

Eckhard Reh  
Chemometrie

## Weitere empfehlenswerte Titel



*Analytik.*  
*Daten, Formeln, Übungsaufgaben*  
Küster, Thiel, 2016  
ISBN 978-3-11-041495-0, e-ISBN 978-3-11-041496-7



*Trennungsmethoden der Analytischen Chemie*  
Bock, Nießner, 2014  
ISBN 978-3-11-026544-6, e-ISBN 978-3-11-026637-5



*Allgemeine und Anorganische Chemie.*  
*11. Auflage*  
Riedel, Meyer, 2013  
ISBN 978-3-11-026919-2, e-ISBN 978-3-11-027013-6



*Grundlagen der Organischen Chemie.*  
*5. Auflage*  
Buddrus, Schmidt, 2014  
ISBN 978-3-11-030559-3, e-ISBN 978-3-11-033105-9



*Physikalische Chemie.*  
*Für die Bachelorprüfung*  
Motschmann, Hofmann, 2014  
ISBN 978-3-11-034877-4, e-ISBN 978-3-11-034878-1

Eckhard Reh

# Chemometrie

---

Grundlagen der Statistik, Numerischen Mathematik  
und Softwareanwendung in der Chemie

**DE GRUYTER**

**Autor**

Prof. Dr. Eckhard Reh  
Technische Hochschule Bingen  
Berlinstr. 109  
55411 Bingen  
e.reh@th-bingen.de

ISBN 978-3-11-045100-9  
e-ISBN (PDF) 978-3-11-045103-0  
e-ISBN (EPUB) 978-3-11-045107-8

**Library of Congress Cataloging-in-Publication Data**

A CIP catalog record for this book has been applied for at the Library of Congress.

**Bibliografische Information der Deutschen Nationalbibliothek**

Die Deutsche Nationalbibliothek verzeichnet diese Publikation in der Deutschen Nationalbibliografie; detaillierte bibliografische Daten sind im Internet über <http://dnb.dnb.de> abrufbar.

© 2017 Walter de Gruyter GmbH, Berlin/Boston  
Satz: PTP-Berlin, Protago-TEX-Production GmbH, Berlin  
Druck und Bindung: CPI books GmbH, Leck  
☉ Gedruckt auf säurefreiem Papier  
Printed in Germany

[www.degruyter.com](http://www.degruyter.com)

---

ad maiorem dei gloriam



## Vorwort

Auf Grund leistungsstarker, kostengünstiger Prozessoren und nahezu unbegrenzter Speicherkapazität hat die Anwendung der EDV in der Chemie einen rasanten Aufschwung erlebt. Auch hier ist die Verwendung effizienter Algorithmen zur täglichen Praxis geworden.

Daher ist es insbesondere in diesem Bereich wichtig, die entsprechenden Grundlagen zu verstehen für einen validierten Einsatz der unterschiedlichsten Tools und Software-Produkte. Beispielhaft wird dies beim Einsatz von Maximum-Entropie-Methoden in der Massenspektrometrie deutlich. Die Resultate sind teilweise frappierend, um so mehr muss hinterfragt werden, welche Relevanz die erzeugten Ergebnisse haben. In diesem Fall werden z.B. in unserem Institut MaxEnt-Resultate nur angegeben, wenn sie mit der konventionellen Dekonvolution zumindest qualitativ bestätigt werden können.

Der Impuls, eine Monographie zur Chemometrie zu verfassen, liegt in der sehr geringen Zahl deutschsprachiger Werke zu dieser Thematik. Während es im angelsächsischen Raum sehr gute Bücher gibt (Brereton, R. G., *Data Analysis für Laboratory and Chemical Plant*, John Wiley oder Massart et al., *Handbook of Chemometrics and Qualimetrics*, Elsevier) besteht im deutschsprachigen Raum ein gravierender Mangel an Fachliteratur. Da die Thematik auf der Grund der numerischen und statistischen Inhalte nicht einfach ist, erschwert eine englischsprachige Behandlung das Verständnis und die Anwendung zumeist. Dies musste auch der Autor erleben bei einem mehrwöchigen Chemometrie-Kurs an der Universität Bristol (Leitung Prof. Brereton). Diese Erfahrung war eine wesentliche Triebkraft, ein deutschsprachiges Werk zur Chemometrie zur Verfügung zu stellen.

Darin liegt eventuell auch der Grund, dass das Fach Chemometrie im Chemie-Studium an den meisten deutschen Universitäten nicht vertreten ist oder das manch etablierter Chemiker sich mit dem Einsatz der Chemometrie schwer tut. Dabei sind die meisten Chemometrie-Themen in unterschiedlichen Bereichen (Physikalische, Organische, Analytische Chemie ..., Chemische Prozesstechnik, Biotechnologie ..., Lebensmittel-Technik, Umwelt-Technik, Klinische Diagnostik) sehr effizient einzusetzen. Wie wichtig wäre z.B. die Anwendung der statistischen Versuchsplanung oder der diversen Optimierungs-Strategien (Simplex, RSM) für meine Diplom- oder Promotionsarbeit, wären diese zum damaligen Zeitpunkt bekannt und einsetzbar gewesen.

