

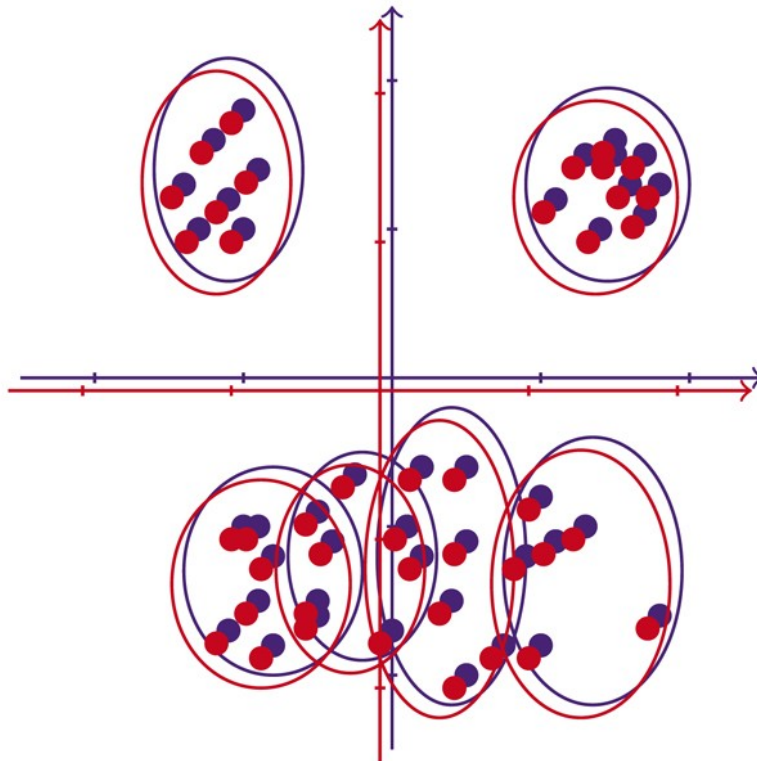


Johann Bacher | Andreas Pöge | Knut Wenzig

Clusteranalyse

Anwendungsorientierte Einführung
in Klassifikationsverfahren

3. Auflage





Clusteranalyse

Anwendungsorientierte Einführung
in Klassifikationsverfahren

von

Universitätsprofessor

Dr. Johann Bacher

Johannes-Kepler-Universität, Linz

Akademischer Rat

Dr. Andreas Pöge

Universität Bielefeld

Diplom-Sozialwirt

Knut Wenzig

Nationales Bildungspanel, Bamberg

3., ergänzte, vollständig überarbeitete und
neu gestaltete Auflage

Oldenbourg Verlag München

www.clusteranalyse.net
3.aufgabe@clusteranalyse.net

Bibliografische Information der Deutschen Nationalbibliothek

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

© 2010 Oldenbourg Wissenschaftsverlag GmbH
Rosenheimer Straße 145, D-81671 München
Telefon: (089) 45051-0
oldenbourg.de

Das Werk einschließlich aller Abbildungen ist urheberrechtlich geschützt. Jede Verwertung außerhalb der Grenzen des Urheberrechtsgesetzes ist ohne Zustimmung des Verlages unzulässig und strafbar. Das gilt insbesondere für Vervielfältigungen, Übersetzungen, Mikroverfilmungen und die Einspeicherung und Bearbeitung in elektronischen Systemen.

Lektorat: Rainer Berger
Herstellung: Anna Grosser
Coverentwurf: Kochan & Partner, München
Umschlagillustration: Cornelia Horn
Gedruckt auf säure- und chlorfreiem Papier
Gesamtherstellung: Grafik + Druck GmbH, München

ISBN 978-3-486-58457-8

Vorwort zur dritten Auflage

Seit der Veröffentlichung der zweiten Auflage sind fast 15 Jahre vergangen, und im Bereich der Clusteranalyse hat es viele Weiterentwicklungen gegeben – die meisten sicherlich im Bereich der probabilistischen Verfahren. Im neuen Autorenteam wurden daher – neben der Umstellung auf L^AT_EX – die beiden vorausgehenden Auflagen grundlegend überarbeitet, um den aktuellen Entwicklungen Rechnung zu tragen. Beispielhaft genannt seien: Eine Erweiterung der Einleitung um eine Sammlung von Beispielen aus der Forschungspraxis, eine taxative Nennung und Beschreibung von Kriterien für eine gute Klassifikation sowie in Teil I eine klare begriffliche Abgrenzung von Hauptkomponentenmethode und Faktorenanalyse. Ausführlich behandelt werden Datenkonstellationen, für welche die Faktorenanalyse brauchbar oder unbrauchbar ist. Daneben wurde in Teil II das K-Means-Verfahren überarbeitet: Aufgenommen wurde die Methode der multiplen zufälligen Startwerte, die eine entscheidende Weiterentwicklung darstellt und das Problem lokaler Minima weitgehend vermeidet. Dargestellt werden auch Verallgemeinerungen des K-Means-Verfahrens, die andere Distanzfunktionen und Lageparameter nutzen.

