

Wissensentdeckung in Datenbanken

Probabilistische Graphische Modelle II

Nico Piatkowski und Uwe Ligges

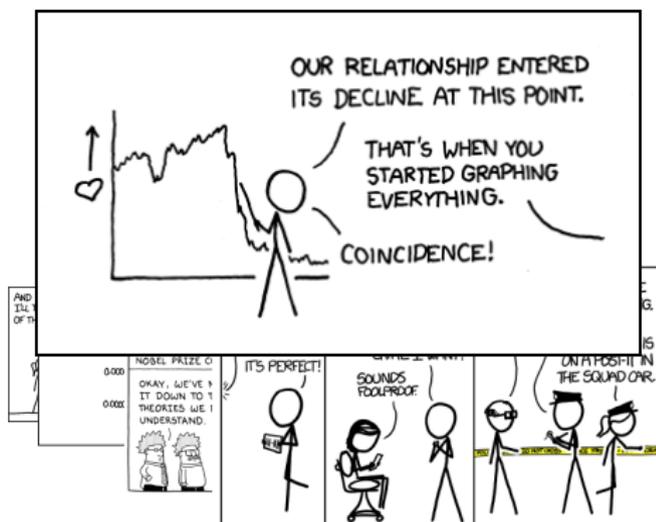
Informatik—Künstliche Intelligenz
Computergestützte Statistik
Technische Universität Dortmund

27.06.2017

Überblick

Was bisher geschah...

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



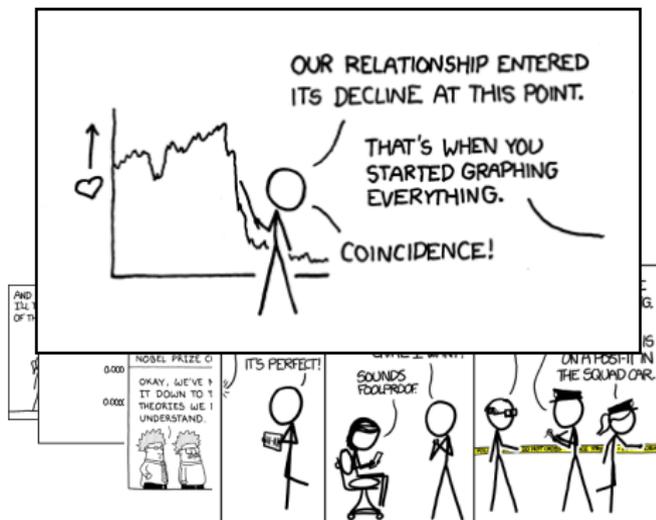
Heute

- Graphische Modelle—Theorie und Algorithmen

Überblick

Was bisher geschah...

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



Heute

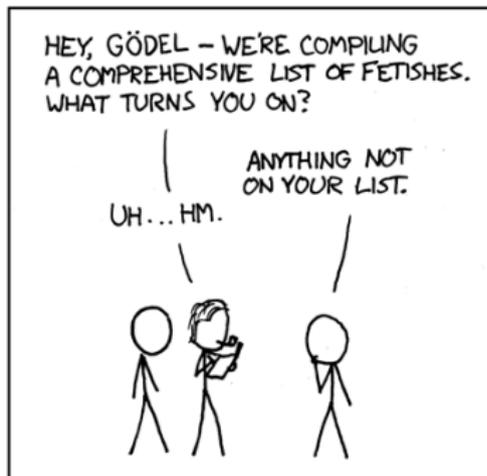
- Graphische Modelle—Theorie und Algorithmen

Überblick

- Suffiziente Statistiken
- Maximum-Entropie
- Gradient
- Randverteilung
- Belief Propagation
- Gibbs Sampling

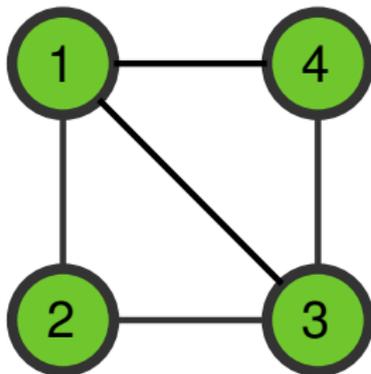
AUTHOR KATHARINE GATES RECENTLY ATTEMPTED TO MAKE A CHART OF ALL SEXUAL FETISHES.

