

# Wissensentdeckung in Datenbanken

## Überanpassung, Häufige Mengen

Nico Piatkowski und Uwe Ligges

Informatik—Künstliche Intelligenz  
Computergestützte Statistik  
Technische Universität Dortmund

09.05.2017

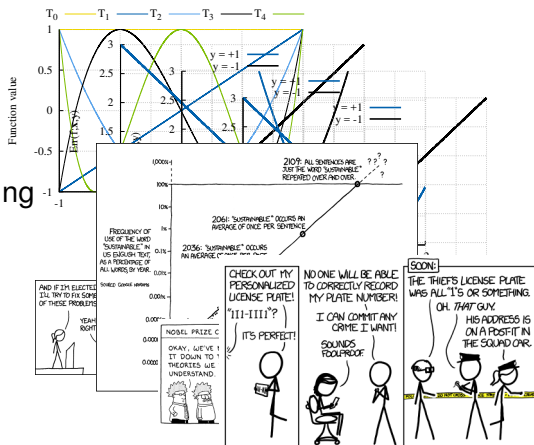
# Überblick

## Was bisher geschah...

- Modellklassen
- Verlustfunktionen
- Numerische Optimierung
- Regularisierung

## Heute

- Überanpassung
- Häufige Mengen



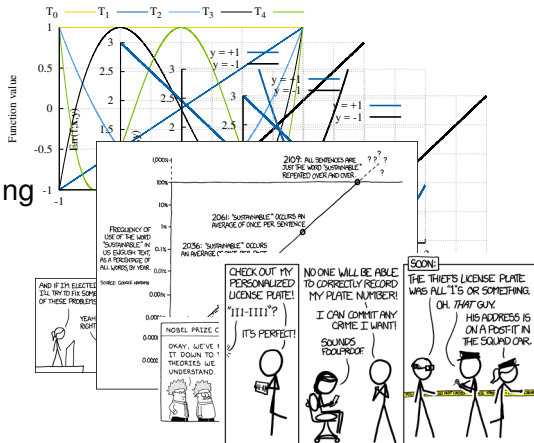
# Überblick

## Was bisher geschah...

- Modellklassen
- Verlustfunktionen
- Numerische Optimierung
- Regularisierung

## Heute

- Überanpassung
- Häufige Mengen





## Fehlertypen

- Lernen des Modells  $f_{\mathcal{D}}$  auf Basis von *Trainingsdaten*  $\mathcal{D}$
- Anwendung des Modells auf *Testdaten*  $\mathcal{T}$
- $\mathcal{D}$  ist Zufallsvariable  $\Rightarrow f_{\mathcal{D}}$  ist Zufallsvariable
- Das **gelernte** Modell soll auf **zufälligen** Testpunkten  $(X, Y)$  möglichst gut funktionieren, d.h.

$$\mathbb{E}[\ell(f_{\mathcal{D}}; (X, Y)) \mid \mathcal{D}]$$

soll minimal sein.

- Aber numerische Optimierung wählt  $f$  so, dass

$$\ell(f; \mathcal{D}) \propto \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \ell(f; (x, y)) = \hat{\mathbb{E}}[\ell(f; (X, Y))]$$

minimiert wird.

## Fehlertypen

- Lernen des Modells  $f_{\mathcal{D}}$  auf Basis von *Trainingsdaten*  $\mathcal{D}$
- Anwendung des Modells auf *Testdaten*  $\mathcal{T}$
- $\mathcal{D}$  ist Zufallsvariable  $\Rightarrow f_{\mathcal{D}}$  ist Zufallsvariable
- Das **gelernte** Modell soll auf **zufälligen** Testpunkten  $(\mathbf{X}, Y)$  möglichst gut funktionieren, d.h.

$$\mathbb{E}[\ell(f_{\mathcal{D}}; (\mathbf{X}, Y)) \mid \mathcal{D}]$$

soll minimal sein.