Teil III wurde um den Latent-GOLD-Ansatz, der die Modellierung von komplexen Clustermodellen ermöglicht, erweitert. Ergänzt wurde dieser Teil um eine Beschreibung von SPSS-TwoStep und AutoClass. Mit AutoClass wird in die derzeit »boomende« Bayes-Statistik eingeführt. Darüber hinaus wird die Performanz von ausgewählter Software verglichen. Im praktisch orientierten Teil IV werden häufig gestellte Anwenderfragen beantwortet und die Klassifikation von Verläufen mittels Optimal Matching, die Bildung von Konsensclustern sowie die formale Gültigkeitsprüfung dargestellt.

In allen Teilen wurden Anwendungsempfehlungen aufgenommen, wobei versucht wurde, einerseits den Praxisanforderungen Rechnung zu tragen und andererseits einfache »Kochbuchrezepte« zu vermeiden. Wie in den bisherigen Auflagen ist das Ziel, eine Darstellung dergestalt zu finden, dass die behandelten Verfahren in ihren Grundlagen von der interessierten Leserin oder dem interessierten Leser nachvollzogen werden können. Für eine wissenschaftliche Analyse erscheint es uns unumgänglich, dass zumindest eine grobe Vorstellung über die Funktionsweise eines Verfahrens vorliegt, bevor ein »Rechnerknopf« am Computer gedrückt wird.

Über Rückmeldungen und Anregungen zu einzelnen Verfahren und Ausführungen freuen wir uns. Unser Dank gilt Arne Bethmann und Heinz Leitgöb für die Mitarbeit an einzelnen Teilen des Buches. Danken möchten wir darüber hinaus unseren Kolleginnen und Kollegen für Hinweise und Korrekturen sowie dem Oldenbourg-Verlag für seine Geduld bei der Abgabe der Druckvorlage. Die Fertigstellung des Manuskripts hat erhebliche Zeit zu Lasten von familiären Verpflichtungen gebunden. Für das Verständnis hierfür herzlichen Dank an alle.

Linz, Bielefeld, Bamberg 2010

Johann Bacher

Andreas Pöge

Knut Wenzig

Vorwort zur ersten und zweiten Auflage

Clusteranalyseverfahren gewinnen in der Forschungspraxis zunehmend an Bedeutung. Sie werden heute in zahlreichen Wissenschaftsdisziplinen zur Lösung von Klassifikationsaufgaben eingesetzt, in den Sozial- und Wirtschaftswissenschaften beispielsweise zur Identifizierung von unterschiedlichen Lebens- und Konsumstilen oder von Wertorientierungstypen.

Die Arbeit gibt – auf Grundlage der neueren Methodenliteratur – eine systematische Einführung, wie bei der Bestimmung von Typen (Clustern) vorzugehen ist, welche Verfahren dabei zur Verfügung stehen und wie die Ergebnisse zu interpretieren sind. Folgende Verfahren werden behandelt: bivariate und multiple Korrespondenzanalyse, nichtmetrische mehrdimensionale Skalierung, Faktorenanalyseverfahren (nominale Faktorenanalyse nach McDonald, R- und Q-Analyse), hierarchisch-agglomerative Verfahren, Repräsentanten-Verfahren, explorative und konfirmatorische K-Means-Verfahren, explorative und konfirmatorische probabilistische Clusteranalyseverfahren (latente Profilanalyse, Analyse latenter Klassen für nominale, ordinale und gemischte Variablen).

Alle Verfahren werden anhand von Beispielen aus der Forschungspraxis behandelt und durchgerechnet. Dabei werden eine Reihe von methodologischen Fragen beantwortet, wie beispielsweise die Frage, welche Konsequenzen eine empirische Standardisierung hat oder wie sich irrelevante Variablen auf die Ergebnisse einer Clusteranalyse auswirken.

