

# BZQ II: Stochastikpraktikum

Block 3: Lineares Modell, Klassifikation, PCA

Randolf Altmeyer

January 9, 2017

- 1 Monte-Carlo-Methoden, Zufallszahlen, statistische Tests
- 2 Nichtparametrische Methoden
- 3 **Lineares Modell, Klassifikation, PCA**
- 4 Markov-Chain-Monte-Carlo-Verfahren
- 5 Simulation stochastischer Prozesse

- 1 Lineares Modell (lineare Regression)
- 2 PCA
- 3 Klassifikation

Literatur:

- Mathias Trabs, Markus Reiß, Moritz Jirak: *Methoden der Statistik* (Skript)

# Das lineare Modell (oder auch: Lineare Regression)

- beobachte  $(x_1, Y_1), \dots, (x_n, Y_n)$ , wobei  $x_k \in \mathbb{R}^p$ ,  $Y_k \in \mathbb{R}$  und

$$Y_k = x_k^\top \beta + \varepsilon_k,$$

$\beta \in \mathbb{R}^p$  unbekannt,  $\varepsilon_k$  iid und  $\mathbb{E}[\varepsilon_k] = 0$

- **Ziel:** finde ein "gutes"  $\hat{\beta}$
- **Idee:**

$$\hat{\beta} = \operatorname{argmin}_{b \in \mathbb{R}^p} \left( \sum_{k=1}^n \|Y_k - x_k^\top b\|^2 \right) = \operatorname{argmin}_{b \in \mathbb{R}^p} \|Y - Xb\|^2,$$

wobei  $Y = (Y_1, \dots, Y_n)^\top \in \mathbb{R}^n$ ,  $X \in \mathbb{R}^{n \times p}$  hat die  $x_k^\top$  als Zeilen

- **explizite Lösung ("kleinster-Quadrate-Schätzer"):**
  - nimm an, dass  $X$  vollen Spaltenrang ( $= p$ ) hat
  - $X\hat{\beta}$  = Bestapproximation von  $Y$  im Spaltenraum von  $X \Rightarrow X\hat{\beta} = \Pi_X Y$
  - kann zeigen, dass  $\Pi_X = X(X^\top X)^{-1}X^\top \Rightarrow \hat{\beta} = (X^\top X)^{-1}X^\top Y$

- Beispiel:

- Polynomregression:  $Y_i = \beta_0 + x\beta_1 + x^2\beta_2 + \dots + x^{p-1}\beta_{p-1} = (1, x, x^2, \dots, x^{p-1})(\beta_0, \beta_1, \dots, \beta_{p-1})^\top$

- **Satz von Gauß-Markov:**

- $\mathbb{E}[\hat{\beta}] = \mathbb{E}[(X^\top X)^{-1}X^\top Y] = (X^\top X)^{-1}X^\top \mathbb{E}[Y] = (X^\top X)^{-1}X^\top X\beta = \beta$

- wenn  $\langle v, \beta \rangle$  geschätzt werden soll für  $v \in \mathbb{R}^p$ , dann ist  $\langle v, \hat{\beta} \rangle$  *best linear unbiased estimator* (BLUE)

- Vorhersagefehler/prediction error:

$$\mathbb{E} \left[ \left( v^\top \hat{\beta} - v^\top \beta \right)^2 \right] = \text{Var} \left( v^\top \hat{\beta} \right) = v^\top \text{Var} \left( \hat{\beta} \right) v,$$

$$\text{Var}(\hat{\beta}) = \mathbb{E} \left[ \left( \hat{\beta} - \beta \right) \left( \hat{\beta} - \beta \right)^\top \right], \text{ Bias ist null!}$$

- aber:

- $\hat{\beta}$  ist nur dann wohldefiniert, wenn  $X$  vollen Spaltenrang hat, d.h. wenn  $p \leq n$
- was tun, wenn  $p \gg n$ ?

