



Haben
Sie
schon mal
geclustert ?
clusteranalyse
@gmx.de



Fragestellung

Definition

Voraussetzung

Methodik

Interpretation

Zentrale Fragestellungen:

Was ist eine Clusteranalyse?

Wie wird eine Clusteranalyse angewendet?

Wann wird eine Clusteranalyse angewendet?



Fragestellung

Definition

Voraussetzung

Methodik

Interpretation

Clusteranalyse = Gruppenbildungsverfahren = eine Vielzahl von Objekten werden zu Gruppen zusammengefasst

zur Etymologie:

engl. Cluster = Haufen, Menge, Ballung

altdt. Kluster = „was dicht und dick zusammensitzt“

(Grimm'sches Wörterbuch)

Fragestellung

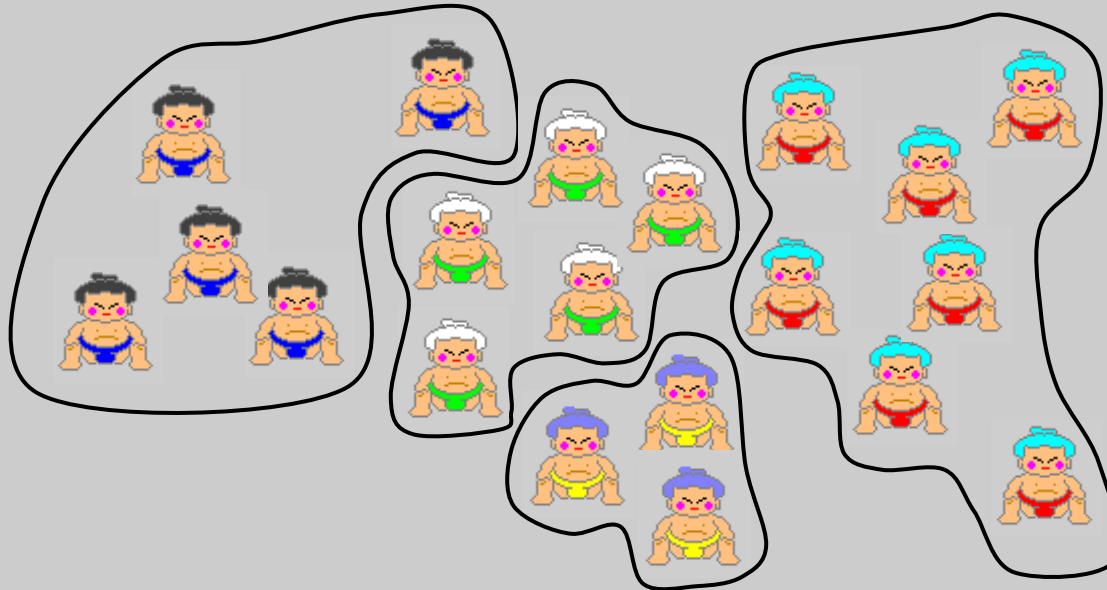
Definition

Voraussetzung

Methodik

Interpretation

„was dicht und dick zusammensitzt“ - Sumo-Ringer



Objekte: Beschreibung durch verschiedene Merkmale unterschiedl. Auspragung

Cluster: Bildung durch Objekte mit ahnlichen Auspragungen

Objekte *innerhalb* einer Gruppe sollen *homogen* sein
Objekte *zwischen* den Gruppen sollen *heterogen* sein

Fragestellung

Definition

Voraussetzung

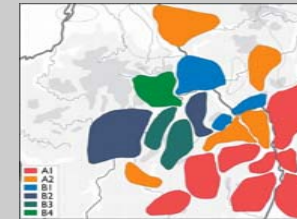
Methodik

Interpretation

- Marketing: Zusammenhang zw. Selbstbild und Wahl einer Automarke



- Archaologie: Kultureller Fingerabdrucke in der Kategorie „Schmuck“ in hallstattzeitlichen Siedlungen im Mittelrheingebiet



- Botanik: Ein pflanzensoziologisches Modell der Schattentoleranz von Baumarten in den Bayerischen Alpen



- Stadttestatistiker: Sozialraumanalyse, Analyse zu Luftverschmutzung u. Larmbelastung, Burgerumfragen, Wirtschaftsraumen, Wahlanalysen etc.