Wie jede Arbeit wäre auch das hier vorliegende Buch nicht ohne die Unterstützung meiner akademischen Lehrer, des Oldenbourg-Verlages, der Teilnehmer an unterschiedlichen Kursen zur Datenanalyse sowie meiner Bekannten und Freunden und meiner Familie zustande gekommen. Ihnen allen möchte ich an dieser Stelle herzlichst danken.

Johann Bacher

Inhalt

1	Einleitung	15
1.1	Zielsetzung clusteranalytischer Verfahren	15
1.2	Homogenität als Grundprinzip der Bildung von Clustern	16
1.3	Clusteranalyseverfahren	18
1.4	Grundlage der Clusterbildung	20
1.5	Konfirmatorische und explorative Clusteranalyse	22
1.6	Anwendungsbeispiele	23
1.7	Modellprüfung und Validierung	27
1.8	Fehleranalyse	28
1.9	Datenanalyse als iterativer Prozess	30
1.10	Computerprogramme	32
I	Unvollständige Clusteranalyseverfahren	35
2	Einleitende Übersicht	37
3	Multiple Korrespondenzanalyse	43
3.1	Ein Anwendungsbeispiel	43
3.1.1	Faktorenanalytische Interpretation	46
3.1.2	Clusteranalytische Interpretation	52
3.2	Das Modell der multiplen Korrespondenzanalyse	57
3.2.1	Berechnung der empirischen Zusammenhangsmatrix \mathbf{G}	58
3.2.2	Berechnung der Eigenwerte, Faktorladungen und Koordinatenwerte	61
3.2.3	Berechnung der Skalenwerte und Interpretation der Koordinaten	63
3.2.4	Unerwünschter Effekt der Reskalierung der Faktorladungen und Rotation der Faktoren	65
3.3	Modellprüfgrößen	67
3.3.1	Signifikanz der Zusammenhangsstruktur	67
3.3.2	Die Zahl maximal möglicher und bedeutsamer Dimensionen	68
3.3.3	Überprüfung der faktorenanalytischen Interpretation	69
3.3.4	Modellprüfgrößen für die clusteranalytische Interpretation	72

3.4	Anwendungsempfehlungen	74
4	Nichtmetrische mehrdimensionale Skalierung	77
4.1	Aufgabenstellung und Ähnlichkeitsmessung	77
4.2	Schätzalgorithmus	80
4.3	Maximale und angemessene Dimensionszahl	87
4.4	Unbekannter Metrikparameter p	90
4.5	Weitere Modellanpassungsgrößen	92
4.6	Freizeitverhalten von Kindern	94
4.6.1	Clusteranalytische Interpretation	96
4.6.2	Faktorenanalytische Interpretation	97
4.6.3	Freizeitaktivitäten und Sozialstruktur	100
4.7	Anwendungsempfehlungen	107
5	Weitere räumliche Darstellungsverfahren	109
5.1	Die bivariate Korrespondenzanalyse	109
5.2	Nominale Faktorenanalyse nach McDonald	117
5.3	Die Hauptkomponenten- und Faktorenanalyse	122
5.3.1	Hauptkomponenten- und R-Faktorenanalyse	122
5.3.2	Die Q-Faktorenanalyse	136
5.4	Anwendungsempfehlungen	143
II	Deterministische Clusteranalyseverfahren	145
6	Einleitende Übersicht	147
6.1	Überlappende und überlappungsfreie Clusterlösungen	147
6.2	Grundvorstellungen über die zu bildenden Cluster	148
6.3	Complete- und Single-Linkage als Basismodelle	150
6.4	Auswahl eines geeigneten Verfahrens	153
6.5	Lösungsschritte einer Klassifikationsaufgabe	156
6.6	Ein Anwendungsbeispiel	156
6.7	Fehlerquellen	165
7	Gewichtung und Transformation von Variablen	175
7.1	Vergleichbarkeit von Klassifikationsmerkmalen	175
7.2	Lösungsstrategien	176
7.3	Theoretische und empirische Standardisierung	177
7.4	Hierarchische Variablen	183
7.5	Gemischte Variablen	184