# Regularisierung

- zurück zur Polynomregression:  
wenn Modellpolynom zu hohem Grad hat, werden die Koeffizienten  $\beta_k$  sehr groß  
⇒ Regularisierung
- $l^2$ -norm (*Ridge-regression*):

$$\hat{\beta}^{RR} = \operatorname{argmin}_{b \in \mathbb{R}^p} \left( \|Y - Xb\|^2 + \lambda \|b\|_{l^2}^2 \right)$$

⇒ kleinere Koeffizienten werden bevorzugt

- $l^0$ -norm (*subset selection*):

$$\hat{\beta}^{SS} = \operatorname{argmin}_{b \in \mathbb{R}^p} \left( \|Y - Xb\|^2 + \lambda \|b\|_{l^0} \right)$$

⇒ Koeffizienten werden oft auf 0 gesetzt (⇒ Variablenselektion!)

- $l^1$ -norm (*lasso = least absolute selection and shrinkage operator*):

$$\hat{\beta}^L = \operatorname{argmin}_{b \in \mathbb{R}^p} \left( \|Y - Xb\|^2 + \lambda \|b\|_{l^1} \right)$$

⇒ kleinere Koeffizienten werden bevorzugt *und* Koeffizienten werden oft auf 0 gesetzt

- **gegeben:**  $X_1, \dots, X_n \in \mathbb{R}^p$ ,  $p$  "groß", *keine* Modellannahme wie im linearen Modell
- **Ziel:** projiziere  $X_k$  auf affinen Unterraum von viel kleinerer Dimension  $q \ll p$
- **Idee:**
  - finde  $\mu \in \mathbb{R}^p$  und Orthogonalprojektion  $\Pi \in \mathbb{R}^{p \times p}$  mit Rang  $q$ , so dass  $\sum_{k=1}^n \|X_k - \mu - \Pi X_k\|^2$  minimal wird
  - oder: finde  $f : \mathbb{R}^q \rightarrow \mathbb{R}^p$  und  $v_1, \dots, v_n \in \mathbb{R}^q$ , so dass  $\sum_{k=1}^n \|X_k - f(v_k)\|^2$  minimal wird, wobei  $f(x) = Ax + \mu$ ,  $A \in \mathbb{R}^{p \times q}$ ,  $A^\top A = I_q$ ,  $\mu \in \mathbb{R}^p$
- durch Ableiten nach  $\mu$  und  $A$  findet man als Lösungen  $\mu = \bar{X} = \frac{1}{n} \sum_{k=1}^n X_k$ ,  $A = (w_1, \dots, w_q)$ ,  $v_k = A^\top (X_k - \mu)$ , wobei  $W = (w_1, \dots, w_p) \in \mathbb{R}^{p \times p}$  und  $X^\top X = W \Lambda W^\top$  die Eigenwertzerlegung von  $X^\top X$  ist (d.h.  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$ )

# PCA (= Hauptkomponentenanalyse)

- $w_k$  = Hauptkomponenten von  $X^T X$
- es gilt für  $f$  oben:  $\sum_{k=1}^n \|X_k - f(v_k)\|^2 = \sum_{k=q+1}^n \lambda_k$   
( $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$ )
- **Achtung:**
  - wenn  $p \gg n$  groß, dann ist die Eigenwertzerlegung von  $X^T X \in \mathbb{R}^{p \times p}$  aufwendig
  - $X^T X$  hat dann Rang  $\min(n, p) = n$
  - besser: Eigenwertzerlegung von  $XX^T = VDV^T \in \mathbb{R}^{n \times n}$
  - dann:  $XX^T v_i = d_i v_i$  für  $V = (v_1, \dots, v_n)$  und  $D = \text{diag}(d_1, \dots, d_n)$   
 $\Rightarrow (X^T X) X^T v_i = d_i X^T v_i \Rightarrow X^T v_i$  ist Eigenvektor von  $X^T X$  zum Eigenwert  $d_i$

- **gegeben:** Daten  $(X_1, Y_1), \dots, (X_n, Y_n)$ ,  $X_k \in \mathbb{R}^p$ ,  $Y_k \in \{0, 1\}$
- **gesucht:** “Klassifikator”  $C : \mathbb{R}^p \rightarrow \{0, 1\}$ , so dass für “neue” unabhängige Samples  $(X, Y)$   $C(X) \approx Y$
- diesmal sind  $X_k$  und  $Y_k$  zufällig
- Beispiel: SPAM-Klassifikation von Emails
  - 1 =SPAM, 0 =HAM
  - $X_k$  = Liste aller Worte in Email  $k$
  - Asymmetrie in SPAM vs. HAM
- Training vs. Testen:
  - teile Datensatz in zwei Teilmengen auf
  - auf dem Ersten *trainiere* den classifier
  - auf dem Zweiten *teste*