Fragestellung

Definition

Voraussetzung

Methodik

Interpretation

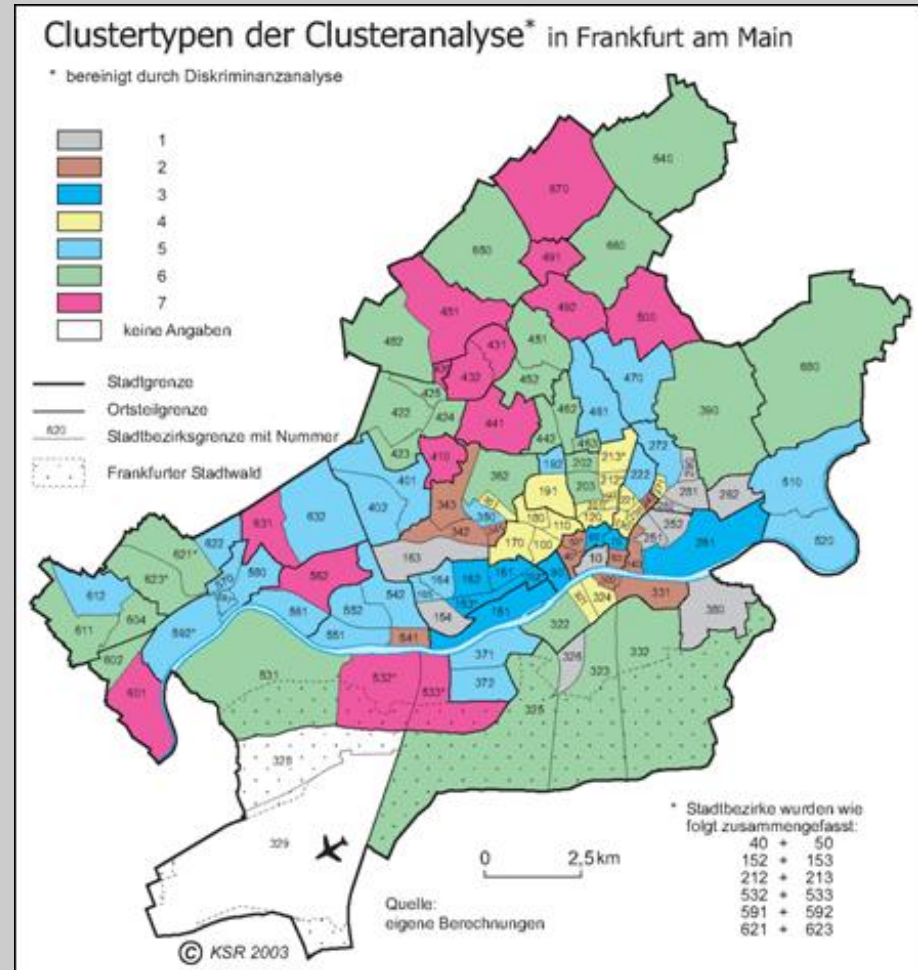
Sozialraumanalyse:

- Cluster = ahnliche Bezirke
- Merkmale sozio-konomische Variablen*:

- unter 6-Jahrig
- ber 65-Jahrig
- Zu- u. Wegzuge
- Auslander/-innen
- Einpersonenhaushalte
- Arbeitslosendichte
- Sozialhilfeempfanger/innen
- Mehrfamilienhuser
- Wahlbeteiligung

* jeweils Anteile

=> Variablenauswahl extrem wichtig



Fragestellung

Definition

Voraussetzung

Methodik

Interpretation

Datenvoraussetzungen - Empfehlungen:

- **kein spezielles Skalenniveau**
- **Standardisierte Merkmale** (z-Transformation)
- **Ausreißer ausschließen** (Verzerrungen)
- **Anzahl der Merkmale**
Keine Begrenzung, aber:
 - nur relevante Variablen einbeziehen - Vorüberlegungen
 - möglichst hoch korrelierenden Variablen ausschließen
 - keine Variablen mit konstanten Ausprägungen bei allen Objekten
- **Anzahl der Objekte**
Keine Begrenzung



Fragestellung

Definition

Voraussetzung

Methodik

Interpretation

Drei Ablaufschritte:

1. Bestimmung der Distanz (Abstand-Differenz) durch Proximitätsmaße
2. Auswahl des Fusionierungsalgorithmuses
3. Bestimmung der Clusteranzahl

Fragestellung

Definition

Voraussetzung

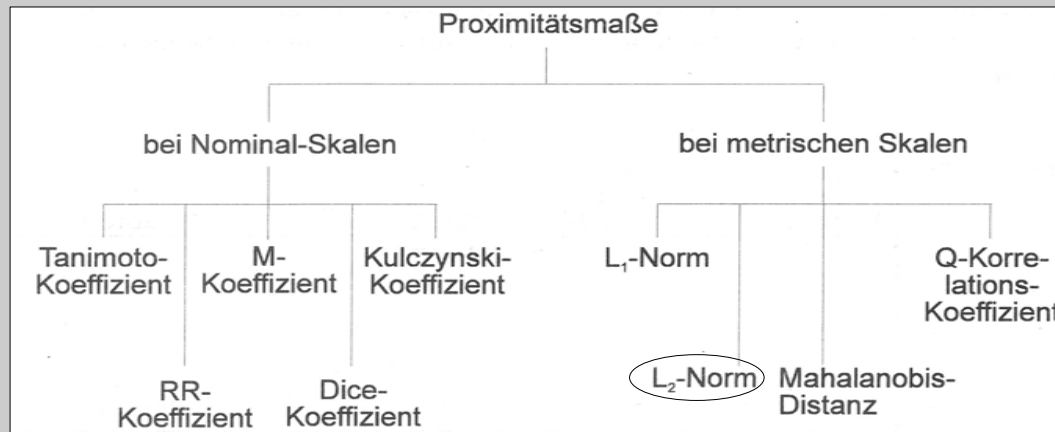
Methodik

Interpretation

Proximitatsma

Proximitatsma: Mazahl zur Quantifizierung des Abstandes der Objekte durch Merkmalswerte

Vielzahl von Maberechnungen - abhangig vom Skalenniveau



**Standard bei Distanzmaen bei metrischen Skalen:
L2-Norm = Quadrierte Euklidische Distanz**

Berechnung: absoluten Differenzwerte werden quadriert und addiert

Fragestellung

Definition

Voraussetzung

Methodik

Interpretation

Fusionierung

Ausgangsdatenmatrix

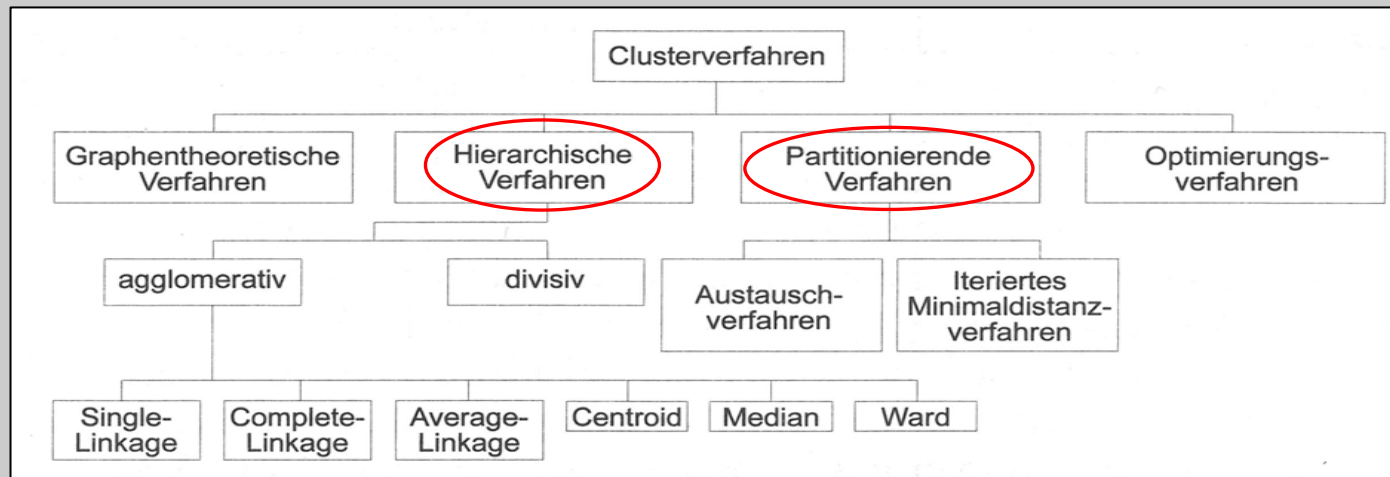
Distanzma

(Quadierte Euklidische Distanz)