- Aber numerische Optimierung wählt  $f$  so, dass

$$\ell(f; \mathcal{D}) \propto \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \ell(f; (x, y)) = \hat{\mathbb{E}}[\ell(f; (\mathbf{X}, Y))]$$

minimiert wird.



## Fehlertypen

- Lernen des Modells  $f_{\mathcal{D}}$  auf Basis von *Trainingsdaten*  $\mathcal{D}$
- Anwendung des Modells auf *Testdaten*  $\mathcal{T}$
- $\mathcal{D}$  ist Zufallsvariable  $\Rightarrow f_{\mathcal{D}}$  ist Zufallsvariable
- Das **gelernte** Modell soll auf **zufälligen** Testpunkten  $(\mathbf{X}, Y)$  möglichst gut funktionieren, d.h.

$$\mathbb{E}[\ell(f_{\mathcal{D}}; (\mathbf{X}, Y)) \mid \mathcal{D}]$$

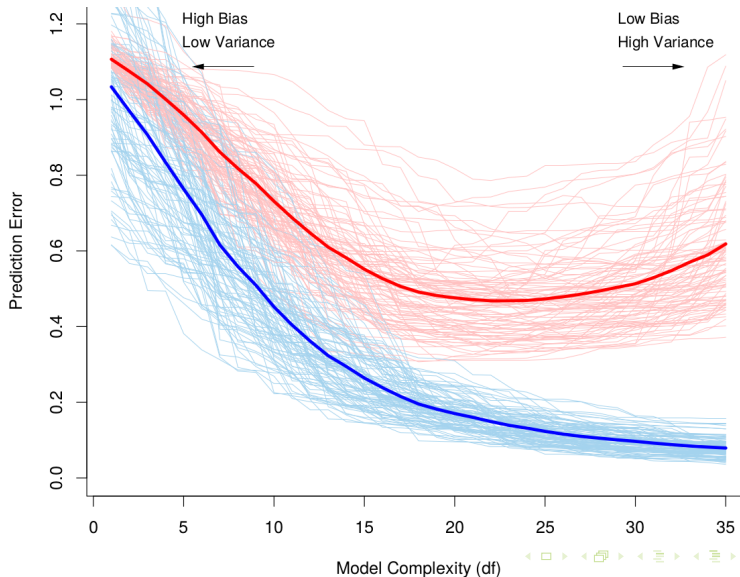
soll minimal sein.

- Aber numerische Optimierung wählt  $f$  so, dass

$$\ell(f; \mathcal{D}) \propto \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, y) \in \mathcal{D}} \ell(f; (\mathbf{x}, y)) = \hat{\mathbb{E}}[\ell(f; (\mathbf{X}, Y))]$$

minimiert wird.

Beispiel: Hellrot =  $\mathbb{E}[\ell(f_{\mathcal{D}}; (\mathbf{X}, Y)) \mid \mathcal{D}]$ , Hellblau =  $\hat{\mathbb{E}}[\ell(f; (\mathbf{X}, Y))]$





# Überanpassung

- Freiheitsgrade im linearen Modell:  $df_{\text{RSS}}(\beta) = d$
- Falls  $d \gg n$  kann sich das Modell “perfekt” an die Daten anpassen
- Modell lernt die Daten “auswendig”  $\equiv$  *Überanpassung*
- Optimum der Verlustfunktion  $\ell(f; \mathcal{D}) = \hat{\mathbb{E}}[\ell(f; (\mathbf{X}, Y))]$  liefert suboptimale Vorhersagen
- Hinzufügen von “Anpassungskosten” soll Überanpassung verhindern





# Überanpassung

- Freiheitsgrade im linearen Modell:  $df_{\text{RSS}}(\beta) = d$
- Falls  $d \gg n$  kann sich das Modell “perfekt” an die Daten anpassen
- Modell lernt die Daten “auswendig”  $\equiv$  *Überanpassung*
- Optimum der Verlustfunktion  $\ell(f; \mathcal{D}) = \hat{\mathbb{E}}[\ell(f; (\mathbf{X}, Y))]$  liefert suboptimale Vorhersagen
- Hinzufügen von “Anpassungskosten” soll Überanpassung verhindern



