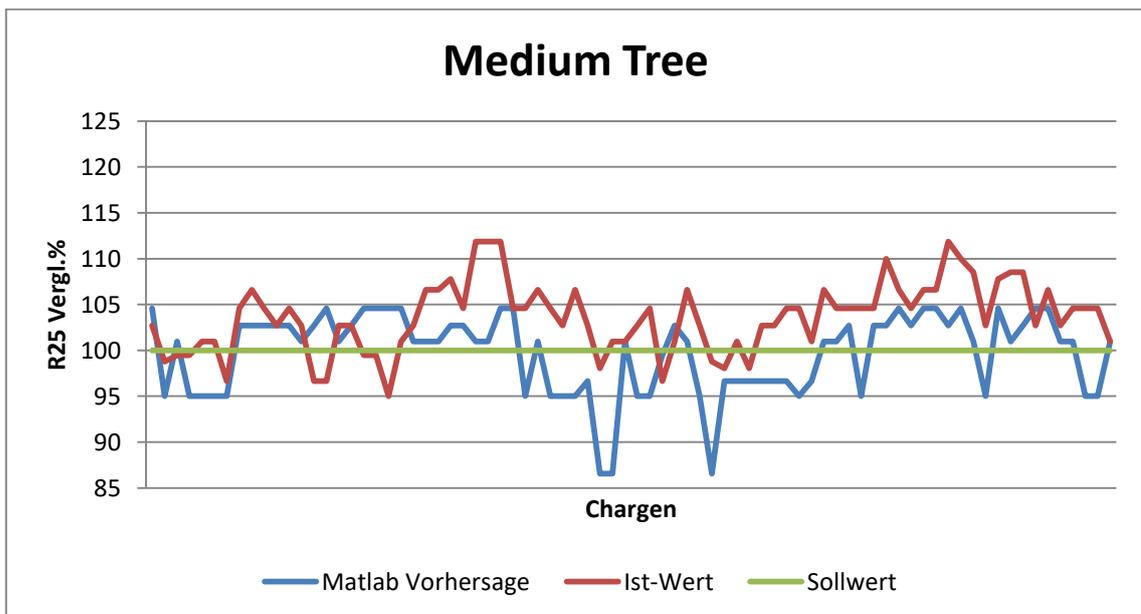


Abbildung 63: MT: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern)<sup>335</sup>

<sup>334</sup> Quelle: eigene Darstellung

<sup>335</sup> Quelle: ebda.

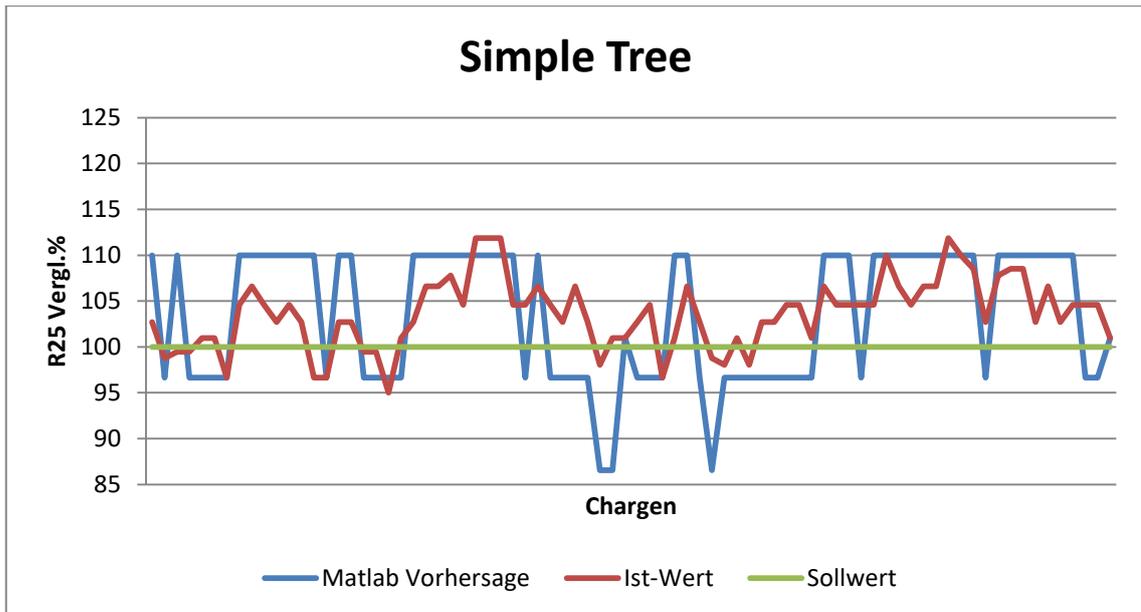


Abbildung 64: ST: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern)<sup>336</sup>

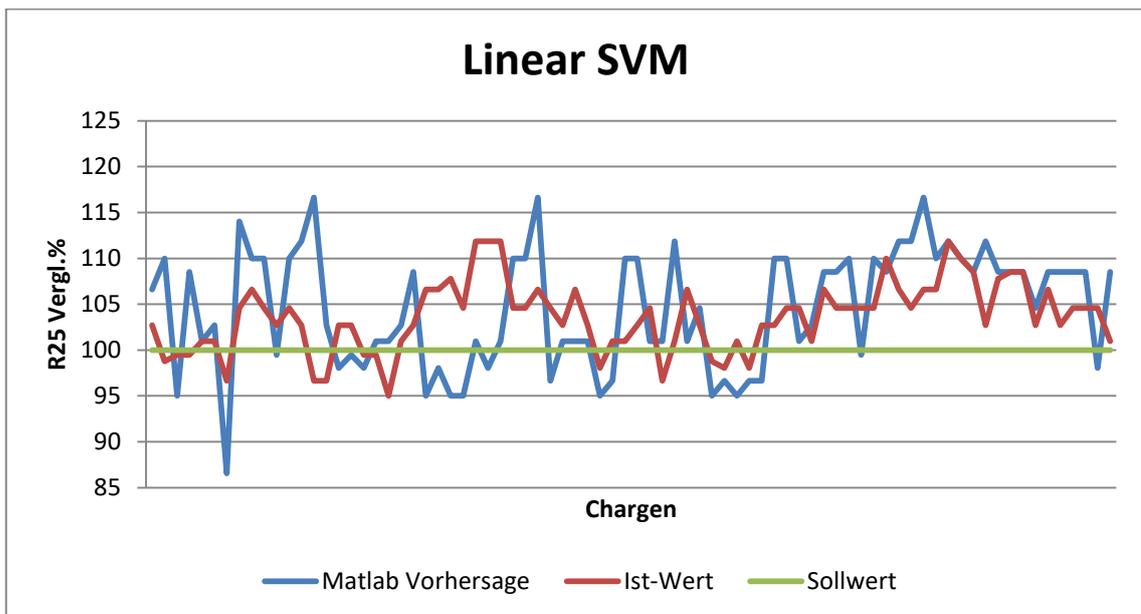


Abbildung 65: LSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern)<sup>337</sup>

<sup>336</sup> Quelle: eigene Darstellung

<sup>337</sup> Quelle: ebda.

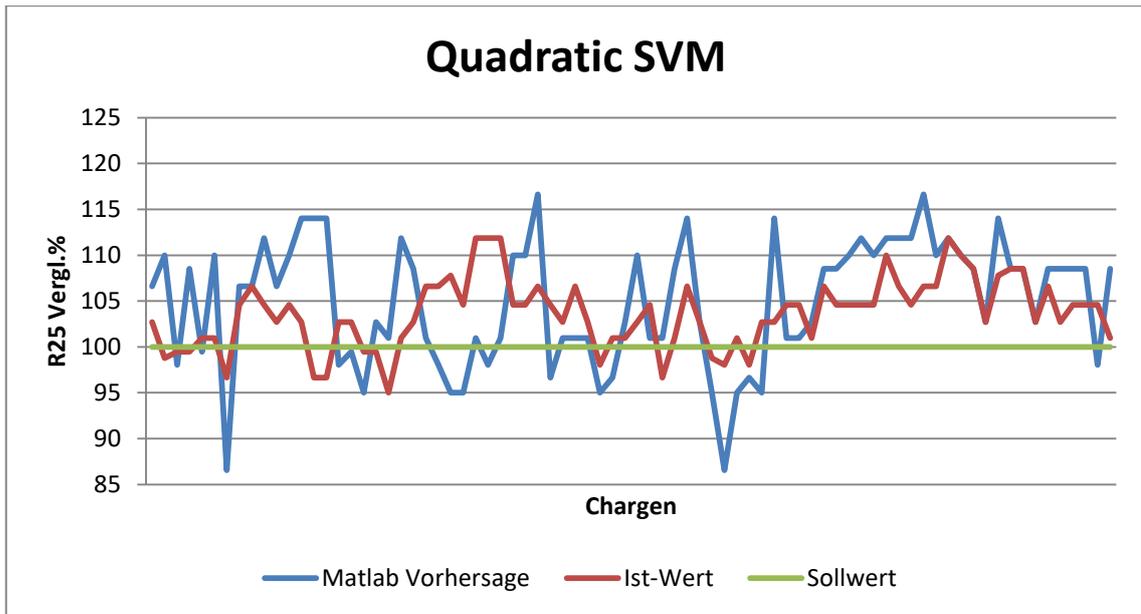


Abbildung 66: QSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern)<sup>338</sup>

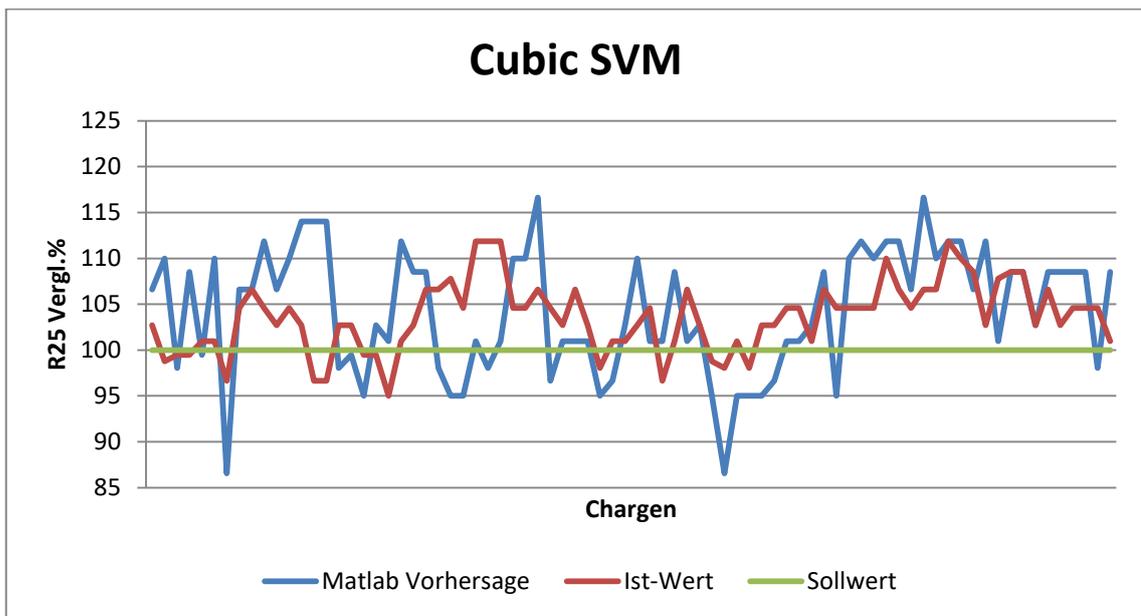


Abbildung 67: CSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern)<sup>339</sup>

<sup>338</sup> Quelle: eigene Darstellung

<sup>339</sup> Quelle: ebda.

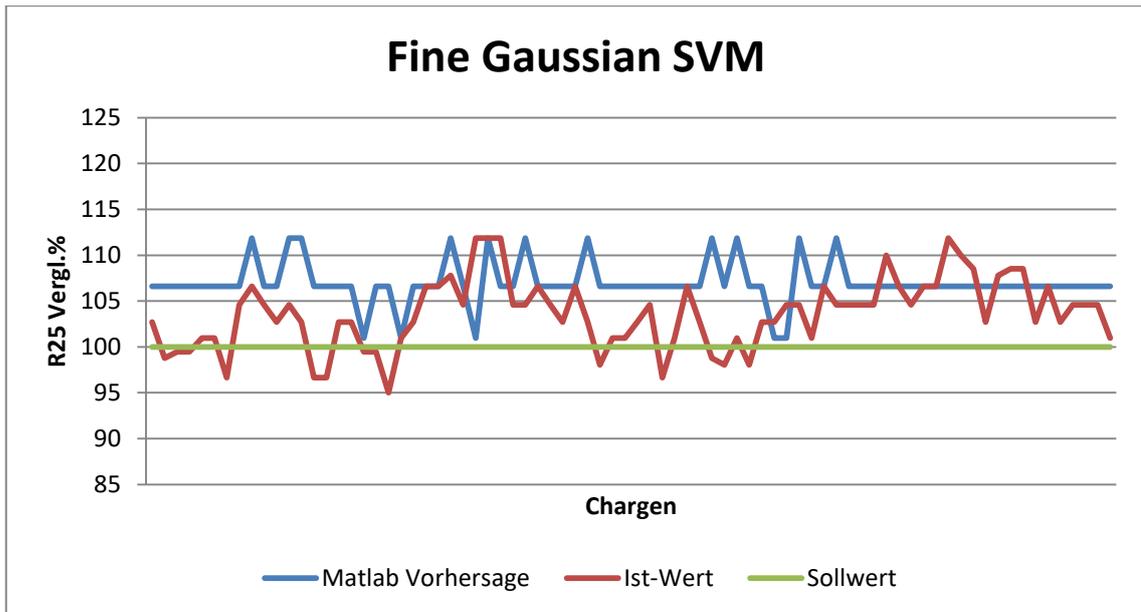


Abbildung 68: FGSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern)<sup>340</sup>

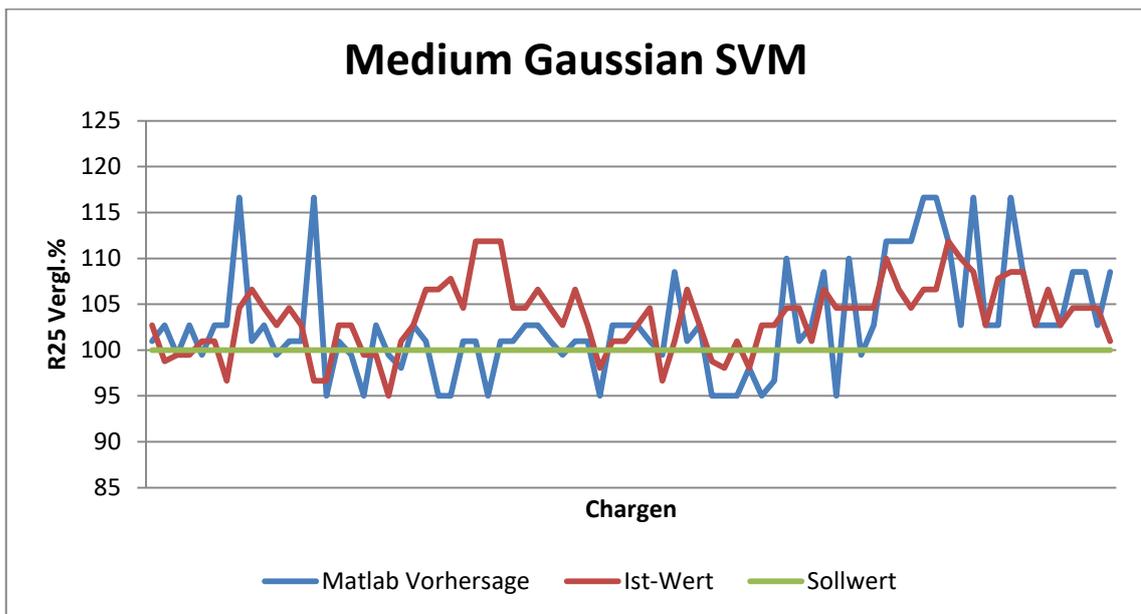


Abbildung 69: MGSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern)<sup>341</sup>

<sup>340</sup> Quelle: eigene Darstellung

<sup>341</sup> Quelle: ebda.

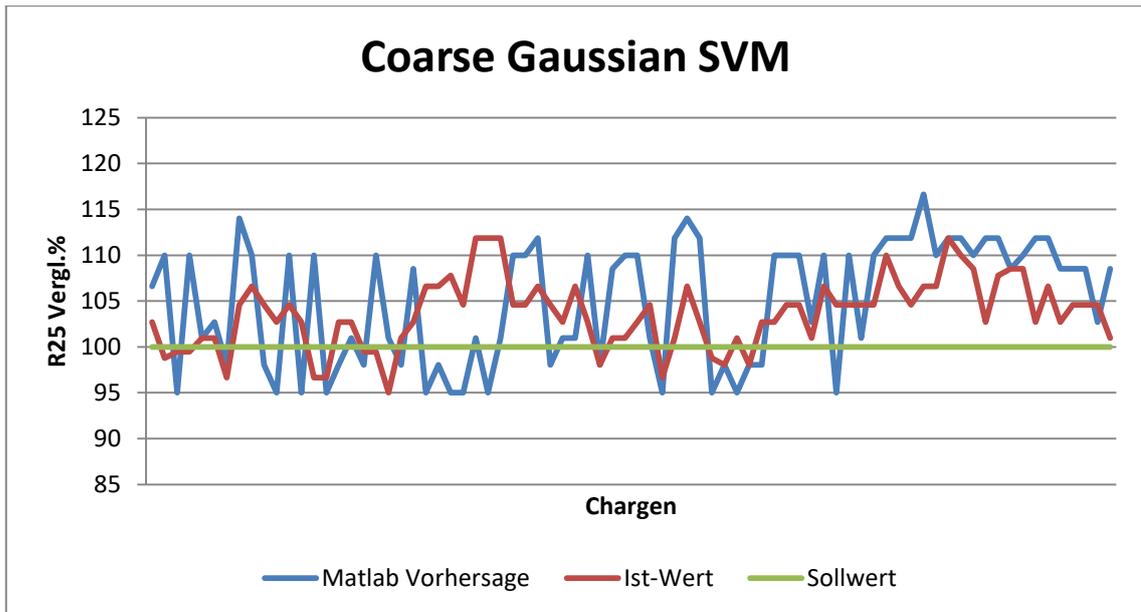


Abbildung 70: CGSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern)<sup>342</sup>

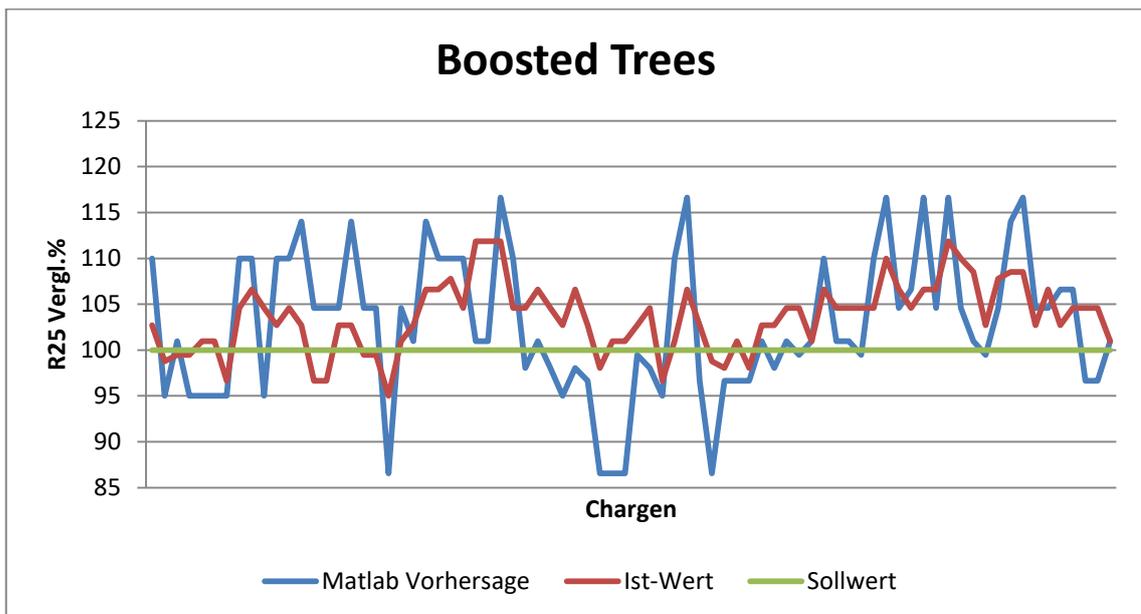


Abbildung 71: BoT: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern)<sup>343</sup>

<sup>342</sup> Quelle: eigene Darstellung

<sup>343</sup> Quelle: ebda.

