

# Wissensentdeckung in Datenbanken

## Strukturlernen, Merkmalsauswahl

Nico Piatkowski und Uwe Ligges

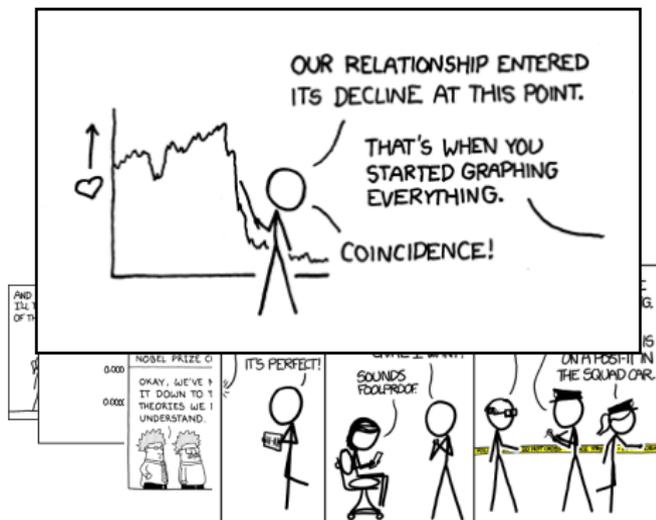
Informatik—Künstliche Intelligenz  
Computergestützte Statistik  
Technische Universität Dortmund

27.06.2017

# Überblick

## Was bisher geschah...

- Modellklassen
- Verlustfunktionen
- Numerische Optimierung
- Regularisierung
- Überanpassung
- SQL, Häufige Mengen
- SVM, xDA, Bäume, ...
- Graphische Modelle



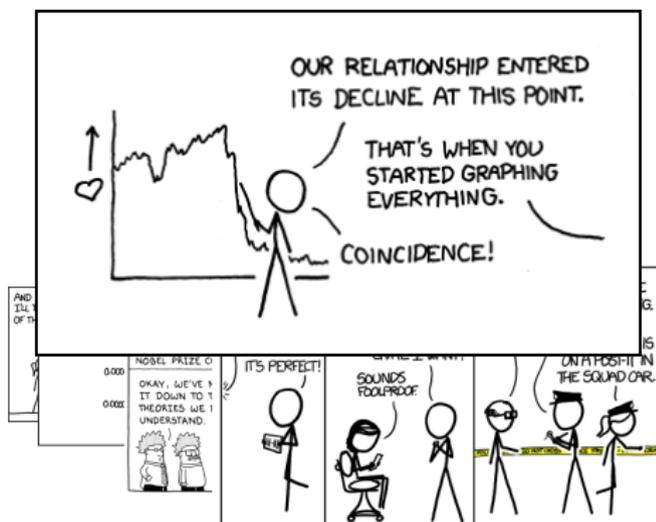
## Heute

- Strukturlernen, Merkmalsauswahl

# Überblick

## Was bisher geschah...

- Modellklassen
- Verlustfunktionen
- Numerische Optimierung
- Regularisierung
- Überanpassung
- SQL, Häufige Mengen
- SVM, xDA, Bäume, ...
- Graphische Modelle



## Heute

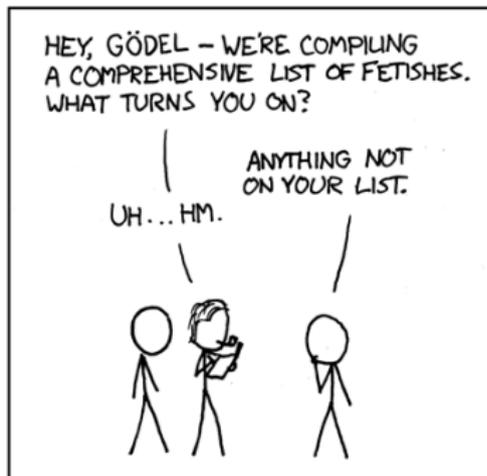
- Strukturlernen, Merkmalsauswahl

# Überblick

- Wiederholung: Belief Propagation, Gibbs Sampling
- Strukturlernen
  - Chow-Liu Bäumen
  - Regularisierung
- Merkmalsauswahl
  - Forward-Selection
  - Backward-Selection
  - Regularisierung

AUTHOR KATHARINE GATES RECENTLY ATTEMPTED TO MAKE A CHART OF ALL SEXUAL FETISHES.

LITTLE DID SHE KNOW THAT RUSSELL AND WHITEHEAD HAD ALREADY FAILED AT THIS SAME TASK.