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



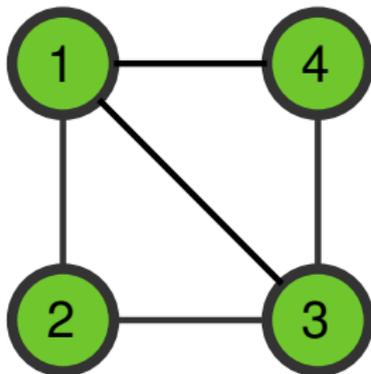
Graph

- $G = (V, E)$ mit Knotenmenge V und Kantenmenge E
- Hier: $V = \{1, 2, 3, 4\}$, $E = \{\{1, 2\}, \{1, 3\}, \{1, 4\}, \{2, 3\}, \{3, 4\}\}$
- Cliques: $\mathcal{C}(G) = V \cup E \cup \{\{1, 2, 3\}, \{1, 3, 4\}\}$



Graph

- $G = (V, E)$ mit Knotenmenge V und Kantenmenge E
- Hier: $V = \{1, 2, 3, 4\}$, $E = \{\{1, 2\}, \{1, 3\}, \{1, 4\}, \{2, 3\}, \{3, 4\}\}$
- Cliques: $\mathcal{C}(G) = V \cup E \cup \{\{1, 2, 3\}, \{1, 3, 4\}\}$



Suffiziente Statistik

Daten \mathcal{D} , Modell mit Parameter β

Funktion ϕ ist eine suffiziente Statistik $\Leftrightarrow \beta \perp\!\!\!\perp \mathcal{D} \mid \phi(\mathcal{D})$

mit $\phi(\mathcal{D}) = \sum_{\mathbf{x} \in \mathcal{D}} \phi(\mathbf{x})$

Für diskrete \mathcal{X} :

ϕ ist immer gegeben durch $\phi_{C=y_C}(\mathbf{x}_C) = \prod_{v \in C} \mathbb{1}_{x_v=y_v}$

mit $C \in \mathcal{C}(G)$, $\mathbf{x}_C \in \mathcal{X}_C$, $\mathbf{x} \in \mathcal{X}$

$\phi_C(\mathbf{x}_C) = (\phi_{C=y_C^1}(\mathbf{x}_C), \phi_{C=y_C^2}(\mathbf{x}_C), \dots) = (\phi_{C=y_C}(\mathbf{x}_C))_{y_C \in \mathcal{X}_C}$

Suffiziente Statistik

Daten \mathcal{D} , Modell mit Parameter β

Funktion ϕ ist eine suffiziente Statistik $\Leftrightarrow \beta \perp\!\!\!\perp \mathcal{D} \mid \phi(\mathcal{D})$

mit $\phi(\mathcal{D}) = \sum_{\mathbf{x} \in \mathcal{D}} \phi(\mathbf{x})$

Für diskrete \mathcal{X} :

ϕ ist immer gegeben durch $\phi_{C=\mathbf{y}_C}(\mathbf{x}_C) = \prod_{v \in C} \mathbb{1}_{\mathbf{x}_v = \mathbf{y}_v}$

mit $C \in \mathcal{C}(G)$, $\mathbf{x}_C \in \mathcal{X}_C$, $\mathbf{x} \in \mathcal{X}$

$\phi_C(\mathbf{x}_C) = (\phi_{C=\mathbf{y}_C^1}(\mathbf{x}_C), \phi_{C=\mathbf{y}_C^2}(\mathbf{x}_C), \dots) = (\phi_{C=\mathbf{y}_C}(\mathbf{x}_C))_{\mathbf{y}_C \in \mathcal{X}_C}$



Beispiel: Suffiziente Statistik





Beispiel: Suffiziente Statistik



Beispiel: Suffiziente Statistik



Beispiel: Suffiziente Statistik





Beispiel: Suffiziente Statistik



Beispiel: Suffiziente Statistik

$$V = \{1, 2, 3, 4\}, E = \{\{1, 2\}, \{1, 3\}, \{1, 4\}, \{2, 3\}, \{3, 4\}\},$$

$$\mathcal{C}(G) = V \cup E \cup \{\{1, 2, 3\}, \{1, 3, 4\}\}$$

$$\phi_C(\mathbf{x}_C) = (\phi_{C=\mathbf{y}_C^1}(\mathbf{x}_C), \phi_{C=\mathbf{y}_C^2}(\mathbf{x}_C), \dots) = (\phi_{C=\mathbf{y}_C}(\mathbf{x}_C))_{\mathbf{y}_C \in \mathcal{X}_C}$$