7.6	Standardisierung von Objekten	188
7.7	Exkurs: Die Problematik einer automatischen Orthogonalisierung	193
8	Unähnlichkeits- und Ähnlichkeitsmaße	195
8.1	Auswahl eines (Un-)Ähnlichkeitsmaßes	196
8.2	Dichotome Variablen	197
8.3	Nominale Variablen	207
8.4	Ordinale Variablen	211
8.5	Quantitative Variablen	219
8.6	A-priori-Prüfung auf Vorhandensein einer Clusterstruktur	224
8.7	Gewichtung von Variablen und Distanzen, Standardisierung von Objekten	226
8.8	Fehlende Werte	228
8.9	Exkurs: Quantifizierung und Konsequenzen der Kategorisierung	230
9	Nächste-Nachbarn- und Mittelwertverfahren	233
9.1	Der Complete-Linkage als Basismodell	233
9.1.1	Der hierarchisch-agglomerative Algorithmus	233
9.1.2	Hierarchische Darstellung von Ähnlichkeitsbeziehungen	237
9.1.3	Maßzahlen zur Bestimmung der Clusterzahl	241
9.1.4	Zufallstestung des Verschmelzungsschemas	245
9.1.5	Maßzahlen zur Beurteilung einer bestimmten Clusterlösung	247
9.2	Der Single-Linkage	251
9.3	Complete-Linkage für überlappende Cluster	255
9.4	Verallgemeinerte Nächste-Nachbarn-Verfahren	259
9.5	Mittelwertverfahren	264
9.6	Anwendungsempfehlungen	274
10	Repräsentanten-Verfahren	277
10.1	Modellansatz	277
10.2	Anwendungsbeispiel	279
10.3	Die Wahl der Schwellenwerte	280
10.4	Weitere Repräsentanten-Verfahren	282
10.5	Anwendungsempfehlungen	283
11	Hierarchische Verfahren zur Konstruktion von Clusterzentren	285
11.1	Modellansätze, Algorithmen und Ward-Verfahren	285
11.2	Bestimmung der Clusterzahl und Modellprüfgrößen	290
11.3	Analyse durchschnittlicher Befragter	290
11.4	Anwendungsempfehlungen	295

12	K-Means-Verfahren	299
12.1	Modellansatz und Algorithmus	299
12.2	Bestimmung der Clusterzahl	305
12.3	Zufallstestung einer bestimmten Clusterlösung	313
12.4	Beschreibung und Interpretation der Cluster	314
12.5	Formale Beschreibung der Cluster	321
12.6	Analyse von Ausreißern	325
12.7	Stabilitätsprüfung	328
12.8	Inhaltliche Validitätsprüfung	332
12.9	Alternative Startwertverfahren	335
12.10	Gemischtes Messniveau	336
12.11	Modifikation des Algorithmus	338
12.12	Verwendung der Mahalanobis-Distanz	339
12.13	Konfirmatorisches K-Means-Verfahren	341
12.14	K-Median- und K-Modus-Verfahren	345
12.15	Anwendungsempfehlungen	347
III	Probabilistische Clusteranalyseverfahren	349
13	Einleitende Übersicht	351
14	Latente Profilanalyse	355
14.1	Modellansatz und Algorithmus	355
14.2	Modellprüfgrößen	362
14.2.1	Bestimmung der Klassenzahl	362
14.2.2	Zufallstestung einer Klassenlösung	366
14.3	Beschreibung und Interpretation einer Klassenlösung	367
14.4	Überlappingsindizes	368
14.5	Überprüfung der Annahme der lokalen Unabhängigkeit	372
14.6	Konfirmatorische latente Profilanalyse	373
15	Analyse latenter Klassen für nominale, ordinale und gemischtskalierte Variablen	377
15.1	Modellansatz und Algorithmus für nominale Variablen	377
15.2	Modellprüfung und Interpretation	382
15.3	Konfirmatorische Analyse	387
15.4	Modellansatz und Algorithmus für ordinale und gemischtskalierte Variablen	390