Vor diesem Hintergrund wurde daher versucht, die teilweise schwierigen, numerischen Grundlagen möglichst verständlich zu präsentieren. Hierzu wurde nicht nur auf eine umfassende Darstellung ohne historische Alternativen oder diverse Optionen verzichtet, wichtig war auch die wiederholte Anwendung der verschiedenen Algorithmen durch einfache Beispiele, die teilweise manuell nachvollzogen werden können.

Es hat sich gezeigt, dass gerade diese Umsetzung für das Verständnis besonders wichtig sein kann.

Trotzdem soll das vorgelegte Buch nicht nur die theoretischen Grundlagen erarbeiten, sondern insbesondere die Anwendung in der Praxis hervorheben. Als Hilfestellung für den Anwender dient daher die direkte Umsetzung mit Hilfe von speziellen Software-Produkten als auch mit dem allgemein etablierten Statistik-Paket R. Das Ziel ist damit, durch verständliche Präsentation der Grundlagen und direkte Umsetzung dem Anwender in der chemischen Praxis eine Hilfestellung für die Routine-Anwendung zu geben.

Vielleicht gelingt es damit auch, der Thematik Chemometrie im deutschsprachigen Raum einen größeren Stellenwert und eine breitere, fundierte Anwendung in der chemischen Praxis zu verleihen.

## Der Autor

Prof. Eckhard Reh studierte Chemie an der Universität Siegen. Er promovierte im Bereich Biochemische Analyse / Klinische Chemie an der Universität des Saarlandes. Fast 10 Jahre leitete er die Gruppe Proteinanalyse im Biochemie Forschungszentrum Tutzing der Fa. Boehringer-Mannheim.

Danach nahm Prof. Reh einen Ruf an die Technische Hochschule Bingen für das Fach Analytische Chemie im Studiengang Biotechnologie an. Seit fast 20 Jahren ist er hier auch Institutsleiter des ZENTRUM PROTEINANALYSE.

# Inhalt

## Vorwort — VII

- 1 Grundlagen der Chemometrie — 1**
  - 1.1 Prinzipien und Disziplinen — 1
  - 1.2 Anwendungsbereiche — 1
  - 1.3 Realisierung — 2
  - 1.4 Zielsetzung des Buchs — 4
  
- 2 Statistische Parameter und Prüfverfahren — 7**
  - 2.1 Einführung — 7
  - 2.2 Deskriptive Statistik — 7
    - 2.2.1 Erstes statistisches Moment — 8
    - 2.2.2 Zweites statistisches Moment — 10
    - 2.2.3 Drittes statistisches Moment — 13
    - 2.2.4 Viertes statistisches Moment — 13
  - 2.3 Prüfmethode — 15
    - 2.3.1 Prüfung der Homogenität — 15
    - 2.3.2 Prüfung der Verteilung der Stichprobe — 16
    - 2.3.3 Trendtest — 18
    - 2.3.4 Zusammenfassung — 19
  - 2.4 Vergleich statistischer Parameter — 19
    - 2.4.1 Vergleich von Mittelwerten — 19
    - 2.4.2 Vergleich von Varianzen — 21
  - 2.5 Literatur — 22
  - 2.6 Übungen — 22
  - 2.7 Softwareanwendung — 22
    - 2.7.1 Einsatz MiniStat — 23
    - 2.7.2 Umsetzung in RStudio/R — 25
  
- 3 Versuchsplanung, Prozessoptimierung — 29**
  - 3.1 Einführung — 29
  - 3.2 Statistische Versuchsplanung — 30
    - 3.2.1 Grundlagen — 30
    - 3.2.2 Aufstellung des Versuchsplans — 32
    - 3.2.3 Auswertung des Versuchsplans — 34
    - 3.2.4 Spezielle Aspekte — 36
    - 3.2.5 Einsatzbereich und Grenzen — 41
  - 3.3 Simplex-Optimierung — 42
    - 3.3.1 Prinzip des Standard-Simplex — 42

3.3.2	Modifizierter Simplex	43
3.3.3	Simplex-Limitierungen	46
3.3.4	Vor- und Nachteile	47
3.4	Response-Surface-Modelling (RSM)	47
3.4.1	Prinzip	47
3.4.2	Untersuchungsdesign	48
3.4.3	Modellberechnung	52
3.4.4	Modellinterpretation	54
3.4.5	Varianzanalyse (ANOVA)	56
3.4.6	Vor- und Nachteile	57
3.5	Literatur	58
3.6	Übungen	58
3.7	Softwareanwendung	59
3.7.1	Einsatz VPlan, SimSoft, RSMSoft	60
3.7.2	Umsetzung in RStudio/R	66
<b>4</b>	<b>Univariate Regression, Kalibration</b>	<b>71</b>
4.1	Einleitung	71
4.2	Regression	72
4.2.1	Modelldesign	72
4.2.2	Lineare Regression	73
4.2.3	Nichtlineare Regression	79
4.2.4	Robuste Regression	81
4.2.5	Prüfung der Adäquatheit des Regressionsmodells	83
4.2.6	Prüfung des Achsenabschnitts	87
4.2.7	Nachweis-, Bestimmungsgrenze	87
4.2.8	Behandlung von Ausnahmen	88
4.3	Kalibrationsmethoden	89
4.3.1	$\vartheta$ -Kalibration	89
4.3.2	$\sigma$ -Kalibration	90
4.3.3	$\delta$ -Kalibration (Standard-Addition)	91
4.4	Literatur	92
4.5	Übungen	92
4.6	Softwareanwendung	94
4.6.1	Einsatz Calib, SimFit	95
4.6.2	Umsetzung in RStudio/R	99
<b>5</b>	<b>Analyse von Messreihen</b>	<b>103</b>
5.1	Digitales Filtern, Glätten	103
5.1.1	Grundlagen	103
5.1.2	Savitzky-Golay-Glättung	103
5.1.3	Kalman-Filter	106

