

# Übersicht

Normalverteilung

Korrelation

Standardverfahren der multivariaten Datenanalyse

- *Vorhersage*
- *Dimensionsreduktion*
- *Klassifizierung*

Versuchsauswertung

Neuere multivariate Verfahren

# Normalverteilung

Bedeutung

Test auf Normalverteilung

- $\chi^2$ -Test
- Kolmogorov-Smirnov mit Lilliefors-Korrektur
- Shapiro-Wilks

Daten-Transformation

- *Box-Cox-Transformation*       $T(x) = \frac{x^\lambda - 1}{\lambda}$       für  $\lambda > 0$

$T(x) = \ln x$       für  $\lambda = 0$

# Korrelation

- Varianz, Kovarianz, Korrelation
- Regression
- Produkt-Moment-Korrelation
- Spearman Rangkorrelation
- Kendalls  $\tau$
- Bestimmtheitsmaß

# Standardverfahren der multivariaten Datenanalyse

	Lineare Regression	Hauptkomponentenanalyse	Korrespondenzanalyse	Clusteranalyse	Diskriminanzanalyse
<b>Zweck:</b>					
Vorhersage	x				
Dimensionsreduktion		x	x	= Ordination	
Klassifizierung				x	x
<b>Eigenschaften:</b>					
nicht-linear					
verteilungsfrei			x	x	
nominal skalierte Var.			x		x

## Nomenklatur - Übersicht

	Multivariate Regression	Hauptkomponenten-analyse	Korrespondenz-analyse	Clusteranalyse	Diskriminanz-analyse
Name der synthetischen Variable	Regressand	Hauptkomponente	-	Clusterzugehörigkeit	Diskriminanzfunktion
Wert der synthetischen Variable	Schätzwert für reale Variable	Faktorwert	Wert in 1./2. Dimension	(Clusterzugehörigkeit)	Diskriminanzwert
durch synth. Variable erklärte Gesamt-Streuung	erklärte Varianz	Anteil der Eigenwerte	Anteil der Eigenwerte (der Trägheit, der Streuung, der Inertia)	(Anteil der Fehlerquadratsumme)	Anteil der Eigenwerte (Diskriminanzanteil)
durch synth. Variable erkl. Streuung der einzelnen Variablen	-	Kommunalität	Anteil der Eigenwerte (der Trägheit, der Streuung, der Inertia)	-	-
Korrelation zwischen realer und synthetischer Variable	partielle Korrelation	Ladung	-	-	Ladung

## Zentrale Begriffe

**Varianz** = mittlere quadratische Abweichung  
 = durch Zahl der Freiheitsgrade geteilte Summe der quadratischen Abweichungen

**Bestimmtheitsmaß** = Anteil der erklärten Varianz  
 = erklärte Varianz / Gesamtvarianz

**F-Wert** = erklärte Varianz / nicht-erklärte Varianz

**Ladung** = Korrelation zwischen beobachteter Variabler und synthetischer Variabler

# ANOVA

Prinzip: Vergleich der

Streuung **zwischen** den Zellen: **erklärte Varianz** (Treatmentvarianz,  
Wechselwirkungsvarianz)

---

Streuung **innerhalb** der Zellen: **nicht-erklärte** Varianz (Fehlervarianz)

= **F-Wert**

## Versuchsauswertung

	<u>ANOVA</u>	<u>Kruskal-Wallis</u>
Voraussetzungen	- Homoskedastizität - Normalverteilung - Unabhängigkeit der Fehlerkomponenten	- symmetrische Streuung um den Median
Erweiterung des	- t-Tests	- Wilcoxon-Tests
Vergleich der	- Mittelwerte	- Mediane
basierend auf	- metrisch skal. Werten	- Ränge
Posthoc-Test	- z.B. Scheffé	

# Konfidenzintervall

= geschätzter Bereich, in dem mit einer gewissen Wahrscheinlichkeit ein bestimmter Wert liegt