16	Latent-GOLD-Ansatz	395
16.1	Allgemeiner Ansatz und Überblick	395
16.2	Modellansatz der Latent-Class-Clusteranalyse	398
16.2.1	Der klassische Ansatz	399
16.2.2	Erweiterung mit Kovariaten	404
16.2.3	Ordinale Indikatorvariablen	407
16.2.4	Kontinuierliche Indikatorvariablen	409
16.2.5	Zählvariablen	409
16.3	Parameterschätzung	411
16.4	Statistiken zur Modellanpassung	416
16.4.1	χ^2 -Statistiken	417
16.4.2	Log-Likelihood-Statistiken	420
16.4.3	Klassifikations-Statistiken	420
16.4.4	Signifikanztests mit parametrischem Bootstrap	422
16.4.5	Bivariate Residuen	424
16.4.6	Beurteilung und Auswahl von Modellen	425
16.5	Ein Anwendungsbeispiel	426
16.5.1	Kontinuierliche Daten (latente Profilanalyse)	426
16.5.2	Hinzunahme von Kovariaten	431
17	Weiterentwicklungen und Modifikationen	439
17.1	AutoClass (<i>gemeinsam mit Arne Bethmann</i>)	439
17.1.1	Modell	441
17.1.2	Schätzverfahren und Bestimmung der Clusterzahl	443
17.1.3	Vergleichsrechnung mit den Denz-Daten	444
17.2	TwoStep-Cluster	446
17.2.1	Allgemeiner Ansatz	446
17.2.2	Ergebnis für die metrischen Denz-Daten und für gemischte Skalenniveaus	450
17.3	Vergleich ausgewählter Software (<i>gemeinsam mit Arne Bethmann</i>)	451
17.3.1	Simulationsmodelle	451
17.3.2	Ergebnisse	452
IV	Spezielle Anwendungsfragen	457
18	Häufig gestellte Anwendungsfragen	459
18.1	Welches Verfahren?	459
18.1.1	Bildung abgeleiteter Variablen	459
18.1.2	Räumliche Darstellung von Objekten oder Variablen	460
18.1.3	Auffinden einer hierarchischen Ähnlichkeitsstruktur	461

18.1.4	Räumliche oder hierarchische Darstellung?	462
18.1.5	Klassifikation von Variablen	463
18.1.6	Klassifikation von Objekten	464
18.2	Verwendung aller Variablen?	466
18.3	Welches Un- bzw. Ähnlichkeitsmaß?	469
18.4	Wie viele Cluster?	470
18.5	Modale Klassenzugehörigkeit oder Zuordnungswahrscheinlichkeiten? . . .	472
19	Klassifikation von Verläufen mittels Optimal Matching	
	<i>Heinz Leitgöb</i>	475
19.1	Einführung	475
19.2	Methodische Grundlagen	476
19.3	Die Hamming-Distanz	479
19.4	Die Levenshtein-Distanz	481
19.5	Ein theoretisches Beispiel	482
19.6	Die Festsetzung der Transformationskosten	485
19.7	Der Analyseablauf	488
19.8	Fazit und Anwendungsempfehlungen	491
20	Formale Gültigkeitsprüfung und Konsenslösungen	493
20.1	Formale Gültigkeitsprüfung	493
20.2	Konsenslösungen	497
20.2.1	Konsensus für Clustermittelwerte	498
20.2.2	Konsensus auf der Basis der Clusterzuordnungen	499
	Literatur	503
	Register	523

1 Einleitung

1.1 Zielsetzung clusteranalytischer Verfahren

Primäres Ziel clusteranalytischer Auswertungsverfahren ist, eine Menge von Klassifikationsobjekten in homogene Gruppen (Klassen, Cluster, Typen) zusammenzufassen – oder kurz ausgedrückt – das *Auffinden einer empirischen Klassifikation* (Gruppeneinteilung, Typologie). Von einer empirischen Klassifikation soll dann gesprochen werden, wenn die Klassifikation auf empirischen Beobachtungen, zum Beispiel auf der Grundlage einer Befragung, basiert. Klassifikationsobjekte (siehe Übersichtstabelle 1.1) können sein: Individuen (Personen, Befragte), Aggregate (Organisationen, Nationen, Berufsgruppen usw.) oder Variablen (Merkmale).

Tab. 1.1: Beispiele für clusteranalytische Auswertungen

Klassifikationsobjekte	Beispiele für eine clusteranalytische Auswertung
Personen	A) Befragte werden aufgrund ihrer sozialstrukturellen Merkmale in (homogene) soziale Schichten zusammengefasst. B) Befragte werden aufgrund ihrer Lebensstile (Freizeitpräferenzen, Musikgeschmack, Wertorientierungen) in homogene Lebensstilgruppen zusammengefasst.
Aggregate	C) Nationen werden aufgrund ihrer demographischen, wirtschaftlichen und/oder sozialen Entwicklung in homogene Gruppen zusammengefasst. D) Berufe werden aufgrund ihrer Tätigkeitsprofile in homogene Gruppen zusammengefasst.
Variablen	E) Freizeitaktivitäten (Variablen) werden aufgrund ihres gemeinsamen Auftretens bei Personen in homogene Gruppen von Variablen zusammengefasst. F) Indikatoren der demographischen, wirtschaftlichen und sozialen Entwicklung werden aufgrund ihrer Korrelationen bei Nationen in homogene (Variablen-)Gruppen zusammengefasst.