- 5.2 Ableitung — 108
- 5.3 Autokorrelation — 109
  - 5.3.1 Grundlagen — 109
  - 5.3.2 Unkorreliertheit — 112
  - 5.3.3 Periodizität — 112
  - 5.3.4 Abfall — 113
  - 5.3.5 Drift — 114
- 5.4 Fourier-Transformation — 114
  - 5.4.1 Aufgabenstellung — 114
  - 5.4.2 Numerische Grundlagen — 116
  - 5.4.3 Signal-Filterung — 120
  - 5.4.4 Faltung — 121
- 5.5 Literatur — 123
- 5.6 Übungen — 123
- 5.7 Softwareanwendung — 124
  - 5.7.1 Einsatz DTrans — 124
  - 5.7.2 Umsetzung in RStudio/R — 126
  
- 6 Signaldekonvolution — 129**
  - 6.1 Einführung — 129
  - 6.2 Peakform-Analyse — 130
    - 6.2.1 Wahl des Modells — 130
    - 6.2.2 Wahl des Algorithmus — 131
    - 6.2.3 Evaluierung — 132
  - 6.3 Fourier-Dekonvolution — 135
    - 6.3.1 Grundlagen — 135
    - 6.3.2 Anwendung in der IR-Spektroskopie — 137
  - 6.4 Maximum-Entropie-Dekonvolution — 138
    - 6.4.1 Grundlagen — 138
    - 6.4.2 Anwendung in der Spektroskopie — 142
  - 6.5 Literatur — 142
  - 6.6 Übungen — 143
  - 6.7 Softwareanwendung — 143
    - 6.7.1 Einsatz PeakCalc — 143
  
- 7 Mustererkennung, Clusteranalyse — 145**
  - 7.1 Einführung — 145
  - 7.2 Grundlagen — 146
    - 7.2.1 Graphische Repräsentation — 147
    - 7.2.2 Konsistenzprüfung — 148
    - 7.2.3 Zentrierung, Skalierung, Normierung — 148
    - 7.2.4 Varianz-Kovarianz-Matrix — 150

7.3	Hauptkomponenten — 152
7.3.1	Hauptkomponentenanalyse ( <i>principal component analysis</i> ) — 153
7.3.2	Zahl signifikanter Hauptkomponenten — 159
7.3.3	Graphische Interpretation — 162
7.4	Mustererkennung — 164
7.4.1	Hierarchische Clusteranalyse — 165
7.4.2	Nicht-hierarchische Clusteranalyse — 171
7.5	Klassifizierung — 176
7.5.1	Methode der $k$ nächsten Nachbarn — 178
7.5.2	SIMCA-Methode — 179
7.6	Literatur — 183
7.7	Übungen — 184
7.8	Softwareanwendung — 184
7.8.1	Einsatz Cluster — 185
7.8.2	Umsetzung in RStudio/R — 187
<b>8</b>	<b>Multivariate Kalibration — 189</b>
8.1	Einführung — 189
8.2	Klassische Kalibration — 191
8.2.1	Direkte Kalibration — 191
8.2.2	Indirekte Kalibration — 193
8.3	Inverse Kalibration — 196
8.3.1	P-Matrix-Verfahren — 197
8.3.2	PCR-Verfahren — 198
8.3.3	PLS-Verfahren — 200
8.4	Validierung der Kalibration — 200
8.4.1	Residuen und RSS-Wert — 201
8.4.2	Kreuzvalidierung und PRESS-Wert — 203
8.4.3	hat-Matrix — 204
8.4.4	Ausreißer — 205
8.5	Literatur — 208
8.6	Übungen — 208
8.7	Softwareanwendung — 209
8.7.1	Einsatz Unscambler <sup>®</sup> — 209
8.7.2	Umsetzung in RStudio/R — 209
<b>9</b>	<b>Softwareanwendung — 211</b>
9.1	Kommerzielle Statistikpakete — 211
9.2	Spezielle Softwarepakete — 211
9.3	RStudio/R — 212
9.4	Literatur — 213

**A Anhang — 215**

**Abkürzungen — 225**

**Stichwortverzeichnis — 227**



# 1 Grundlagen der Chemometrie

## 1.1 Prinzipien und Disziplinen

Die Chemometrie basiert auf den grundlegenden Arbeiten von Kowalski und Wold ab 1971, der Begriff Chemometrics wurde ab 1972 explizit in der Fachliteratur verwendet. Eine bis heute gültige Definition geht ebenfalls auf Kowalski zurück:

Chemometrics is defined as a chemical diszipline, that uses mathematical, statistical and other methods employing formal logic a) to design or select optimal measurement procedures and experiments, b) to provide maximum relevant chemical information by analyzing chemical data.