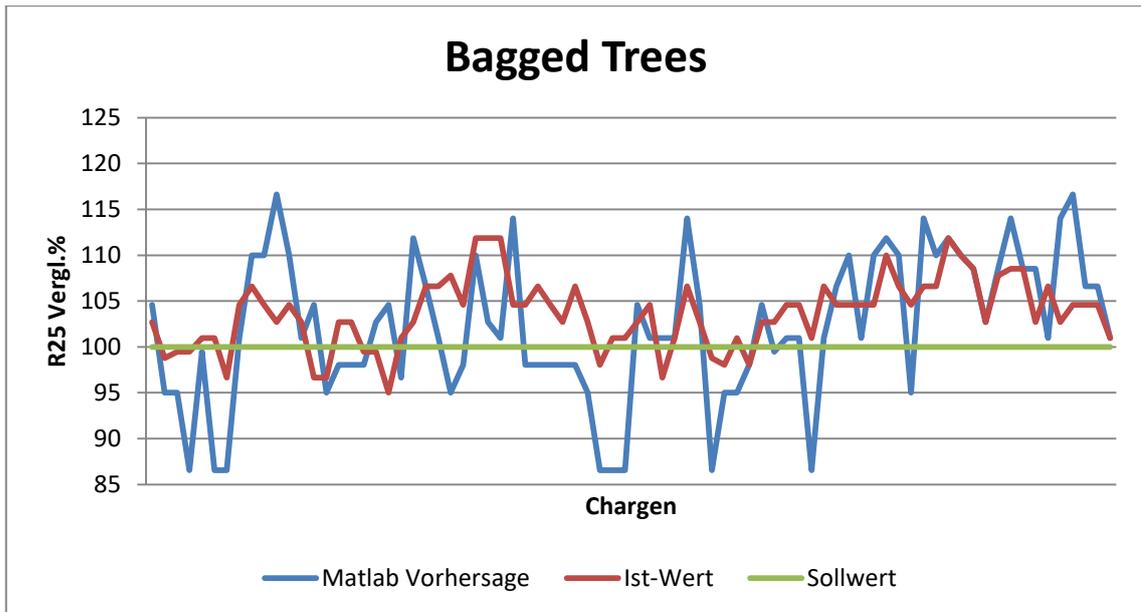


Abbildung 72: BaT: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern)<sup>344</sup>

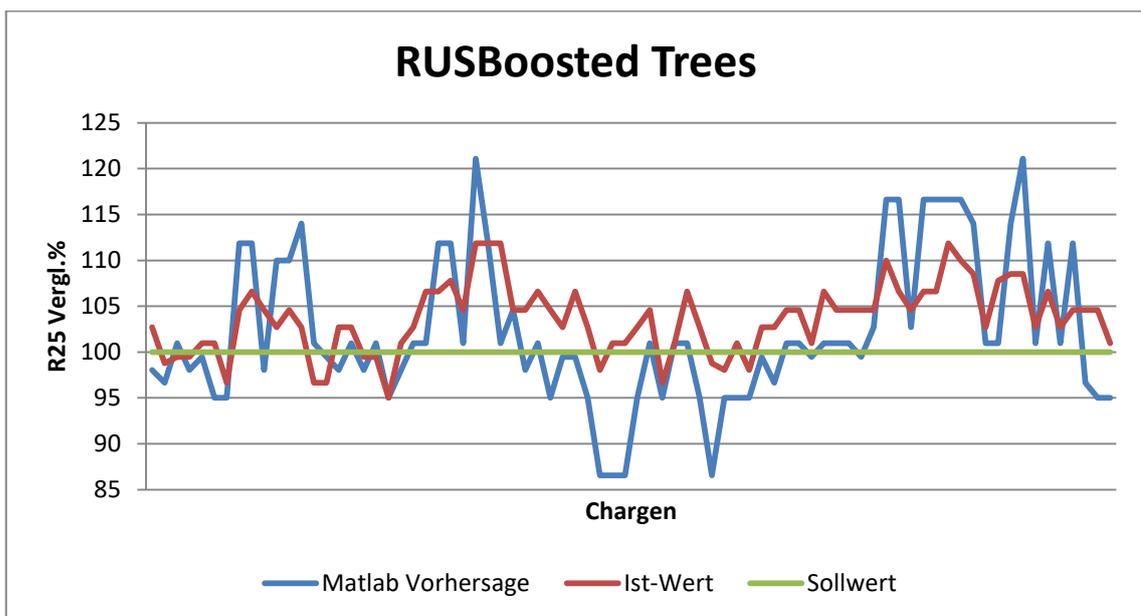


Abbildung 73: RUSBT: Vergleich Vorhersage zu Ist-Wert (VM+NM mit Ausreißern)<sup>345</sup>

<sup>344</sup> Quelle: eigene Darstellung

<sup>345</sup> Quelle: ebda.

**Medium Tree VM+NM mit Ausreißern**

```
1  if R25Widerstand<-1.36532 then node 2 elseif R25Widerstand>=-
1.36532 then node 3 else 101.0526667
2  if Durchfluss<0.144277 then node 4 elseif
Durchfluss>=0.144277 then node 5 else 88.62466667
3  if YAcet<0.238895 then node 6 elseif YAcet>=0.238895 then
node 7 else 102.9246667
4  class = 101.0526667
5  class = 88.62466667
6  if R25Widerstand<0.119414 then node 8 elseif
R25Widerstand>=0.119414 then node 9 else 98.03066667
7  if SchlickerdruckMW_VM<1.51411 then node 10 elseif
SchlickerdruckMW_VM>=1.51411 then node 11 else 106.6053333
8  if DifferenzdruckFilter<1.50558 then node 12 elseif
DifferenzdruckFilter>=1.50558 then node 13 else 96.936
9  if YAcet<-0.57182 then node 14 elseif YAcet>=-0.57182 then
node 15 else 110.0333333
10 if d90NM<-1.49482 then node 16 elseif d90NM>=-1.49482 then
node 17 else 108.566
11 if R25Widerstand<-1.15983 then node 18 elseif
R25Widerstand>=-1.15983 then node 19 else 106.6053333
12 if pHWertSchlicker_VM<-1.38409 then node 20 elseif
pHWertSchlicker_VM>=-1.38409 then node 21 else 96.936
13 class = 104.642
14 if R25Widerstand<0.860772 then node 22 elseif
R25Widerstand>=0.860772 then node 23 else 110.0333333
15 if DurchflussMW_VM<0.282349 then node 24 elseif
DurchflussMW_VM>=0.282349 then node 25 else 101.0526667
16 class = 114.066
17 if R25Widerstand<1.04813 then node 26 elseif
R25Widerstand>=1.04813 then node 27 else 108.566
18 class = 104.642
19 if DifferenzdruckFilter<-0.20924 then node 28 elseif
DifferenzdruckFilter>=-0.20924 then node 29 else 106.6053333
20 class = 99.516
21 if d50VM_VM<-0.978351 then node 30 elseif d50VM_VM>=-
0.978351 then node 31 else 94.87333333
22 if SegmentzeitMAX_VM<0.561374 then node 32 elseif
SegmentzeitMAX_VM>=0.561374 then node 33 else 110.0333333
23 if d10NM<0.345118 then node 34 elseif d10NM>=0.345118 then
node 35 else 111.9966667
24 if DurchflussMW_VM<-0.350973 then node 36 elseif
DurchflussMW_VM>=-0.350973 then node 37 else 101.0526667
25 class = 104.642
26 if UD_VM_1 in {10 11 7} then node 38 elseif UD_VM_1 in {8 9}
then node 39 else 108.566
27 if Granulatfeuchte<-0.635176 then node 40 elseif
Granulatfeuchte>=-0.635176 then node 41 else 116.8253333
28 class = 106.6053333
29 class = 102.9246667
30 class = 96.936
31 class = 94.87333333
32 class = 104.642
33 class = 110.0333333
34 class = 111.9966667
```

```

35 class = 116.8253333
36 class = 101.0526667
37 class = 102.9246667
38 class = 111.9966667
39 class = 106.6053333
40 class = 116.8253333
41 class = 119.52375

```

### Complex Tree VM+NM mit Ausreißern

```

1  if R25Widerstand<-1.36532 then node 2 elseif
R25Widerstand>=-1.36532 then node 3 else 101.0526667
2  if Durchfluss<0.144277 then node 4 elseif
Durchfluss>=0.144277 then node 5 else 88.62466667
3  if YAcet<0.238895 then node 6 elseif YAcet>=0.238895 then
node 7 else 102.9246667
4  class = 101.0526667
5  class = 88.62466667
6  if R25Widerstand<0.119414 then node 8 elseif
R25Widerstand>=0.119414 then node 9 else 98.03066667
7  if SchlickerdruckMW_VM<1.51411 then node 10 elseif
SchlickerdruckMW_VM>=1.51411 then node 11 else 106.6053333
8  if DifferenzdruckFilter<1.50558 then node 12 elseif
DifferenzdruckFilter>=1.50558 then node 13 else 96.936
9  if YAcet<-0.57182 then node 14 elseif YAcet>=-0.57182 then
node 15 else 110.0333333
10 if d90NM<-1.49482 then node 16 elseif d90NM>=-1.49482 then
node 17 else 108.566
11 if R25Widerstand<-1.15983 then node 18 elseif
R25Widerstand>=-1.15983 then node 19 else 106.6053333
12 if pHWertSchlicker_VM<-1.38409 then node 20 elseif
pHWertSchlicker_VM>=-1.38409 then node 21 else 96.936
13 class = 104.642
14 if R25Widerstand<0.860772 then node 22 elseif
R25Widerstand>=0.860772 then node 23 else 110.0333333
15 if DurchflussMW_VM<0.282349 then node 24 elseif
DurchflussMW_VM>=0.282349 then node 25 else 101.0526667
16 class = 114.066
17 if R25Widerstand<1.04813 then node 26 elseif
R25Widerstand>=1.04813 then node 27 else 108.566
18 class = 104.642
19 if DifferenzdruckFilter<-0.20924 then node 28 elseif
DifferenzdruckFilter>=-0.20924 then node 29 else 106.6053333
20 class = 99.516
21 if d50VM_VM<-0.978351 then node 30 elseif d50VM_VM>=-
0.978351 then node 31 else 94.87333333
22 if SegmentzeitMAX_VM<0.561374 then node 32 elseif
SegmentzeitMAX_VM>=0.561374 then node 33 else 110.0333333
23 if d10NM<0.345118 then node 34 elseif d10NM>=0.345118 then
node 35 else 111.9966667
24 if DurchflussMW_VM<-0.350973 then node 36 elseif
DurchflussMW_VM>=-0.350973 then node 37 else 101.0526667
25 class = 104.642
26 if UD_VM_1 in {10 11 7} then node 38 elseif UD_VM_1 in {8 9}
then node 39 else 108.566

```

```
27 if Granulatfeuchte<-0.635176 then node 40 elseif
Granulatfeuchte>=-0.635176 then node 41 else 116.8253333
28 class = 106.6053333
29 class = 102.9246667
30 if d90VM_VM<-1.00417 then node 42 elseif d90VM_VM>=-1.00417
then node 43 else 96.936
31 if R25Widerstand<-0.517188 then node 44 elseif
R25Widerstand>=-0.517188 then node 45 else 94.873333333
32 class = 104.642
33 class = 110.0333333
34 class = 111.9966667
35 class = 116.8253333
36 class = 101.0526667
37 class = 102.9246667
38 if d10GranulatNM<-1.40239 then node 46 elseif
d10GranulatNM>=-1.40239 then node 47 else 111.9966667
39 if d10NM<-0.693852 then node 48 elseif d10NM>=-0.693852 then
node 49 else 106.6053333
40 class = 116.8253333
41 class = 119.52375
42 if d90NM<-0.555518 then node 50 elseif d90NM>=-0.555518 then
node 51 else 96.936
43 class = 101.0526667
44 if UD_VM_1 in {11 8 9} then node 52 elseif UD_VM_1 in {10 7}
then node 53 else 94.873333333
45 if d10NM<-0.407239 then node 54 elseif d10NM>=-0.407239 then
node 55 else 98.03066667
46 class = 119.52375
47 if d90NM<-0.636609 then node 56 elseif d90NM>=-0.636609 then
node 57 else 111.9966667
48 class = 106.6053333
49 if SchlickerdichteMIN_VM<0.0450895 then node 58 elseif
SchlickerdichteMIN_VM>=0.0450895 then node 59 else 108.566
50 class = 98.03066667
51 class = 96.936
52 if Gluehverlust<-0.852915 then node 60 elseif
Gluehverlust>=-0.852915 then node 61 else 94.873333333
53 class = 96.936
54 class = 94.873333333
55 if d90NM<0.0188779 then node 62 elseif d90NM>=0.0188779 then
node 63 else 99.516
56 class = 111.9966667
57 if d10GranulatNM<-0.437824 then node 64 elseif
d10GranulatNM>=-0.437824 then node 65 else 102.9246667
58 class = 101.0526667
59 class = 108.566
60 class = 96.936
61 if d50NM<-0.602723 then node 66 elseif d50NM>=-0.602723 then
node 67 else 94.873333333
62 class = 98.03066667
63 if Symm3MW<0.403742 then node 68 elseif Symm3MW>=0.403742
then node 69 else 99.516
64 class = 110.0333333
65 if DurchflussMIN_VM<0.152816 then node 70 elseif
DurchflussMIN_VM>=0.152816 then node 71 else 102.9246667
66 class = 88.62466667
```

```

67 if DifferenzdruckFilter<1.16262 then node 72 elseif
DifferenzdruckFilter>=1.16262 then node 73 else 94.87333333
68 class = 99.516
69 class = 101.0526667
70 if d90NM<-0.187229 then node 74 elseif d90NM>=-0.187229 then
node 75 else 102.9246667
71 if Tb<0.774428 then node 76 elseif Tb>=0.774428 then node 77
else 104.642
72 class = 94.87333333
73 class = 98.03066667
74 class = 108.566
75 class = 102.9246667
76 class = 114.066
77 class = 104.642

```

**Tabelle 35: Ergebnis CL und Testdaten: VM+NM ohne Ausreißer<sup>346</sup>**

	Genauigkeit bei Training	Genauigkeit bei Anwendung auf Testdaten	Grenzfälle bei Anwendung auf Testdaten	Mittelwert der absoluten Differenz zwischen der vorhergesagten und Ist-Klasse
Complex Tree	16,5%	11,5%	1,3%	2,36
Medium Tree	16,1%	14,1%	1,3%	2,06
Simple Tree	11,9%	9,0%	1,3%	2,71
Linear SVM	17,0%	9,0%	0,0%	3,13
Quadratic SVM	18,8%	10,3%	0,0%	3,14
Cubic SVM	19,3%	9,0%	0,0%	3,37
Fine Gaussian SVM	15,1%	10,3%	0,0%	2,72
Medium Gaussian SVM	19,7%	10,3%	0,0%	3,14
Coarse Gaussian SVM	17,9%	2,6%	0,0%	3,59
Boosted Trees	17,9%	14,1%	1,3%	2,31
Bagged Trees	18,3%	10,3%	0,0%	2,92
RUSBoosted Trees	19,3%	9,0%	1,3%	2,88

<sup>346</sup> Quelle: eigene Darstellung

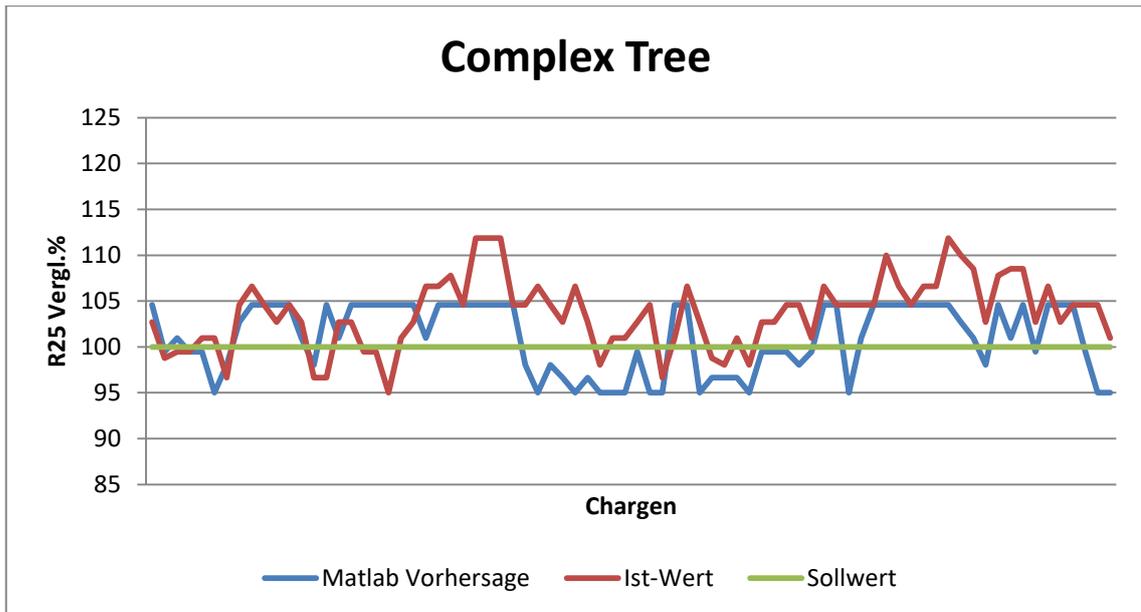


Abbildung 74: CT: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer)<sup>347</sup>

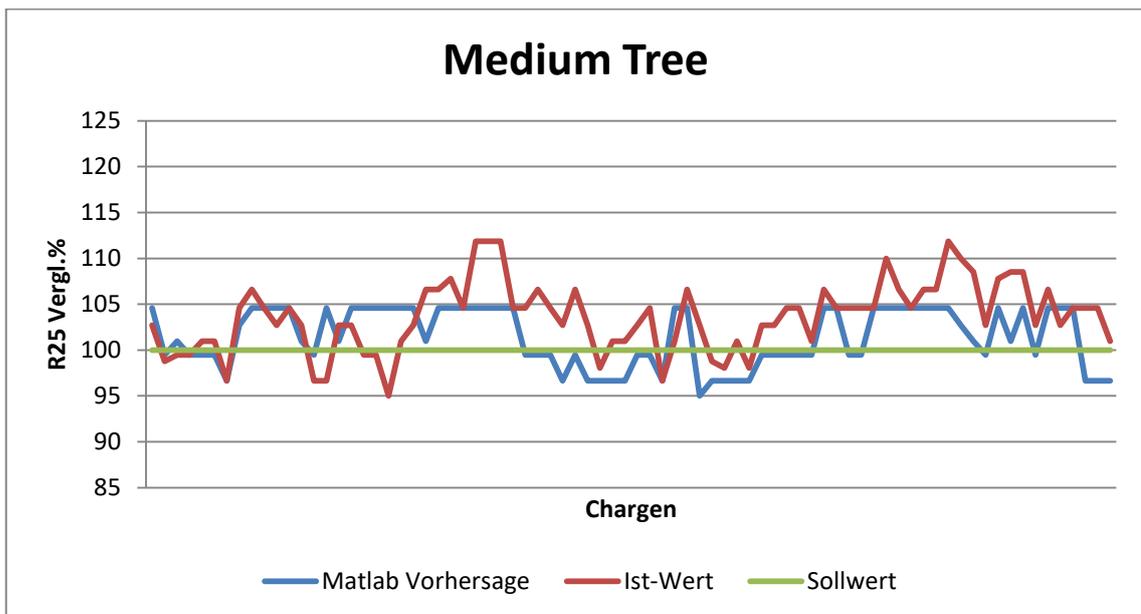


Abbildung 75: MT: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer)<sup>348</sup>

<sup>347</sup> Quelle: eigene Darstellung

<sup>348</sup> Quelle: ebda.