# Modellkomplexität

Allgemein für lineare Modelle mit  $\beta \in \mathbb{R}^d$  und  $y = f(\mathbf{x}) + \varepsilon$

$$df(f) = \frac{1}{\sigma_\varepsilon^2} \sum_{(\mathbf{x}, y) \in \mathcal{D}} \mathbb{C}[f(\mathbf{x}), y]$$

Idee: bestrafe Überanpassung mittels *Regularisierung*

$R: \mathcal{M} \rightarrow \mathbb{R}$

$$f_{\mathcal{D}} = \arg \min_{f \in \mathcal{M}} \ell(f; \mathcal{D}) + \lambda R(f)$$

Der Parameter  $\lambda > 0$  bestimmt den Einfluss der Regularisierung.

Oft:  $R(f) = \|f\|_q^q$



# Modellkomplexität

Allgemein für lineare Modelle mit  $\beta \in \mathbb{R}^d$  und  $y = f(\mathbf{x}) + \varepsilon$

$$df(f) = \frac{1}{\sigma_\varepsilon^2} \sum_{(\mathbf{x}, y) \in \mathcal{D}} \mathbb{C}[f(\mathbf{x}), y]$$

Idee: bestrafe Überanpassung mittels *Regularisierung*

$R: \mathcal{M} \rightarrow \mathbb{R}$

$$f_{\mathcal{D}} = \arg \min_{f \in \mathcal{M}} \ell(f; \mathcal{D}) + \lambda R(f)$$

Der Parameter  $\lambda > 0$  bestimmt den Einfluss der Regularisierung.

Oft:  $R(f) = \|f\|_q^q$



# Modellkomplexität

Allgemein für lineare Modelle mit  $\beta \in \mathbb{R}^d$  und  $y = f(\mathbf{x}) + \varepsilon$

$$df(f) = \frac{1}{\sigma_\varepsilon^2} \sum_{(\mathbf{x}, y) \in \mathcal{D}} \mathbb{C}[f(\mathbf{x}), y]$$

Idee: bestrafe Überanpassung mittels *Regularisierung*

$R: \mathcal{M} \rightarrow \mathbb{R}$

$$f_{\mathcal{D}} = \arg \min_{f \in \mathcal{M}} \ell(f; \mathcal{D}) + \lambda R(f)$$

Der Parameter  $\lambda > 0$  bestimmt den Einfluss der Regularisierung.

Oft:  $R(f) = \|f\|_q^q$



## Regularisierung kontrolliert Modellkomplexität

Freiheitsgrade für  $R_{l_2}(f) = \|\beta\|_2^2$  und  $R_{l_1}(f) = \|\beta\|_1$ :

$$\text{df}_{\text{RSS}+l_2}(\beta) = \text{trace}[\mathbf{D}(\mathbf{D}^\top \mathbf{D} + \lambda \mathbf{I})^{-1} \mathbf{D}^\top]$$

$$\text{df}_{\text{RSS}+l_1}(\beta) \approx \sum_{i=1}^d \mathbb{1}_{\beta_i \neq 0}$$

mit Datenmatrix  $\mathbf{D} = (\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N)^\top$  und Einheitsmatrix  $\mathbf{I}$ .

Beispiel:

- $f(\mathbf{x}) = \langle \beta, \mathbf{x} \rangle$
- $\ell(f; \mathcal{D}) = \text{RSS}(f; \mathcal{D}) = \sum_{(x,y) \in \mathcal{D}} (y - f(\mathbf{x}))^2$
- $\lambda = 0.1, \mathcal{D} = \{(0.25, 0.5)\}$



## Regularisierung kontrolliert Modellkomplexität

Freiheitsgrade für  $R_{l_2}(f) = \|\beta\|_2^2$  und  $R_{l_1}(f) = \|\beta\|_1$ :

$$df_{\text{RSS}+l_2}(\beta) = \text{trace}[\mathbf{D}(\mathbf{D}^\top \mathbf{D} + \lambda \mathbf{I})^{-1} \mathbf{D}^\top]$$

