

Funktionale Hauptkomponentenanalyse

Philipp Lintl

12. Januar 2018

Inhaltsverzeichnis

- ① Motivation Hauptkomponentenanalyse
- ② Hauptkomponentenanalyse
 - Für multivariate Daten
 - Für funktionale Daten
- ③ Anwendung auf Datensatz
 - Anzahl benötigter Hauptkomponenten
 - Visualisierung der Hauptkomponenten
 - Glättung der Hauptkomponenten

Hauptkomponentenanalyse

- umfangreiche Datensätze zu strukturieren, zu vereinfachen und zu veranschaulichen
- Grundidee: Datenreduktion bzw. Dimensionsreduktion
- betrachtete Daten auf möglichst wenige Hauptkomponenten reduzieren, ohne zu großen Informationsverlust
- → Hauptkomponenten sollen einen möglichst großen Teil der Varianz der Daten erklären
- ⇒ funktional: wichtigste Arten der Variabilität / Strukturen in Daten

Beispiel

- Variablen x_1, x_2
- gesucht:
Linearkombination
 $\alpha^T \mathbf{x} = \alpha_1 x_1 + \alpha_2 x_2$
mit maximaler
Varianz
- HK1: Gerade, bei der
die Summe der
Fehlerquadrate
minimal ist (rote
Linien)
- HK2: zu HK1
orthogonal, wieder
maximierte Varianz

Für multivariate Daten

- Datenmatrix: $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$, $\mathbf{X} \in \mathbb{R}^{N \times p}$
- N Beobachtungen, p Variablen
- zentrierte Daten: $\tilde{x}_{ij} = x_{ij} - \frac{1}{N} \sum_{i=1}^N x_{ij}$
- Varianz: $\widehat{Var}(\tilde{x}_j) = \frac{1}{N-1} \sum_{i=1}^N \tilde{x}_{ij}^2$
- Kovarianz: $\widehat{Cov}(\tilde{x}_j, \tilde{x}_k) = \frac{1}{N-1} \sum_{i=1}^N \tilde{x}_{ij} \tilde{x}_{ik}$
- Kovarianzmatrix: $\mathbf{V} = \frac{1}{N-1} \mathbf{X}^T \mathbf{X}$, $\mathbf{V} \in \mathbb{R}^{p \times p}$

- Idee: Varianz in $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_p)$ durch unkorrelierte Variablen $\boldsymbol{\xi} = (\xi_1, \dots, \xi_p)^T$ beschreiben
- $\boldsymbol{\xi}$ als Linearkombination der originalen Variablen und den **Gewichten** $\boldsymbol{\phi}$

$$\xi_1 = \sum_{j=1}^p \phi_{j1} x_j$$
$$\xi_1 = \boldsymbol{\phi}_1^T \mathbf{x}$$

- Schrittweise finden der derjenigen **Hauptkomponenten** $\boldsymbol{\phi}_i$, die Varianz der $\boldsymbol{\xi}_i$ **maximieren**

- 1. Hauptkomponente: Der Gewichtsvektor $\phi_1 = (\phi_{11}, \dots, \phi_{p1})^T$ für den

$$\xi_{i1} = \sum_j \phi_{j1} x_{ij} = \phi_1^T \mathbf{x}_i, \quad i = 1, \dots, N, j = 1, \dots, p$$

- Varianz **maximal**: $Var(\xi_1) = \frac{1}{N-1} \sum_i \xi_{i1}^2 \rightarrow \max$
- \Rightarrow erklärt stärkste Art der Variabilität in den Daten
- Eindeutigkeit durch Normierung: $\|\phi_1\|^2 = \sum_j \phi_{j1}^2 = 1$

- m-te HK: Gewicht ϕ_m mit

$$\text{Var}(\xi_m) \rightarrow \max, \quad \|\phi_m\|^2 = 1$$

- und m-1 zusätzlichen Bedingungen

$$\langle \phi_k, \phi_m \rangle = \phi_k^T \phi_m = 0, \quad k < m.$$

- zueinander orthogonale Hauptkomponenten \rightarrow jede HK erklärt Neues
- zu jedem Schritt Varianzmaximierung: insgesamt sinkende erklärte Varianz
- \Rightarrow Oft Großteil der Varianz schon durch wenige HK