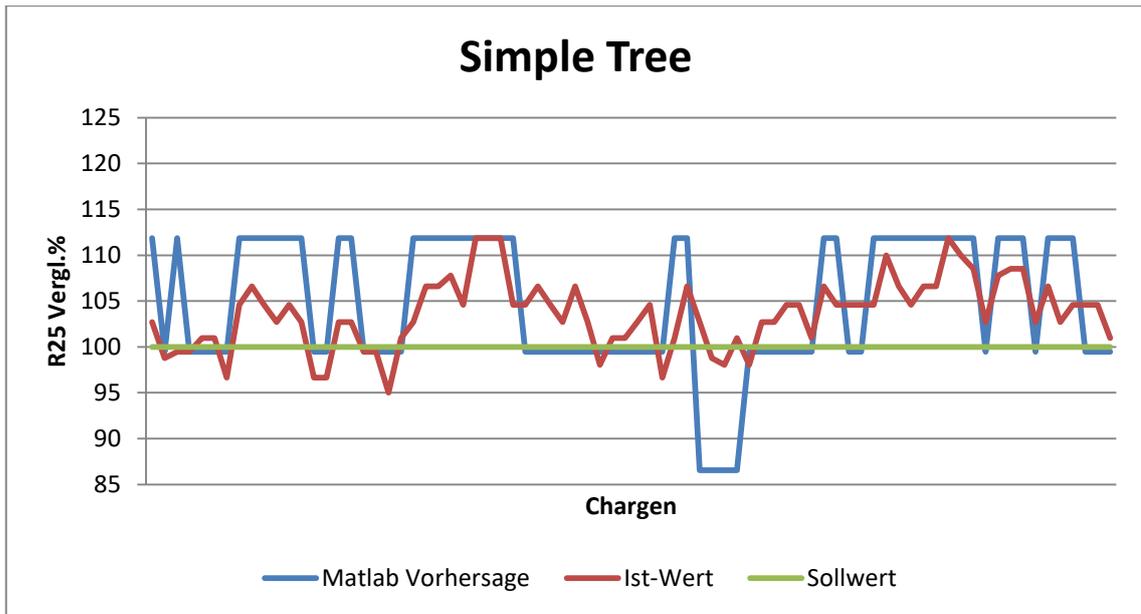


Abbildung 76: ST: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer)<sup>349</sup>

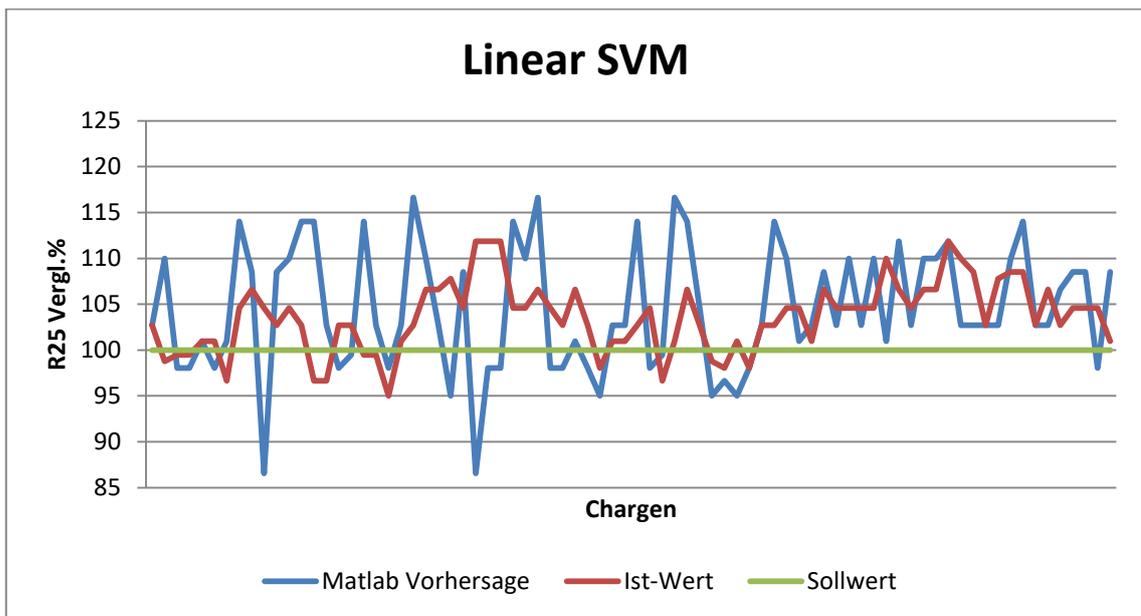


Abbildung 77: LSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer)<sup>350</sup>

<sup>349</sup> Quelle: eigene Darstellung

<sup>350</sup> Quelle: ebda.

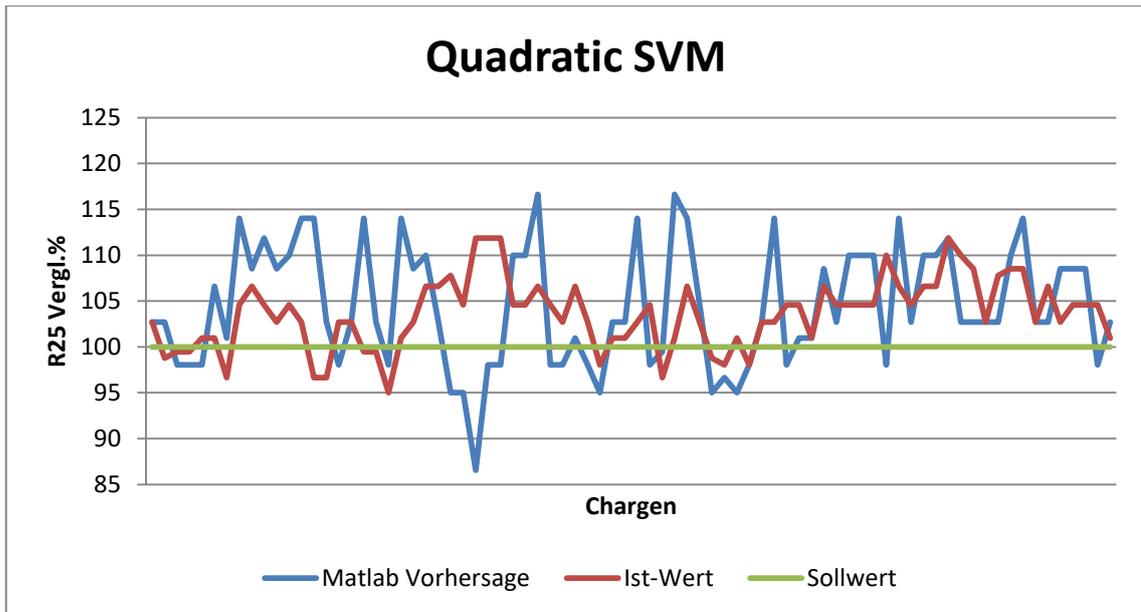


Abbildung 78: QSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer)<sup>351</sup>

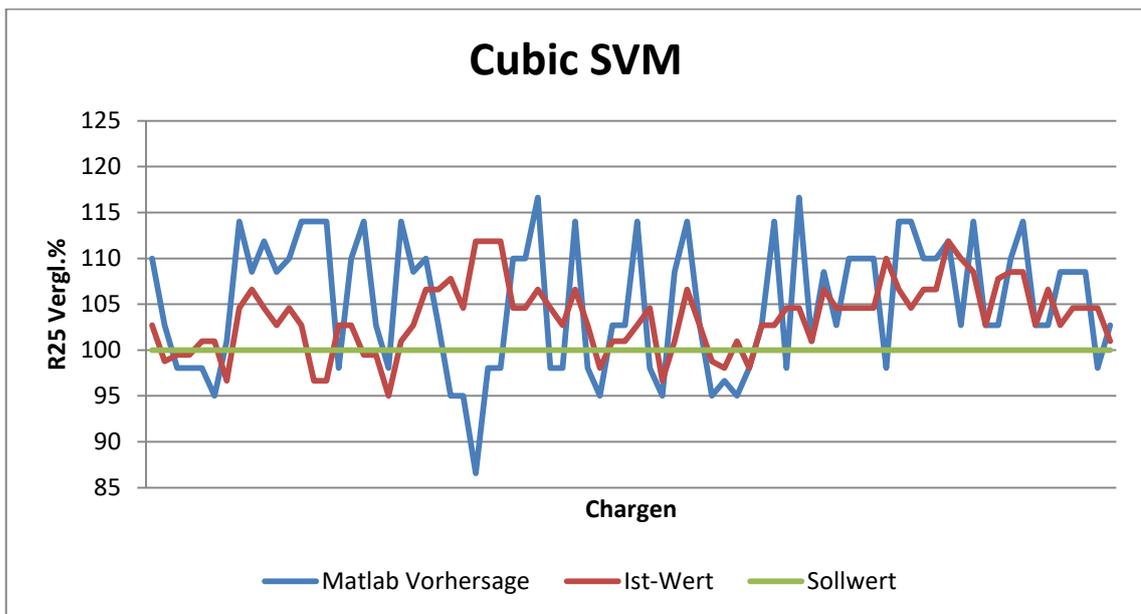


Abbildung 79: CSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer)<sup>352</sup>

<sup>351</sup> Quelle: eigene Darstellung

<sup>352</sup> Quelle: ebda.

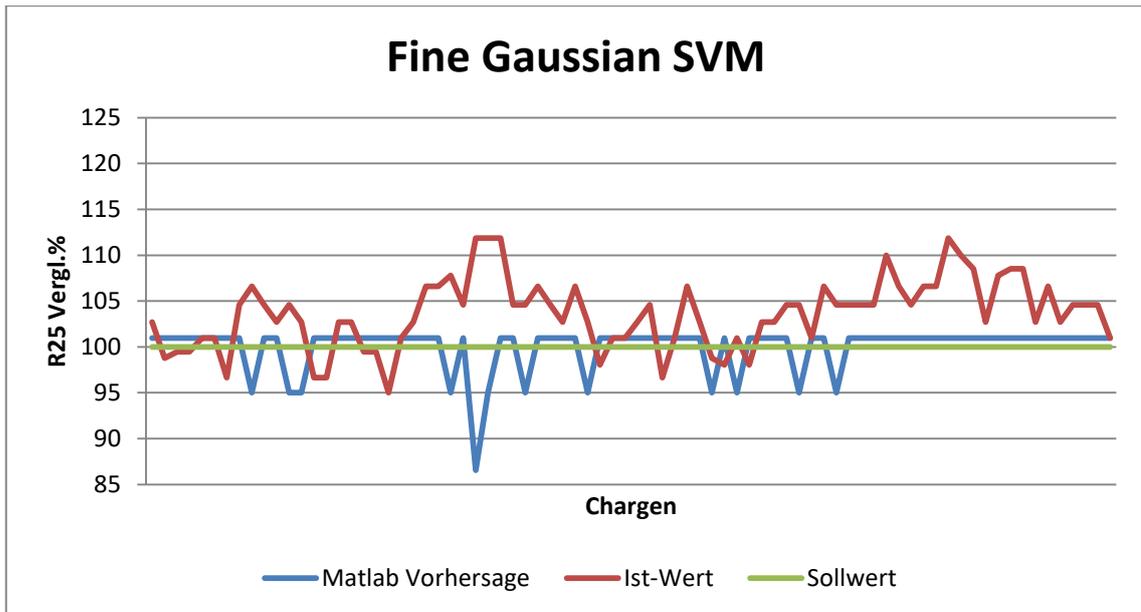


Abbildung 80: FGSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer)<sup>353</sup>

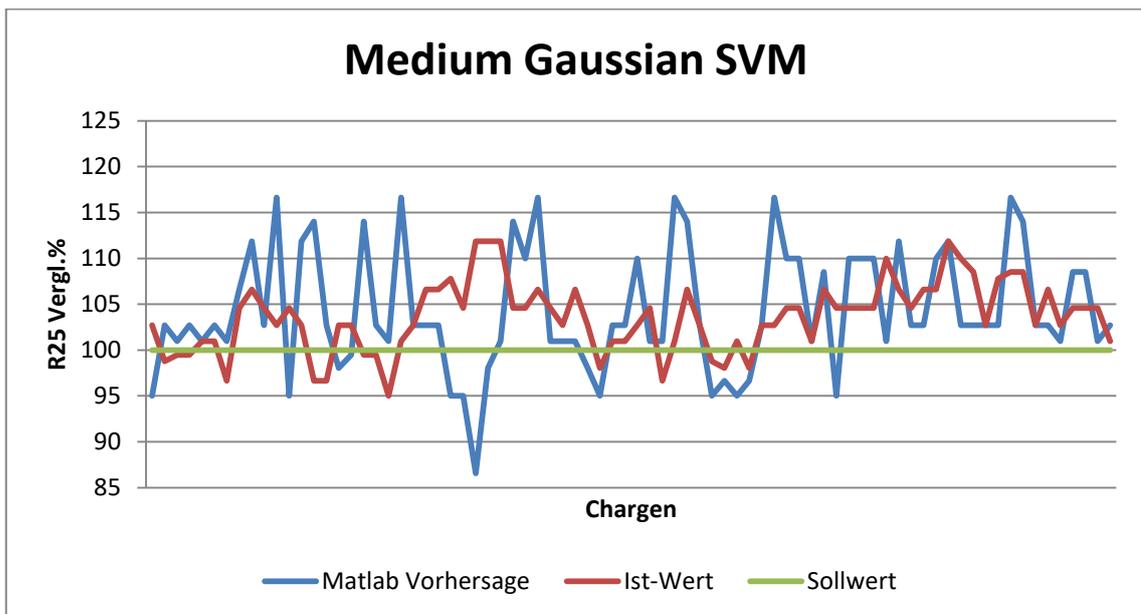


Abbildung 81: MGSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer)<sup>354</sup>

<sup>353</sup> Quelle: eigene Darstellung

<sup>354</sup> Quelle: ebda.

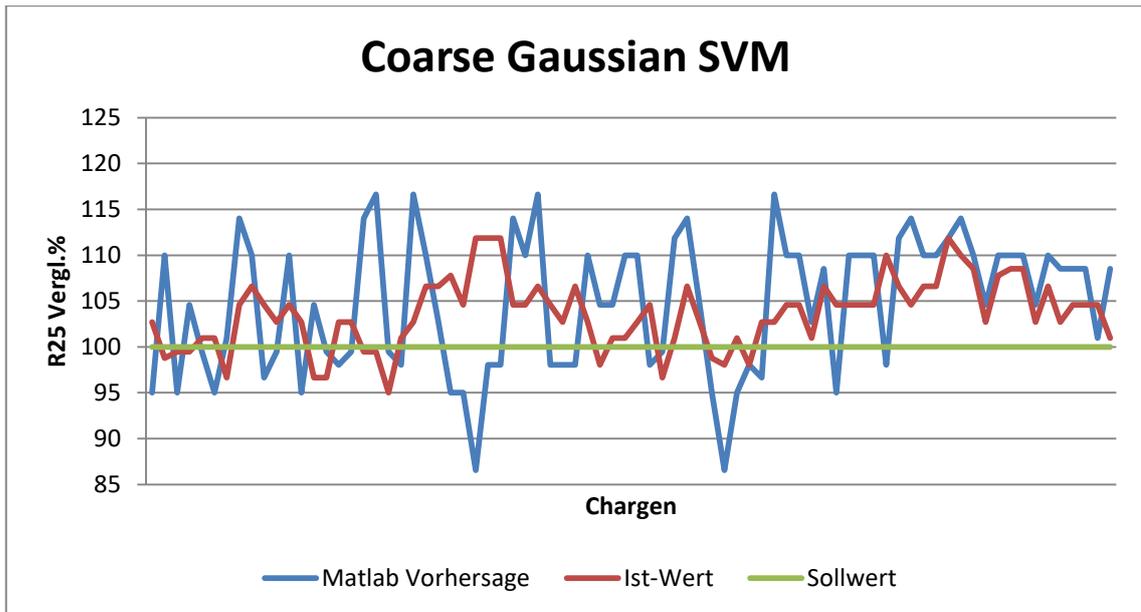


Abbildung 82: CGSVM: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer)<sup>355</sup>

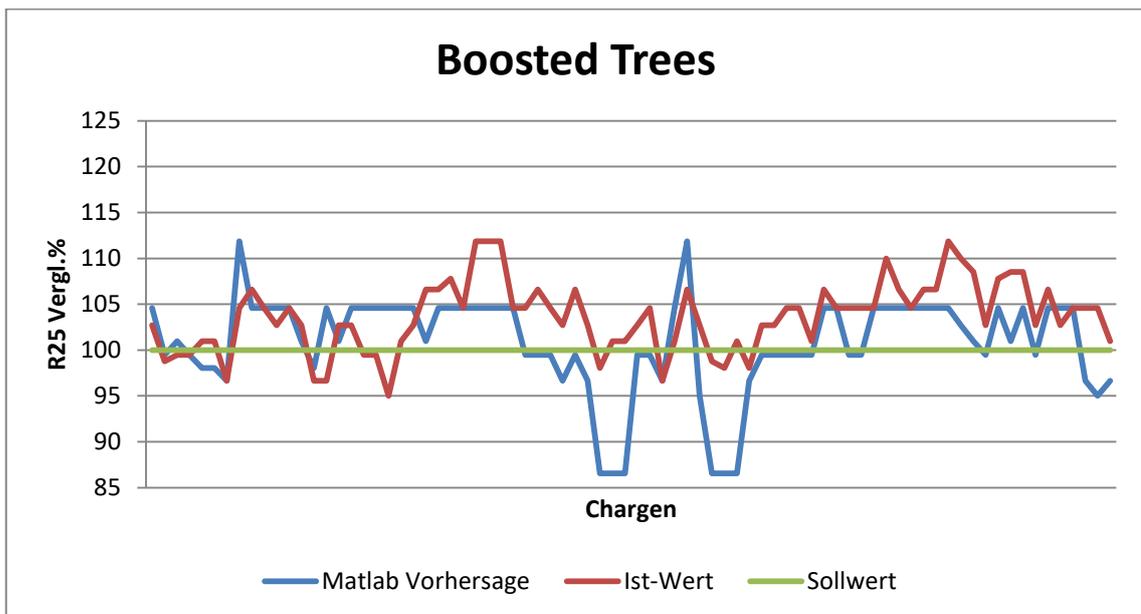


Abbildung 83: BoT: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer)<sup>356</sup>

<sup>355</sup> Quelle: eigene Darstellung

<sup>356</sup> Quelle: ebda.