Distanzmatrix

Fusionierung

mit Hilfe von Cluster-Algorithmen:



Fragestellung

Definition

Voraussetzung

Methodik

Interpretation

Fusionierung

**Partitionierendes
Verfahren:**

1. **Vorgabe einer Anfangspartition**
2. **jedes Objekt kann im Prozess jederzeit verschoben werden**
3. **Festlegung der Clusteranzahl im vorhinein**

**Hierarchisch-
agglomerativ
Verfahren:**

1. **Keine Vorgabe – Start mit feinsten Partition**
(jedes Objekt ist ein Cluster)
2. **Objekte mit der geringsten Distanz werden verbunden, später Gruppen**
3. **„Durchlaufen“ zu einem Großcluster**

Fragestellung

Definition

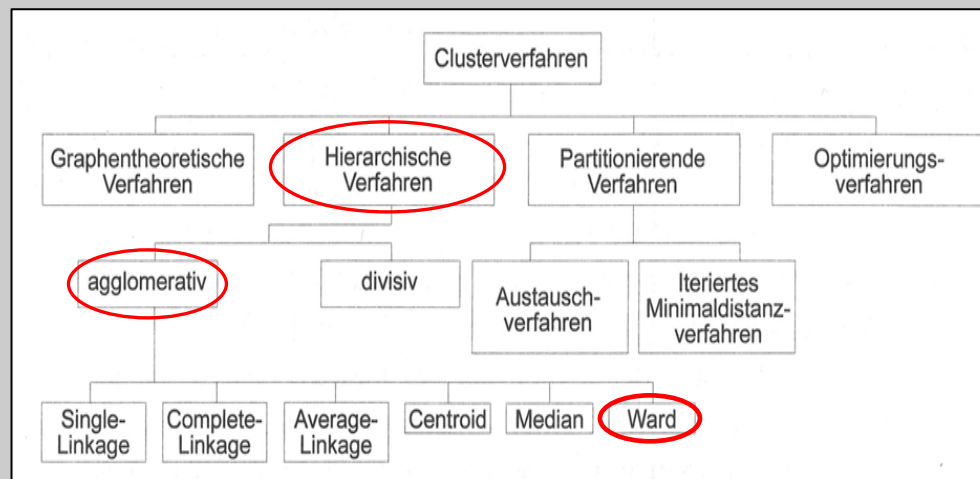
Voraussetzung

Methodik

Interpretation

Fusionierung

Hufigste Anwendung:



Ward Verfahren:

Bildung von homogeneren Clustern –

„Vereinigt diejenigen Objekte, die die Fehlerquadratsumme (Varianz/Streuung) am wenigsten erhohen“



Fragestellung

Definition

Voraussetzung

Methodik

Interpretation

Clusteranzahl

Moglichkeiten zur Bestimmung der Clusteranzahl:

1. Fehlerquadratsumme
2. Elbow-Kriterium
3. Dendrogramm

Zur Erinnerung! : aggl. Verf. nach Ward geht von der kleinsten Partition aus und endet bei einem Grocluster!

Fragestellung

Definition

Voraussetzung

Methodik

Interpretation

Clusteranzahl

1. Fehlerquadratsumme

= gibt die Varianzveränderung an, bei großen Sprüngen werden heterogene Cluster zusammengeführt

Schritt	Zusammengeführte Cluster		Koeffizienten	Erstes Vorkommen des Clusters		Nächster Schritt	Diff. CL1:CL2	Clusteranzahl
	Cluster 1	Cluster 2		Cluster 1	Cluster 2			
100	25	62	267,35	99	91	103	21,8	10
101	14	18	291,56	94	97	105	24,2	9
102	3	13	317,16	96	0	108	25,6	8
103	25	47	346,46	100	74	107	29,3	7
104	2	7	387,24	93	89	106	40,8	6
105	14	65	428,42	101	98	107	41,2	5
106	1	2	473,85	95	104	108	45,4	4
107	14	25	599,89	105	103	109	126,0	3
108	1	3	738,35	106	102	109	138,5	2
109	1	14	981,00	108	107	0	242,7	1