Lösung des Optimierungsproblems

- wird Optimierungsproblem zu

$$\begin{aligned} \max Var(\xi) &= \max \frac{1}{N-1} \sum (\phi^T \mathbf{x})^2 = \max \frac{1}{N-1} \phi^T \mathbf{X}^T \mathbf{X} \phi \\ &= \max \phi^T \mathbf{V} \phi \quad \text{mit NB: } \|\phi\|^2 = 1, \end{aligned}$$

- \Rightarrow Suche nach größten Eigenwerten aus

$$\mathbf{V} \phi = \lambda \phi$$

- mit (λ_m, ϕ_m) Eigenwert-Eigenvektor Paare, $Var(\xi_m) = \lambda_m$ m-größter Eigenwert

Hauptkomponentenanalyse für funktionale Daten

Für funktionale Daten

- Zufallsstichprobe reellwertiger Funktionen $x_1(t), \dots, x_N(t)$ auf Intervall $\mathcal{T} = [0, T]$
- Individuell: Realisierungen eines eindimensionalen stochastischen Prozesses $X = X(t)$
- Wieder zentrierte Daten: $x_i(t) = \tilde{x}_i(t) - \frac{1}{N} \sum_{j=1}^N \tilde{x}_j(t)$
- unendlich dimensionale funktionale Daten \xrightarrow{FPCA} endlich dimensionale Darstellung

	PCA	FPCA
Variablen	$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_p),$ $\mathbf{x}_i = (x_{1i}, \dots, x_{Ni}),$ $i = 1, \dots, p$	$(x_1(t), \dots, x_N(t)),$ $t \in [0, T]$
Daten	Vektoren $\in \mathbb{R}^p$	Kurven $\in L_2(\mathcal{T})$
Mittelwert	$\mu = \mathbb{E}(X)$	$\mu(t) = \mathbb{E}(X(t))$
Kovarianz	$\text{Cov}(\mathbf{x}_j, \mathbf{x}_k) = \mathbf{V}_{jk}$	$\text{Cov}(x(s), x(t)) = \gamma(s, t)$

- Kovarianzfunktion:

$$\hat{\gamma}(s, t) = \frac{1}{N-1} \sum_{i=1}^N x_i(s)x_i(t)$$

- Multivariate Linearkombination:

$$\langle \phi, \mathbf{x} \rangle = \sum_{j=1}^p \phi_j \mathbf{x}_j \quad \xrightarrow{\text{wird zu}} \quad \langle \phi, \mathbf{x} \rangle = \int_{\mathcal{T}} \phi(t) x(t) dt$$

- \Rightarrow Gewichts-/Hauptkomponentenfunktionen $\phi(t)$
- \Rightarrow Hauptkomponentenscores

$$\xi_{i1} = \langle \phi_1, x_i \rangle = \int \phi_1(t) x_i(t) dt.$$

- *Karhunen-Loeve Erweiterung* zentrierter Daten:

$$x(t) = \sum_{j=1}^{\infty} \xi_j \phi_j(t)$$

- Dimensionsreduktion, wenn

$$x(t) \approx \sum_{j=1}^k \xi_j \phi_j(t)$$

- gute Approximation der unendlichen Summe

- gleiche schrittweise Prozedur:
- 1. $\phi_1(t)$ aus:

$$\text{Var}(\xi_1) = \frac{1}{N-1} \sum_{i=1}^N (\xi_{i1})^2 \rightarrow \max$$

- Bedingung:

$$\|\phi_1\|^2 = \langle \phi_1, \phi_1 \rangle = \int \phi_1(t)^2 dt = 1$$

- 2. weitere HK durch maximale Varianz der m-ten HK
- und zusätzlich

$$\int \phi_k(t) \phi_m(t) dt = 0, k < m$$

- Gesuchte Gewichtsfunktionen $\phi(t)$ lösen:
- **funktionales Eigenwertproblem:**

$$\int \hat{\gamma}(s, t)\phi(t)dt = \lambda\phi(s)$$

- Kovarianzoperator Γ einer Funktion ϕ :

$$\Gamma\phi(s) = \int \hat{\gamma}(s, t)\phi(t)dt$$

- Eigenwertproblem also wieder der Form:

$$\Gamma\phi = \lambda\phi$$

- ϕ nun Eigenfunktionen und $V(\xi_m) = \lambda_m$ m-größter Eigenwert von Γ .