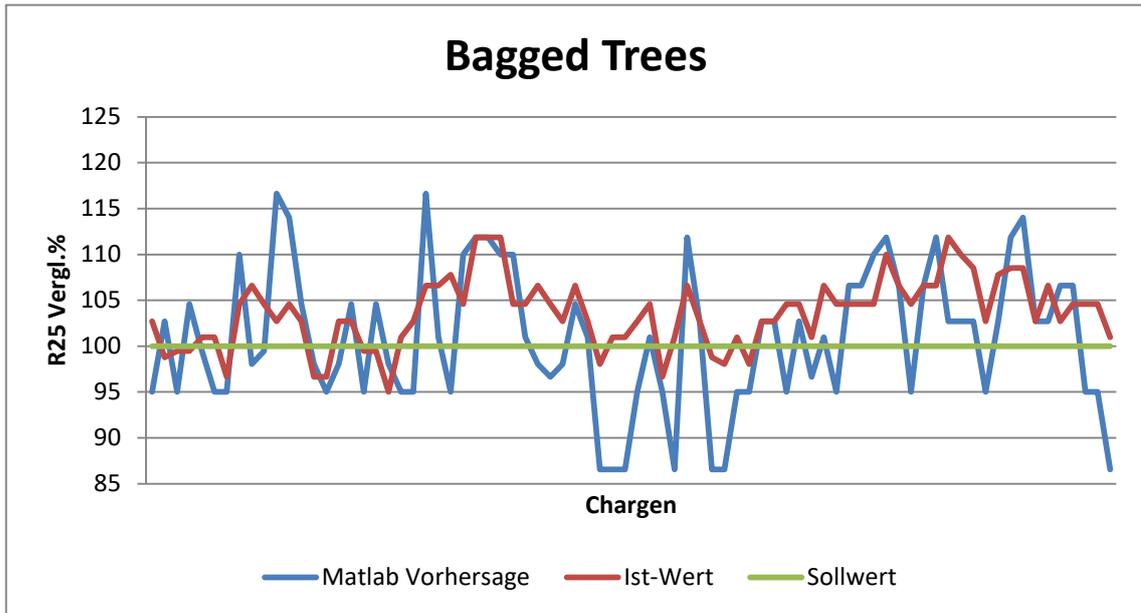


Abbildung 84: BaT: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer)<sup>357</sup>

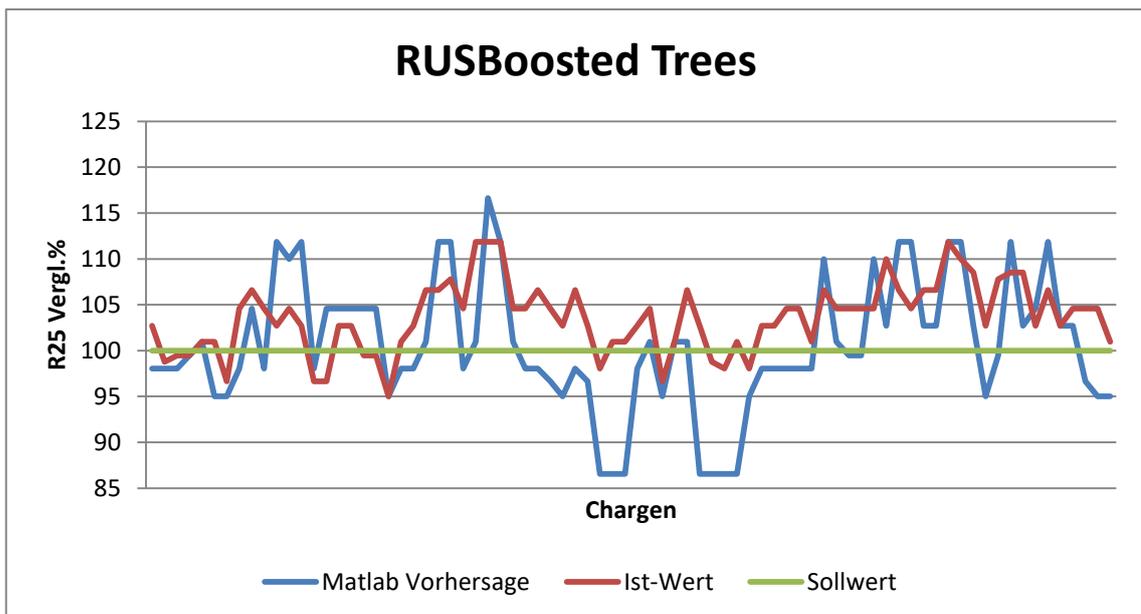


Abbildung 85: RUSBT: Vergleich Vorhersage zu Ist-Wert (VM+NM ohne Ausreißer)<sup>358</sup>

<sup>357</sup> Quelle: eigene Darstellung

<sup>358</sup> Quelle: ebda.

**Medium Tree VM+NM ohne Ausreißer**

```
1  if YAcet<0.492441 then node 2 elseif YAcet>=0.492441 then
node 3 else 88.625
2  if R25Widerstand<0.200899 then node 4 elseif
R25Widerstand>=0.200899 then node 5 else 88.625
3  if d90NM<-1.81552 then node 6 elseif d90NM>=-1.81552 then
node 7 else 106.61
4  if d10NM<-1.05918 then node 8 elseif d10NM>=-1.05918 then
node 9 else 88.625
5  if YAcet<-0.697275 then node 10 elseif YAcet>=-0.697275 then
node 11 else 112
6  class = 114.07
7  if SchlickerdruckMW_VM<1.51411 then node 12 elseif
SchlickerdruckMW_VM>=1.51411 then node 13 else 106.61
8  if TRmin<0.338175 then node 14 elseif TRmin>=0.338175 then
node 15 else 88.625
9  if DifferenzdruckFilter<1.50558 then node 16 elseif
DifferenzdruckFilter>=1.50558 then node 17 else 99.516
10 if R25Widerstand<1.07384 then node 18 elseif
R25Widerstand>=1.07384 then node 19 else 110.03
11 if DurchflussMW_VM<-0.859258 then node 20 elseif
DurchflussMW_VM>=-0.859258 then node 21 else 101.05
12 if UD_VM_1 in {10 11 7} then node 22 elseif UD_VM_1 in {8 9}
then node 23 else 108.57
13 if R25Widerstand<-1.37386 then node 24 elseif
R25Widerstand>=-1.37386 then node 25 else 106.61
14 if d90GranulatNM<0.0729912 then node 26 elseif
d90GranulatNM>=0.0729912 then node 27 else 88.625
15 class = 94.873
16 if R25Widerstand<-0.582763 then node 28 elseif
R25Widerstand>=-0.582763 then node 29 else 99.516
17 class = 104.64
18 if SegmentzeitMAX_VM<0.599055 then node 30 elseif
SegmentzeitMAX_VM>=0.599055 then node 31 else 110.03
19 if d10NM<0.408982 then node 32 elseif d10NM>=0.408982 then
node 33 else 112
20 class = 101.05
21 if Schlickerdichte<-0.0478531 then node 34 elseif
Schlickerdichte>=-0.0478531 then node 35 else 102.92
22 if Tb<0.202637 then node 36 elseif Tb>=0.202637 then node 37
else 116.83
23 if d10NM<-0.741741 then node 38 elseif d10NM>=-0.741741 then
node 39 else 106.61
24 class = 104.64
25 if DifferenzdruckFilter<-0.20924 then node 40 elseif
DifferenzdruckFilter>=-0.20924 then node 41 else 106.61
26 class = 88.625
27 class = 96.936
28 class = 96.936
29 class = 99.516
30 class = 104.64
31 class = 110.03
32 class = 112
33 class = 116.83
34 class = 102.92
```

```
35 class = 104.64
36 class = 116.83
37 class = 112
38 class = 106.61
39 class = 108.57
40 class = 106.61
41 class = 102.92
```

### Complex Tree VM+NM ohne Ausreißer

```
1  if YAcet<0.492441 then node 2 elseif YAcet>=0.492441 then
node 3 else 88.625
    2  if R25Widerstand<0.200899 then node 4 elseif
R25Widerstand>=0.200899 then node 5 else 88.625
    3  if d90NM<-1.81552 then node 6 elseif d90NM>=-1.81552 then
node 7 else 106.61
    4  if d10NM<-1.05918 then node 8 elseif d10NM>=-1.05918 then
node 9 else 88.625
    5  if YAcet<-0.697275 then node 10 elseif YAcet>=-0.697275 then
node 11 else 112
    6  class = 114.07
    7  if SchlickerdruckMW_VM<1.51411 then node 12 elseif
SchlickerdruckMW_VM>=1.51411 then node 13 else 106.61
    8  if TRmin<0.338175 then node 14 elseif TRmin>=0.338175 then
node 15 else 88.625
    9  if DifferenzdruckFilter<1.50558 then node 16 elseif
DifferenzdruckFilter>=1.50558 then node 17 else 99.516
    10  if R25Widerstand<1.07384 then node 18 elseif
R25Widerstand>=1.07384 then node 19 else 110.03
    11  if DurchflussMW_VM<-0.859258 then node 20 elseif
DurchflussMW_VM>=-0.859258 then node 21 else 101.05
    12  if UD_VM_1 in {10 11 7} then node 22 elseif UD_VM_1 in {8 9}
then node 23 else 108.57
    13  if R25Widerstand<-1.37386 then node 24 elseif
R25Widerstand>=-1.37386 then node 25 else 106.61
    14  if d90GranulatNM<0.0729912 then node 26 elseif
d90GranulatNM>=0.0729912 then node 27 else 88.625
    15  class = 94.873
    16  if R25Widerstand<-0.582763 then node 28 elseif
R25Widerstand>=-0.582763 then node 29 else 99.516
    17  class = 104.64
    18  if SegmentzeitMAX_VM<0.599055 then node 30 elseif
SegmentzeitMAX_VM>=0.599055 then node 31 else 110.03
    19  if d10NM<0.408982 then node 32 elseif d10NM>=0.408982 then
node 33 else 112
    20  class = 101.05
    21  if Schlickerdichte<-0.0478531 then node 34 elseif
Schlickerdichte>=-0.0478531 then node 35 else 102.92
    22  if Tb<0.202637 then node 36 elseif Tb>=0.202637 then node 37
else 116.83
    23  if d10NM<-0.741741 then node 38 elseif d10NM>=-0.741741 then
node 39 else 106.61
    24  class = 104.64
    25  if DifferenzdruckFilter<-0.20924 then node 40 elseif
DifferenzdruckFilter>=-0.20924 then node 41 else 106.61
```

```
26  if Durchfluss<-0.448149 then node 42 elseif Durchfluss>=-
0.448149 then node 43 else 88.625
27  class = 96.936
28  if DurchflussStabw_VM<0.554187 then node 44 elseif
DurchflussStabw_VM>=0.554187 then node 45 else 96.936
29  if Heizgruppel1_VM<-1.23287 then node 46 elseif
Heizgruppel1_VM>=-1.23287 then node 47 else 99.516
30  class = 104.64
31  class = 110.03
32  class = 112
33  class = 116.83
34  class = 102.92
35  class = 104.64
36  if SchuettdichteGranulat_VM<-0.643787 then node 48 elseif
SchuettdichteGranulat_VM>=-0.643787 then node 49 else 116.83
37  if pHWertSchlicker_VM<-0.000115404 then node 50 elseif
pHWertSchlicker_VM>=-0.000115404 then node 51 else 112
38  class = 106.61
39  if TRmin<-0.243261 then node 52 elseif TRmin>=-0.243261 then
node 53 else 108.57
40  class = 106.61
41  class = 102.92
42  class = 98.031
43  class = 88.625
44  if R25Widerstand<-1.04899 then node 54 elseif
R25Widerstand>=-1.04899 then node 55 else 96.936
45  class = 94.873
46  class = 101.05
47  if Ringtemperatur_VM<-0.335616 then node 56 elseif
Ringtemperatur_VM>=-0.335616 then node 57 else 99.516
48  class = 119.52
49  if Durchfluss<-0.986719 then node 58 elseif Durchfluss>=-
0.986719 then node 59 else 116.83
50  if d50NM<0.797421 then node 60 elseif d50NM>=0.797421 then
node 61 else 104.64
51  if Gluehverlust<-0.623386 then node 62 elseif
Gluehverlust>=-0.623386 then node 63 else 112
52  class = 101.05
53  class = 108.57
54  class = 94.873
55  if DifferenzdruckFilter<0.991135 then node 64 elseif
DifferenzdruckFilter>=0.991135 then node 65 else 96.936
56  class = 99.516
57  if Ablufttemperatur_VM<0.313699 then node 66 elseif
Ablufttemperatur_VM>=0.313699 then node 67 else 98.031
58  class = 114.07
59  if SchlickertempStabw_VM<-0.120338 then node 68 elseif
SchlickertempStabw_VM>=-0.120338 then node 69 else 116.83
60  class = 104.64
61  class = 110.03
62  class = 102.92
63  class = 112
64  if DurchflussMW_VM<-0.72149 then node 70 elseif
DurchflussMW_VM>=-0.72149 then node 71 else 96.936
65  class = 98.031
66  class = 98.031
```

67 class = 94.873  
68 class = 116.83  
69 class = 119.52  
70 class = 99.516  
71 class = 96.936

Tabelle 36: Statistische Maße für Šumperk (mit Ausreißern)<sup>359</sup>

Bezeichnung	$\tilde{s}^2$	$\tilde{s}$	$\gamma$	$g_m$	Korrelationskoeffizient mit MW
MW	1,33E+01	3,65E+00	4,54	1,52	1,00
Kugelmenge_DL	8,87E-02	2,98E-01	-4,30	-0,19	0,22
Durchfluss_DL	7,32E-01	8,56E-01	-3,30	-0,38	-0,03
Schlickertemperatur_DL	5,86E+00	2,42E+00	-3,44	-0,37	0,07
d10 (NM)_DL	9,59E-05	9,79E-03	-3,41	-0,34	0,27
d50 (NM)_DL	3,16E-04	1,78E-02	-2,43	-0,99	-0,08
d90 (NM)_DL	1,16E-02	1,08E-01	-2,87	-0,44	-0,15
Auslaufzeit_DL	2,27E+00	1,51E+00	9,48	2,81	-0,09
Schuettdichte Granulat_DL	1,27E-04	1,13E-02	-1,78	-0,61	0,16
Granulatfeuchte_DL	4,94E-05	7,03E-03	-2,88	0,04	0,01
Differenzdruck Filter_DL	6,84E-04	2,61E-02	-2,34	-0,27	-0,04
d10 Granulat (NM)_DL	5,91E+00	2,43E+00	-2,99	0,31	-0,30
d50 Granulat (NM)_DL	1,93E+01	4,39E+00	-3,06	0,22	-0,24
d90 Granulat (NM)_DL	2,25E+01	4,74E+00	-2,71	0,52	-0,18
Schlickerdichte_DL	1,18E-04	1,09E-02	-2,40	-0,01	-0,11
pH Wert Schlicker_DL	2,62E-02	1,62E-01	-3,58	-0,30	0,04
Gluehverlust_DL	7,33E-04	2,71E-02	-3,44	0,43	0,18
Sinterdichte_DL	9,80E-05	9,90E-03	-3,56	-0,40	0,01
R25 Widerstand_DL	3,42E+00	1,85E+00	-2,71	0,55	0,18
R25 Vergl,%_DL	4,07E+01	6,38E+00	-3,93	-0,15	0,39
R25_DL	1,94E+00	1,39E+00	-3,34	-0,02	0,19
Rmin_DL	1,15E+00	1,07E+00	-3,21	0,01	0,23
T-Rmin_DL	3,42E+00	1,85E+00	-2,70	-0,56	-0,15
Tb_DL	3,50E-01	5,92E-01	-2,52	-0,48	-0,25
Tcorr_DL	2,96E+00	1,72E+00	-2,25	-0,62	-0,18
Gleitwert_DL	3,61E-04	1,90E-02	-1,19	-1,10	0,29
Gruenteilfestigkeit_DL	3,46E-01	5,88E-01	-3,38	0,04	-0,25
Kohaesion_DL	8,30E-03	9,11E-02	-1,85	-0,60	0,16
Pressdruck_DL	1,26E+01	3,55E+00	-2,80	0,76	-0,28
Radiale Auffederung_DL	3,50E-07	5,92E-04	-3,62	0,26	-0,19
Schuettwinkel_DL	4,79E-01	6,92E-01	-3,23	0,29	0,19
SPHT3 MW_DL	2,16E-06	1,47E-03	-2,87	-0,05	-0,14
Symm3 MW_DL	1,19E-06	1,09E-03	-2,26	-0,66	0,10
b I3 MW_DL	1,02E-05	3,19E-03	-2,97	-0,15	0,09
Contraspum					
K1012_Menge1_NM_Gehalt_DL	8,19E-10	2,86E-05	11,85	3,98	0,16
Methocell A4C_NM_Gehalt_DL	1,93E-13	4,39E-07	-3,85	0,14	0,08
BaO_NM_Gehalt_DL	1,33E-03	3,64E-02	-2,39	0,22	0,00
CaO_NM_Gehalt_DL	2,48E-04	1,57E-02	-3,25	0,68	0,00
MnO_NM_Gehalt_DL	7,59E-08	2,76E-04	-3,58	-0,85	-0,18
Na2O_NM_Gehalt_DL	1,73E-07	4,15E-04	-4,38	-0,57	-0,10
PbO_NM_Gehalt_DL	5,18E-04	2,28E-02	-1,99	-0,42	-0,20
SiO2_NM_Gehalt_DL	1,35E-05	3,67E-03	-3,47	-0,48	-0,10
SrO_NM_Gehalt_DL	5,72E-05	7,56E-03	3,77	2,39	0,44
TiO2_NM_Gehalt_DL	6,94E-04	2,63E-02	-1,83	0,48	0,07

Fortsetzung von Tabelle 36: Statistische Maße für Šumperk (mit Ausreißern)<sup>360</sup><sup>359</sup> Quelle: eigene Darstellung