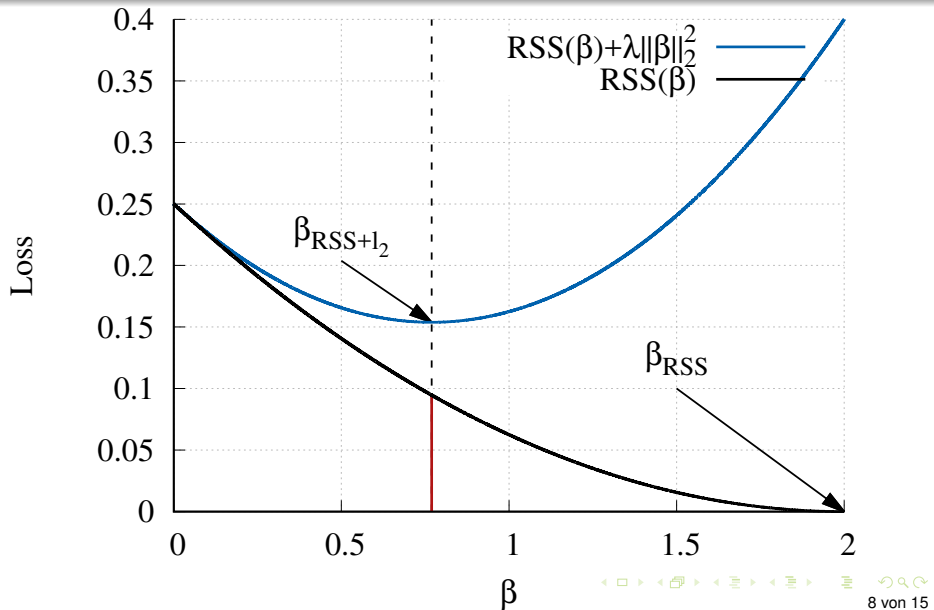
$$df_{\text{RSS}+l_1}(\beta) \approx \sum_{i=1}^d \mathbb{1}_{\beta_i \neq 0}$$

mit Datenmatrix  $\mathbf{D} = (\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N)^\top$  und Einheitsmatrix  $\mathbf{I}$ .

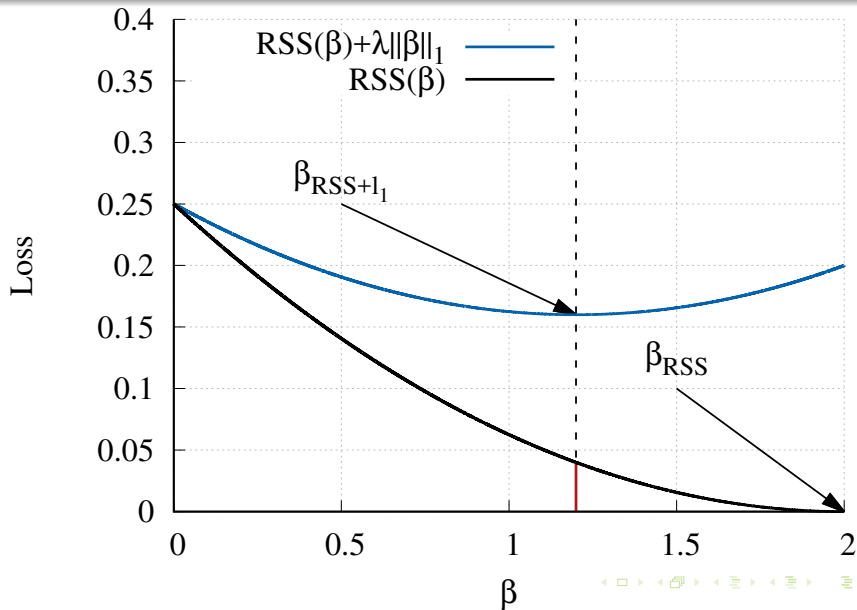
Beispiel:

- $f(\mathbf{x}) = \langle \beta, \mathbf{x} \rangle$
- $\ell(f; \mathcal{D}) = \text{RSS}(f; \mathcal{D}) = \sum_{(\mathbf{x}, y) \in \mathcal{D}} (y - f(\mathbf{x}))^2$
- $\lambda = 0.1, \mathcal{D} = \{(0.25, 0.5)\}$

## Beispiel: $l_2$ -Regularisierung



# Beispiel: $l_1$ -Regularisierung







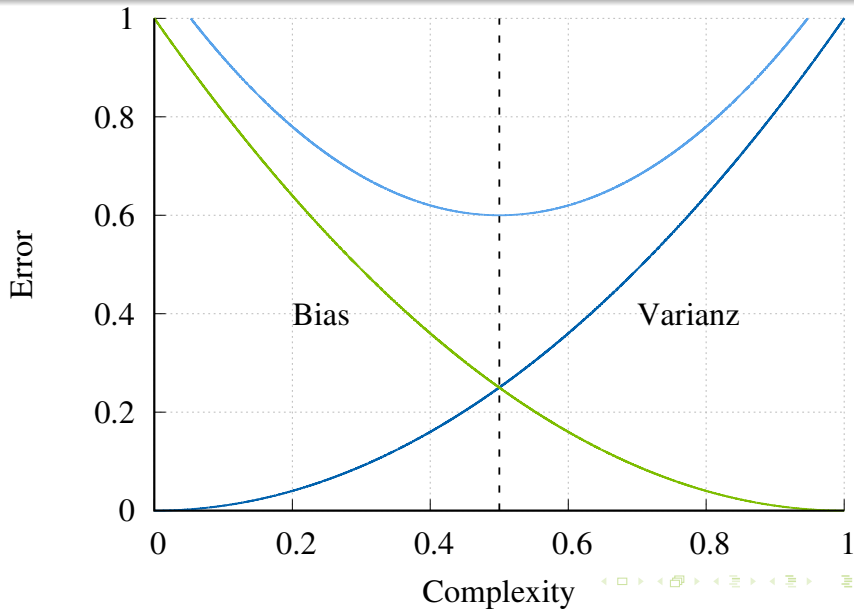
## Verzerrung und Varianz

Falls  $Y = f(\mathbf{x}) + \varepsilon$  und  $\mathbb{E}[\varepsilon] = 0$ , dann

$$\begin{aligned}\mathbb{E}[\text{RSS}(f_{\mathcal{D}}; (\mathbf{x}, Y))] &= \mathbb{E}[(Y - f_{\mathcal{D}}(\mathbf{x}))^2] \\ &= \mathbb{E}[Y^2 - 2Y f_{\mathcal{D}}(\mathbf{x}) + f_{\mathcal{D}}(\mathbf{x})^2] \\ &= \mathbb{E}[(f(\mathbf{x}) + \varepsilon)^2] - \mathbb{E}[2Y f_{\mathcal{D}}(\mathbf{x})] + \mathbb{E}[f_{\mathcal{D}}(\mathbf{x})^2] \\ &= \mathbb{E}[f(\mathbf{x})^2 + 2f(\mathbf{x})\varepsilon + \varepsilon^2] - \mathbb{E}[2Y f_{\mathcal{D}}(\mathbf{x})] + \mathbb{E}[f_{\mathcal{D}}(\mathbf{x})^2] \\ &= f(\mathbf{x})^2 + \sigma_{\varepsilon}^2 - \mathbb{E}[2Y f_{\mathcal{D}}(\mathbf{x})] + \mathbb{V}[f_{\mathcal{D}}(\mathbf{x})] + \mathbb{E}[f_{\mathcal{D}}(\mathbf{x})]^2 \\ &= \sigma_{\varepsilon}^2 + \mathbb{V}[f_{\mathcal{D}}(\mathbf{x})] + \mathbb{B}[f_{\mathcal{D}}(\mathbf{x})]^2\end{aligned}$$

mit  $\mathbb{B}[f_{\mathcal{D}}(\mathbf{x})] = f(\mathbf{x}) - \mathbb{E}[f_{\mathcal{D}}(\mathbf{x})]$  und alle Erwartungswerte sind bedingt auf  $\mathbf{x}$ .