- Unterschied zum multivariaten Fall: Anzahl mgl. Eigenwert-Eigenfunktionspaare
- theoretisch: $\max \# \text{Eigenfunktionen} = \# \text{Funktionswerte } x(t)$
 \Rightarrow unbegrenzt
- In der Praxis: Basisdarstellung der Funktionen $x_i(t)$:

$$\hat{x}_i(t) = \sum_{m=1}^k c_{im} v_m(t)$$

gemäß bekannter Basisfunktionen $v_m(t)$ (Spline, Fourier,...)

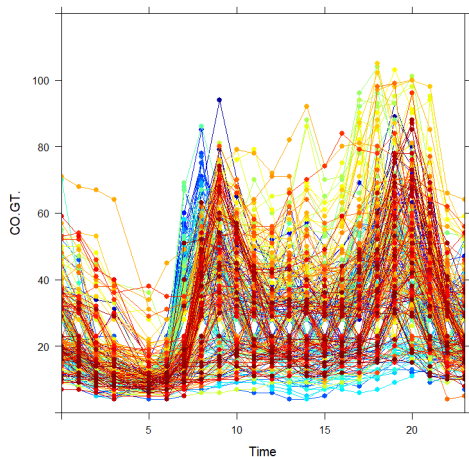
Lösung Eigenwertproblem: Diskretisierung

- diskretisieren gefitteter $x_i \in L^2(\mathcal{T})$ auf Gitter mit K Werten gleichen Abstands
- neue Datenmatrix $\mathbf{X} \in \mathbb{R}^{N \times K}$
- \Rightarrow multivariate Hauptkomponentenanalyse
- neues Eigenwertproblem \Rightarrow neue Eigenvektoren
Rücktransformation \rightarrow Funktionen

Anwendung auf Luftverschmutzungsdaten

Datensatz

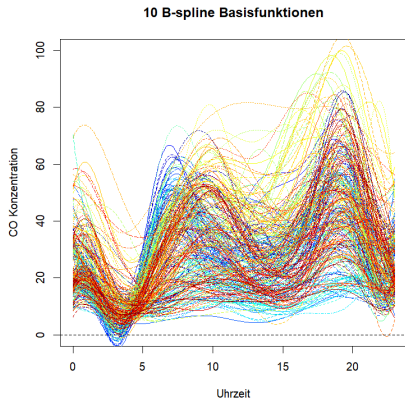
- Luftverschmutzung einer italienischen Stadt
- stündliche Mittelwerte CO-Konzentration [mg/m^3]
- Zeitraum von 03.2004 - 04.2005
- eine Beobachtung: 0-23 Uhr (4 Uhr oft NA -> entfernt)
- 282 Beobachtungen zur Analyse
- Quelle: University of California, Irvine



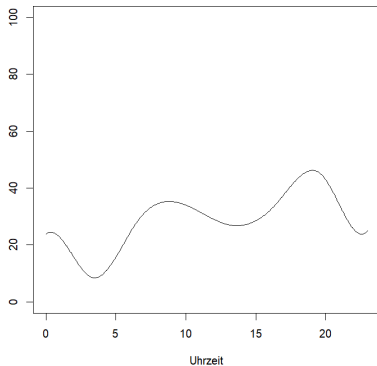
- bimodal (ca. 9 Uhr und 20 Uhr)
 - An welchen Stellen liegt nun größte Variabilität in den Daten vor?
- ⇒ Hauptkomponenten