Bezeichnung	$\bar{s}^2$	$\bar{s}$	$\gamma$	$g_m$	Korrelationskoeffizient mit MW
Y2O3_NM_Gehalt_DL	6,07E-07	7,79E-04	-1,81	0,71	-0,03
Pb_Me_DL	7,75E-10	2,78E-05	1,10	1,31	-0,14
BaCO3_DL	4,88E-09	6,98E-05	-2,86	-0,23	-0,13
SrCO3_DL	2,92E-09	5,41E-05	-0,07	1,48	0,13
CaCO3_DL	2,09E-10	1,44E-05	-2,89	-0,21	0,04
Pb3O4_DL	7,75E-10	2,78E-05	1,10	1,31	-0,14
Y-Acet_DL	8,73E-07	9,35E-04	-4,44	0,01	0,16
TiO2_DL	4,73E-09	6,88E-05	-0,77	-0,03	0,07
A40_Me_DL	9,18E-14	3,03E-07	-4,62	-0,59	-0,26
Optapi_DL	1,36E-14	1,17E-07	-4,62	0,59	0,26
Kugelmenge_VM_DL	4,30E-07	6,56E-04	52,02	-7,55	-0,58
Siebueckstand feucht_VM_DL	2,46E+03	4,96E+01	-3,87	0,22	0,05
d10 (VM)_VM_DL	5,66E-05	7,52E-03	-3,19	0,12	-0,08
d50 (VM)_VM_DL	1,80E-03	4,24E-02	-3,74	0,06	0,00
d90 (VM)_VM_DL	1,99E-02	1,41E-01	-3,66	0,02	0,11
pH Wert Schlicker_VM_DL	1,16E-02	1,08E-01	-3,94	-0,10	-0,07
Schuettdichte Granulat_VM_DL	1,01E-04	1,00E-02	-3,38	0,23	-0,17
Granulatfeuchte_VM_DL	5,51E-05	7,42E-03	1,46	1,77	0,04
Differenzdruck Filter_VM_DL	3,95E-03	6,28E-02	-3,61	-0,36	-0,04
Ablufttemperatur_VM_DL	3,29E+01	5,74E+00	-3,77	-0,13	-0,02
Gewoelbeunterdruck_VM_DL	6,91E-04	2,63E-02	-3,65	0,20	-0,18
Heizgruppe 7_VM_DL	2,98E+01	5,46E+00	2,68	-1,00	-0,06
Heizgruppe 9_VM_DL	1,77E+01	4,21E+00	3,34	-2,11	0,02
Heizgruppe 11_VM_DL	1,48E+01	3,84E+00	4,27	-2,33	0,01
Heizgruppe 13_VM_DL	2,35E+01	4,84E+00	3,73	-0,10	-0,07
Differenzdruck MAX_VM_DL	3,81E-04	1,95E-02	5,99	-1,36	-0,04
Differenzdruck MIN_VM_DL	3,22E-04	1,79E-02	6,00	-2,02	-0,03
Differenzdruck MW_VM_DL	3,04E-04	1,74E-02	8,63	-2,13	-0,04
Differenzdruck Stabw_VM_DL	3,28E-06	1,81E-03	1,53	1,92	0,06
Durchfluss MAX_VM_DL	3,25E-02	1,80E-01	4,30	2,49	-0,22
Durchfluss MIN_VM_DL	5,89E-01	7,68E-01	45,58	-6,91	0,11
Durchfluss MW_VM_DL	2,13E-02	1,46E-01	31,66	-5,34	0,13
Durchfluss Stabw_VM_DL	6,34E-02	2,52E-01	49,56	7,30	-0,07
Leistung MAX_VM_DL	4,79E+05	6,92E+02	-3,72	0,62	0,17
Leistung MIN_VM_DL	3,63E+05	6,03E+02	-3,26	0,86	0,13
Leistung MW_VM_DL	4,07E+05	6,38E+02	-3,57	0,72	0,14
Leistung Stabw_VM_DL	2,65E+03	5,15E+01	5,98	2,49	0,31
Schlickerdichte MAX_VM_DL	1,49E-05	3,86E-03	2,99	0,46	-0,04
Schlickerdichte MIN_VM_DL	6,45E-03	8,03E-02	50,77	-7,43	0,06
Schlickerdichte MW_VM_DL	1,33E-04	1,15E-02	42,26	-6,57	0,03
Schlickerdichte Stabw_VM_DL	4,33E-04	2,08E-02	50,81	7,43	-0,05
Schlickerdruck MAX_VM_DL	6,65E-03	8,16E-02	-3,67	0,10	-0,07
Schlickerdruck MIN_VM_DL	3,50E-03	5,91E-02	-2,78	0,40	-0,01

<sup>360</sup> Quelle: eigene Darstellung

Fortsetzung von Tabelle 36: Statistische Maße für Šumperk (mit Ausreißern)<sup>361</sup>

Bezeichnung	$\tilde{s}^2$	$\tilde{s}$	$\gamma$	$g_m$	Korrelationskoeffizient mit MW
Schlickerdruck MW_VM_DL	4,35E-03	6,59E-02	-3,14	0,17	-0,04
Schlickerdruck Stabw_VM_DL	1,45E-04	1,21E-02	10,42	3,07	-0,10
Schlickertemp, MAX_VM_DL	1,86E+00	1,36E+00	9,13	2,96	0,15
Schlickertemp, MIN_VM_DL	9,84E-01	9,92E-01	-1,73	1,37	0,19
Schlickertemp, MW_VM_DL	9,52E-01	9,76E-01	-2,02	1,25	0,19
Schlickertemp, Stabw_VM_DL	3,44E-02	1,85E-01	29,85	5,29	0,05
Anzahl Segmente_VM_DL	1,41E-01	3,76E-01	-1,19	1,51	0,29
Segmentzeit MAX_VM_DL	1,22E+06	1,10E+03	2,11	-1,87	-0,26
Segmentzeit MIN_VM_DL	1,03E+07	3,21E+03	-3,77	-0,68	-0,21
MW Dauer je Segment_VM_DL	4,10E+06	2,02E+03	-3,64	-0,63	-0,27
Stabw Segmentzeit_VM_DL	4,54E+14	2,13E+07	-0,82	-1,35	-0,28
Gesamtzeit Segmente_VM_DL	1,01E+06	1,01E+03	5,57	-2,28	-0,24
Ringtemperatur_VM_DL	5,41E+00	2,33E+00	-3,78	-0,10	0,24

<sup>361</sup> Quelle: eigene Darstellung

Tabelle 37: MSA für einzelne Attribute (Šumperk mit Ausreißern)<sup>362</sup>

Attribute	Durchgang			
	1	2	3	4
Kugelmenge_DL	nan	0,61	0,72	0,78
Durchfluss_DL	nan			
Schlickertemperatur_DL	nan	0,64	0,70	0,71
d10 (NM)_DL	nan	0,58	0,53	0,62
d50 (NM)_DL	nan			
d90 (NM)_DL	nan	0,32		
Auslaufzeit_DL	nan	0,46		
Schuettdichte Granulat_DL	nan			
Granulatfeuchte_DL	nan			
Differenzdruck Filter_DL	nan	0,24		
d10 Granulat (NM)_DL	nan	0,51	0,75	0,73
d50 Granulat (NM)_DL	nan			
d90 Granulat (NM)_DL	nan			
Schlickerdichte_DL	nan			
pH Wert Schlicker_DL	nan			
Gluehverlust_DL	nan	0,48		
Sinterdichte_DL	nan			
R25 Widerstand_DL	nan			
R25 Vergl,%_DL	nan			
R25_DL	nan			
Rmin_DL	nan			
T-Rmin_DL	nan			
Tb_DL	nan			
Tcorr_DL	nan			
Gleitwert_DL	nan			
Gruenteilfestigkeit_DL	nan			
Kohaesion_DL	nan			
Pressdruck_DL	nan	0,60	0,81	0,88
Radiale Auffederung_DL	nan			
Schuettwinkel_DL	nan	0,44		
SPHT3 MW_DL	nan	0,48		
Symm3 MW_DL	nan	0,20		
b-13 MW_DL	nan	0,16		
Contraspum				
K1012_Menge1_NM_Gehalt_DL	nan	0,27		
Methocell A4C_NM_Gehalt_DL	nan	0,15		
BaO_NM_Gehalt_DL	nan	0,26		
CaO_NM_Gehalt_DL	nan			
MnO_NM_Gehalt_DL	nan			
Na2O_NM_Gehalt_DL	nan			
PbO_NM_Gehalt_DL	nan			
SiO2_NM_Gehalt_DL	nan			
SrO_NM_Gehalt_DL	nan			
TiO2_NM_Gehalt_DL	nan			
Y2O3_NM_Gehalt_DL	nan			
Pb_Me_DL	nan			

<sup>362</sup> Quelle: eigene Darstellung

Fortsetzung von Tabelle 37: MSA für einzelne Attribute (Šumperk mit Ausreißern)<sup>363</sup>

Attribute	Durchgang			
	1	2	3	4
BaCO3_DL	nan	0,23		
SrCO3_DL	nan			
CaCO3_DL	nan			
Pb3O4_DL	nan	0,17		
Y-Acet_DL	nan	0,70	0,82	0,81
TiO2_DL	nan			
A40_Me_DL	nan			
Optapi_DL	nan	0,58	0,76	0,76
Kugelmenge_VM_DL	nan	0,23		
Siebrueckstand feucht_VM_DL	nan			
d10 (VM)_VM_DL	nan			
d50 (VM)_VM_DL	nan	0,32		
d90 (VM)_VM_DL	nan	0,37		
pH Wert Schlicker_VM_DL	nan			
Schuettdichte Granulat_VM_DL	nan	0,40		
Granulatfeuchte_VM_DL	nan			
Differenzdruck Filter_VM_DL	nan			
Ablufttemperatur_VM_DL	nan	0,15		
Gewoelbeunterdruck_VM_DL	nan	0,24		
Heizgruppe 7_VM_DL	nan			
Heizgruppe 9_VM_DL	nan			
Heizgruppe 11_VM_DL	nan			
Heizgruppe 13_VM_DL	nan			
Differenzdruck MAX_VM_DL	nan	0,32		
Differenzdruck MIN_VM_DL	nan			
Differenzdruck MW_VM_DL	nan			
Differenzdruck Stabw_VM_DL	nan	0,39		
Durchfluss MAX_VM_DL	nan			
Durchfluss MIN_VM_DL	nan	0,80	0,34	
Durchfluss MW_VM_DL	nan			
Durchfluss Stabw_VM_DL	nan			
Leistung MAX_VM_DL	nan			
Leistung MIN_VM_DL	nan			
Leistung MW_VM_DL	nan			
Leistung Stabw_VM_DL	nan	0,32		
Schlickerdichte MAX_VM_DL	nan			
Schlickerdichte MIN_VM_DL	nan	0,39		
Schlickerdichte MW_VM_DL	nan			
Schlickerdichte Stabw_VM_DL	nan			
Schlickerdruck MAX_VM_DL	nan			
Schlickerdruck MIN_VM_DL	nan	0,27		
Schlickerdruck MW_VM_DL	nan			
Schlickerdruck Stabw_VM_DL	nan	0,22		
Schlickertemp, MAX_VM_DL	nan	0,37		
Schlickertemp, MIN_VM_DL	nan			
Schlickertemp, MW_VM_DL	nan	0,31		
Schlickertemp, Stabw_VM_DL	nan	0,39		

<sup>363</sup> Quelle: eigene Darstellung

Fortsetzung von Tabelle 37: MSA für einzelne Attribute (Šumperk mit Ausreißern)<sup>364</sup>

Attribute	Durchgang			
	1	2	3	4
Anzahl Segmente VM_DL	nan	0,60	0,62	0,62
Segmentzeit MAX_VM_DL	nan			
Segmentzeit MIN_VM_DL	nan	0,57	0,61	0,63
MW Dauer je Segment_VM_DL	nan			
Stabw Segmentzeit_VM_DL	nan			
Gesamtzeit Segmente_VM_DL	nan			
Ringtemperatur_VM_DL	nan			

Tabelle 38: Ergebnis CL und Testdaten: Šumperk mit Ausreißern<sup>365</sup>

	Genauigkeit bei Training	Genauigkeit bei Anwendung auf Testdaten	Grenzfälle bei Anwendung auf Testdaten	Mittelwert der absoluten Differenz zwischen der vorhergesagten und Ist-Klasse
Complex Tree	25,0%	24,0%	0,0%	1,94
Medium Tree	25,0%	24,0%	0,0%	1,94
Simple Tree	23,3%	12,0%	4,0%	2,38
Linear SVM	15,0%	24,0%	4,0%	2,30
Quadratic SVM	21,7%	20,0%	4,0%	2,34
Cubic SVM	21,7%	12,0%	4,0%	2,62
Fine Gaussian SVM	18,3%	8,0%	0,0%	3,10
Medium Gaussian SVM	20,0%	0,0%	0,0%	3,30
Coarse Gaussian SVM	15,0%	8,0%	0,0%	3,42
Boosted Trees	23,3%	20,0%	4,0%	2,18
Bagged Trees	15,0%	8,0%	4,0%	2,98
RUSBoosted Trees	19,3%	8,0%	0,0%	2,06

<sup>364</sup> Quelle: eigene Darstellung<sup>365</sup> Quelle: ebda.

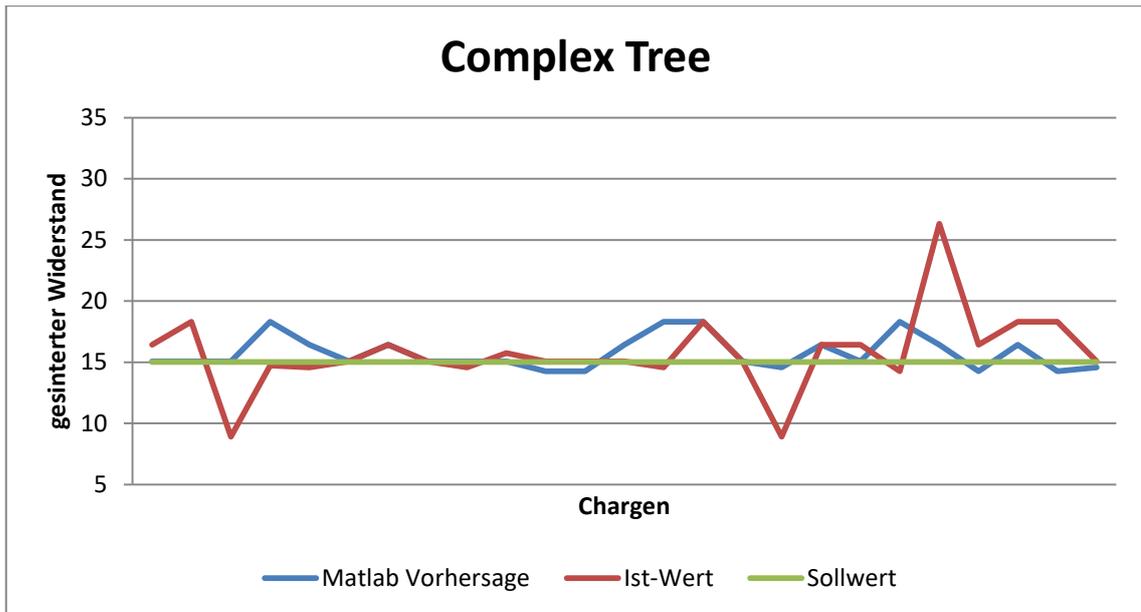


Abbildung 86: CT: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern)<sup>366</sup>

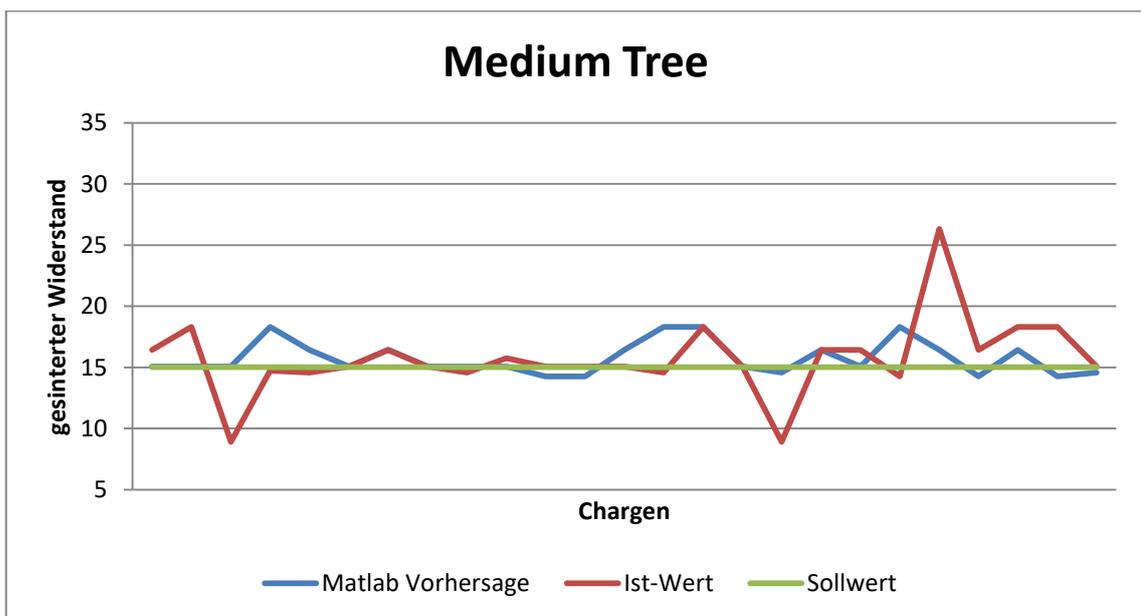


Abbildung 87: MT: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern)<sup>367</sup>

<sup>366</sup> Quelle: eigene Darstellung

<sup>367</sup> Quelle: ebda.