Fragestellung

Definition

Voraussetzung

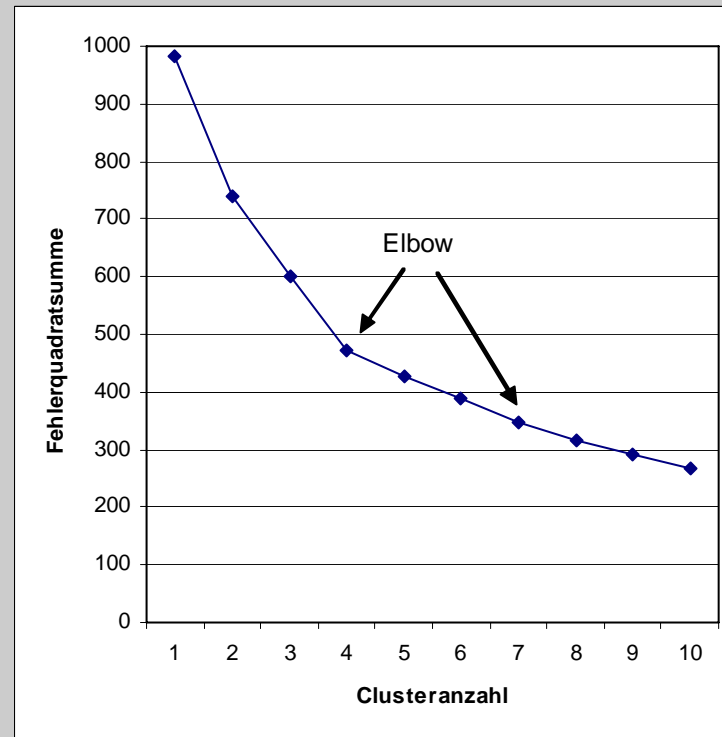
Methodik

Interpretation

Clusteranzahl

2. Elbow-Kriterium

= Abtragung der Fehlerquadratsumme in ein Diagramm



Fragestellung

Definition

Voraussetzung

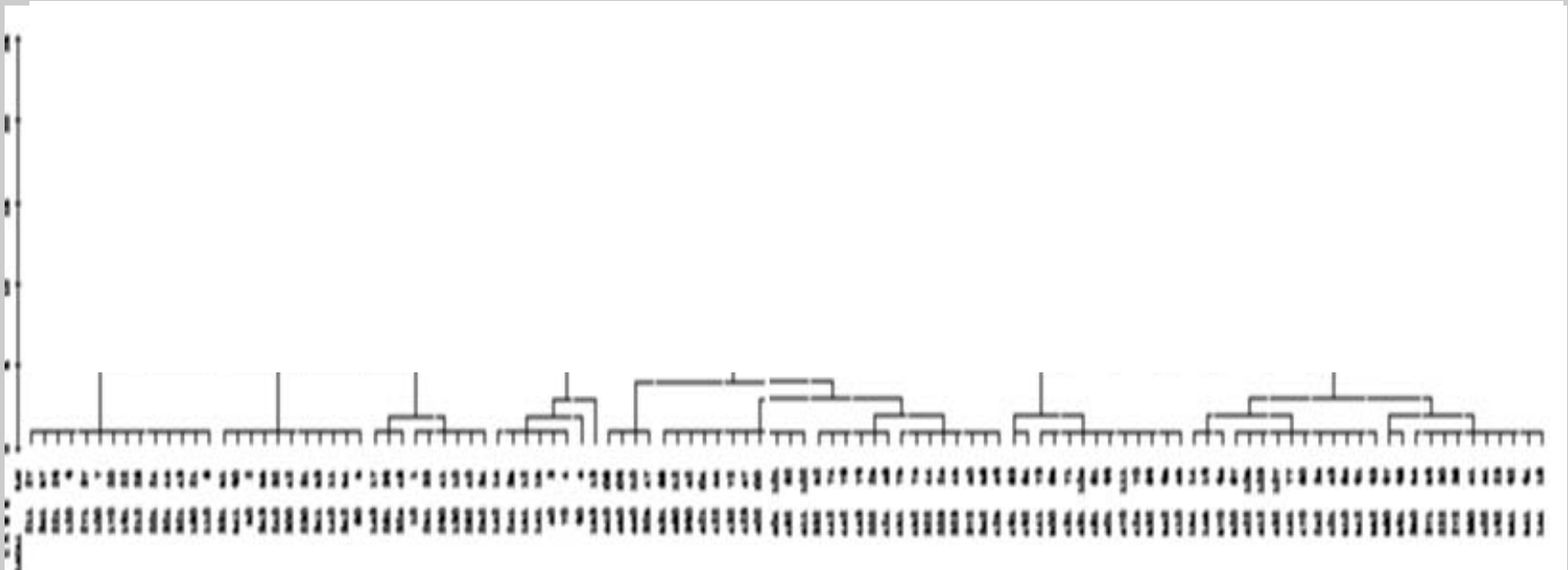
Methodik

Interpretation

Clusteranzahl

3. Dendrogramm

= Abtragung der Fusionsierungsschritte



Festlegung Clusteranzahl => Clustermatrix

Fragestellung

Definition