Visualisierung

- Daten gemäß 10 Basisfunktionen



- Mittelwertsfunktion $\hat{\mu}(t)$



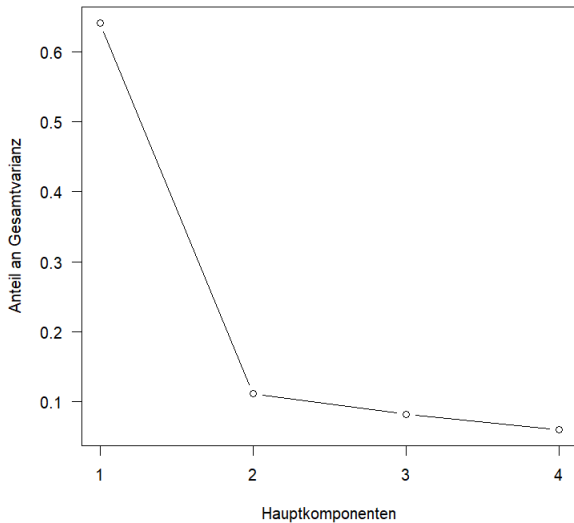
Wahl der Anzahl an Hauptkomponenten

- intuitiv: wähle so viele HK, bis $t\%$ der Gesamtvarianz in den Daten erklärt
- t meist zwischen 70% und 100%
- wegen $Var(\xi_m) = \lambda_m$, gilt $\sum_k^\infty \lambda_k = p$, mit Varianz der k -ten HK λ_k und Gesamtvarianz in den Daten p
- daher entfallen auf ersten m Hauptkomponenten

$$t_m = 100 \frac{\sum_{k=1}^m \lambda_k}{p}$$

- Sobald $t_k \geq t$ Anzahl k an Hauptkomponenten gefunden

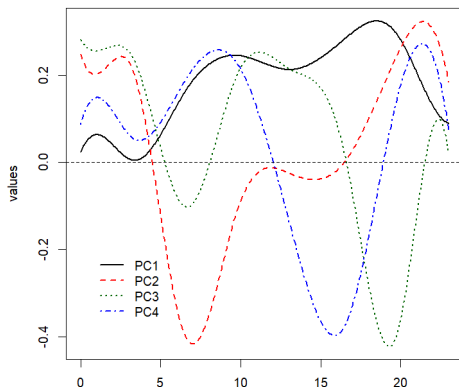
Screeplot



- HK 1: meiste Varianz
Scree: Nur die 1. HK betrachten
- Zur Veranschaulichung: 4 Hauptkomponenten

Visualisierung: Gewichtsfunktionen

- Gewichtsfunktionen ϕ_1, \dots, ϕ_4



- HK 1: Ähnlich zu $\hat{\mu}(t)$
6 - 21 Uhr stark gewichtet
→ größte Variabilität
Luftverschmutzungsverlauf
- HK 2:
Nacht- & Abendstunden positiv
→ zweite Art der Variabilität
Verschmutzung Nachts und Abends

- *Karhunen-Loeve* nicht zentrierter Daten:

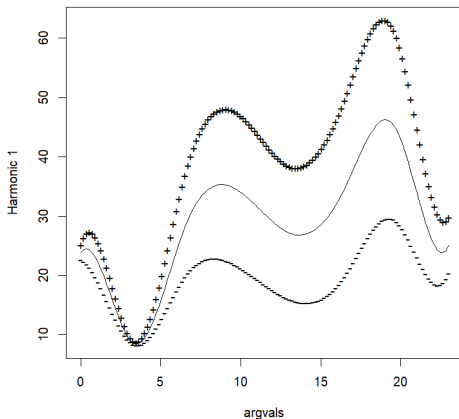
$$x(t) \approx \hat{\mu}(t) + \sum_{j=1}^k \xi_j \phi_j(t)$$

- Idee: Neue Funktion \tilde{x} , mit Score-Vektor: $\tilde{\xi} = (\pm c, 0, \dots, 0)$
- Die *Karhunen-Loeve* Darstellung davon:

$$\tilde{x}(t) = \hat{\mu}(t) + \sum_{j=1}^k \tilde{\xi}_j \phi_j(t) = \hat{\mu}(t) \pm c \phi_1(t)$$

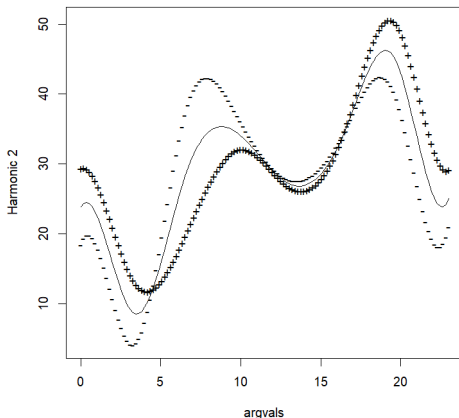
Visualisierung: Auswirkung auf Mittelwert