- ein “guter” classifier minimiert das *0-1-Risiko*

$$\mathbb{P}(C(X) \neq Y) = \mathbb{E}[|\mathbf{1}(C(X) \in \{1\}) - \mathbf{1}(Y \in \{1\})|]$$

- es gilt:

$$\begin{aligned}\mathbb{P}(C(X) \neq Y) &= 1 - \mathbb{P}(C(X) = Y) \\ &= 1 - \int \mathbb{P}(C(x) = Y | X = x) \mathbb{P}^X(dx) \\ &= 1 - \int \left( \mathbf{1}(C(x) = 1) \mathbb{P}(Y = 1 | X = x) \right. \\ &\quad \left. + \mathbf{1}(C(x) = 0) \mathbb{P}(Y = 0 | X = x) \right) \mathbb{P}^X(dx)\end{aligned}$$

⇒ maximiere punktweise den Integranden bezüglich  $x$ , d.h. wähle  $C(x) = k$ , wenn  $\mathbb{P}(Y = k | X = x)$  maximal

- schätze  $\mathbb{P}(Y = 1|X = x)$  aus den Daten
- **Satz von Bayes** (für diskrete Zufallsvariablen  $X$  und  $\mathbb{P}(X = x) \neq 0$ ):

$$\mathbb{P}(Y = 1|X = x) = \frac{\mathbb{P}(X = x|Y = 1)\mathbb{P}(Y = 1)}{\mathbb{P}(X = x)}$$

$$\Rightarrow \mathbb{P}(X = x|Y = 1) \approx \frac{\sum_{k=1}^n \mathbf{1}(X_k=x, Y_k=1)}{\sum_{k=1}^n \mathbf{1}(Y_k=1)}$$

- Beispiel: SPAM-Klassifikation
  - $X_k$  = Liste aller Worte in Email  $k$
  - Problem:  $X_k$  sehr hochdimensional

$\Rightarrow$  *naive Bayes*: Wir nehmen an, dass alle Features, bedingt auf  $Y$ , unabhängig sind, d.h.

$$\begin{aligned}\mathbb{P}(X = x|Y = 1) &= \mathbb{P}\left(X^{(1)} = x_1, \dots, X^{(p)} = x_p \mid Y = 1\right) \\ &= \prod_{k=1}^p \mathbb{P}\left(X^{(k)} = x_k \mid Y = 1\right)\end{aligned}$$

# Naive Bayes classifier

- schätze nun  $\mathbb{P}(X^{(k)} = x_k | Y = 1) \approx \frac{\sum_{j=1}^n \mathbf{1}(X_j^{(k)} = x_k, Y_j = 1)}{\sum_{j=1}^n \mathbf{1}(Y_j = 1)} := \Phi_{x_k|1}$ ,  
ähnlich  $\Phi_{x_k|0}$
- für  $x \in \mathbb{R}^p$ :
  - bestimme  $\Gamma_1 = \frac{1}{n} \sum_{j=1}^n \mathbf{1}(Y_j = 1)$ ,  $\Gamma_0 = 1 - \Gamma_1$
  - bestimme  $\Phi_{x_k|1}$ ,  $\Phi_{x_k|0}$
  - berechne  $\Phi_{x,1} := \prod_{k=1}^p \Phi_{x_k|1}$ ,  $\Phi_{x,0} := \prod_{k=1}^p \Phi_{x_k|0}$
  - definiere den classifier

$$C(x) = \mathbf{1}(\Phi_{x,1} \cdot \Gamma_1 > \Phi_{x,0} \cdot \Gamma_0)$$

- normalerweise verwende stattdessen

$$C(x) = \mathbf{1}\left(\sum_{k=1}^p \log(\Phi_{x_k|1}) + \log \Gamma_1 > \sum_{k=1}^p \log(\Phi_{x_k|0}) + \log \Gamma_0\right)$$