Die Themenbereiche, die neben den chemischen Aspekten in der Chemometrie Eingang finden, sind statistische und mathematische Methoden, Algorithmen der Informatik und Datenstrukturen für Bibliotheken im Hinblick auf eine Anwendung in der Chemie. Anders als in anderen chemischen Disziplinen wie z. B. in der organischen Chemie, wo hunderte organischer Reaktionen zur Anwendung kommen, sind es in der Chemometrie nur wenig Basisprinzipien, auf die im Wesentlichen zurück gegriffen werden kann. Dazu gehören Prinzipien, z. B. aus der

- Stochastik (Basis-Statistik, Fehlerfortpflanzung, Versuchsplanung, ...)
- Datenanalyse (Messreihenbehandlung, Hauptkomponenten, Clusteranalyse, ...)
- Modellierung (Korrelationsanalyse, Regression, Variablentransformation, ...)
- Systemtheorie (Informationstheorie, Bildverarbeitung, Künstliche Intelligenz, ...)

Bekanntermaßen ist nicht bei allen Chemikern der Einsatz der Mathematik und Statistik sehr positiv besetzt. Daher soll hervorgehoben werden, dass Mathematik und Statistik nicht die Basisdisziplinen der Chemometrie sind, sondern wichtige Werkzeuge liefern. Trotzdem kann ein grundlegendes Verständnis für diese Werkzeuge in der Chemometrie nicht ignoriert werden. Hauptbestandteil bleibt aber die Umsetzung für eine Anwendung in der Chemie und damit eine Thematik, die viele Chemiker in unterschiedlichsten Bereichen betrifft.

Besonders hervorzuheben ist, dass die Chemometrie keine Spezialdisziplin der Analytischen Chemie ist, wie zu Beginn propagiert wurde. Dies zeigt sich bereits in den Anfängen an einem der Begründer, Wold, aus dem Institut für Organische Chemie der Universität Umea.

## 1.2 Anwendungsbereiche

Als grundlegende Aufgabenfelder der Chemometrie können formuliert werden

- Planung, Optimierung chemischer Prozesse oder Experimente (Versuchsplanung, Simplex-, RSM-Optimierung, ...)

- Auswertung von Resultaten z. B. in Prozess-Kontrolle oder experimentellem Verlauf (Messreihenanalyse, Datentransformation, ...)
- Interpretation und Validierung von Resultaten chemischer Experimente (Regression, Teststatistik, Varianzanalyse, ...)

Diese Anwendungsfelder sind in unterschiedlichsten Disziplinen angesiedelt, z. B.

- Physikalische Chemie, Organische Chemie, Analytische Chemie, ...
- Chemische Prozesstechnik, Biotechnologie, ...
- Lebensmitteltechnik, Umwelttechnik, Klinische Diagnostik, ...

Entsprechend wird die weitreichende Anwendung der Chemometrie deutlich, leider hat sich dies in vielen Bereichen noch nicht ausreichend durchgesetzt.

Schon bei einfachen Fragen, z. B.

- wann besser der Median statt des Mittelwertes einzusetzen ist,
- dass aufgrund der RSM-Limitierung die Simplex-Optimierung verwendet wird oder
- ob eine Zentrierung oder Autoskalierung der Messwerte vor der Auswertung erfolgen soll,

in allen solchen Fällen ist eine fundierte, chemometrische Beurteilung essentiell.

Die besondere Bedeutung der Chemometrie wird heutzutage aber durch die Vieldimensionalität und Komplexität aktueller chemischer Fragestellungen deutlich. Auf Grund der fortgeschrittenen Automatisierung, effizienten Messtechniken, hoher Rechnerleistung und nahezu unbegrenzter Speicher-Kapazität sind die anfallenden Daten oft nicht mehr konventionell zu behandeln. Die zu analysierenden Resultate sind oft nur durch multivariate Auswerte- und Interpretationsmethoden zu bearbeiten, um die inneren Strukturen und komplexen Zusammenhänge zu erkennen.

### 1.3 Realisierung

Aktuelle chemometrische Methoden sind zumeist nicht mehr manuell oder unter Einsatz eines Taschenrechners durchzuführen. Für eine effektive Anwendung der Chemometrie ist daher eine entsprechende Softwarelösung unabdingbar. Auf Grund der Leistungsfähigkeit der aktuellen Prozessoren sind Hardware als auch Software keine limitierenden Aspekte. Kontrovers wird dagegen der Einsatz der unterschiedlichen Software-Tools beschrieben.

Einige Autoren plädieren für einen Einsatz der weitverbreiteten Tabellenkalkulation, die limitierte Funktionalität und Leistungsfähigkeit soll durch den Einsatz von Makros bzw. dem Einbinden von eigens entwickelten Visual-Basic-Routinen kompensiert werden. Dabei sollte jedoch nicht aus den Augen verloren werden, dass eine Benutzung einer Tabelle einer Tabellenkalkulation immer neu die Problematik einer Fehlbedienung bedingt. Auch ist mit Visual Basic nicht eine leistungsfähige, struktu-

rierte oder objektorientierte Sprache für komplexere Anwendungen gegeben. Einfachere Realisierungen oder erste Testauswertung können sicher auf diese Weise realisiert werden, eine validierte Anwendung in der Routine oder chemischen Produktion ist damit vermutlich nicht empfehlenswert.