## Beispiele:

Schätzung des      wenn

Erwartungswerts    - Normalverteilung  
                          - Varianz bekannt

$$x_{KI} = \bar{x} \pm \frac{\sigma}{\sqrt{n}} \cdot z\left(1 - \frac{\alpha}{2}\right)$$

Erwartungswerts    - Normalverteilung  
                          - Varianz unbekannt

$$x_{KI} = \bar{x} \pm \frac{s}{\sqrt{n}} \cdot t\left(1 - \frac{\alpha}{2}; n - 1\right)$$

Erwartungswerts    - Verteilung unbekannt  
                          - Varianz unbekannt  
                          -  $n > 50$

$$x_{KI} = \bar{x} \pm \frac{s}{\sqrt{n}} \cdot z\left(1 - \frac{\alpha}{2}\right)$$

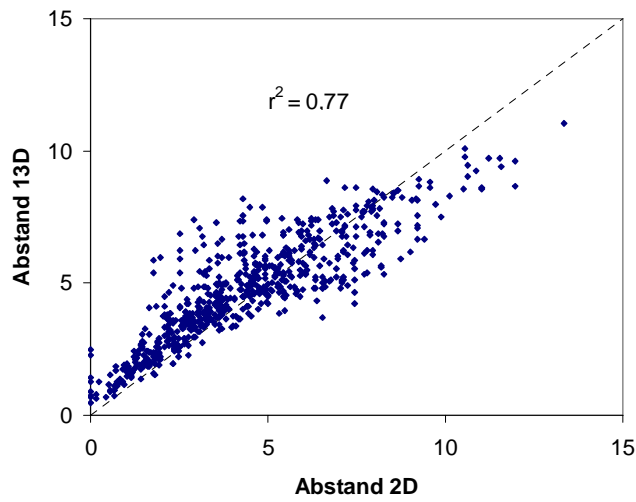
$z(\cdot)$ : Quantil der Normalverteilung