$$\mathcal{X}_1 = \{1, 2, 3, 4, 5\}, \mathcal{X}_2 = \{-, \text{♫}\},$$

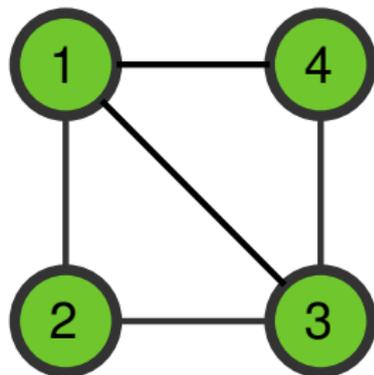
$$\mathcal{X}_3 = \{-, \text{Punk}, \text{Pop}, \dots\}, \mathcal{X}_4 = \{\text{■}, \text{■}, \text{■}\}$$

$$\mathbf{x} = (2, -, -, \text{■})$$

$$\mathbf{x}' = (1, -, \text{Pop}, \text{■})$$

$$\mathbf{x}'' = (1, \text{♫}, \text{Punk}, \text{■})$$

...



Beispiel: Suffiziente Statistik

$$V = \{1, 2, 3, 4\}, E = \{\{1, 2\}, \{1, 3\}, \{1, 4\}, \{2, 3\}, \{3, 4\}\}, \\
 \mathcal{C}(G) = V \cup E \cup \{\{1, 2, 3\}, \{1, 3, 4\}\}$$

$$\phi_C(\mathbf{x}_C) = (\phi_{C=\mathbf{y}_C^1}(\mathbf{x}_C), \phi_{C=\mathbf{y}_C^2}(\mathbf{x}_C), \dots) = (\phi_{C=\mathbf{y}_C}(\mathbf{x}_C))_{\mathbf{y}_C \in \mathcal{X}_C}$$

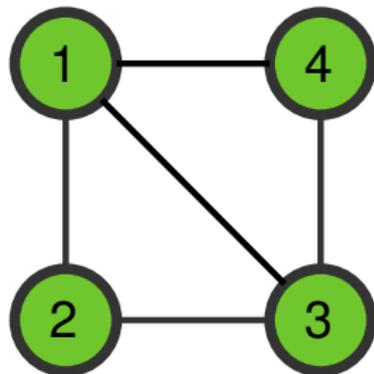
$$\mathcal{X}_1 = \{1, 2, 3, 4, 5\}, \mathcal{X}_2 = \{-, \text{♫}\}, \\
 \mathcal{X}_3 = \{-, \text{Punk, Pop, \dots}\}, \mathcal{X}_4 = \{\text{■, ■, ■}\}$$

$$\mathbf{x} = (2, -, -, \text{■})$$

$$\mathbf{x}' = (1, -, \text{Pop}, \text{■})$$

$$\mathbf{x}'' = (1, \text{♫}, \text{Punk}, \text{■})$$

...



Beispiel: Suffiziente Statistik

$$V = \{1, 2, 3, 4\}, E = \{\{1, 2\}, \{1, 3\}, \{1, 4\}, \{2, 3\}, \{3, 4\}\}, \\
 \mathcal{C}(G) = V \cup E \cup \{\{1, 2, 3\}, \{1, 3, 4\}\}$$

$$\phi_C(\mathbf{x}_C) = (\phi_{C=\mathbf{y}_C^1}(\mathbf{x}_C), \phi_{C=\mathbf{y}_C^2}(\mathbf{x}_C), \dots) = (\phi_{C=\mathbf{y}_C}(\mathbf{x}_C))_{\mathbf{y}_C \in \mathcal{X}_C}$$

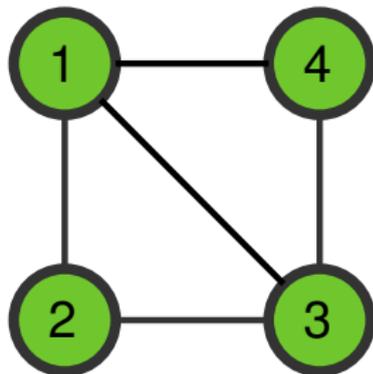
$$\mathcal{X}_1 = \{1, 2, 3, 4, 5\}, \mathcal{X}_2 = \{-, \text{♫}\}, \\
 \mathcal{X}_3 = \{-, \text{Punk, Pop, \dots}\}, \mathcal{X}_4 = \{\text{■, ■, ■}\}$$

$$\mathbf{x} = (2, -, -, \text{■})$$

$$\mathbf{x}' = (1, -, \text{Pop}, \text{■})$$

$$\mathbf{x}'' = (1, \text{♫}, \text{Punk}, \text{■})$$

...



