

# Wissensentdeckung in Datenbanken

## Organisation und Überblick

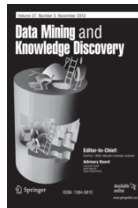
Nico Piatkowski und Uwe Ligges

Informatik—Künstliche Intelligenz  
Computergestützte Statistik  
Technische Universität Dortmund

18.04.2017

# Fakten

- Team
  - Vorlesung: Uwe Ligges, Nico Piatkowski
  - Übung: Sarah Schnackenberg, Sebastian Buschjäger
- 28 Termine (13/13/2)
  - Dienstags — 10:15 Uhr bis 12:00 Uhr
  - Donnerstags — 14:15 Uhr bis 16:00 Uhr
- Feiertage
  - 25.05. (Christi Himmelfahrt)
  - 15.06. (Fronleichnam)



## Inhalt

- 18.4. L,P Übersicht, Einführung
- 20.4. L Statistik 1
- 25.4. L Statistik 2
- 27.4. L Stichproben, Versuchsplanung, Datenvorverarbeitung
- 02.5. Optimierung, Modellklassen, Lineares Modell,  
P Bias-Varianz, Overfitting 1
- 04.5. Optimierung, Modellklassen, Lineares Modell,  
P Bias-Varianz, Overfitting 2
- 09.5. P Data Cube, Frequent sets, Apriori, FPgrowth 1
- 11.5. P Data Cube, Frequent sets, Apriori, FPgrowth 2
- 16.5. L Einführung Klassifikation, BayesRegel,  
Logistische Regression
- 18.5. L kNN, Ähnlichkeitsmaße, Modellselektion
- 23.5. L Resampling, Klassifikationsbeurteilung 2

# Inhalt

- 30.5. P SVM
- 01.6. L Diskriminanzanalyse (LDA) 1
- 06.6. L Diskriminanzanalyse (LDA, QDA, RDA) 2
- 08.6. L von Entscheidungsbäumen (CART) zu Wäldern
- 13.6. L Ensemble Methoden (Bagging, Boosting)
- 20.6. L Stetige Modelle Hauptkomponentenanalyse)
- 22.6. P Graphische Modelle 1
- 27.6. P Graphische Modelle 2
- 29.6. P Nicht-glatte und stochastische Optimierung
- 04.7. P Merkmalsselektion, Struktur Lernen, Regularisierung



## Inhalt

- 06.7. P Clustering, k-Means, Gaussian Mixture,  
Latent Dirichlet Allocation 1
- 11.7. P Clustering, k-Means, Gaussian Mixture,  
Latent Dirichlet Allocation 2
- 13.7. L Hierarchisches Clustern; Zeitreihen 1
- 18.7. L Zeitreihen 2
- 20.7. P Künstliche Neuronale Netze 1
- 25.7. P Künstliche Neuronale Netze 2
- 27.7. L,P Zusammenfassung; Rückblick



## Fragen und Kommentare

### **Unsere Bitte: Mitdenken, kommentieren und Fragen stellen!**

Die unterschiedliche Terminologie mag zunächst verwirren:  
*Merkmal, Feature, Variable, Parameter* im Sprachgebrauch in  
Informatik, Mathematik und Statistik.



## Beispiele für KDD / Data Mining

- Telekommunikationstechnik
  - riesige Datenbestände, Prozesskontrolle, Zeitreihenanalyse, Zuverlässigkeitsanalysen
- CallCenter
  - Warteschlangen, Klassifikation des Problems mit wenig Information
- Kundenkartenanalyse
  - riesige Datenbestände (und große Anzahl Variablen), Vorhersage von Kaufverhalten und Warengruppenzuordnung
- Genomik
  - riesige Datenbestände (und große Anzahl Variablen), viel Rauschen, Klassifikation zur Wissensentdeckung