$s$ : Standardabw. der Stichprobe

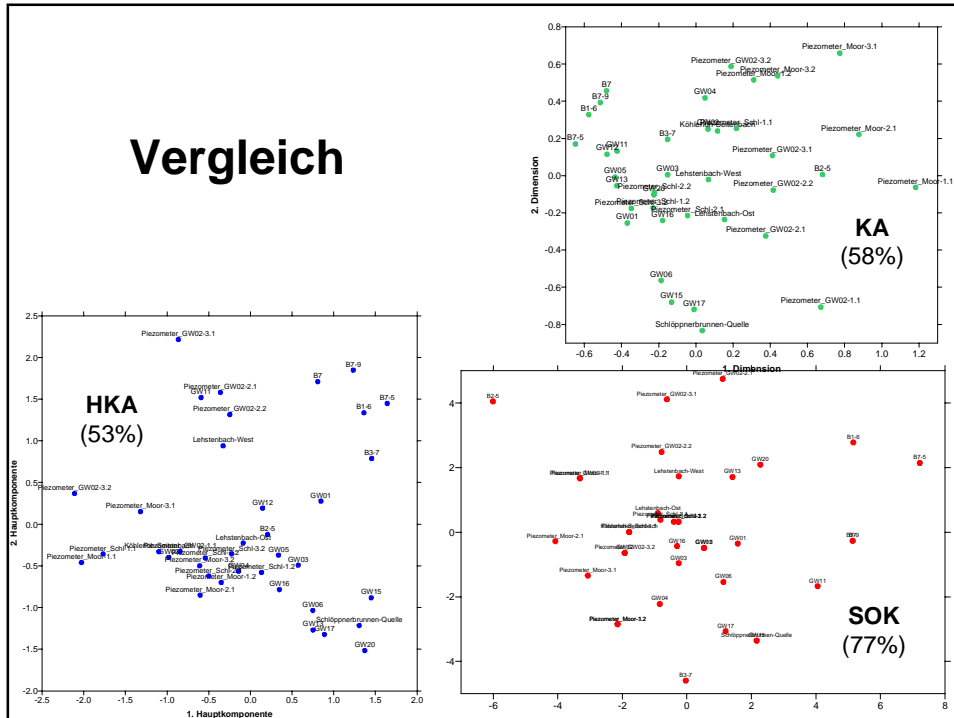
# Neuere multivariate Verfahren

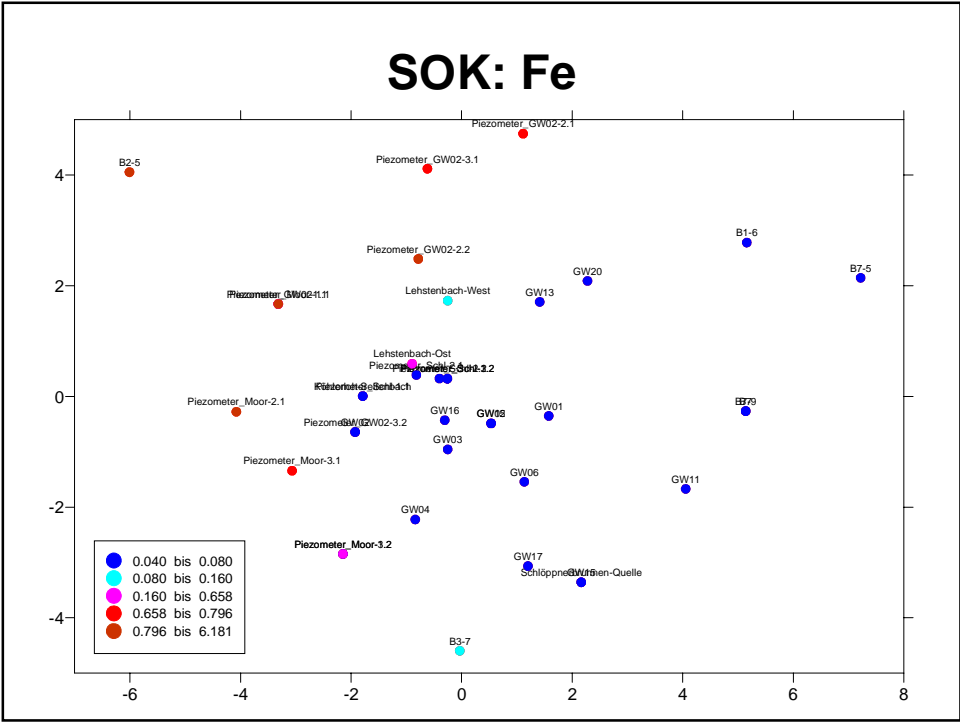
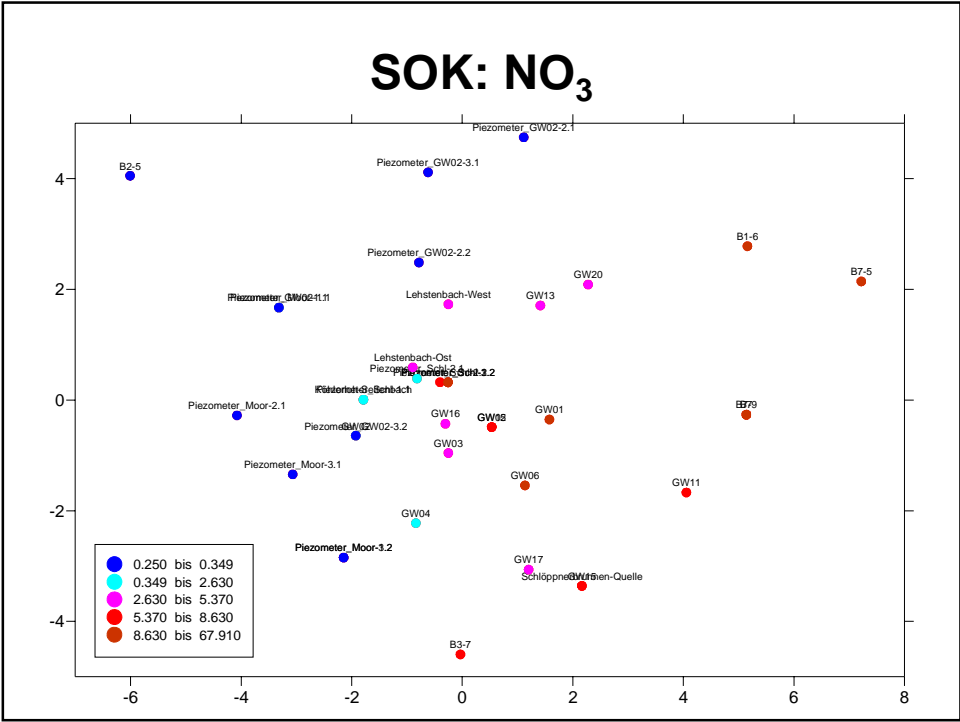
- Mehrschichtperzeptron
- Independent Component Analysis,  
Non-linear Principal Component Analysis
- Selbst-organisierende Karte, Lernende Vektor-  
Quantisierung