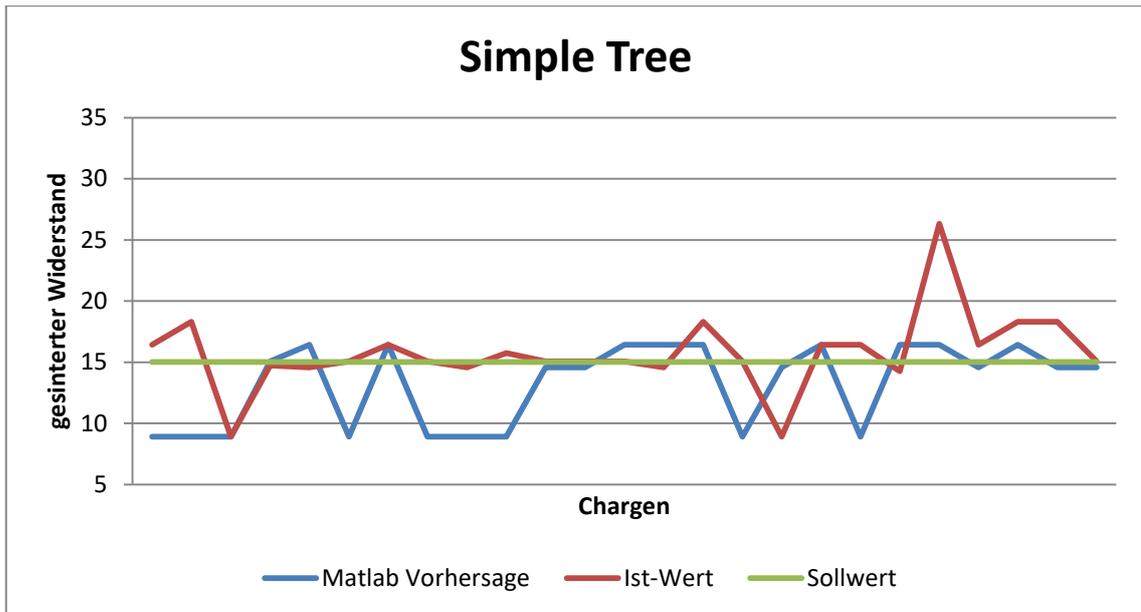


Abbildung 88: ST: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern)<sup>368</sup>

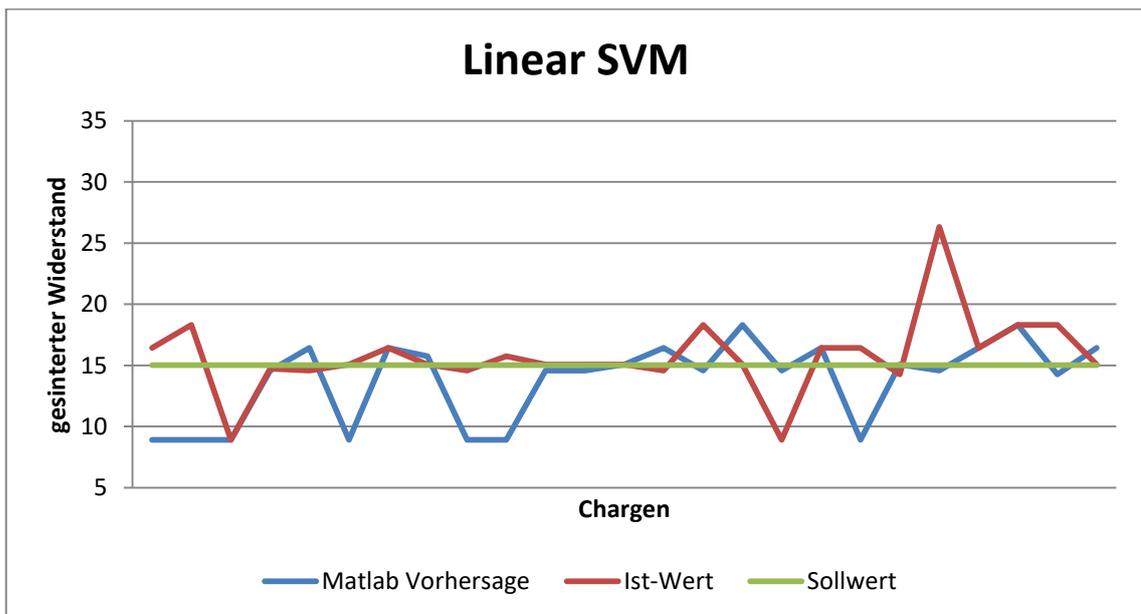


Abbildung 89: LSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern)<sup>369</sup>

<sup>368</sup> Quelle: eigene Darstellung

<sup>369</sup> Quelle: ebda.

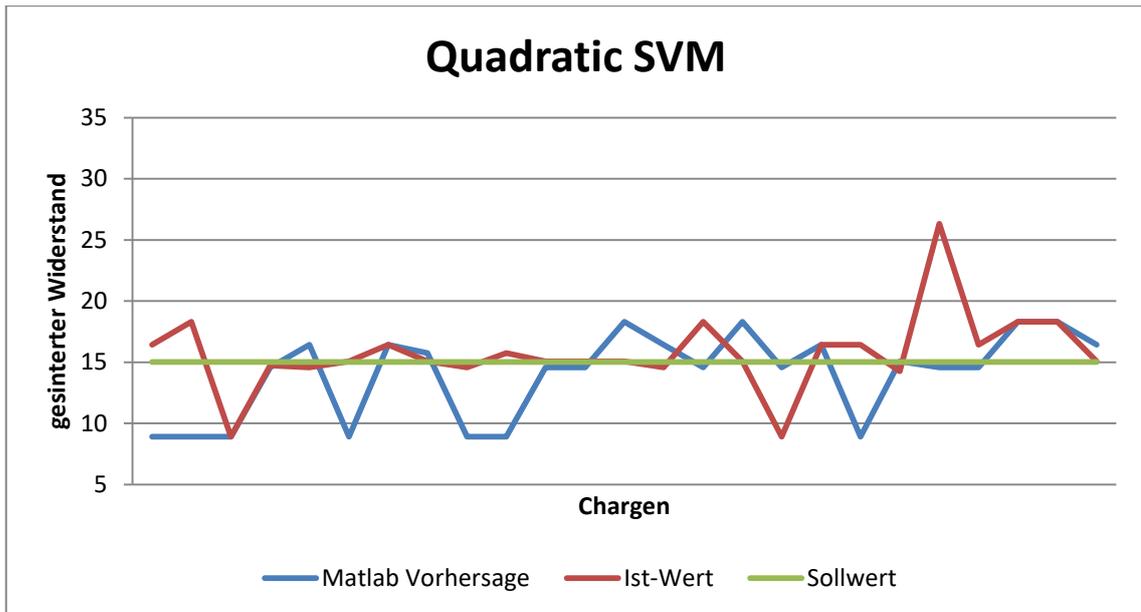


Abbildung 90: QSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern)<sup>370</sup>

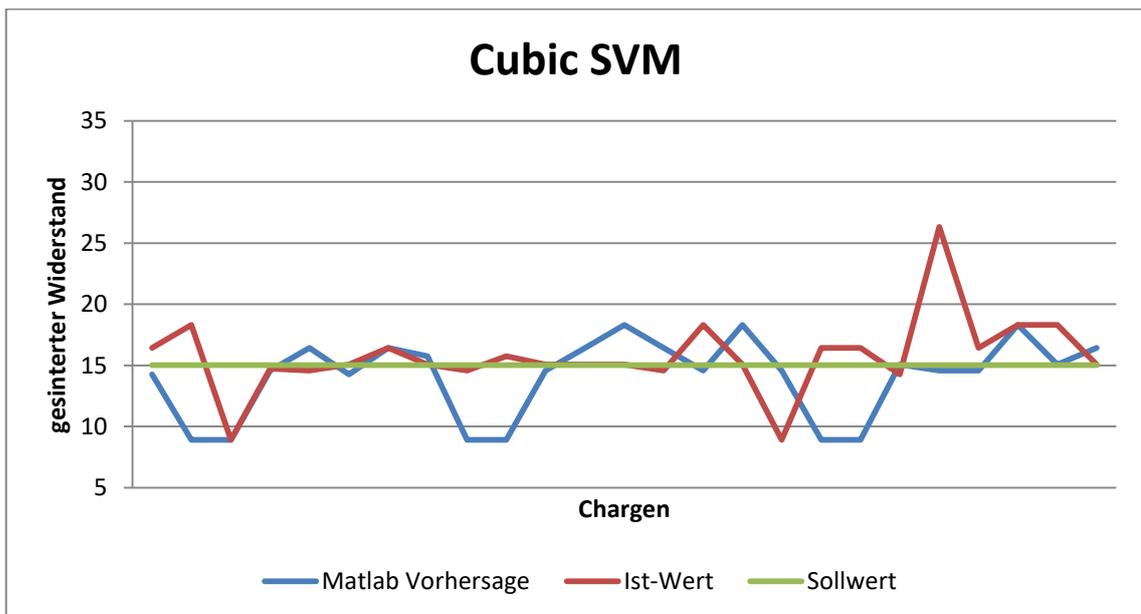


Abbildung 91: CSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern)<sup>371</sup>

<sup>370</sup> Quelle: eigene Darstellung

<sup>371</sup> Quelle: ebda.

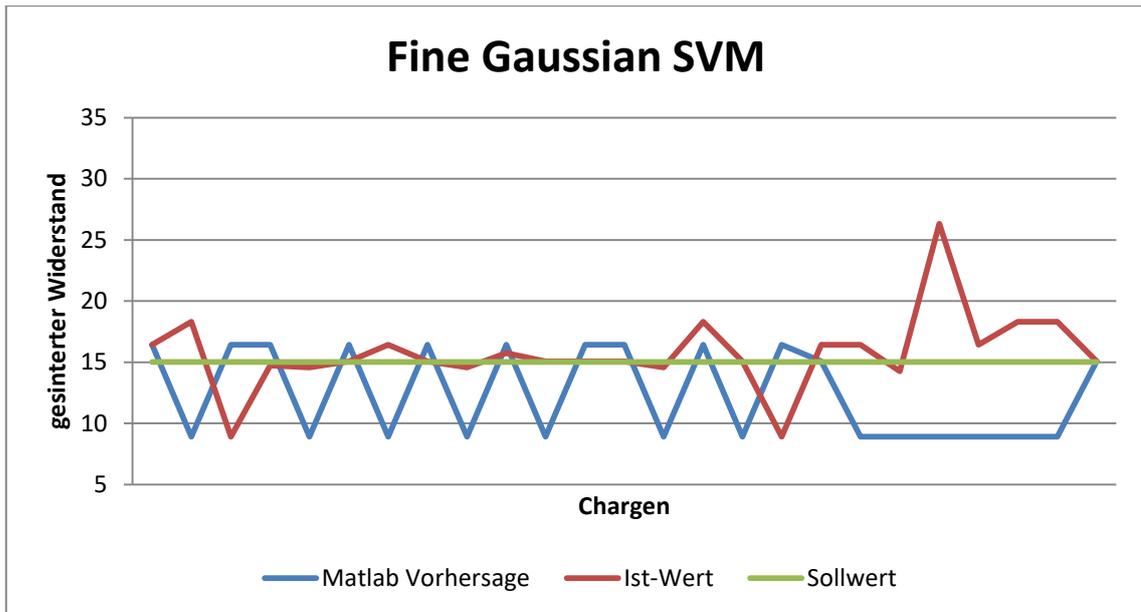


Abbildung 92: FGSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern)<sup>372</sup>

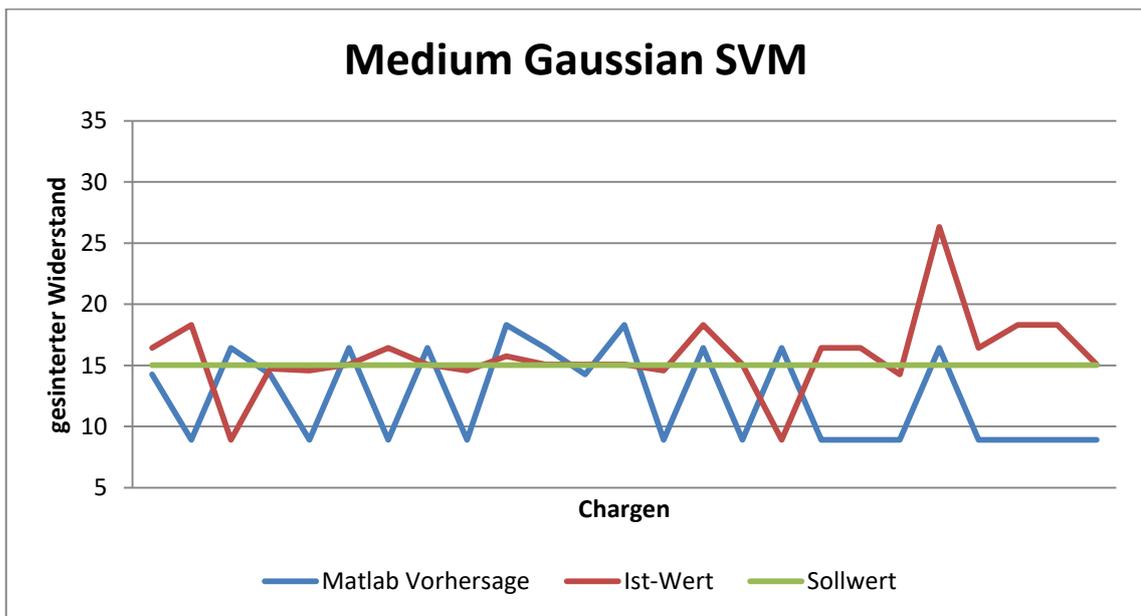


Abbildung 93: MGSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern)<sup>373</sup>

<sup>372</sup> Quelle: eigene Darstellung

<sup>373</sup> Quelle: ebda.

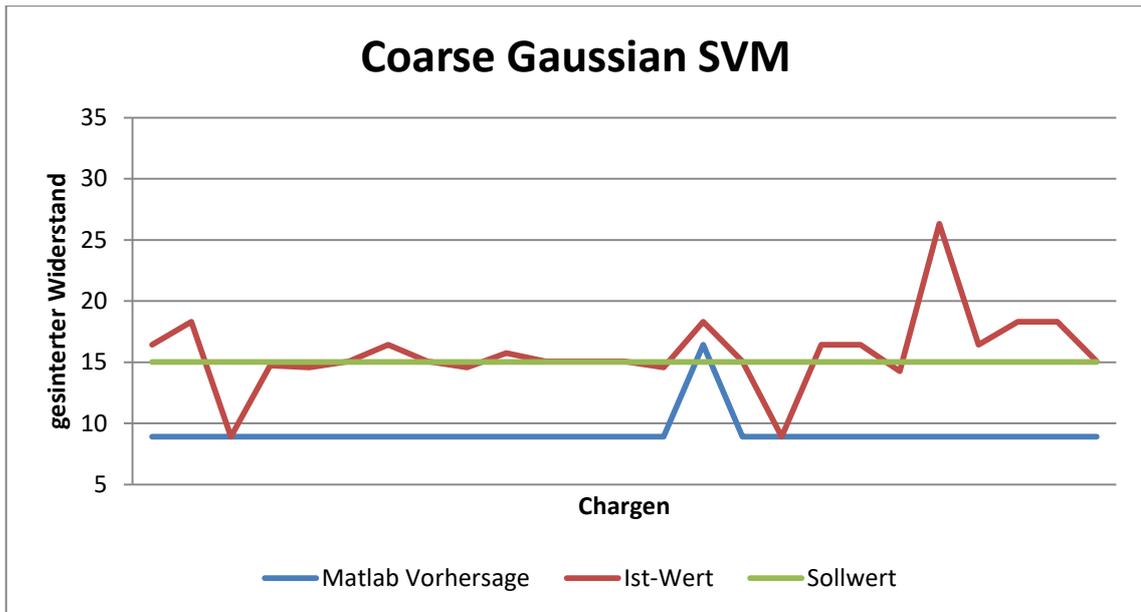


Abbildung 94: CGSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern)<sup>374</sup>

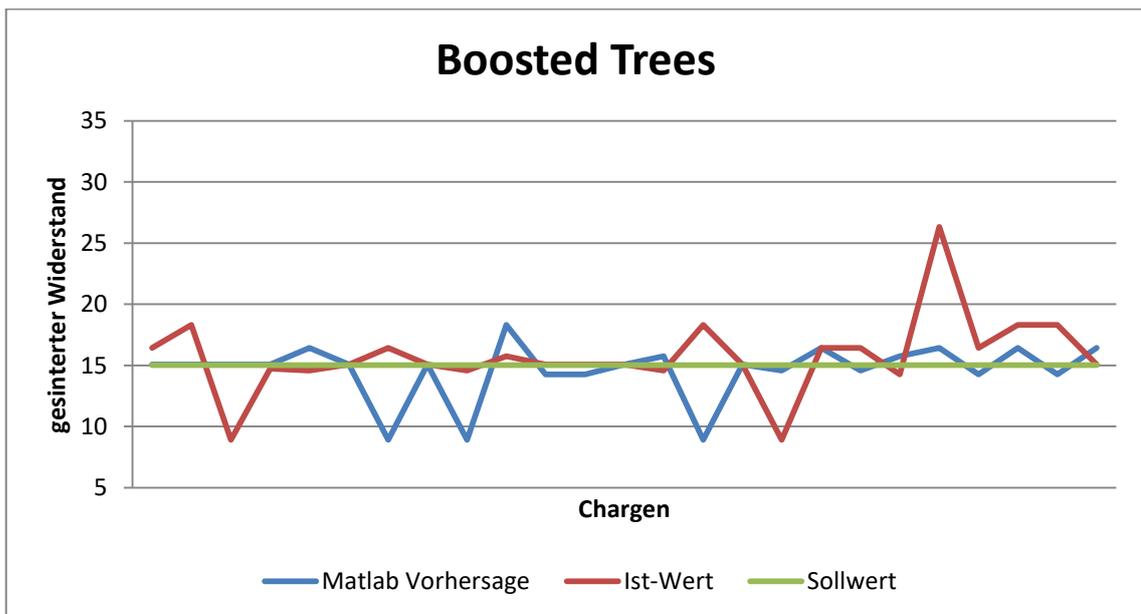


Abbildung 95: BoT: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern)<sup>375</sup>

<sup>374</sup> Quelle: eigene Darstellung

<sup>375</sup> Quelle: ebda.

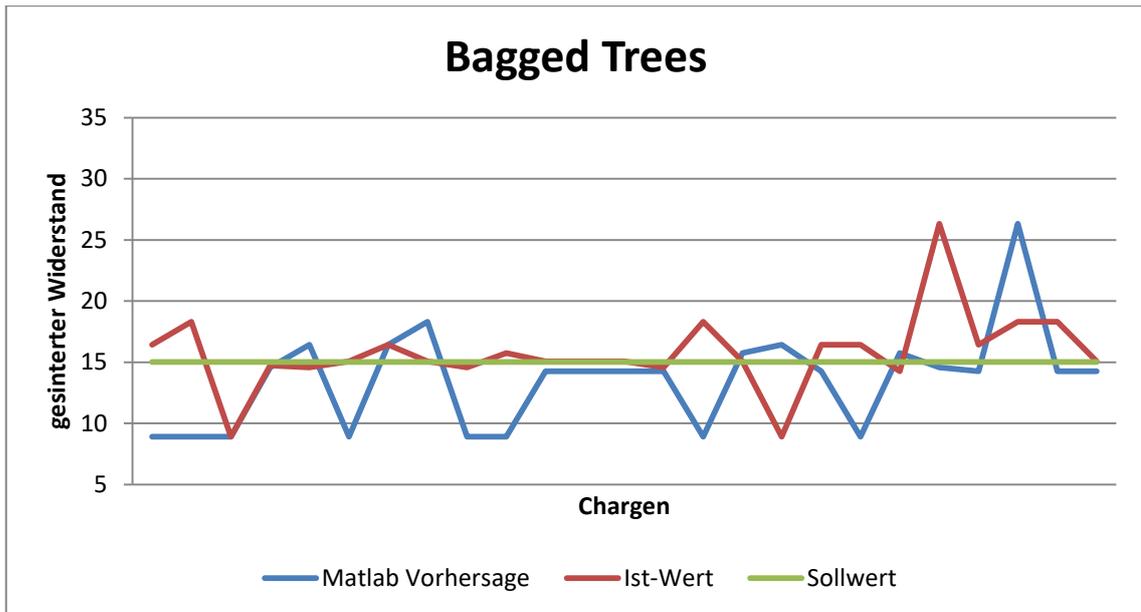


Abbildung 96: BaT: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern)<sup>376</sup>

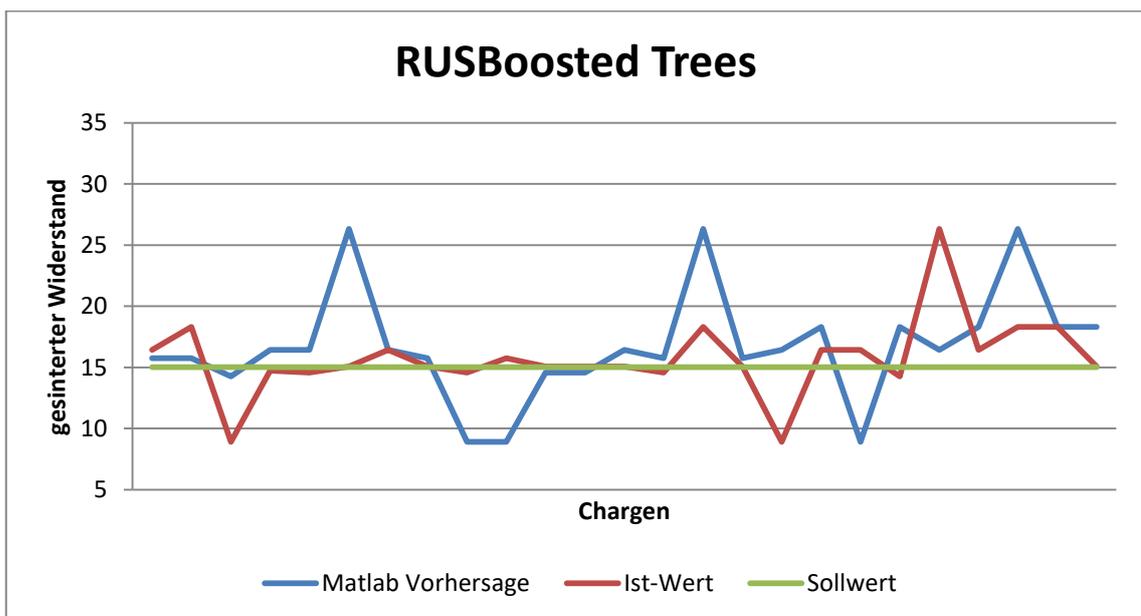


Abbildung 97: RUSBT: Vergleich Vorhersage zu Ist-Wert (Šumperk mit Ausreißern)<sup>377</sup>

<sup>376</sup> Quelle: eigene Darstellung

<sup>377</sup> Quelle: ebda.