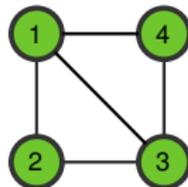
Beispiel: Suffiziente Statistik II

$$C = \{1, 2, 3\}, \mathcal{X}_C = \mathcal{X}_1 \times \mathcal{X}_2 \times \mathcal{X}_3$$

$$\mathcal{X}_1 = \{1, 2, 3, 4, 5\}, \mathcal{X}_2 = \{-, \text{♫}\},$$

$$\mathcal{X}_3 = \{-, \text{Punk}, \text{Pop}, \text{Rock}, \text{Schlager}\}$$

$$\mathcal{X}_C = \{(1, -, -), (2, -, -), \dots, (5, -, -), (1, \text{♫}, -), \dots, (5, \text{♫}, -), (1, \text{♫}, \text{Punk}), \dots, (5, \text{♫}, \text{Schlager})\}$$



$$\phi_{\{1,2,3\}}(5, \text{♫}, \text{Schlager}) = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix}$$

$$\phi_{\{1,2,3\}}(2, \text{♫}, -) = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \end{pmatrix}$$

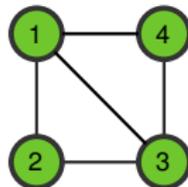
Beispiel: Suffiziente Statistik II

$$C = \{1, 2, 3\}, \mathcal{X}_C = \mathcal{X}_1 \times \mathcal{X}_2 \times \mathcal{X}_3$$

$$\mathcal{X}_1 = \{1, 2, 3, 4, 5\}, \mathcal{X}_2 = \{-, \text{♫}\},$$

$$\mathcal{X}_3 = \{-, \text{Punk}, \text{Pop}, \text{Rock}, \text{Schlager}\}$$

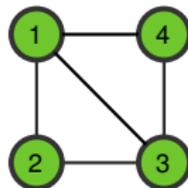
$$\mathcal{X}_C = \{(1, -, -), (2, -, -), \dots, (5, -, -), (1, \text{♫}, -), \dots, (5, \text{♫}, -), (1, \text{♫}, \text{Punk}), \dots, (5, \text{♫}, \text{Schlager})\}$$



$$\phi_{\{1,2,3\}}(5, \text{♫}, \text{Schlager}) = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix} \quad \phi_{\{1,2,3\}}(2, \text{♫}, -) = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \end{pmatrix}$$

Exponentialfamilie

$$\begin{aligned} \mathbb{P}_{\beta}(\mathbf{x}) &= \frac{1}{Z(\beta)} \exp(\langle \beta_{\{1,2,3\}}, \phi_{\{1,2,3\}}(\mathbf{x}_{\{1,2,3\}}) \rangle) \\ &\quad \exp(\langle \beta_{\{1,3,4\}}, \phi_{\{1,3,4\}}(\mathbf{x}_{\{1,3,4\}}) \rangle) \\ &= \exp(\langle \beta, \phi(\mathbf{x}) \rangle - A(\beta)) \end{aligned}$$



$$A(\beta) = \log Z(\beta) = \log \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle)$$



Maximum Entropie Prinzip: Aufgabe

Gegeben: Daten \mathcal{X} (Realisierungen einer ZV \mathbf{X})
und beliebige Funktion $f : \mathcal{X} \rightarrow \mathbb{R}^d$