Voraussetzung

Methodik

Interpretation

Drei Ablaufschritte:

1. Bestimmung der **Distanz** durch Proximitätsmaße:
Quadrierte Euklidische Distanz
2. Auswahl des **Fusionierungsalgorithmuses**:
hierarchisch, agglomerativ nach Ward
3. Bestimmung der **Clusteranzahl**:
Fehlerquadratsumme, Elbow und Dendrogramm

Fragestellung

Definition

Voraussetzung

Methodik

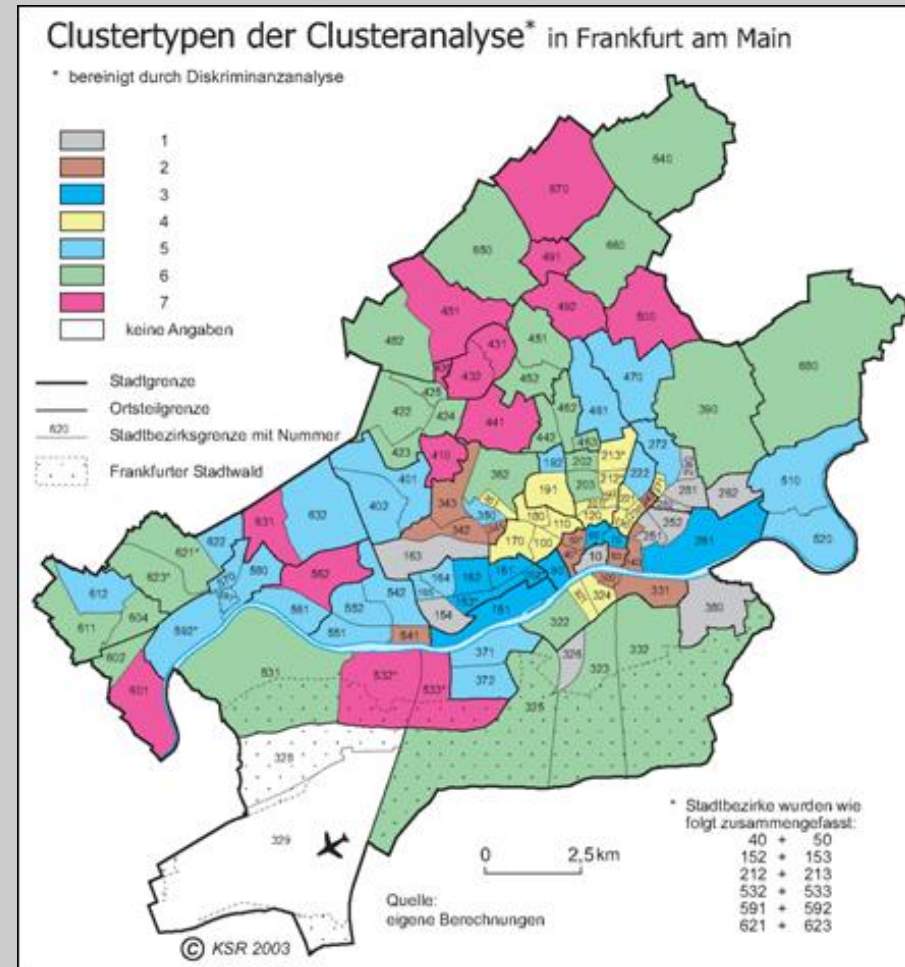
Interpretation

Hilfestellung fur die Interpretation:

- absoluten Merkmalsauspragungen
- Karte

Fur Spezialisten:

- F- Werte (Varianz/Streuung)
- t-Wert (Standardabweichung)



Fragestellung

Definition

Voraussetzung

Methodik

Interpretation

- **F-Werte (Varianz): Homogenitat einer Gruppe;**
- **t-Werte (Standardabweichung): Uberbewertung/Unterbewertung einer Variablen)**

Cluster	Arbeitslosendichte	Anteil d. Sozialhilfeempfanger/-innen an der Bevolkerung	Wahlbeteiligung b.d. Bundestagswahl	Anteil d. Einpersonenhaushalte a.d. Privathaush.	Anteil der Bev. im Alter v. u. 6 Jahren a.d. Bevolkerung	Anteil d. Mehrfamilienhuser an den Wohngebuden	Anteil d. Zu- u. Wegzuge an der Bevolkerung	Anteil der ausl. Bev. an der Bevolkerung	Anteil der Bev. v. 65 Jahren und alter a.d. Bevolkerung
1	0,35	0,40	0,40	0,13	0,19	0,32	0,15	0,27	0,76
2	0,42	0,35	0,35	0,34	0,18	0,28	0,43	0,10	0,58
3	1,38	0,60	0,60	0,34	1,88	0,07	1,12	0,28	0,66
4	0,38	0,15	0,15	0,06	0,09	0,07	0,11	0,08	0,21
5	0,41	0,34	0,34	0,30	0,49	0,35	0,33	0,21	0,40
6	0,50	0,39	0,39	0,71	0,42	0,50	0,31	0,19	0,89
7	0,23	0,18	0,18	0,12	0,48	0,23	0,09	0,06	0,28

Cluster	Arbeitslosendichte	Anteil d. Sozialhilfeempfanger/-innen an der Bevolkerung	Wahlbeteiligung b.d. Bundestagswahl	Anteil d. Einpersonenhaushalte a.d. Privathaush.	Anteil der Bev. im Alter v. u. 6 Jahren a.d. Bevolkerung	Anteil d. Mehrfamilienhuser an den Wohngebuden	Anteil d. Zu- u. Wegzuge an der Bevolkerung	Anteil der ausl. Bev. an der Bevolkerung	Anteil der Bev. v. 65 Jahren und alter a.d. Bevolkerung
1	0,55	0,37	0,02	0,44	-0,81	0,60	-0,50	-0,30	1,18
2	-0,09	0,25	-0,76	0,90	-0,86	0,77	1,21	0,97	-0,74
3	1,57	1,16	-2,19	1,07	-0,47	0,98	2,31	2,42	-1,18
4	-0,82	-0,83	0,64	1,15	-0,92	1,12	0,22	-0,22	-0,75
5	0,63	0,48	-0,46	-0,28	0,67	0,00	-0,13	0,45	-0,04
6	-0,80	-0,92	0,83	-0,61	0,09	-0,91	-0,61	-0,97	0,65
7	0,14	0,77	0,24	-1,25	1,25	-1,11	-0,65	-0,49	-0,06



Problembehandlung

Alternative Methoden

Fazit

- **Vorüberlegungen sind wichtig:**
 - Probleme bei der Auswahl der Variablen
 - Korrelationen
 - Ausreißer
- Festlegung von Grenzwerten – aber welche Grenzwerte sind gültig?
- Bestimmung der Clusteranzahl – ein Glücksspiel?
- Zuweisung bei hierarchisch-agglomerativen Verfahren
(nicht revidierbar, dafür Tendenz, gleich große Gruppen zu bilden, Ausreißer erkennbarer)