**Medium Tree Šumperk mit Ausreißern**

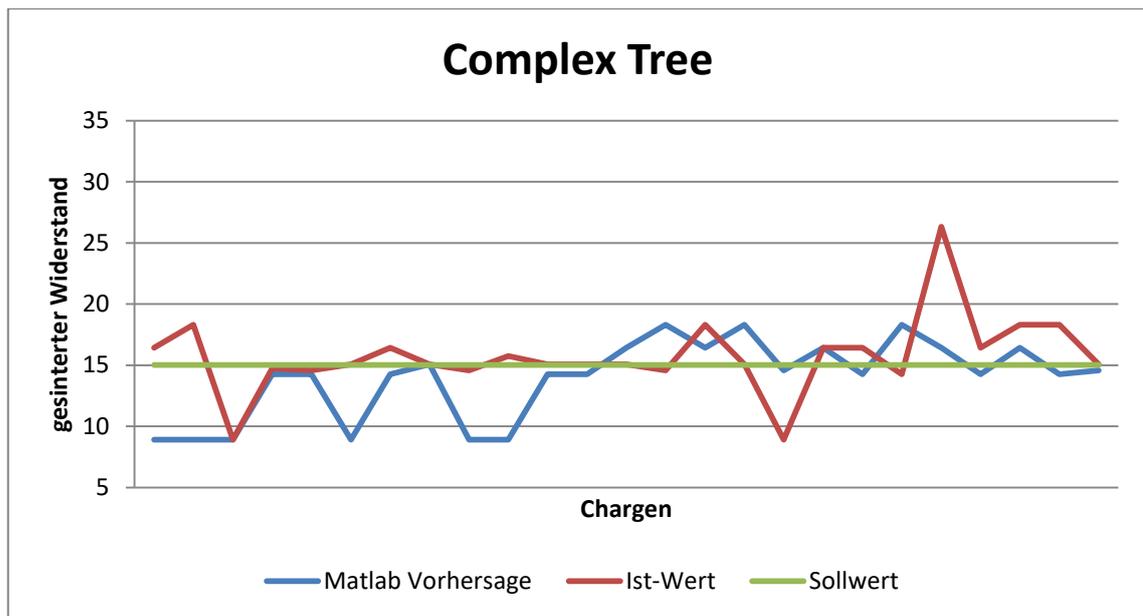
```
1  if SchlickerdichteMIN_VM_DL<0.191284 then node 2 elseif
SchlickerdichteMIN_VM_DL>=0.191284 then node 3 else 11.664
2  if SchlickertempMAX_VM_DL<-0.262925 then node 4 elseif
SchlickertempMAX_VM_DL>=-0.262925 then node 5 else 16.5
3  if pHWertSchlicker_VM_DL<-0.550633 then node 6 elseif
pHWertSchlicker_VM_DL>=-0.550633 then node 7 else 11.664
4  if d90VM_VM_DL<-0.773474 then node 8 elseif d90VM_VM_DL>=-
0.773474 then node 9 else 15.089
5  if Symm3MW_DL<-0.307843 then node 10 elseif Symm3MW_DL>=-
0.307843 then node 11 else 16.5
6  class = 15.68
7  if RadialeAuffederung_DL<-0.793705 then node 12 elseif
RadialeAuffederung_DL>=-0.793705 then node 13 else 11.664
8  class = 15.089
9  class = 18.336
10 if d10GranulatNM_DL<-0.147212 then node 14 elseif
d10GranulatNM_DL>=-0.147212 then node 15 else 14.589
11 if Ringtemperatur_VM_DL<0.336911 then node 16 elseif
Ringtemperatur_VM_DL>=0.336911 then node 17 else 16.5
12 class = 25.907
13 if SchlickerdichteMIN_VM_DL<0.251001 then node 18 elseif
SchlickerdichteMIN_VM_DL>=0.251001 then node 19 else 11.664
14 class = 14.241
15 if Heizgruppe11_VM_DL<0.646222 then node 20 elseif
Heizgruppe11_VM_DL>=0.646222 then node 21 else 14.589
16 class = 16.5
17 class = 18.336
18 class = 11.664
19 class = 15.089
20 class = 14.589
21 class = 14.241
```

**Complex Tree**

```
1  if SchlickerdichteMIN_VM_DL<0.191284 then node 2 elseif
SchlickerdichteMIN_VM_DL>=0.191284 then node 3 else 11.664
2  if SchlickertempMAX_VM_DL<-0.262925 then node 4 elseif
SchlickertempMAX_VM_DL>=-0.262925 then node 5 else 16.5
3  if pHWertSchlicker_VM_DL<-0.550633 then node 6 elseif
pHWertSchlicker_VM_DL>=-0.550633 then node 7 else 11.664
4  if d90VM_VM_DL<-0.773474 then node 8 elseif d90VM_VM_DL>=-
0.773474 then node 9 else 15.089
5  if Symm3MW_DL<-0.307843 then node 10 elseif Symm3MW_DL>=-
0.307843 then node 11 else 16.5
6  class = 15.68
7  if RadialeAuffederung_DL<-0.793705 then node 12 elseif
RadialeAuffederung_DL>=-0.793705 then node 13 else 11.664
8  class = 15.089
9  class = 18.336
10 if d10GranulatNM_DL<-0.147212 then node 14 elseif
d10GranulatNM_DL>=-0.147212 then node 15 else 14.589
11 if Ringtemperatur_VM_DL<0.336911 then node 16 elseif
Ringtemperatur_VM_DL>=0.336911 then node 17 else 16.5
12 class = 25.907
13 if SchlickerdichteMIN_VM_DL<0.251001 then node 18 elseif
SchlickerdichteMIN_VM_DL>=0.251001 then node 19 else 11.664
14 class = 14.241
15 if Heizgruppe11_VM_DL<0.646222 then node 20 elseif
Heizgruppe11_VM_DL>=0.646222 then node 21 else 14.589
16 class = 16.5
17 class = 18.336
18 class = 11.664
19 class = 15.089
20 class = 14.589
21 class = 14.241
```

Tabelle 39: Ergebnis CL und Testdaten: Šumperk ohne Ausreißer<sup>378</sup>

	Genauigkeit bei Training	Genauigkeit bei Anwendung auf Testdaten	Grenzfälle bei Anwendung auf Testdaten	Mittelwert der absoluten Differenz zwischen der vorhergesagten und Ist-Klasse
Complex Tree	25,0%	12,0%	0,0%	2,58
Medium Tree	25,0%	12,0%	0,0%	2,58
Simple Tree	23,3%	20,0%	4,0%	2,02
Linear SVM	15,0%	8,0%	0,0%	2,74
Quadratic SVM	21,7%	24,0%	0,0%	2,14
Cubic SVM	21,7%	24,0%	0,0%	1,82
Fine Gaussian SVM	18,3%	16,0%	4,0%	2,82
Medium Gaussian SVM	20,0%	4,0%	0,0%	3,02
Coarse Gaussian SVM	15,0%	8,0%	0,0%	3,62
Boosted Trees	23,3%	12,0%	0,0%	2,34
Bagged Trees	15,0%	16,0%	0,0%	2,58
RUSBoosted Trees	19,3%	4,0%	4,0%	2,50

Abbildung 98: CT: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer)<sup>379</sup><sup>378</sup> Quelle: eigene Darstellung<sup>379</sup> Quelle: ebda.

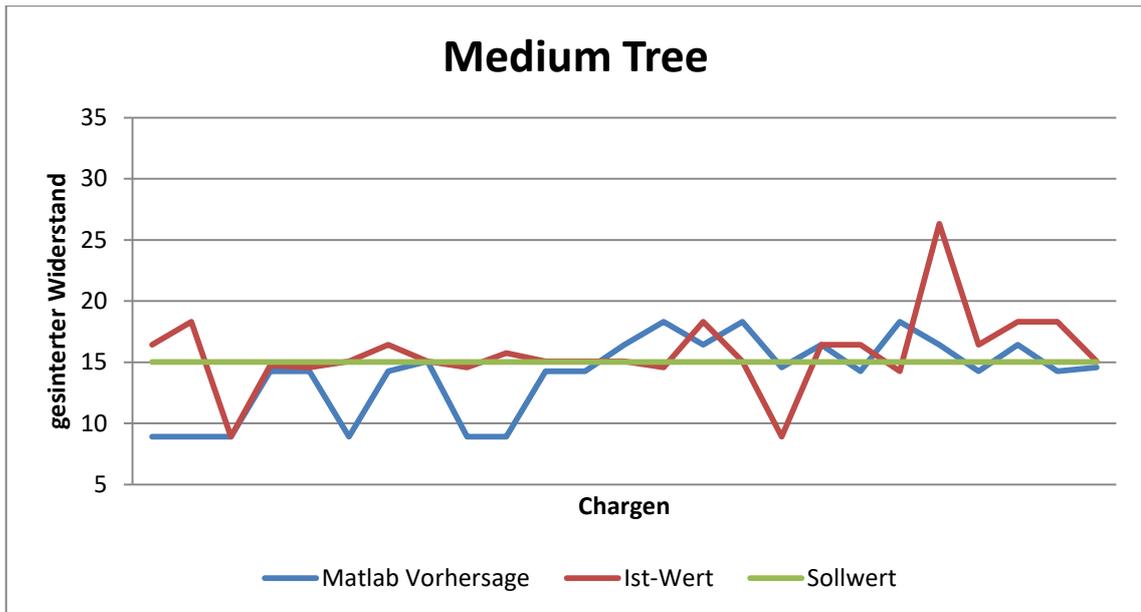


Abbildung 99: MT: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer)<sup>380</sup>

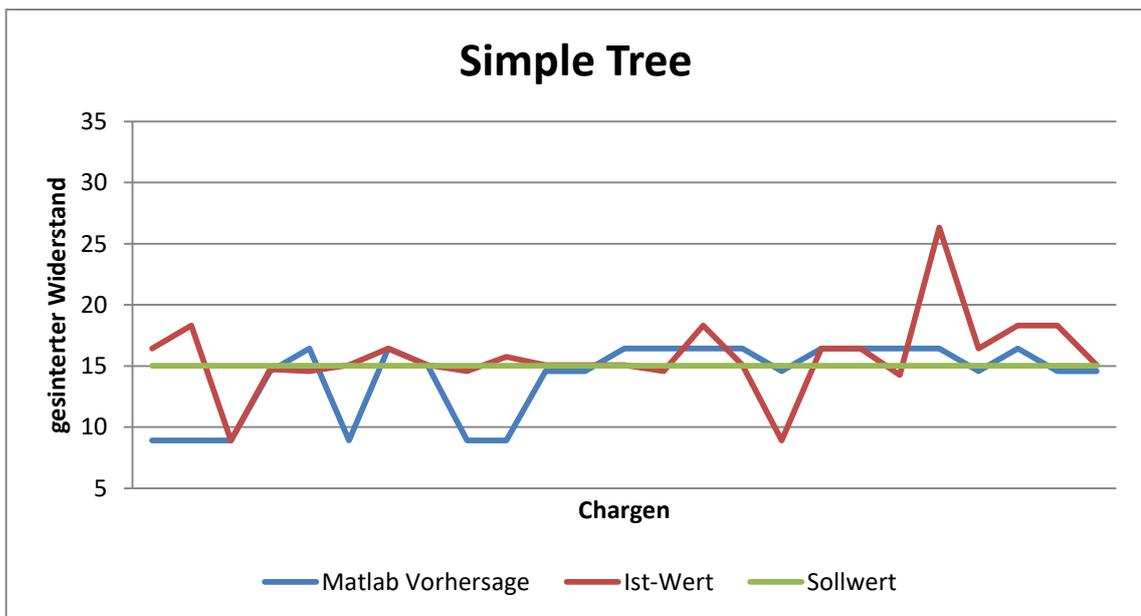


Abbildung 100: ST: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer)<sup>381</sup>

<sup>380</sup> Quelle: eigene Darstellung

<sup>381</sup> Quelle: ebda.

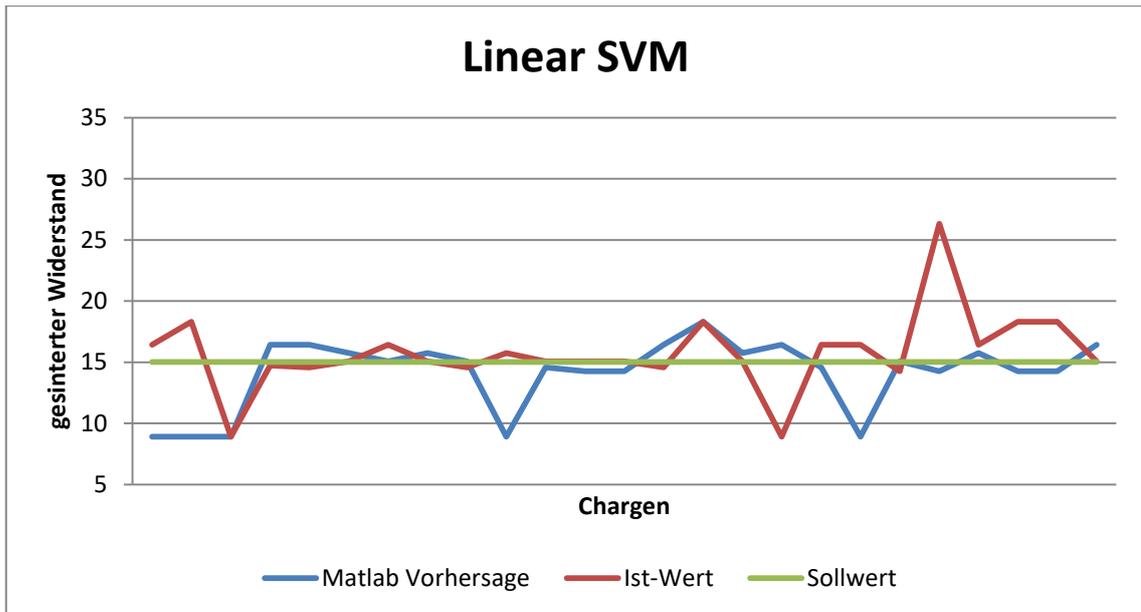


Abbildung 101: LSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer)<sup>382</sup>

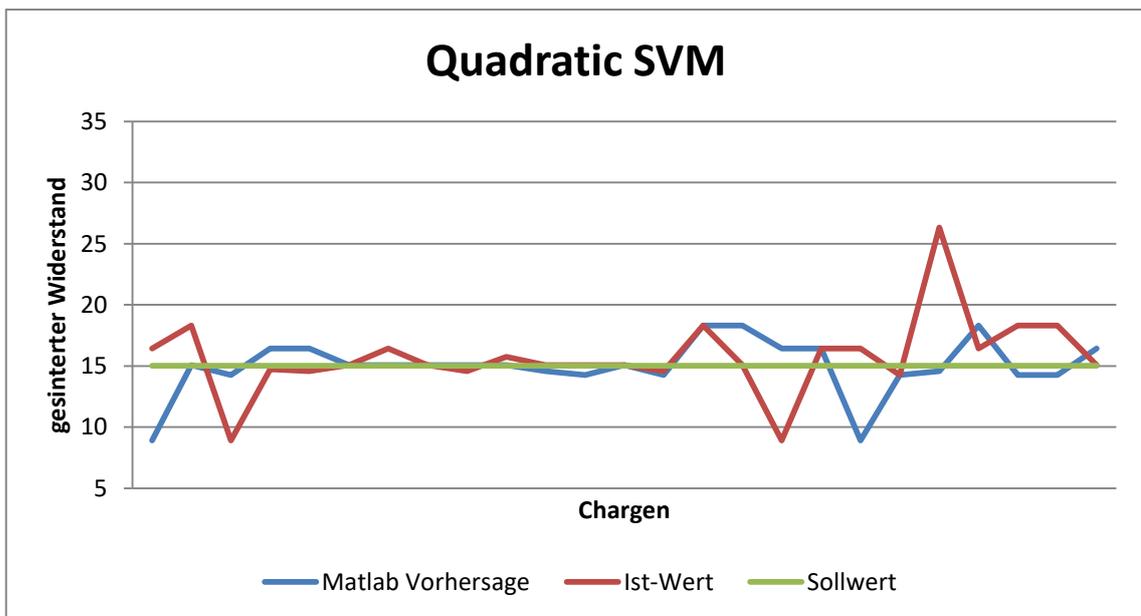


Abbildung 102: QSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer)<sup>383</sup>

<sup>382</sup> Quelle: eigene Darstellung

<sup>383</sup> Quelle: ebda.