Andere Chemometrie-anwender favorisieren den Einsatz von Matlab. Die Leistungsfähigkeit dieses Softwarewerkzeugs ist sicherlich in vielen Anwendungsbereichen überzeugend. Es sollte aber nicht übersehen werden, dass Matlab ein Tool ist, das unterschiedliche Auswerterroutinen zur Verfügung stellt. Der qualifizierte Einsatz und die validierte Anwendung bleibt dem versierten Matlab-Spezialisten vorbehalten. Für eine sporadische Anwendung durch den praktischen Chemiker ist dieses Werkzeug möglicherweise nicht optimal.

Ein ähnliches Werkzeug ist R, das ursprünglich als Statistiksprache realisiert wurde. Zusammen mit einer graphischen Oberfläche (z. B. RStudio) ist dies in Verbindung mit vorgegebenen Skripten auch für eine Anwendung durch den Praktiker im Labor gut handhabbar. Die nötigen Skripte sollten von den R-Spezialisten entwickelt und validiert sein, mittlerweile sind bereits viele tausend Pakete verfügbar auch für die unterschiedlichsten, chemometrischen Aufgabenstellungen. Wesentlicher Vorteil ist, dass R/RStudio kostenfrei auf allen Betriebssystemen verfügbar ist ebenso wie die meisten Skripte. Nachteilig ist, dass die Basis von R durch Anwendung in der Statistik dominiert ist, was sich oft in der Interpretation und Darstellung der Resultate auswirkt. Zudem sind R-Skripte keine Programmiersprache, um strukturiert oder objektorientiert komplexe Softwarelösungen zu realisieren.

Einige der relevanten Chemometriethemen werden durch etablierte Statistikpakete gut abgedeckt, auch speziellere Aufgaben wie die Erstellung eines statistischen Versuchsplans oder der Hauptkomponentenregression können hiermit erfolgreich durchgeführt werden. Viele der speziellen, chemometrischen Aufgaben sind hier aber meist nicht implementiert. Zudem sind die leistungsfähigen Statistikpakete nicht nur kostenintensiv, in der Anwendung sind sie oftmals auch nicht einfach. Insbesondere für eine schnelle, sporadische Routineanwendung im Labor sind sie nicht gut geeignet.

Für einige Aufgabenstellungen ist auch spezielle, kommerzielle Software verfügbar, leider trifft dies nicht für alle relevanten Chemometriethemen zu. Die kommerziellen Pakete decken die avisierte Aufgabenstellung zumeist sehr gut ab, zeichnen sich aber teilweise durch eine sehr detaillierte Funktionalität aus, die in manchen Fällen nicht notwendig und für eine schnelle, einfache Anwendung in der Praxis hinderlich ist. Darüber hinaus sollte den implementierten Algorithmen besondere Beachtung geschenkt werden.

## 1.4 Zielsetzung des Buchs

Es stehen mittlerweile einige sehr gute Monographien zur Verfügung, die sowohl die Grundlagen als auch fortgeschrittene, aktuelle Aspekte der Chemometrie behandeln. Diese sind für den erfahrenen Chemometrieanwender uneingeschränkt zu empfehlen.

Leider kommen nahezu alle diese Bücher aus dem angelsächsischen Raum und sind in englischer Sprache verfasst. Dies macht die Thematik für den Einsteiger zu meist nicht einfacher, zumal eine englische Behandlung der nicht immer favorisierten mathematischen oder statistischen Grundlagen nicht für eine Verbesserung sorgt. Es gibt aktuell nur sehr wenige, deutschsprachige Chemometriebücher, so wird das Werk von Otto nicht mehr in deutscher Sprache aufgelegt. Auch eine deutsche Übersetzung z. B. des didaktisch und inhaltlich sehr guten Lehrbuchs von Brereton ist nicht verfügbar.

Dieser Mangel an diversen deutschsprachigen Chemometriefachbüchern limitiert die Verbreitung und Anwendung dieser Disziplin im deutschsprachigen Raum doch deutlich. Gerade für den Einsteiger aus der chemischen Praxis wäre eine deutsche Einführung in die wichtigsten Grundlagen der Chemometrie eine willkommene Basis für eine Anwendung bei der eigenen Problemstellung. Das vorliegende Buch soll die Brücke schlagen zwischen den Grundlagen der Chemometrie und dem noch unerfahrenen Anwender in der chemischen Praxis.

Folgende Aspekte wurden daher explizit berücksichtigt:

- es werden nur ausgewählte, wichtige Themen behandelt,
- mathematische, statistische Aspekte kommen nur vor, soweit dies zwingend nötig ist,
- besonderen Wert wird auf eine verständliche Darstellung der Grundlagen gelegt,
- Beispiele mit einfacher Umsetzung ohne Hilfsmittel verbessern das Verständnis,
- zusätzliche oder alternative Details sind aus Gründen der Übersichtlichkeit ausgelassen,
- für die Anwendung wird entsprechende Software und die Realisierung mit R vorgestellt.