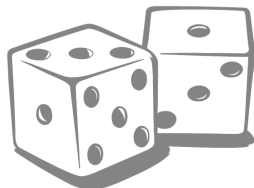
# Statistische Grundlagen

- Maß, Dichte, Erwartungswert, Zufallsvariable, ...

$$\mathbb{P}(\mathbf{X} \in S) = \int_S \frac{d\mathbb{P} \circ \mathbf{X}^{-1}}{d\nu} d\nu = \int_S p d\nu$$

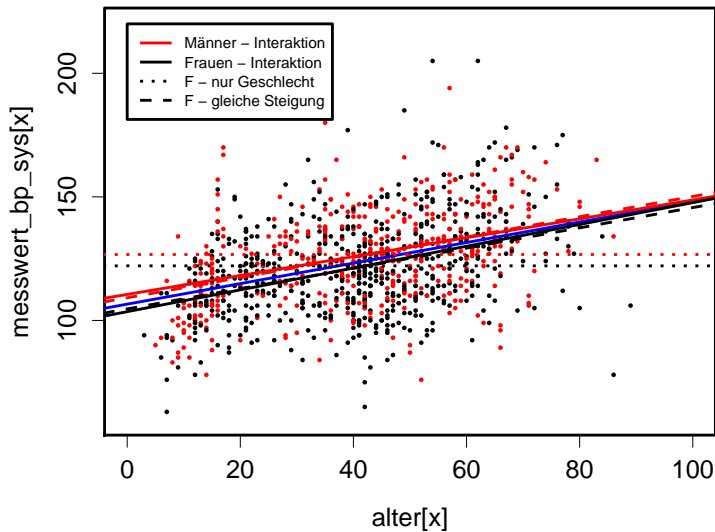
$$p(\mathbf{X} = \mathbf{x} \mid \mathbf{Y} = \mathbf{y}) = \frac{p(\mathbf{X} = \mathbf{x}, \mathbf{Y} = \mathbf{y})}{p_{\mathbf{Y}}(\mathbf{y})}$$

$$\mathbb{E}[\phi(\mathbf{X})] = \int_{\mathcal{X}} \phi d\nu$$





# Regression



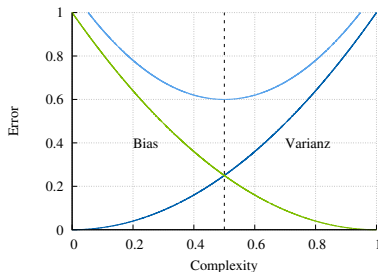
# Maschinelles Lernen

- Optimierung, Modell, Verlustfunktion, Overfitting, ...

$$\mathbb{E}[(y - \hat{f}(\mathbf{x}))^2] = \mathbb{B}[\hat{f}(\mathbf{x})^2] + \mathbb{V}[\hat{f}(\mathbf{x})] + \sigma^2$$

$$\beta^{t+1} = \beta^t - \eta_t \nabla \ell(\beta^t; \mathcal{D})$$

$$\hat{f} = \min_{f \in \mathcal{F}} \ell(f; \mathcal{D})$$



# Frequent Set Mining

- Datenbanken, Häufige Mengen, A-Priori, FP-Trees, ...

SELECT \* FROM transactions WHERE ...

$$p(A, B) \leq \min\{p(A), p(B)\}$$

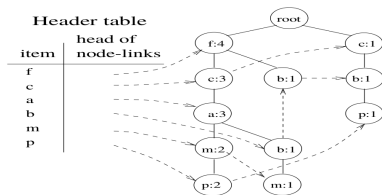


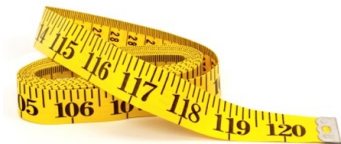
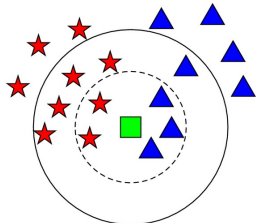
Figure 1: The FP-tree in Example 1.