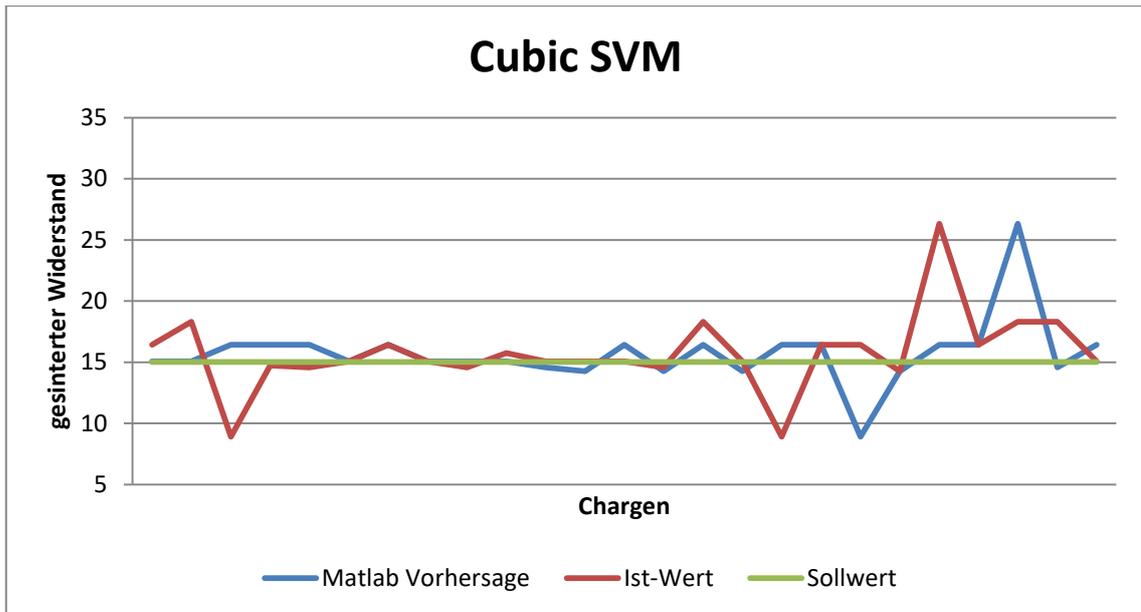


Abbildung 103: CSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer)<sup>384</sup>

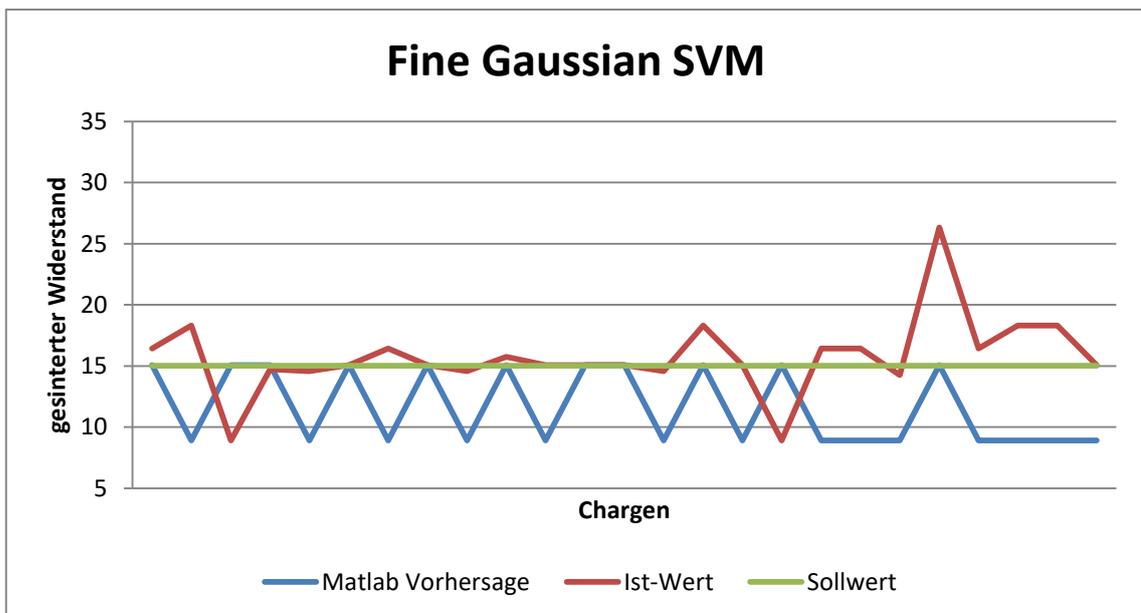


Abbildung 104: FGSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer)<sup>385</sup>

<sup>384</sup> Quelle: eigene Darstellung

<sup>385</sup> Quelle: ebda.

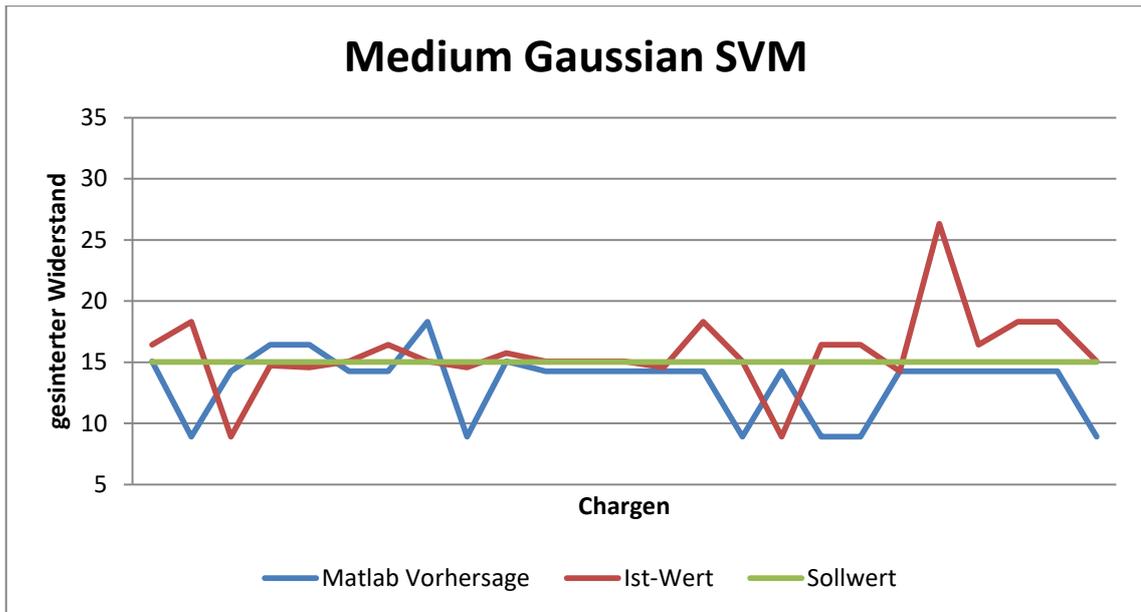


Abbildung 105: MGSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer)<sup>386</sup>

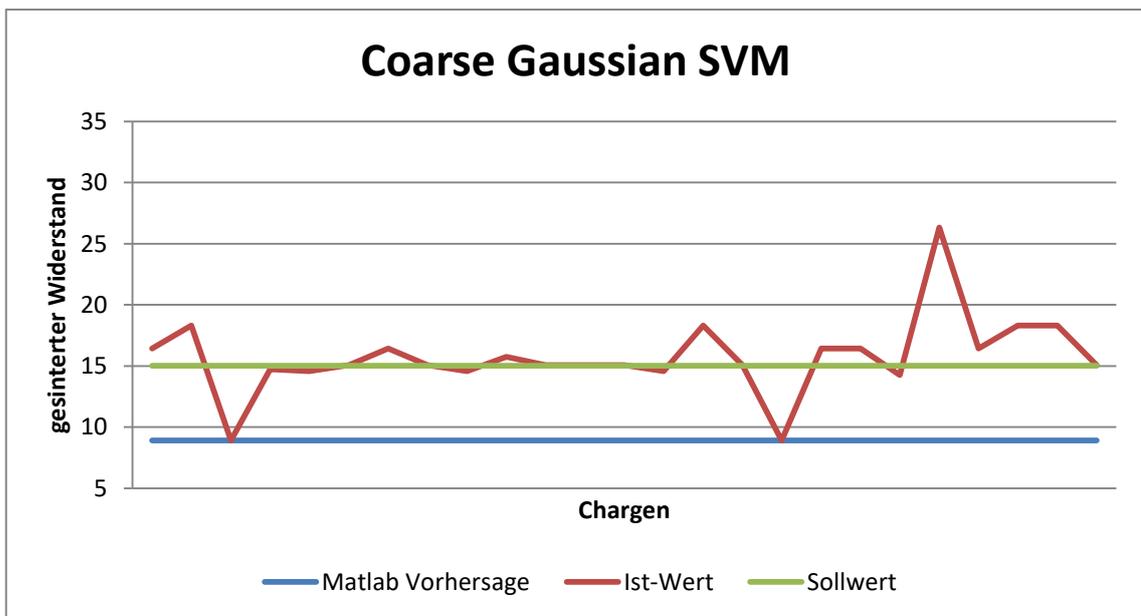


Abbildung 106: CGSVM: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer)<sup>387</sup>

<sup>386</sup> Quelle: eigene Darstellung

<sup>387</sup> Quelle: ebda.

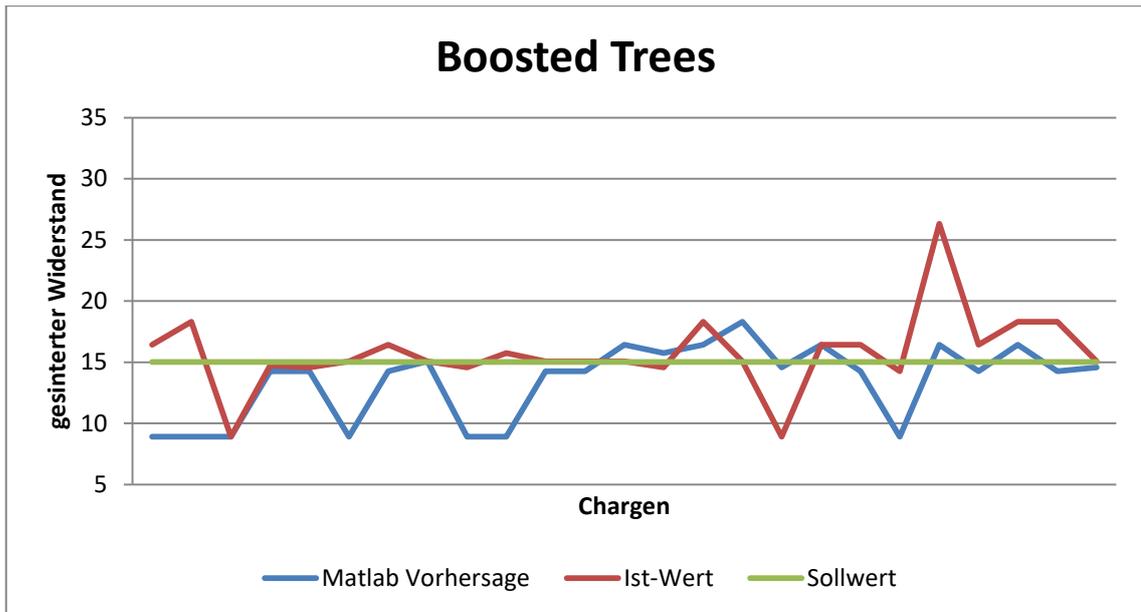


Abbildung 107: BoT: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer)<sup>388</sup>

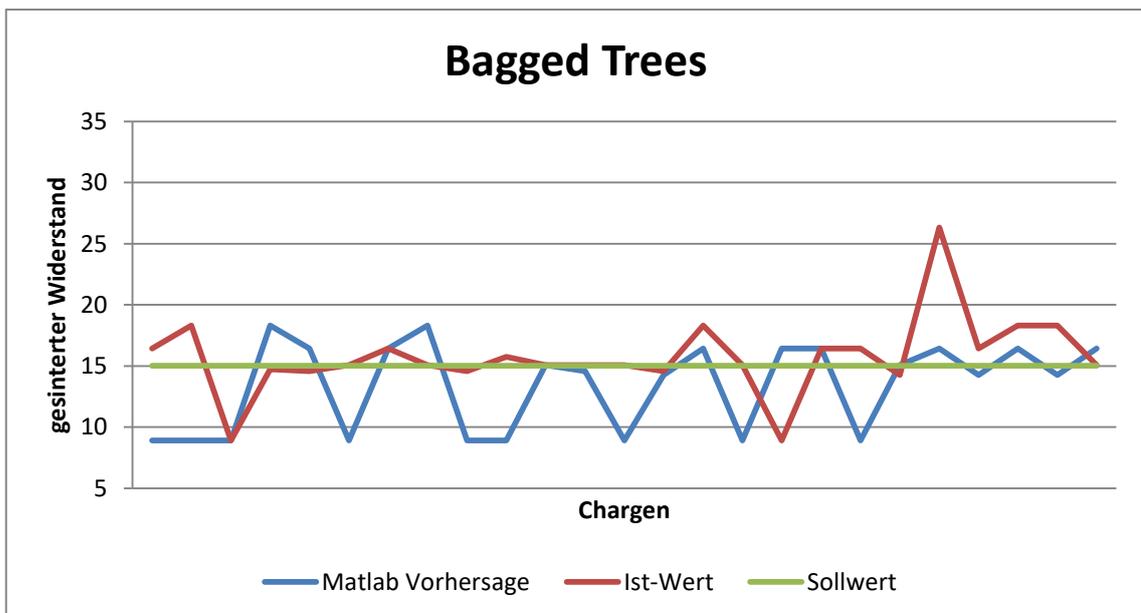


Abbildung 108: BaT: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer)<sup>389</sup>

<sup>388</sup> Quelle: eigene Darstellung

<sup>389</sup> Quelle: ebda.

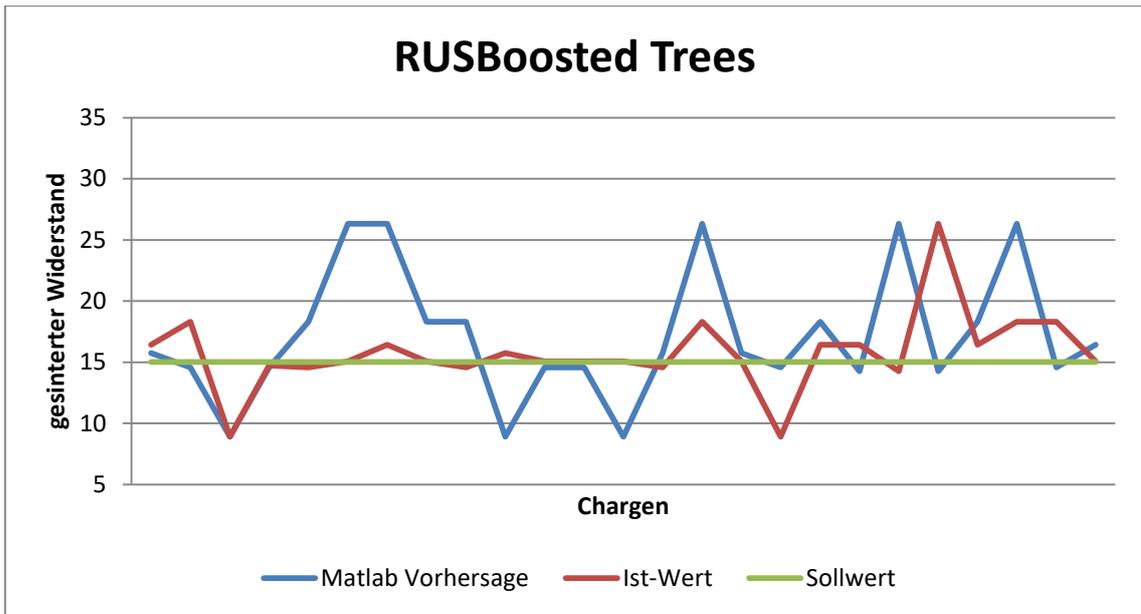


Abbildung 109: RUSBT: Vergleich Vorhersage zu Ist-Wert (Šumperk ohne Ausreißer)<sup>390</sup>

<sup>390</sup> Quelle: eigene Darstellung

**Medium Tree Šumperk ohne Ausreißer**

```
1  if SchlickerdichteMIN_VM_DL<0.571801 then node 2 elseif
SchlickerdichteMIN_VM_DL>=0.571801 then node 3 else 11.664
2  if Symm3MW_DL<-0.435851 then node 4 elseif Symm3MW_DL>=-
0.435851 then node 5 else 16.5
3  if pHWertSchlicker_VM_DL<-0.550633 then node 6 elseif
pHWertSchlicker_VM_DL>=-0.550633 then node 7 else 11.664
4  if d10GranulatNM_DL<-0.147212 then node 8 elseif
d10GranulatNM_DL>=-0.147212 then node 9 else 14.589
5  if SchlickertempMAX_VM_DL<0.0256806 then node 10 elseif
SchlickertempMAX_VM_DL>=0.0256806 then node 11 else 16.5
6  class = 15.68
7  if SchlickertempStabw_VM_DL<-1.26292 then node 12 elseif
SchlickertempStabw_VM_DL>=-1.26292 then node 13 else 11.664
8  class = 14.241
9  if Heizgruppe11_VM_DL<0.660181 then node 14 elseif
Heizgruppe11_VM_DL>=0.660181 then node 15 else 14.589
10 if d90VM_VM_DL<-0.773474 then node 16 elseif d90VM_VM_DL>=-
0.773474 then node 17 else 15.089
11 if Symm3MW_DL<0.44651 then node 18 elseif
Symm3MW_DL>=0.44651 then node 19 else 16.5
12 class = 15.089
13 if RadialeAuffederung_DL<-0.793705 then node 20 elseif
RadialeAuffederung_DL>=-0.793705 then node 21 else 11.664
14 class = 14.589
15 class = 14.241
16 class = 15.089
17 if Heizgruppe13_VM_DL<0.605687 then node 22 elseif
Heizgruppe13_VM_DL>=0.605687 then node 23 else 18.336
18 class = 16.5
19 class = 18.336
20 class = 25.907
21 if Ofen in {19 20 22 25 BV} then node 24 elseif Ofen=24 then
node 25 else 11.664
22 class = 14.241
23 class = 18.336
24 class = 11.664
25 class = 18.336
```

**Complex Tree Šumperk ohne Ausreißer**

```
1  if SchlickerdichteMIN_VM_DL<0.571801 then node 2 elseif
SchlickerdichteMIN_VM_DL>=0.571801 then node 3 else 11.664
2  if Symm3MW_DL<-0.435851 then node 4 elseif Symm3MW_DL>=-
0.435851 then node 5 else 16.5
3  if pHWertSchlicker_VM_DL<-0.550633 then node 6 elseif
pHWertSchlicker_VM_DL>=-0.550633 then node 7 else 11.664
4  if d10GranulatNM_DL<-0.147212 then node 8 elseif
d10GranulatNM_DL>=-0.147212 then node 9 else 14.589
5  if SchlickertempMAX_VM_DL<0.0256806 then node 10 elseif
SchlickertempMAX_VM_DL>=0.0256806 then node 11 else 16.5
6  class = 15.68
7  if SchlickertempStabw_VM_DL<-1.26292 then node 12 elseif
SchlickertempStabw_VM_DL>=-1.26292 then node 13 else 11.664
8  class = 14.241
9  if Heizgruppe11_VM_DL<0.660181 then node 14 elseif
Heizgruppe11_VM_DL>=0.660181 then node 15 else 14.589
10 if d90VM_VM_DL<-0.773474 then node 16 elseif d90VM_VM_DL>=-
0.773474 then node 17 else 15.089
11 if Symm3MW_DL<0.44651 then node 18 elseif
Symm3MW_DL>=0.44651 then node 19 else 16.5
12 class = 15.089
13 if RadialeAuffederung_DL<-0.793705 then node 20 elseif
RadialeAuffederung_DL>=-0.793705 then node 21 else 11.664
14 class = 14.589
15 class = 14.241
16 class = 15.089
17 if Heizgruppe13_VM_DL<0.605687 then node 22 elseif
Heizgruppe13_VM_DL>=0.605687 then node 23 else 18.336
18 class = 16.5
19 class = 18.336
20 class = 25.907
21 if Ofen in {19 20 22 25 BV} then node 24 elseif Ofen=24 then
node 25 else 11.664
22 class = 14.241
23 class = 18.336
24 class = 11.664
25 class = 18.336
```