Das Buch konzentriert sich auf die wichtigsten Themengebiete der Chemometrie wie

- Statistische Kennzahlen und Prüfverfahren
- Versuchsplanung, Prozessoptimierung
- Univariate Regression und Kalibration
- Analyse von Messreihen
- Signaldekonvolution
- Mustererkennung, Clusteranalyse
- Multikomponentenanalyse

Es ist zu erwarten, dass die Ausführungen in dem vorgelegten Buch für den versierten Chemometrieexperten oft nicht tiefgründig, evtl. nicht ausführlich genug sein mögen. Einige aktuelle Themen wurden auch explizit nicht behandelt.

Dabei sollte aber berücksichtigt werden, dass die Zielgruppe nicht der Chemometriespezialist sondern der Chemometrieinsteiger ist. Auf diese Weise wird als weitere Zielsetzung vielleicht auch eine größere Verbreitung der Chemometrie in neuen, diversen Anwendungsfeldern der Chemie erreicht.



## 2 Statistische Parameter und Prüfverfahren

### 2.1 Einführung

Die deskriptive Statistik ist mit den beschreibenden Parametern wie Mittelwert und Standardabweichung weithin bekannt. Nicht immer wird jedoch berücksichtigt, dass hierzu Bedingungen erfüllt sein sollten, wie zum Beispiel das Vorliegen einer symmetrischen, eingipfligen Dichtefunktion der Grundgesamtheit (z. B. Gauss-Dichtefunktion) oder die Homogenität der Stichprobe (keine Ausreißer).

Es wurde z. B. die Reaktionsausbeute im Laufe eines Tages wiederholt ermittelt. Die Messwerte zeigen keine gravierenden Besonderheiten, eventuell kann eine geringfügige Steigerung der Ausbeute im Laufe der Zeit konstatiert werden. Liegt hier möglicherweise ein Trend vor, da die Reaktionstemperatur im Laufe des Tages, von den kühlen Morgenstunden bis zur Hitze des Nachmittags, einen Einfluss hatte? Welche Aussage hat in diesem Fall der Mittelwert?

Schon zur Berechnung der grundlegenden, statistischen Parameter sollte die Stichprobe hinsichtlich Homogenität, symmetrischer Dichtefunktion, signifikanter Schiefe oder Trend beurteilt werden. Wurde bei einer dieser Prüfungen die entsprechende Nullhypothese<sup>1</sup> nicht bestätigt, sollte die Angabe von Mittelwert oder Standardabweichung hinterfragt werden. Die alternative, nichtparametrische Statistik umgeht diese Problematik, liefert aber mit ihren robusten Parametern (z. B. Median, Interquartilsdispersionskoeffizient) nicht die gleichen Aussagen.

Eine valide Bewertung einer Stichprobe kann zumeist nicht manuell erfolgen. Auch eine Tabellenkalkulation hilft in den meisten Fällen nicht weiter bzw. dessen Anwendung ist kompliziert und fehleranfällig. Daher sind Softwarepakete wie Mini-Stat, Statistica<sup>®</sup> oder SPSS<sup>®</sup> zu empfehlen, die eine möglichst einfache Anwendung der wichtigsten Routinen erlauben sollten. Entsprechend erfüllen die folgenden Kapitel nicht den Anspruch einer umfassenden Behandlung der deskriptiven Statistik, sondern der wichtigsten Aufgabenstellungen im naturwissenschaftlich-technischen Bereich.

### 2.2 Deskriptive Statistik

Die deskriptive Statistik beschreibt eine Stichprobe unkorrelierter Werte mit ihrem Zufallsfehler durch ihre statistischen Momente, wie z. B. Mittelwert, Standardabweichung, Schiefe und Excess.

---

<sup>1</sup> Nullhypothese: Annahme (positiv formuliert), die durch einen statistischen Test bestätigt werden soll, z. B. Verifizierung, Wert 8 der Nitratresultate (Beispiel 2.1) gehört zur Grundgesamtheit (kein Ausreißer).

Auch für messtechnisch aufwendige Untersuchungen sollte die Stichprobe einen Umfang  $n$  von etwa 10 aufweisen. Dies scheint ein guter Kompromiss zwischen einer zuverlässigen, statistischen Aussage und dem Material-/Zeitaufwand der Messung zu sein.

Die diversen Aspekte werden im Folgenden anhand von Nitratanalysen verdeutlicht.

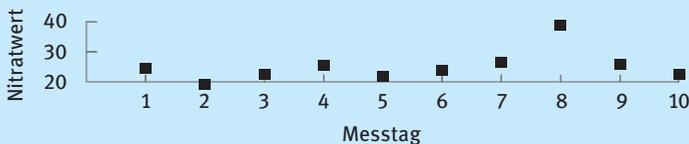


**Bsp. 2.1:** 2-wöchige Nitratmessungen in einer Kläranlage.

Die Resultate im Verlauf folgender Messtage dienen als Stichprobe:

24,32 / 19,97 / 22,45 / 25,50 / 21,68 / 23,77 / 26,51 / 39,00 / 25,86 / 22,51

Unten stehende Abbildung stellt die Messwerte graphisch dar.



### 2.2.1 Erstes statistisches Moment

#### Mittelwert

Für diskrete Stichproben<sup>2</sup> mit kleinem Stichprobenumfang  $n$  ist das arithmetische Mittel ein anerkannter Lageschätzwert mit

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i . \quad (2.1)$$

Für naturwissenschaftlich-technische Untersuchungen hat dieser Erwartungswert eine besondere Bedeutung. Er ist definiert als der „richtige Messwert“, der dem unbekanntem „wahren Wert“ am ehesten entspricht.