# Klassifikation

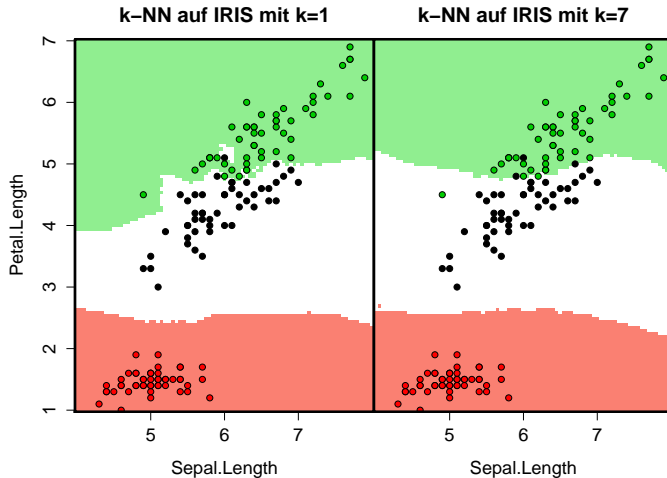
- Logistische Regression,  $k$ -NN, Distanzmaße, .....

$$p(y \mid \mathbf{X} = \mathbf{x}) = \frac{1}{1 + \exp(-(\beta_0 + \langle \beta, \mathbf{x} \rangle))}$$

$$\text{dist}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2 = \sqrt{\sum_{i=1}^d (\mathbf{x}_i - \mathbf{y}_i)^2}$$



# Klassifikation



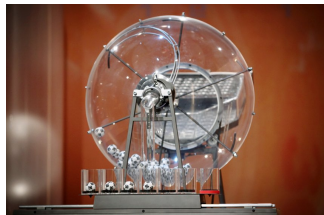
# Modelle und Sampling

- Modellselektion, Sampling, Klassifikationsgüte, ...

$$\text{BIC} = d \log(N) - 2\ell(\beta^*; \mathcal{D})$$

$$F_1 = \frac{2 \times \text{PREC} \times \text{REC}}{\text{PREC} + \text{REC}}$$

$$\mathbf{x} \sim \mathbb{P}$$



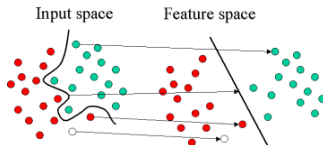
# Stützvektormethode

- Hyperebene, Hinge-Loss, Merkmalsraum, Kernel, ...

$$f(\mathbf{x}) = b + \langle \boldsymbol{\beta}, \phi(\mathbf{x}) \rangle$$

$$K_{\text{Gauss}}(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle = \exp\left(-\frac{1}{\gamma} \|\mathbf{x} - \mathbf{y}\|_2^2\right)$$

$$\ell_{\text{Hinge}}(y^*, \mathbf{x}, \boldsymbol{\beta}) = \max\{0, 1 - y^* \times (b + \langle \boldsymbol{\beta}, \phi(\mathbf{x}) \rangle)\}$$





# Diskriminanzanalyse

- Linear, Quadratic, Regularized, ...

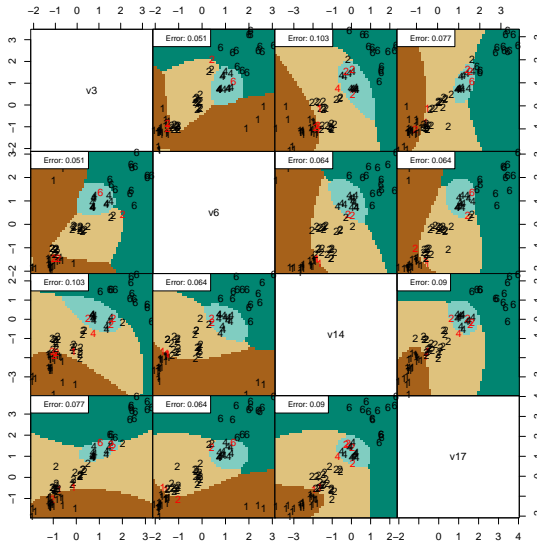
$$h_i^L(x) := (\Sigma^{-1} \mu_i)' x - 0.5 \mu_i' \Sigma^{-1} \mu_i + \ln(\pi_i)$$