PCA function 1 (Percentage of variability 64.1)



- HK 1: genereller Luftverschmutzungsverlauf
- Beobachtung mit hohem ξ_1 :
Hohe Verschmutzung (überdurchschnittlich) von 6-21 Uhr, besonders an Gipfeln
- niedriges ξ_1 :
geringe Luftverschmutzung (unterdurchschnittlich)
- \Rightarrow Größte Art der Variabilität (64.1%)

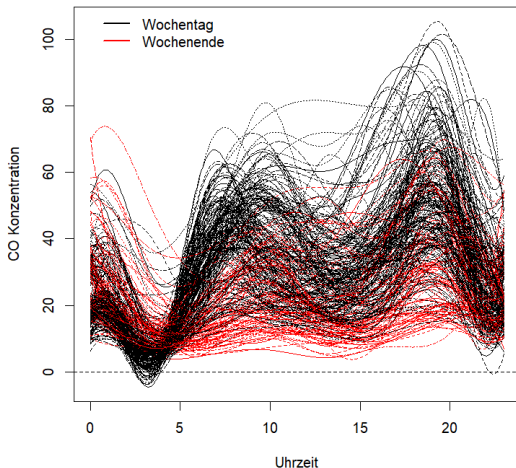
PCA function 2 (Percentage of variability 11.2)



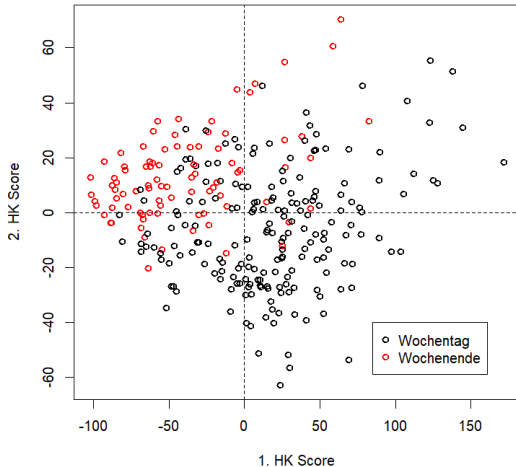
- HK 2 Variabilität nachts / spät
- Beobachtung mit hohem ξ_2 :
Hohe Verschmutzung 0-5 und 16-23 Uhr, sonst niedrig
- niedriges ξ_2 :
hohes CO zwischen 5 und 15 Uhr, sonst niedrig
- \Rightarrow Zweitgrößte Variation (11.2%)

Visualisierung der Hauptkomponenten: Scores

Beobachtungen nach Wochentag /-ende



- Nachts-/Abendstunden teilweise höhere Belastung
- über den Tag verteilt allerdings unterdurchschnittlich
- erwartbar: kleines ξ_1 und hohes ξ_2

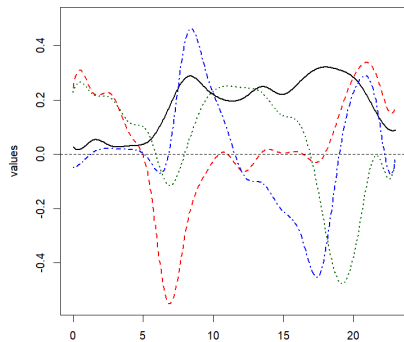
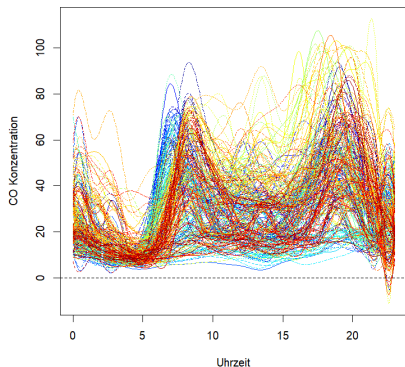


- Wochenendbeobachtungen vor allem links oben
- also niedrige generelle Verschmutzung
- höhere Belastung früh, spät und unterdurchschnittlich bei erstem Gipfel