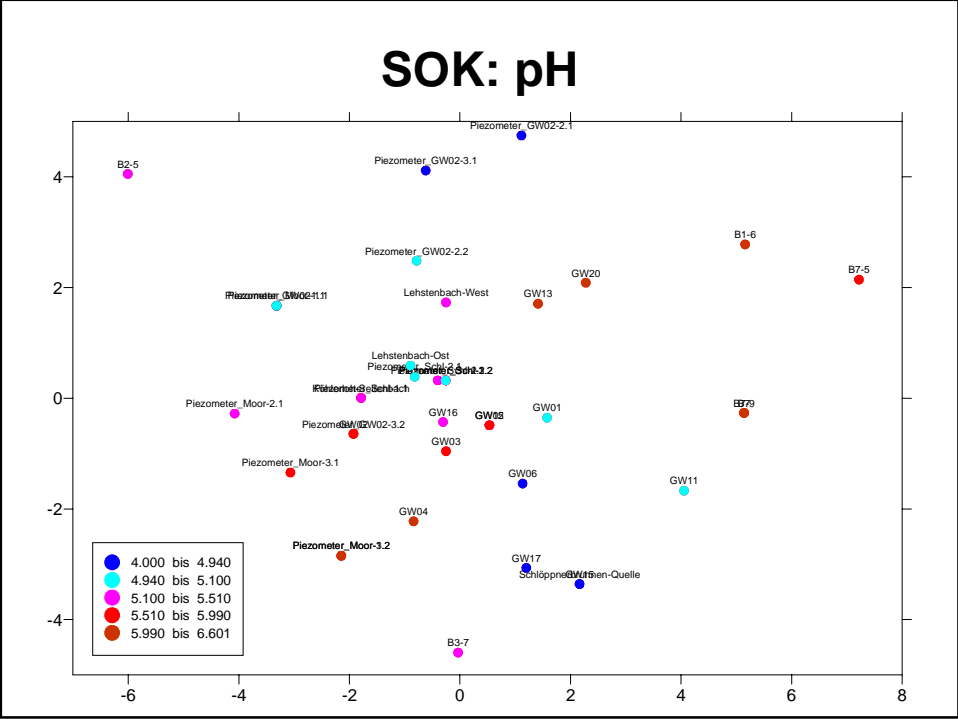
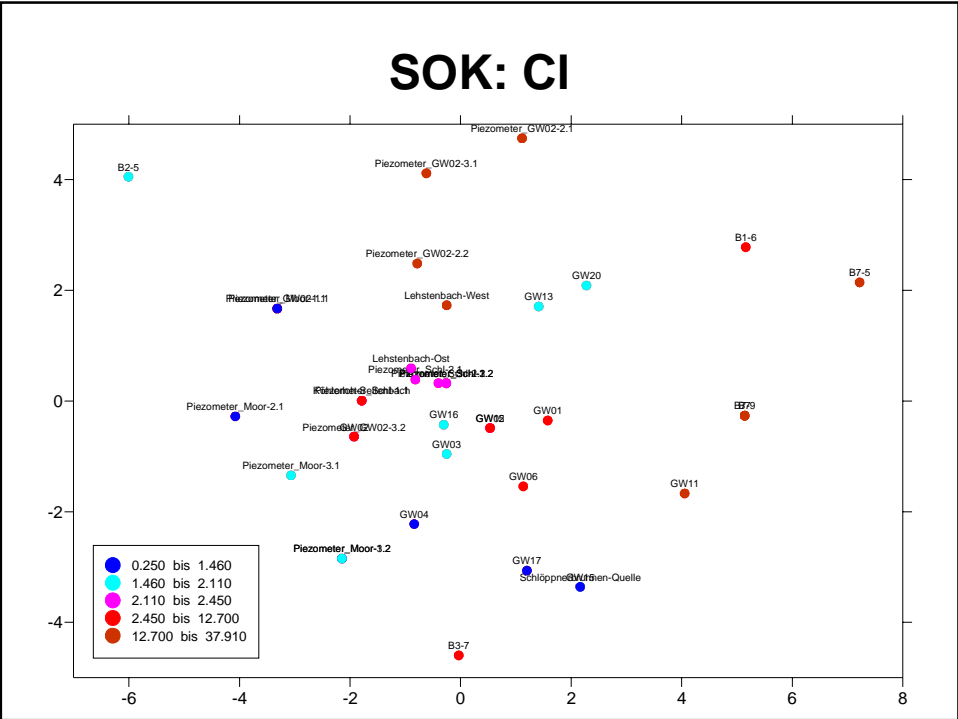
# SOK: Abstände