$$h_i^Q(x) := -0.5(x - \mu_i)' \Sigma_i^{-1} (x - \mu_i) + \ln(\pi_i) - 0.5 \ln(\det(\Sigma_i))$$

$$\hat{\Sigma}_i(\delta, \lambda) := (1 - \lambda) \hat{\Sigma}_i(\delta) + \frac{\lambda \cdot \text{tr}[\hat{\Sigma}_i(\delta)]}{p} I$$

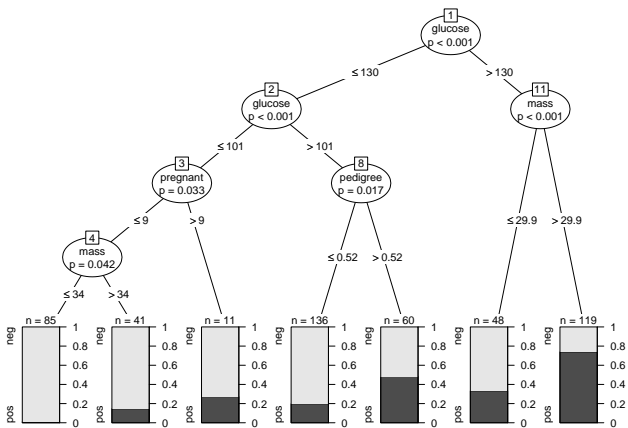


# Diskriminanzanalyse



# Ensembles

- Bäume, Bagging, Boosting, Radom Forests ...



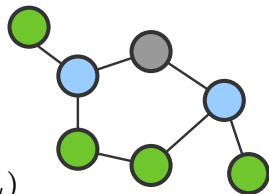
# Graphische Modelle

- Exponentialfamilien, Markov Random Fields, Belief-Propagation, ...

$$p_{\theta}(\mathbf{x}) = \exp(\langle \theta, \phi(\mathbf{x}) \rangle - A(\theta))$$

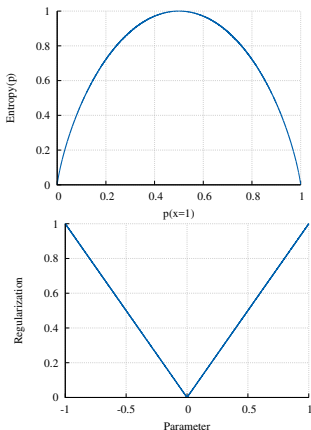
$$p_{\theta}(\mathbf{x}) = \frac{1}{Z} \prod_{v \in V} \psi_v(\mathbf{x}_v) \prod_{(v,u) \in E} \psi_{vu}(\mathbf{x}_v, \mathbf{x}_u)$$

$$\mathbf{m}_{u \rightarrow v}(x) = \sum_{y \in \mathcal{X}_u} \psi_u(y) \psi_{uv}(x, y) \prod_{w \in \mathcal{N}_u \setminus \{v\}} \mathbf{m}_{w \rightarrow u}(y)$$



# Strukturlernen

- Merkmalsselektion, Regularisierung, nicht-glatte Optimierung, ...



$$I(v, u) = \sum_{x,y} p_{vu}(x, y) \log \frac{p_{vu}(x, y)}{p_v(x)p_u(y)}$$

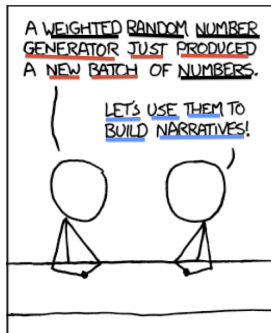
$$\min_{\beta \in \mathbb{R}^d} \ell(\beta; \mathcal{D}) + \lambda R(\beta)$$

$$\text{prox}_{\lambda \|\cdot\|_1}(\beta)_i = \begin{cases} \beta_i - \lambda & , \beta_i \geq +\lambda \\ 0 & , |\beta_i| < \lambda \\ \beta_i + \lambda & , \beta_i \leq -\lambda \end{cases}$$