Glättung der Hauptkomponenten

- Für raue Daten: 20 anstatt 10 Basisfunktionen
→ raue HK → schlecht interpretierbar

20 Basisfunktionen



- Glattheitsanforderung für weitere Verwendung der Hauptkomponenten
- 2 Ansätze: Daten vor HKA glätten (Splines) vs. Hauptkomponenten glätten
- mittels **Penalisierungsterm** in der Hauptkomponentenanalyse
- Einführung Rauheitsmaß:

$$PEN_2(\phi) = \|D^2\phi\|^2 = \int \phi''(t)^2 dt$$

- 2. Ableitung kontrolliert Krümmung (Rauheit)
- zuvor: $\max_{\|\phi_m\|^2=1} \text{Var}(\xi_m)$ gelöst durch λ, ϕ

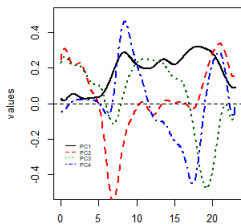
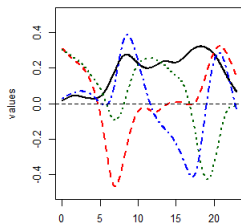
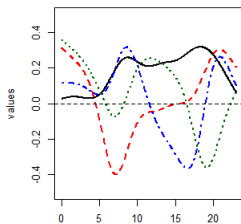
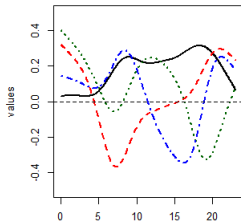
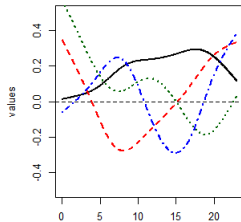
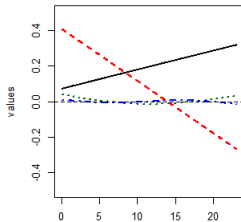
- jetzt: zusätzliche Berücksichtigung der Rauheit \Rightarrow Einführung Glattheitsparameter $\lambda \geq 0$

- **penalisierte Varianz**

$$PCAPSV(\xi) = \frac{\text{Var}(\int \phi x_i)}{\|\phi\|^2 + \lambda \text{PEN}_2(\phi)}$$

- $\lambda \rightarrow 0 : PCAPSV(\xi) \rightarrow \frac{\text{var}(\int \phi x_i)}{1}$, gleiche Hauptkomponente wie zuvor
- $\lambda \rightarrow \infty$, ergibt Konstante $\phi = a$ im periodischen Fall oder $\phi = a + bt$ im nichtperiodischen Fall
- Für optimales λ : Leave-one-out Kreuzvalidierung

Anwendung der Glättung

 $\lambda = 0$  $\lambda = 0.01$  $\lambda = 0.5$  $\lambda = 1$  $\lambda = 10$  $\lambda = 100000$ 

Zusammenfassung

- PCA: beobachtete Variablen $\xrightarrow{\text{Orthogonaltransformation}}$ linear unabhängige Hauptkomponenten (Dimensionsreduktion ohne zu großen Informationsverlust)
- FPCA: Explorative Methode zur Erkennung von Mustern und Variationsquellen in funktionalen Daten
- Anwendung: Meiste Variabilität in den Tagesverläufen (mit Peaks bei Hauptverkehrszeiten)
- Unterschied Werkstage / Wochenendtage
- Glättung durch Basenwahl oder Rauheitsmaß (ähnliche Ergebnisse)