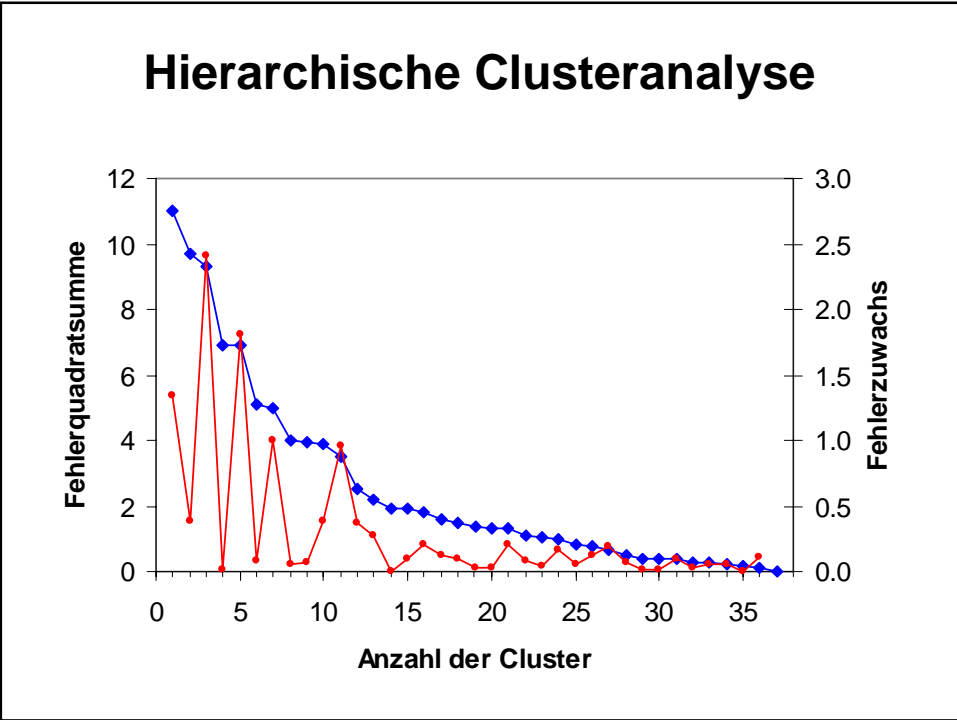
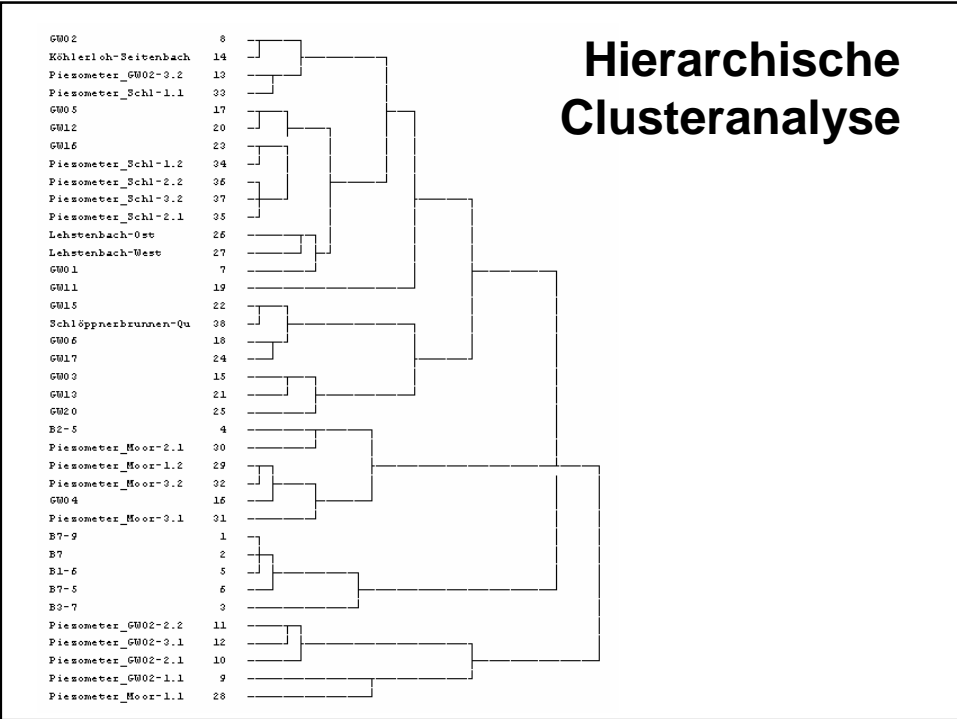
# Vergleich