## Belief Propagation

**Ziel:** Berechnung von  $\mathbb{E}_{\beta}[\phi(\mathbf{X})_{C=\mathbf{y}}] = \mathbb{P}(\mathbf{X}_C = \mathbf{y})$

$$Z(\beta) = \sum_{\mathbf{x}_1 \in \mathcal{X}_1} \sum_{\mathbf{x}_2 \in \mathcal{X}_2} \cdots \sum_{\mathbf{x}_n \in \mathcal{X}_n} \prod_{uv \in E} \underbrace{\exp(\langle \beta_{uv}, \phi_{uv}(\mathbf{x}_{uv}) \rangle)}_{\psi_{uv}(\mathbf{x}_{uv})}$$

Distributivität ausnutzen...

$$m_{v \rightarrow u}(x) = \sum_{y \in \mathcal{X}_v} \psi_{uv}(yx) \prod_{w \in \mathcal{N}(v) \setminus \{u\}} m_{w \rightarrow v}(y)$$

$$Z(\beta) = \sum_{y \in \mathcal{X}_v} \prod_{w \in \mathcal{N}(v)} m_{w \rightarrow v}(y)$$

$$\mathbb{P}(\mathbf{X}_{uv} = xy) = \frac{\psi_{uv}(yx)}{Z(\beta)} \prod_{w \in \mathcal{N}(v) \setminus \{u\}} m_{w \rightarrow v}(y) \prod_{w \in \mathcal{N}(u) \setminus \{v\}} m_{w \rightarrow u}(y)$$

## Belief Propagation

**Ziel:** Berechnung von  $\mathbb{E}_{\beta}[\phi(\mathbf{X})_{C=\mathbf{y}}] = \mathbb{P}(\mathbf{X}_C = \mathbf{y})$

$$Z(\beta) = \sum_{\mathbf{x}_1 \in \mathcal{X}_1} \sum_{\mathbf{x}_2 \in \mathcal{X}_2} \cdots \sum_{\mathbf{x}_n \in \mathcal{X}_n} \prod_{uv \in E} \underbrace{\exp(\langle \beta_{uv}, \phi_{uv}(\mathbf{x}_{uv}) \rangle)}_{\psi_{uv}(\mathbf{x}_{uv})}$$

Distributivität ausnutzen...

$$m_{v \rightarrow u}(x) = \sum_{y \in \mathcal{X}_v} \psi_{uv}(yx) \prod_{w \in \mathcal{N}(v) \setminus \{u\}} m_{w \rightarrow v}(y)$$

$$Z(\beta) = \sum_{y \in \mathcal{X}_v} \prod_{w \in \mathcal{N}(v)} m_{w \rightarrow v}(y)$$

$$\mathbb{P}(\mathbf{X}_{uv} = xy) = \frac{\psi_{uv}(yx)}{Z(\beta)} \prod_{w \in \mathcal{N}(v) \setminus \{u\}} m_{w \rightarrow v}(y) \prod_{w \in \mathcal{N}(u) \setminus \{v\}} m_{w \rightarrow u}(y)$$

## Belief Propagation

**Ziel:** Berechnung von  $\mathbb{E}_{\beta}[\phi(\mathbf{X})_{C=\mathbf{y}}] = \mathbb{P}(\mathbf{X}_C = \mathbf{y})$

$$Z(\beta) = \sum_{\mathbf{x}_1 \in \mathcal{X}_1} \sum_{\mathbf{x}_2 \in \mathcal{X}_2} \cdots \sum_{\mathbf{x}_n \in \mathcal{X}_n} \prod_{uv \in E} \underbrace{\exp(\langle \beta_{uv}, \phi_{uv}(\mathbf{x}_{uv}) \rangle)}_{\psi_{uv}(\mathbf{x}_{uv})}$$

Distributivität ausnutzen...

$$m_{v \rightarrow u}(x) = \sum_{y \in \mathcal{X}_v} \psi_{uv}(yx) \prod_{w \in \mathcal{N}(v) \setminus \{u\}} m_{w \rightarrow v}(y)$$

$$Z(\beta) = \sum_{y \in \mathcal{X}_v} \prod_{w \in \mathcal{N}(v)} m_{w \rightarrow v}(y)$$

$$\mathbb{P}(\mathbf{X}_{uv} = xy) = \frac{\psi_{uv}(yx)}{Z(\beta)} \prod_{w \in \mathcal{N}(v) \setminus \{u\}} m_{w \rightarrow v}(y) \prod_{w \in \mathcal{N}(u) \setminus \{v\}} m_{w \rightarrow u}(y)$$

## Gibbs Sampling: Algorithmus

- 1 Erzeuge  $x$  zufällig Gleichverteilt (das entspricht **NICHT**  $\mathbb{P}_\beta$ !)
- 2 Besuche jeden Knoten  $v \in V$  und weise gemäß  $\mathbb{P}_v(x \mid \mathbf{x}_{\mathcal{N}(v)})$  neuen Wert zu
- 3 Wiederhole Schritt 2 so oft wie möglich

Man kann zeigen: Nach endlicher Anzahl von Schritten ist  $x$  ein echtes Sample aus  $\mathbb{P}_\beta$ !