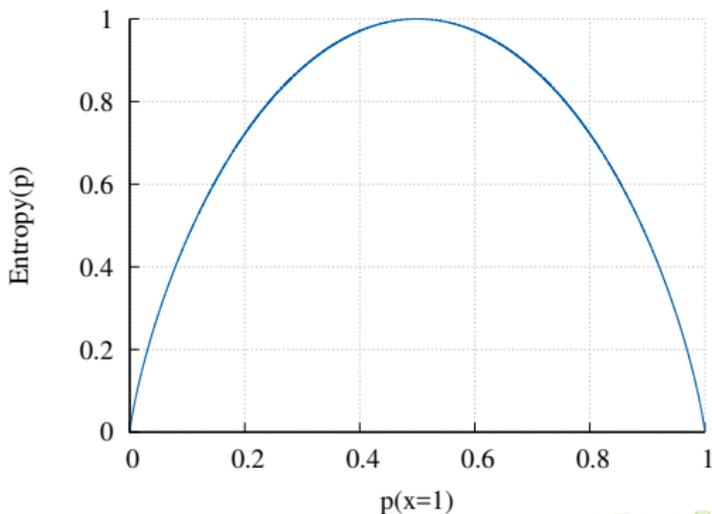
Gesucht: \mathbb{P} mit $\mathbb{E}_{\mathbb{P}}[f(\mathbf{X})] = \tilde{\mathbb{E}}_{\mathcal{D}}[f(\mathbf{X})] = \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x})$

Problem: Viele \mathbb{P} haben die gesuchte Eigenschaft!!

Maximum Entropie Prinzip: Intuition

Entropie \mathcal{H} einer Zufallsvariable X mit \mathbb{P}

$$\mathcal{H}(X) = - \sum_{x \in \mathcal{X}} \mathbb{P}(x) \log \mathbb{P}(x)$$

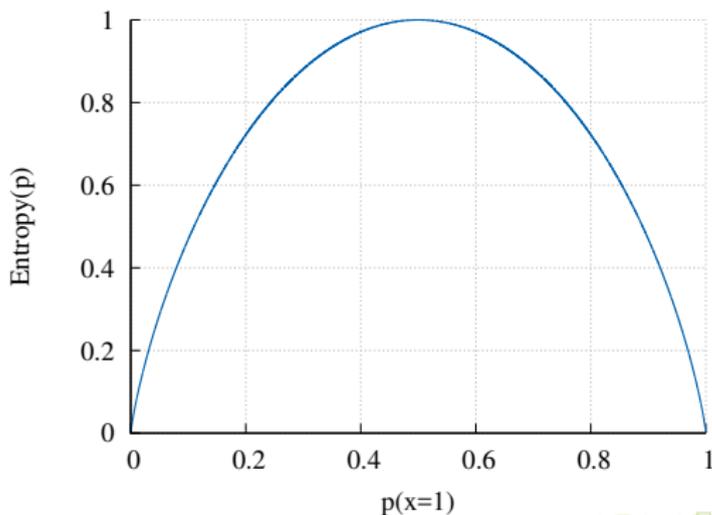


Maximum Entropie Prinzip: Intuition

Sei \mathcal{P} der Raum aller Wahrscheinlichkeitsfunktionen

$$\max_{\mathbb{P} \in \mathcal{P}} \mathcal{H}(\mathbb{P})$$

s.t. $\mathbb{E}_{\mathbb{P}}[f(\mathbf{X})] = \tilde{\mathbb{E}}_{\mathcal{D}}[f(\mathbf{X})] \leftarrow d \text{ Nebenbedingungen}$



Maximum Entropie Prinzip: Lösung

Umformulieren in Lagrange Funktion

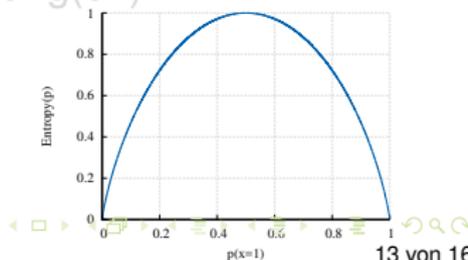
$$\max_{\mathbb{P} \in \mathcal{P}} \mathcal{H}(\mathbb{P}) + \sum_{i=1}^d \lambda_i \mathbb{E}_{\mathbb{P}}[f(\mathbf{X})]_i - \tilde{\mathbb{E}}_{\mathcal{D}}[f(\mathbf{X})]_i$$

ableiten (nach \mathbb{P} !) und $= 0$ setzen liefert

$$\mathbb{P}(\mathbf{x}) = \exp(\langle \boldsymbol{\lambda}, f(\mathbf{x}) \rangle - A(\boldsymbol{\lambda}))$$