Formal unterscheiden sich die sechs Beispiele der Übersicht darin, dass in den Beispielen A bis D die Zeilen einer Datenmatrix die Klassifikationsobjekte bilden. In den Beispielen E bis F sind dies dagegen die Spalten einer Datenmatrix. Im ersten Fall soll von einer *objektorientierten* Datenanalyse gesprochen werden, im letzten Fall von einer *variablenorientierten* Datenanalyse. Die Bezeichnung Klassifikationsobjekte kann sich auf die Spalten oder Zeilen einer Datenmatrix beziehen und bezeichnet jene Einheiten (Personen, Aggregate oder Variablen), die geclustert werden. Die untersuchten Objekte (Zeilen) müssen dabei keinesfalls mit den Erhebungseinheiten identisch sein. So zum Beispiel können durch Aggregation über eine oder mehrere Variablen (zum Beispiel Beruf, regionale oder nationale Zugehörigkeit usw.) »neue« Objekte (Aggregate) erzeugt werden, die dann geclustert werden.

1.2 Homogenität als Grundprinzip der Bildung von Clustern

Jeder Clusterbildung liegt – unabhängig von Unterschieden im Detail – die Grundvorstellung der *Homogenität* bzw. von »homogenen« Gruppen zugrunde (Kozelka 1982, S. 6; Sodeur 1974, S. 118–124; u. a.). Mit dem Begriff *homogene Gruppe* sind folgende Vorstellungen verbunden:

1. Die Klassifikationsobjekte, die einer homogenen Gruppe angehören, sollen untereinander ähnlich sein. Es soll *Homogenität innerhalb der Cluster* vorliegen.
2. Die Klassifikationsobjekte, die unterschiedlichen homogenen Gruppen angehören, sollen verschieden sein. Es soll *Heterogenität zwischen den Clustern* vorliegen.

Diese beiden Grundvorstellungen werden in der Literatur unterschiedlich bezeichnet. Cormack (1971) spricht von *externer Isolierung* (Heterogenität zwischen den Clustern) und *interner Kohäsion* (Homogenität innerhalb der Cluster). Everitt (1980, S. 60) und andere sprechen von »natürlichen« Clustern, wenn beide Forderungen erfüllt sind, und beschreiben Cluster als dichte Punktwolken in einem p -dimensionalen Raum, die durch Regionen mit einer geringen Dichte voneinander getrennt sind.¹ Sind die beiden Grundvorstellungen nicht erfüllt, ist es wenig sinnvoll, eine Klassifikation durchzuführen. Abbildung 1.1a verdeutlicht diesen Sachverhalt. Die untersuchten Klassifikationsobjekte bilden in den beiden Variablen X und Y eine große, relativ geschlossene Punktwolke. Eine Aussage der Art, den Daten liegen K Cluster (zum Beispiel $K = 3$) zugrunde, ist nicht sinnvoll. In der Abbildung 1.1b dagegen sind drei Cluster zu erkennen. Sie sind

¹ Zur Vorstellung von »natürlichen« Clustern siehe auch Aldenderfer und Blashfield (1984, S. 33–34) oder Vogel (1975, S. 16–17).