Der „richtige Wert“ ist deshalb nicht gleich zu setzen mit dem „wahren Wert“, da jedes Messverfahren prinzipiell mit einem systematischen Fehler behaftet sein kann, der statistisch nicht bestimmbar ist. Diese Deutung des „richtigen Werts“ trifft aber

<sup>2</sup> Mittelwert  $\mu$  einer Grundgesamtheit ( $n = \infty$ ) ist definiert zu  $\mu = \int x f(x) dx$  mit  $f(x)$ : Wahrscheinlichkeitsdichte Gauss-Funktion

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp \frac{-(x-\mu)^2}{2\sigma^2} .$$

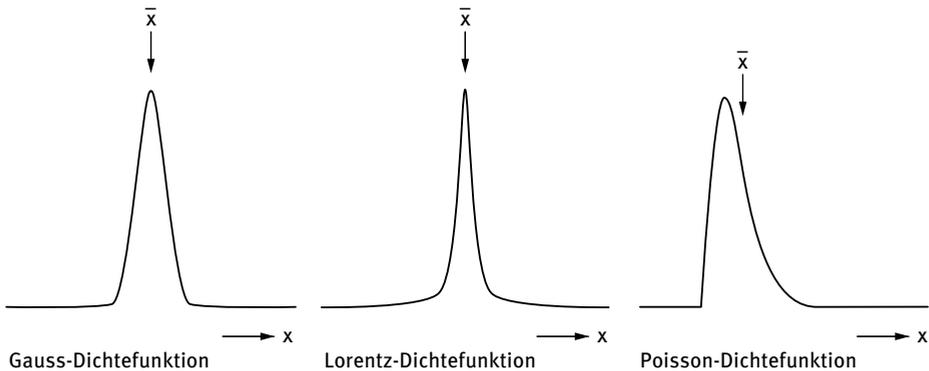


Abb. 2.1: Mittelwert zufallsbehafteter Werte bei Gauss-, Lorentz- bzw. Poisson-Dichtefunktion.

nur zu bei einer eingipfligen, symmetrischen Wahrscheinlichkeitsfunktion<sup>3</sup>, da hier dieser auch der Wert mit der höchsten Wahrscheinlichkeit ist.

Abbildung 2.1 verdeutlicht dies an drei verschiedenen Dichtefunktionen.

Annahmen zur Verwendung des Mittelwertes sind eingipflige, symmetrische Dichtefunktionen (z. B. Gauss-, Lorentz-Funktion), Homogenität (keine Ausreißer), kein Trend. Dies trifft für die Poisson-Dichtefunktion nicht zu, hier ist der arithmetische Mittelwert  $\bar{x}$  nicht der Wert mit der höchsten Wahrscheinlichkeit.

### Median

Der Median (mittleres Quartil)  $\tilde{x}$  erfordert keine Voraussetzungen (nicht parametrisch) und ist invariant gegenüber Ausreißern. Er ist definiert als zentraler Wert einer der Größe nach sortierten Stichprobe mit  $n$  Werten, d. h. 50 % der Werte der Stichprobe liegen oberhalb und 50 % unterhalb des Medians.

$$\tilde{x} = \begin{cases} x_{(n+1)/2} & \text{für } n \text{ ungeradzahlig} \\ (x_{n/2} + x_{n/2+1})/2 & \text{für } n \text{ geradzahlig} \end{cases} \quad (2.2)$$

mit  $x_n$ :  $n$ -ter-Wert aus der Reihe der sortierten Stichprobe.

Das heißt, bei einer der Größe nach sortierten Stichprobe mit 11 Werten, entspricht der Median dem 6. Wert, bei einer Stichprobe mit 10 Werten ist der Median der Mittelwert aus dem 5. und 6. Wert.

Der Median hat nicht die gleiche Aussage wie der Mittelwert, so dass er verwendet werden sollte, wenn der Mittelwert nicht anwendbar ist, z. B. da keine symmetrisch eingipflige Dichtefunktion der Stichprobe konstatiert werden kann.

<sup>3</sup> Grundgesamtheit: Häufigkeits-, Wahrscheinlichkeitsdichtefunktion; diskrete Stichprobe: Wahrscheinlichkeits-, Dichtefunktion.

### 2.2.2 Zweites statistisches Moment

#### Varianz

Das zweite statistische Moment macht eine Aussage über die Unsicherheit mit der ein Wert aufgrund des Zufallsfehlers behaftet ist (Streuungsmaß) und wird meist als Varianz<sup>4</sup>  $\sigma^2$  angegeben. Die Voraussetzungen entsprechen denen des ersten statistischen Moments.