# Clusteranalyse

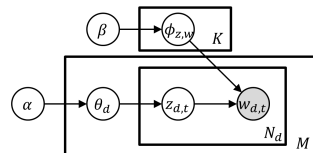
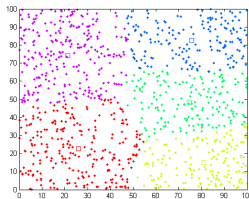
- $k$ -Means, Gaussian Mixture, Latent Dirichlet Allocation, ...

$$\min_{C \subseteq \mathbb{R}^d, |C|=k} \sum_{\mathbf{x} \in \mathcal{D}} \min_{\mathbf{c} \in C} \|\mathbf{x} - \mathbf{c}\|_2$$



ALL SPORTS COMMENTARY

$$p(\mathbf{w}_i^d | \mathbf{z}_i^d, \phi) p(\mathbf{z}_i^d | \theta^d) p(\theta^d | \alpha) p(\phi | \beta)$$





# Zeitreihenanalyse

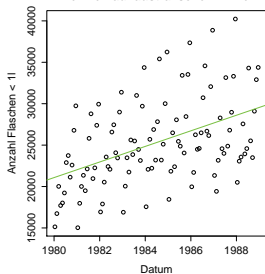
- Glättung, von Autoregression bis SARIMA, Schwingungen

$$y_t = \beta_1 + \beta_2 y_{t-1} + \dots + \beta_{p+1} y_{t-p} + \epsilon_t - \gamma_1 \epsilon_{t-1} - \dots - \gamma_q \epsilon_{t-q}$$

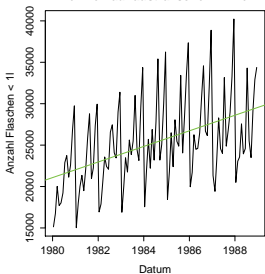
$$y_t = \beta_1 + \sum_{k=1}^K (\beta_{2k} \cos(2\pi \lambda_k t) + \beta_{2k+1} \sin(2\pi \lambda_k t)) + \epsilon_t$$

# Zeitreihenanalyse

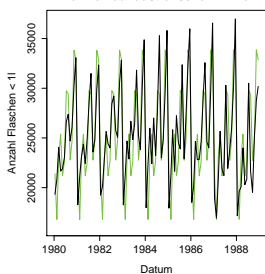
Weinverkauf australischer Winzer



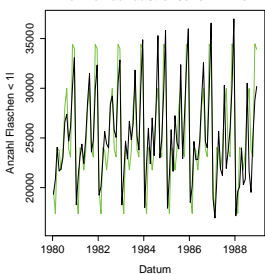
Weinverkauf australischer Winzer



Weinverkauf australischer Winzer



Weinverkauf australischer Winzer



# Künstliche Neuronale Netze

- Versteckte Schichten, Back-Propagation, LRU, Embeddings, ...

$$\frac{\partial \text{Err}(\mathbf{x}, \mathbf{y}^*, \boldsymbol{\beta})}{\partial \beta_{l,i,j}} = \frac{\partial \text{Err}(\mathbf{x}, \mathbf{y}^*, \boldsymbol{\beta})}{\partial \text{Out}^\top} \frac{\partial \text{Out}}{\partial \text{Net}_n^\top} \cdots \frac{\partial \text{Net}_{l+1}}{\partial \text{Net}_l^\top} \left( \frac{\partial \text{Net}_l}{\partial \beta_{l,i,j}} \right)^\top$$



Sample of cats & dogs images from Kaggle Dataset