## Verzerrung und Varianz (II)





# Datenbanken und Häufige Mengen



# Datenbanken und SQL

Jetzt:  $\mathcal{M} \subset \mathcal{D}$  bzw.  $\mathcal{M} \subset \mathcal{X}$

- Relationale Datenbanken  $\equiv$  Menge von Tabellen
- Relationales Datenbankmanagementsystem erlaubt Anfrage und Manipulation von Daten mittels Structured Query Language (SQL)
- Beispielhafte Anfragen  $Q$ :
  - `SELECT DISTINCT x1,x2 FROM data WHERE ...`
  - `SELECT * FROM ... ORDER BY x`
  - `SELECT AVG(x) FROM ... GROUP BY y`
  - `SELECT MIN(x) FROM ... GROUP BY y`
  - `SELECT MAX(x) FROM ... GROUP BY y`
  - `SELECT COUNT(x) AS c FROM ... HAVING c>10`



# Datenbanken und SQL

Jetzt:  $\mathcal{M} \subset \mathcal{D}$  bzw.  $\mathcal{M} \subset \mathcal{X}$

- Relationale Datenbanken  $\equiv$  Menge von Tabellen
- Relationales Datenbankmanagementsystem erlaubt Anfrage und Manipulation von Daten mittels Structured Query Language (SQL)
- Beispielhafte Anfragen  $Q$ :
  - `SELECT DISTINCT x1,x2 FROM data WHERE ...`
  - `SELECT * FROM ... ORDER BY x`
  - `SELECT AVG(x) FROM ... GROUP BY y`
  - `SELECT MIN(x) FROM ... GROUP BY y`
  - `SELECT MAX(x) FROM ... GROUP BY y`
  - `SELECT COUNT(x) AS c FROM ... HAVING c>10`



# Anfragen und Mengen

Resultat einer Datenbankanfrage  $Q$  ist eine neue Tabelle  $D$

- Annahme: Datenbank besteht aus Binärdaten, oder wird mittels SQL Anfragen konvertiert
- Die Einträge von  $D$  heißen dann *Transaktionen*
- Transaktion  $t \in D$  entspricht Indikatorvektor einer Menge

001010011101000110010  $\equiv \{x_3, x_5, x_8, x_9, x_{10}, x_{12}, x_{16}, x_{17}, x_{20}\}$

- Die Elemente  $x_i$  der Menge nennt man auch *Items*





# Anfragen und Mengen

Resultat einer Datenbankanfrage  $Q$  ist eine neue Tabelle  $D$

- Annahme: Datenbank besteht aus Binärdaten, oder wird mittels SQL Anfragen konvertiert
- Die Einträge von  $D$  heißen dann *Transaktionen*
- Transaktion  $t \in D$  entspricht Indikatorvektor einer Menge

$$001010011101000110010 \equiv \{x_3, x_5, x_8, x_9, x_{10}, x_{12}, x_{16}, x_{17}, x_{20}\}$$

- Die Elemente  $x_i$  der Menge nennt man auch *Items*



## Ausblick: Häufige Mengen und Apriori

### Frequent Itemset Mining:

- Bestimme alle Mengen von Items die in mindestens  $s$ -prozent aller Transaktionen vorkommen
  - Solche Menge heißen *häufig*
- $2^n$  mögliche häufige Mengen!
- Aber: Eine Menge kann nur dann häufig sein, wenn alle ihre Teilmengen häufig sind
- *Apriori*: Bottom-Up Algorithmus über Teilmengenverband zur Berechnung aller häufigen Mengen