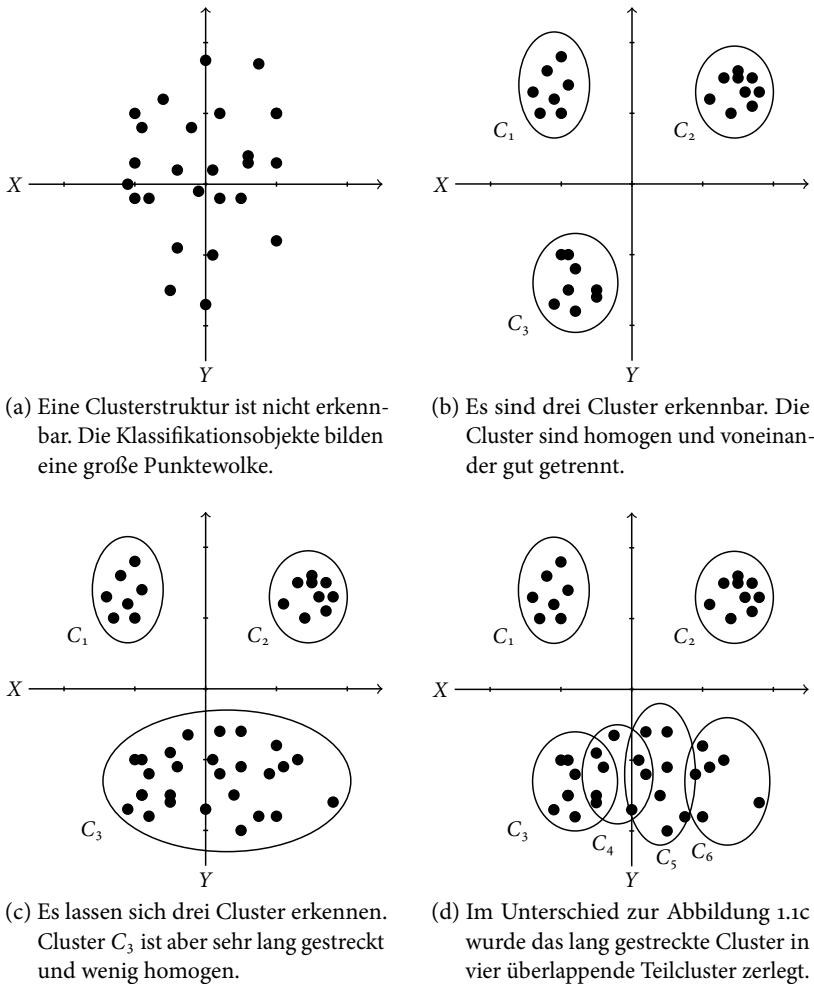


Abb. 1.1: Datenkonstellationen mit erkennbarer und nicht erkennbarer Clusterstruktur

in sich homogen und voneinander verschieden. In der Abbildung 1.1c ist zwar die Vorstellung der Heterogenität zwischen den Clustern erfüllt, Cluster 3 erfüllt aber in einem geringeren Ausmaß die Homogenitätsbedingung innerhalb der Cluster, da das Cluster sehr lang gestreckt ist. Abhängig von dem Gewicht beider Grundvorstellungen bei der Clusterbildung wird man sich entweder für drei Cluster entscheiden oder aus dem lang gestreckten Cluster mehrere überlappende Cluster bilden (siehe Abbildung 1.1d).

Neben diesen beiden Grundvorstellungen werden aus forschungspraktischen, aber auch aus inhaltlichen Gründen weitere Anforderungen an die gesuchte Klassifikation gestellt

(Schlosser 1976, S. 186; Sodeur 1974, S. 125–129), so zum Beispiel, dass die Zahl der Cluster möglichst klein sein sollte (Schlosser 1976, S. 186). Zusammenfassend lassen sich folgende Anforderungen bzw. Kriterien auflisten, die häufig von einer guten Clusterlösung verlangt werden:

1. Die Cluster sollen in sich homogen sein. Objekte, die einem Cluster angehören, sollen zueinander ähnlich sein.
2. Die Cluster sollen voneinander isoliert sein. Objekte, die verschiedenen Clustern angehören, sollen sich voneinander unterscheiden.
3. Die Cluster sollen den Daten gut angepasst sein. Die Klassifikation soll in der Lage sein, die Variation in den Daten zu erklären.
4. Die Cluster sollen stabil sein. Geringfügige Änderungen in den Daten oder im Verfahren sollen in keinen gravierenden Änderungen der Ergebnisse resultieren.
5. Die Cluster sollen inhaltlich gut interpretierbar sein. Den Clustern sollen inhaltlich sinnvolle Namen gegeben werden können. Im Idealfall sollen die Namen aus einer Theorie abgeleitet werden.
6. Die Cluster sollen (inhaltlich) valide sein. Die Cluster sollen mit externen Variablen korrelieren, von denen bekannt ist, dass sie im Zusammenhang mit den Typen stehen, die aber nicht in die Bildung der Cluster eingehen.

Gefordert wird mitunter ferner:

7. Die Zahl der Cluster soll klein und damit überschaubar sein. Angenommen wird, dass dies die inhaltliche Interpretierbarkeit (Kriterium 5) erleichtert und die Stabilität erhöht (Kriterium 4).
8. Die Cluster selbst sollen eine gewisse Mindestgröße haben. Dies soll zur Stabilität (Kriterium 4) beitragen.