**Dann:** Nutze den Algorithmus um “viele” (so viele wie möglich) Samples zu erzeugen und berechne  $\mathbb{P}(X_{uv} = xy)$  (für alle Kanten  $\{v, u\} \in E$ ) durch “abzählen”



## Gibbs Sampling: Algorithmus

- 1 Erzeuge  $x$  zufällig Gleichverteilt (das entspricht **NICHT**  $\mathbb{P}_\beta$ !)
- 2 Besuche jeden Knoten  $v \in V$  und weise gemäß  $\mathbb{P}_v(x \mid \mathbf{x}_{\mathcal{N}(v)})$  neuen Wert zu
- 3 Wiederhole Schritt 2 so oft wie möglich

Man kann zeigen: Nach endlicher Anzahl von Schritten ist  $x$  ein echtes Sample aus  $\mathbb{P}_\beta$ !

**Dann:** Nutze den Algorithmus um “viele” (so viele wie möglich) Samples zu erzeugen und berechne  $\mathbb{P}(\mathbf{X}_{uv} = xy)$  (für alle Kanten  $\{v, u\} \in E$ ) durch “abzählen”



## Chow-Liu Bäume

Minimierung der Distanz zwischen optimalem Graph und  
“bestem” Baum  $T$

Hier: Distanz gemessen durch Kullback-Leiber Divergenz

$$\text{KL}(\mathbb{P}^*, \mathbb{P}_T) = \sum_{\mathbf{x} \in \mathcal{X}} \mathbb{P}^*(\mathbf{x}) \log \frac{\mathbb{P}^*(\mathbf{x})}{\mathbb{P}_T(\mathbf{x})}$$

Kann umgeformt werden zu

$$\text{KL}(\mathbb{P}^*, \mathbb{P}_T) = -\mathcal{H}(\mathbb{P}^*) + \sum_{v \in V} \mathcal{H}(\mathbb{P}_v^*) - \underbrace{\sum_{vu \in E(T)} I(\mathbf{X}_v, \mathbf{X}_u)}_{\text{Maximaler Spannbaum!}}$$

mit  $I(\mathbf{X}_v, \mathbf{X}_u) = \text{KL}(\mathbb{P}_{vu}, \mathbb{P}_v \mathbb{P}_u)$  (Allgemeines Maß für  
Unabhängigkeit!)



## Chow-Liu Bäume

Minimierung der Distanz zwischen optimalem Graph und  
“bestem” Baum  $T$

Hier: Distanz gemessen durch Kullback-Leiber Divergenz

$$\text{KL}(\mathbb{P}^*, \mathbb{P}_T) = \sum_{\mathbf{x} \in \mathcal{X}} \mathbb{P}^*(\mathbf{x}) \log \frac{\mathbb{P}^*(\mathbf{x})}{\mathbb{P}_T(\mathbf{x})}$$

Kann umgeformt werden zu

$$\text{KL}(\mathbb{P}^*, \mathbb{P}_T) = -\mathcal{H}(\mathbb{P}^*) + \sum_{v \in V} \mathcal{H}(\mathbb{P}_v^*) - \underbrace{\sum_{vu \in E(T)} I(\mathbf{X}_v, \mathbf{X}_u)}_{\text{Maximaler Spannbaum!}}$$

mit  $I(\mathbf{X}_v, \mathbf{X}_u) = \text{KL}(\mathbb{P}_{vu}, \mathbb{P}_v \mathbb{P}_u)$  (Allgemeines Maß für  
Unabhängigkeit!)

## Graphen mittels Regularisierung

**Baobachtung:** Sind ist kompletter Parametervektor  $\beta_{vu}$  einer Kante = 0, so hat diese Kante keinen Einfluss auf  $\mathbb{P}(\mathbf{X} = \mathbf{x})!$

**Idee:** Minimiere  $\ell(\beta; \mathcal{D}) + \lambda \|\beta\|_1$

$\|\cdot\|_1$  nicht differenzierbar!! → Proximaler Gradientenabstieg (nächste Woche)

**Spoiler:**

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