Für Stichproben mit kleinem Stichprobenumfang ist die Standardabweichung  $s$  der etablierte Schätzwert bezogen auf einen Einzelwert der Stichprobe<sup>5</sup>  $s_E$

$$s_E = \sqrt{\frac{\sum(x_i - \bar{x})^2}{n - 1}} \quad (2.3)$$

bzw.  $s_M$  bezogen auf den Mittelwert der Stichprobe (Standardfehler)

$$s_M = \sqrt{\frac{\sum(x_i - \bar{x})^2}{n(n - 1)}} \quad (2.4)$$

#### Relative Standardabweichung $s_{rel}$

Die relative Standardabweichung wird oft als Schätzwert für die Reproduzierbarkeit eines Messverfahrens herangezogen, da sie unabhängig von den absoluten Größen der Messwerte ist und einen allgemeinen Vergleich erlaubt. Es gilt:

$$s_{rel} = \frac{s}{\bar{x}} 100 [\%] \quad (2.5)$$

So sollte die relative Standardabweichung (Reproduzierbarkeit in der Messserie) für ein quantitatives Analysenverfahren kleiner als 3 % sein (siehe Vertrauensbereich).

#### Vertrauensbereich

Der Vertrauensbereich  $\Delta\bar{x}$  (Konfidenzintervall) gibt an, in welchem Bereich um den Mittelwert ein Messwert bei einem gegebenen Signifikanzniveau  $\alpha$  zu erwarten ist. Für den Mittelwert und der entsprechenden Standardabweichung einer Stichprobe ist der

<sup>4</sup> Die Varianz  $\sigma^2$  der Grundgesamtheit (Stichprobenumfang  $n = \infty$ ) ist definiert zu

$$\sigma^2 = \int (x - \mu)^2 f(x) dx .$$

<sup>5</sup> Im Folgenden wird die Notation  $s$  verwendet, wenn sich auf die Standardabweichung  $s_E$  bezogen wird

Vertrauensbereich definiert zu

$$\Delta\bar{x} = \frac{t(\alpha, f)s}{\sqrt{n}} \quad (2.6)$$

mit  $t(\alpha, f)$ : Student- $t$ -Faktor<sup>6</sup> (zweiseitig<sup>7</sup>, vgl. Tabelle A.1),  $\alpha$ : Signifikanzniveau ( $\alpha$ : 0,05 d. h. Wahrscheinlichkeit der Nullhypothese  $P = 0,95$ )<sup>8</sup>,  $f$ : Freiheitsgrade =  $n - 1$ .

In der physikalisch-chemischen Praxis erfolgt bei den meisten Untersuchungen keine Messwertwiederholung, die eine zuverlässige Schätzung der Streuung erlaubt. Zumeist werden z. B. nur zwei Wiederholungsmessungen (Doppelbestimmung) vorgenommen. In vielen Fällen ist jedoch aus einer Verfahrensvalidierung ein relevanter Schätzwert für die Standardabweichung  $s$  gegeben.

In solchen Fällen wird in Gleichung (2.6) dieser Schätzwert für die Standardabweichung eingesetzt, mit dem entsprechenden Freiheitsgrad  $f$  aus der Validierung. Im Nenner von Gleichung (2.6) wird dann die aktuelle Wiederholungszahl  $n_a$  eingetragen (z. B.  $n_a = 2$  für eine Doppelbestimmung).

Der Vertrauensbereich für den Mittelwert der aktuellen Untersuchung ist damit nicht nur abhängig von der Streuung des Messverfahrens, sondern auch von der Zahl der aktuellen Wiederholungsmessungen. Je mehr Messungen für die aktuelle Untersuchung gemacht wurden, umso sicherer wird die Aussage, desto kleiner der Vertrauensbereich. Abbildung 2.2 verdeutlicht den Einfluss von  $f$  und  $n_a$  auf den Vertrauensbereich.

Abbildung 2.2 zeigt, dass der Vertrauensbereich für eine Doppelbestimmung bei  $\pm 5\%$  liegt ( $f = 7$ ,  $s_{\text{rel}} = 3\%$ ). Aus einer Streuung (Reproduzierbarkeit) von  $6\%$ , würde ein Vertrauensbereich von etwa  $\pm 10\%$  resultieren, was hinsichtlich der Angabe der Nachkommastellen bei  $\bar{x}$  bedacht werden sollte.

In der Regel sollte ein Messresultat als Mittelwert mit Vertrauensbereich angegeben werden  $\bar{x} \pm \Delta\bar{x}$ , damit die Aussagekraft des Ergebnisses beurteilt werden kann.

<sup>6</sup> Bei Stichproben mit kleinerem Stichprobenumfang  $n$  erfolgen Prüfungen nicht auf Grundlage der Gauss-Verteilung sondern abgeleiteter Verteilungen, hier der Student- $t$ -Verteilung.

<sup>7</sup> Einseitige Fragestellung liegt z. B. vor wenn gefragt ist, ob z. B. ein Wert größer als ein Grenzwert ist, eine zweiseitige Fragestellung, wenn der Wert größer als eine Unter- und kleiner als eine Obergrenze sein soll.

<sup>8</sup> Formal ist die statistische Notation für den Student- $t$ -Wert bei einseitiger Fragestellung  $t(\alpha, f)$  bzw. bei zweiseitiger Fragestellung  $t(\alpha/2, f)$ . Aus Praktikabilitätsgründen wird im Folgenden generell das gewählte Signifikanzniveau verwendet, z. B.  $\alpha = 0,05$  oder  $0,01$ . Entsprechend sind in Tabelle A.1 zwei Tabellen für die Student- $t$ -Verteilung aufgelistet, für die einseitige oder zweiseitige Fragestellung.