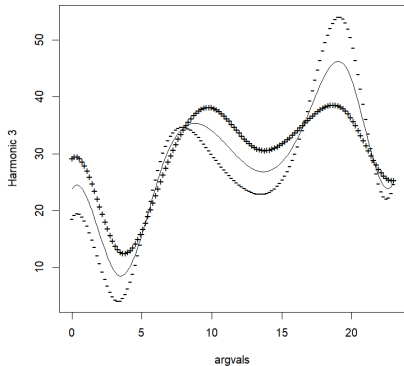
Referenzen

- J. O. Ramsay and B. W. Silverman. Functional Data Analysis. Springer, 2005.
- I. Jolliffe. Principal Component Analysis. Springer, 2 edition, 2002.
- J. Ramsay and B. Silverman. Applied Functional Data Analysis: Methods and Case Studies. Springer, 2002.
- Ramsay, J. O., Wickham, H., Graves, S., and Hooker, G. (2011). fda: Functional Data Analysis. R package version 2.2.6.
- J.-L.Wang, J.-M. Chiou, and H.-G. Müller. Functional data analysis. Annual Review
- Datensatz: <https://archive.ics.uci.edu/ml/datasets/Air+Quality#>

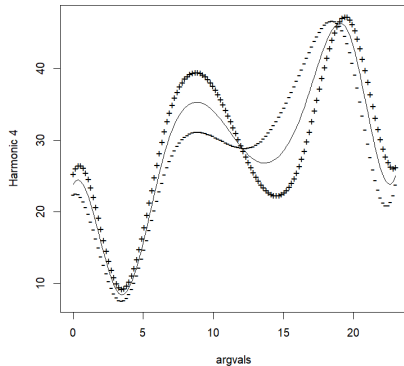
Anhangsfolien

- Hauptkomponenten 3,4

PCA function 3 (Percentage of variability 8.3)



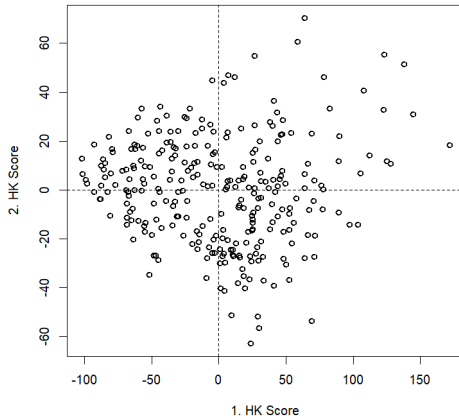
PCA function 4 (Percentage of variability 6.1)



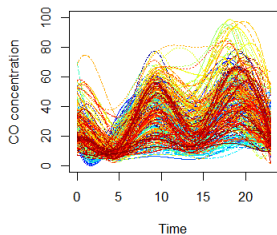
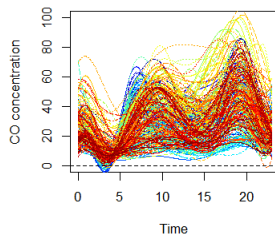
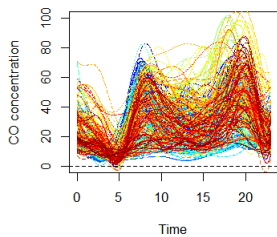
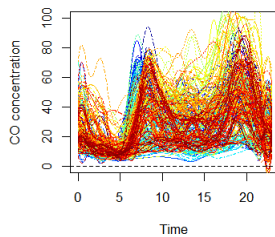
- Interpretation immer schwieriger

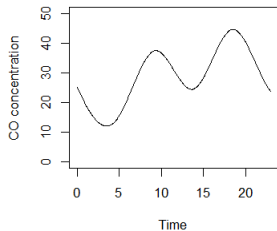
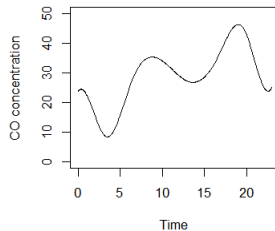
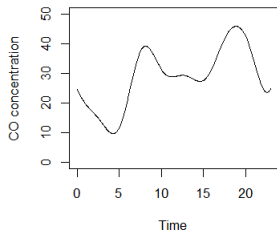
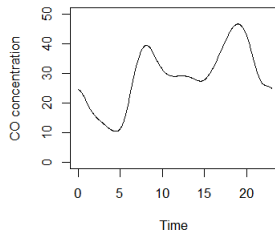
Visualisierung der Hauptkomponenten: Scores

● Hauptkomponentenscores: ξ_1 vs. ξ_2



- tendenziell gleich auf Bereiche verteilt
- großes $\xi_1 \rightarrow$ großes ξ_2
- Kurven mit hoher Verschmutzung auch früh und spät stark verschmutzt
- \rightarrow eher schwer interpretierbar
- mgl. Unterschiede zwischen Wochenende/Arbeitstage