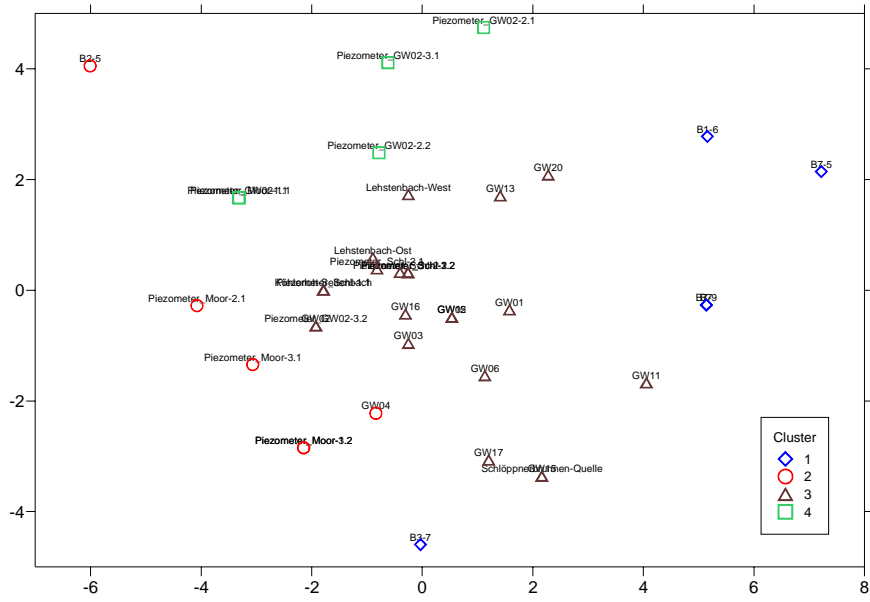




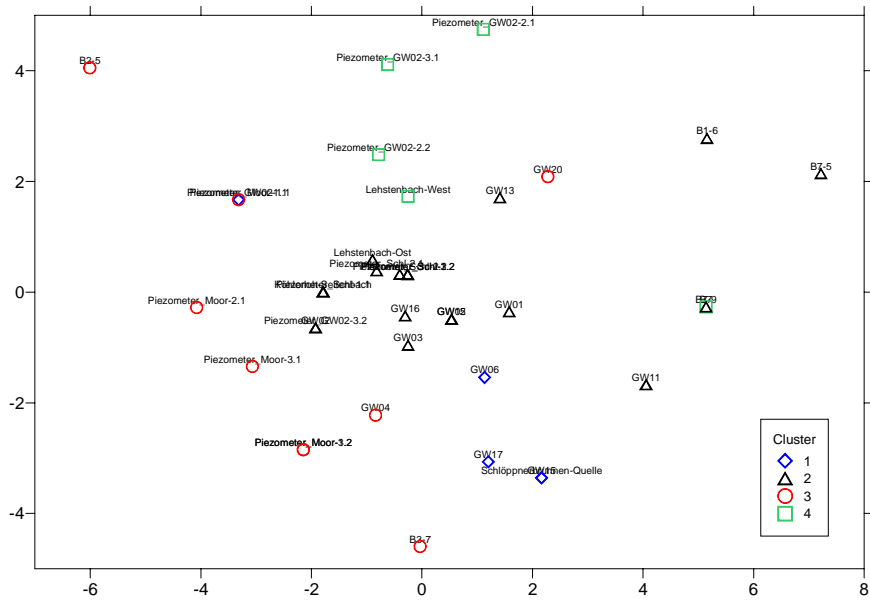




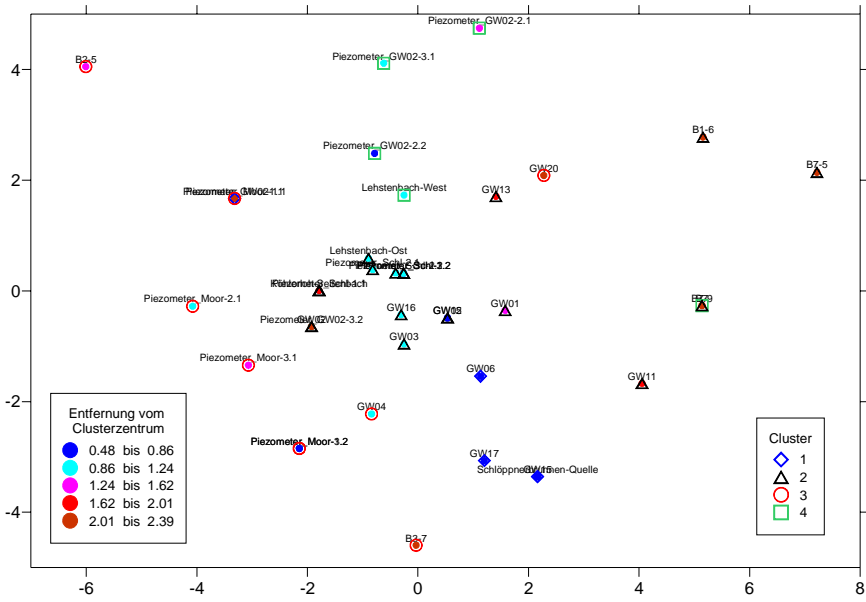
# Hierarchische Clusteranalyse



# Clusterzentrenanalyse



# Clusterzentrenanalyse

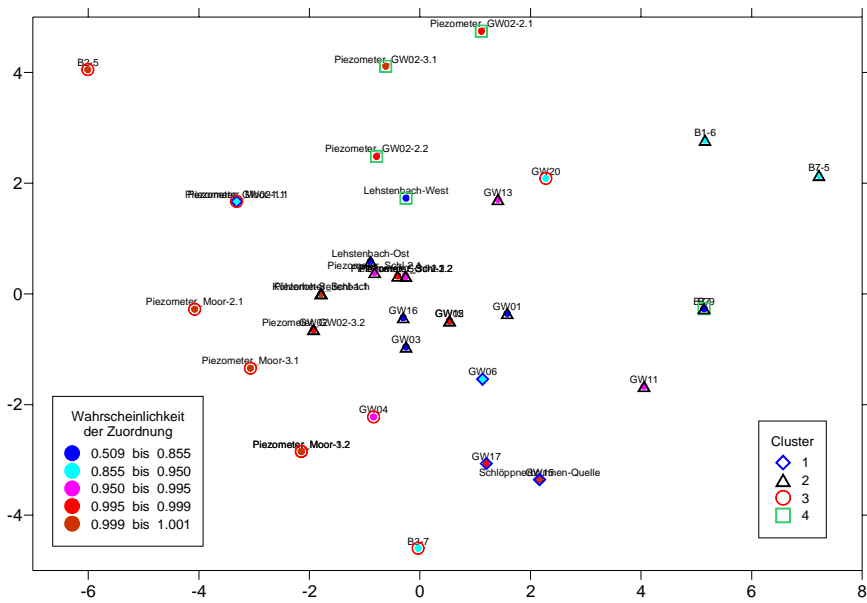


# Diskriminanzanalyse

Klassifizierungsergebnisse<sup>c</sup>

	Cluster-Nr. des Falls	Vorhergesagte Gruppenzugehörigkeit				Gesamt	
		1	2	3	4		
Original	Anzahl	1	5	0	0	0	5
		2	0	19	0	0	19
		3	0	0	9	0	9
		4	0	0	0	5	5
	%	1	100.0	.0	.0	.0	100.0
		2	.0	100.0	.0	.0	100.0
		3	.0	.0	100.0	.0	100.0
		4	.0	.0	.0	100.0	100.0
Kreuzvalidiert	Anzahl	1	5	0	0	0	5
		2	1	17	0	1	19
		3	0	1	8	0	9
		4	0	1	0	4	5
	%	1	100.0	.0	.0	.0	100.0
		2	5.3	89.5	.0	5.3	100.0
		3	.0	11.1	88.9	.0	100.0
		4	.0	20.0	.0	80.0	100.0

## Diskriminanzanalyse



## Multivariate Analyse: Take-Home Message

1. Die meisten Standardmethoden basieren auf **linearen Zusammenhängen** im Datensatz.
2. In der Regel gibt es verschiedene **Maßzahlen für die Güte des Verfahrens**, die unbedingt beachtet werden sollten (ein großes  $r^2$  alleine sagt noch nicht viel aus!).
3. Multivariate Verfahren werden überwiegend eingesetzt, um
  - einzelne Werte **vorherzusagen** (Regression)
  - die **Dimension** des Datensatzes zu **reduzieren** (Prozessanalyse, Visualisierung)
  - **Gruppen** zu identifizieren (klassifizieren).
4. Die Verfahren verlangen i.d.R. **±subjektive Entscheidungen** des Anwenders, die zu begründen sind.
5. Die Ergebnisse der hier vorgestellten **Verfahren sind nicht unabhängig** voneinander.

# Diplomarbeit

Wo: Bayer. Geologisches Landesamt

Was: Multivariate Statistik, Geostatistik

Warum: EU-Wasserrahmenrichtlinie und  
Hydrogeologische Landesaufnahme