Problembehandlung

Alternative Methoden

Fazit

- **Zur Überprüfung der Klassenzugehörigkeit:
Diskriminanzanalyse**
- **Vorüberlegungen sind wichtig:
Konkretisierung der Problemstellung der Untersuchung
Verfahren:**
 - **Faktorenanalyse** (Reduzierung der Merkmale)
 - **Klassische Raumanalyse** (feste Grenzwerte)

Problembehandlung

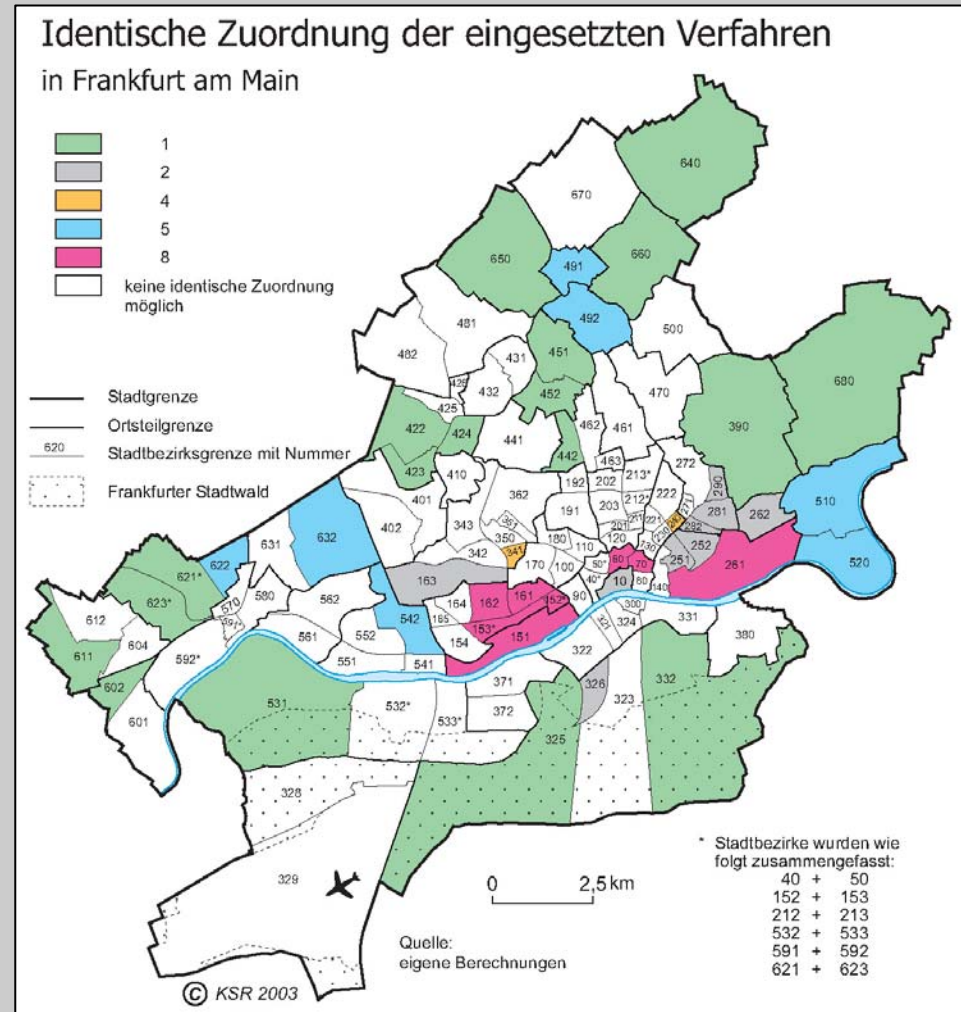
Alternative Methoden

Fazit

Ergebnisvergleich:

**Konkrante Zuordnung –
Klassische Raumanalyse:**

- 43 Bezirke von
110 Bezirken waren identisch
(39,1 %)
- Bezirke mit eindeutiger
„inhaltlicher Aussage“



Ergebnis:

- Auswahl der Analyseverfahren muss sich an konkreten Fragestellungen orientieren
- methodische Annahmen, Voraussetzungen und Bedingungen mussen erlautert werden
- Clusteranalyse mit unterschiedlichen Programmen durchfuhrbar

=> Clusteranalyse ist ein geeignetes Verfahren fur Stadtestatistiker



Sind Sie nun von der

Clusteranalyse

überzeugt ! ?