8 Basisfunctions**10 Basisfunctions****12 Basisfunctions****20 Basisfunctions**

8 Basisfunctions**10 Basisfunctions****12 Basisfunctions****20 Basisfunctions**

- Definition $\tilde{\phi}$ aus k Werten $\phi(s_j)$

- Dann gilt approximativ:

$$V\phi(s_j) = \int \gamma(s_j, s)\phi(s)ds \approx \frac{T}{n} \sum \gamma(s_j, s_k)\tilde{\phi}_k,$$

- mit Elementen der Kovarianzmatrix \mathbf{V} : $\gamma(s_j, s_k)$
- diskrete Form des funktionalen Eigenwertproblems:

$$\frac{T}{k} \mathbf{V} \tilde{\phi} = \lambda \tilde{\phi}$$

- unter $\frac{T}{k} \|\tilde{\phi}\|^2 = 1$
- gilt für die diskrete Approximierung: $\tilde{\phi} = \frac{T}{k}^{-\frac{1}{2}} \mathbf{u}$
- Funktion ϕ dann aus $\tilde{\phi}$ durch geeignete Interpolation

- Lösung: Singulärwertzerlegung: $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{W}^T$
 - \mathbf{U} ist $N \times q$ und orthogonal: $\mathbf{U}^T \mathbf{U} = \mathbf{I}_q$;
 - \mathbf{D} ist eine $q \times q$ Diagonalmatrix mit $\text{diag}(\mathbf{D}) = d_1 \geq \dots \geq d_q \geq 0$;
 - \mathbf{W} ist $k \times q$ und auch orthogonal: $\mathbf{W}^T \mathbf{W} = \mathbf{I}_q$
- symmetrische Matrix $\mathbf{V} \Rightarrow d_i$ beinhalten alle nichtnegativen Eigenwerte von \mathbf{X}
- Auswirkungen auf Kovarianzoperator \mathbf{V} :

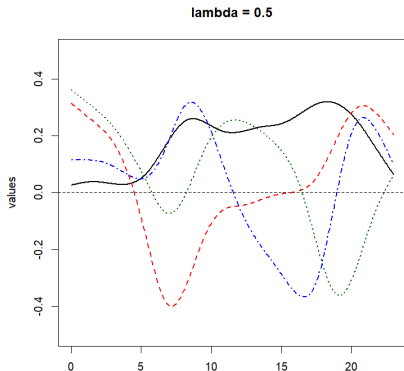
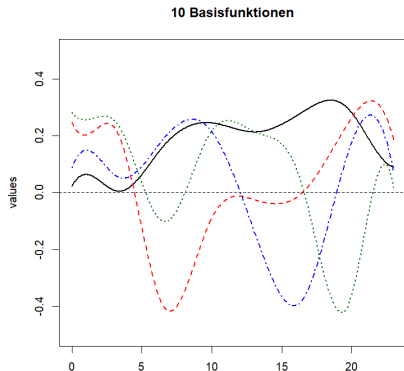
$$\mathbf{NV} = \mathbf{X}^T \mathbf{X} = (\mathbf{W}\mathbf{D}^T \mathbf{U}^T)(\mathbf{U}\mathbf{D}\mathbf{W}^T) = \mathbf{W}\mathbf{D}^2 \mathbf{W}^T$$
- Eigenwerte von \mathbf{V} dann $\text{diag}(\mathbf{D})^2$; Eigenvektoren in Spalten von \mathbf{W}

- maximieren von $\mathbf{PCAPSV}(\phi_j)$ bzgl.
- $\|\phi_j\|^2 = 1$ und einer modifizierten Orthogonalitätsbedingung:

$$\int \phi_j(t)\phi_k(t)dt + \int D^2\phi_j(t)D^2\phi_k(t)dt = 0, \quad k = 1, \dots, j-1$$

- ergibt j-te Hauptkomponente ϕ_j

- Unterschied zur Glättung vor Hauptkomponentenanalyse



- Anfangsbereich verschieden, danach ähnlich
- Für optimales λ : Leave-one-out Kreuzvalidierung