1.3 Clusteranalyseverfahren

Zum Auffinden von Clustern wurde eine Vielzahl von Verfahren entwickelt, für die unterschiedliche Einteilungen und Zusammenfassungen vorgeschlagen wurden. In der vorliegenden Arbeit werden drei große Verfahrensgruppen unterschieden: *unvollständige Clusteranalyseverfahren*, *deterministische Clusteranalyseverfahren* und *probabilistische Clusteranalyseverfahren*. Grundlage der Differenzierung ist die Zuordnung der Klassifikationsobjekte zu den Clustern:

- *Unvollständige Clusteranalyseverfahren* (siehe Teil I): Diese Verfahren werden in der Literatur auch als geometrische Methoden (Gordon 1981, S. 80–120), als Repräsentations-

oder Projektionsverfahren (Jain und Dubes 1988; Opitz 1980) bezeichnet. In dieser Arbeit wurde die Bezeichnung »unvollständige Clusteranalyseverfahren« gewählt, da die Bildung von Clustern und die Zuordnung der Klassifikationsobjekte zu den Clustern von der Anwenderin bei der Interpretation der räumlichen Darstellung vorgenommen werden muss. Die unvollständigen Clusteranalyseverfahren selbst führen nur zu einer räumlichen Darstellung. Die Verfahren können auch dazu verwendet werden, abgeleitete Variablen (zum Beispiel Faktorwerte) zu bilden oder eine metrische Ähnlichkeits- oder Unähnlichkeitsmatrix (zum Beispiel mittels mehrdimensionaler Skalierung) zu schätzen, die anschließend mittels eines hierarchischen Clusteranalyseverfahrens untersucht wird.

- *Deterministische Clusteranalyseverfahren* (siehe Teil II): Die Klassifikationsobjekte werden mit einer Wahrscheinlichkeit von 1 einem oder mehreren Clustern zugeordnet. Es lassen sich zwei Verfahrensgruppen unterscheiden: *hierarchische Verfahren* und *partitionierende Verfahren*. Bei den hierarchischen Verfahren erfolgt die Clusterbildung schrittweise: Bei den sogenannten *hierarchisch-agglomerativen Verfahren* werden aus n Objekten zunächst n Cluster gebildet, aus den n Clustern durch Zusammenfassung der beiden ähnlichsten Cluster dann $n - 1$ Cluster, aus diesen wiederum $n - 2$ Cluster usw. Bei zehn Objekten werden also zunächst zehn Cluster gebildet, aus diesen zehn dann neun, aus den neun dann acht usw. Bei den *divisiven Verfahren* wird umgekehrt vorgegangen: Die n Objekte bilden ein einziges großes Cluster, dieses wird dann in zwei Cluster aufgespaltet, die zwei Cluster dann in drei usw. Bei den *partitionierenden Verfahren* muss dagegen eine Clusterzahl vorgegeben werden. Die Objekte werden dann den Clustern so zugeordnet, dass ein bestimmtes Kriterium maximiert bzw. minimiert wird. Das bekannteste partitionierende Verfahren ist das K-Means-Verfahren. Es ordnet die n Objekte K Clustern so zu, dass die Varianz in den Clustern minimiert wird.
- *Probabilistische Clusteranalyseverfahren* (siehe Teil III): Die Klassifikationsobjekte werden den Clustern nicht deterministisch mit einer Wahrscheinlichkeit von 1 oder 0, sondern mit einer dazwischen liegenden Wahrscheinlichkeit zugeordnet. Beispielsweise gehört ein Objekt g mit einer Wahrscheinlichkeit von 0,6 dem Cluster 1, mit einer Wahrscheinlichkeit von 0,2 dem Cluster 2 und mit einer Wahrscheinlichkeit von 0,1 dem Cluster 3 usw. an.

Die Beziehung zwischen unvollständigen, deterministischen und probabilistischen Clusteranalyseverfahren können wir uns wie folgt vorstellen: Eine graphische Bildung von Clustern, wie sie bei den unvollständigen Clusteranalyseverfahren durch den Anwender vorgenommen wird, ist nur in einem ein-, zwei- oder maximal dreidimensionalen Raum bei einer kleinen Klassifikationsobjektmenge möglich. Liegt ein höherdimensionaler Raum oder eine größere Klassifikationsobjektmenge vor, sind formale Verfahren zur Clusterbildung erforderlich. Hinsichtlich der Zuordnung der Klassifikationsobjekte