Also: Parameter von Exponentialfamilien sind
Lagrange-Multiplikatoren der Nebenbedingung(en)

$$\mathbb{E}_{\mathbb{P}}[f(\mathbf{X})] = \tilde{\mathbb{E}}_{\mathcal{D}}[f(\mathbf{X})]$$



Maximum Entropie Prinzip: Lösung

Umformulieren in Lagrange Funktion

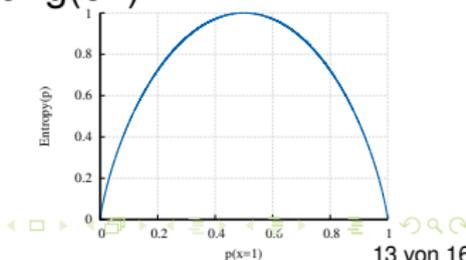
$$\max_{\mathbb{P} \in \mathcal{P}} \mathcal{H}(\mathbb{P}) + \sum_{i=1}^d \lambda_i \mathbb{E}_{\mathbb{P}}[f(\mathbf{X})]_i - \tilde{\mathbb{E}}_{\mathcal{D}}[f(\mathbf{X})]_i$$

ableiten (nach \mathbb{P} !) und $= 0$ setzen liefert

$$\mathbb{P}(\mathbf{x}) = \exp(\langle \boldsymbol{\lambda}, f(\mathbf{x}) \rangle - A(\boldsymbol{\lambda}))$$

Also: Parameter von Exponentialfamilien sind
Lagrange-Multiplikatoren der Nebenbedingung(en)

$$\mathbb{E}_{\mathbb{P}}[f(\mathbf{X})] = \tilde{\mathbb{E}}_{\mathcal{D}}[f(\mathbf{X})]$$





Parameterlernen durch Gradientenabstieg

Gegeben Datensatz \mathcal{D} , Funktion ϕ (binär) Verlustfunktion:
Negative mittlere log-Likelihood

$$\begin{aligned}\ell(\beta, \mathcal{D}) &= -\frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \log \mathbb{P}_{\beta}(\mathbf{x}) \\ &= -\frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \langle \beta, \phi(\mathbf{x}) \rangle + A(\beta) \\ &= -\langle \beta, \tilde{\mu} \rangle + A(\beta)\end{aligned}$$

Partielle Ableitung:

$$\frac{\partial \ell(\beta, \mathcal{D})}{\partial \beta_i} = -\tilde{\mu}_i + \frac{\partial}{\partial \beta_i} A(\beta) = \hat{\mu}_i - \tilde{\mu}_i$$

Parameterlernen durch Gradientenabstieg

Gegeben Datensatz \mathcal{D} , Funktion ϕ (binär) Verlustfunktion:
 Negative mittlere log-Likelihood

$$\begin{aligned}
 \ell(\boldsymbol{\beta}, \mathcal{D}) &= -\frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \log \mathbb{P}_{\boldsymbol{\beta}}(\mathbf{x}) \\
 &= -\frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \langle \boldsymbol{\beta}, \phi(\mathbf{x}) \rangle + A(\boldsymbol{\beta}) \\
 &= -\langle \boldsymbol{\beta}, \tilde{\boldsymbol{\mu}} \rangle + A(\boldsymbol{\beta})
 \end{aligned}$$

Partielle Ableitung:

$$\frac{\partial \ell(\boldsymbol{\beta}, \mathcal{D})}{\partial \beta_i} = -\tilde{\mu}_i + \frac{\partial}{\partial \beta_i} A(\boldsymbol{\beta}) = \hat{\mu}_i - \tilde{\mu}_i$$

Parameterlernen durch Gradientenabstieg

Gegeben Datensatz \mathcal{D} , Funktion ϕ (binär) Verlustfunktion:
 Negative mittlere log-Likelihood

$$\begin{aligned}
 \ell(\boldsymbol{\beta}, \mathcal{D}) &= -\frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \log \mathbb{P}_{\boldsymbol{\beta}}(\mathbf{x}) \\
 &= -\frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \langle \boldsymbol{\beta}, \phi(\mathbf{x}) \rangle + A(\boldsymbol{\beta}) \\
 &= -\langle \boldsymbol{\beta}, \tilde{\boldsymbol{\mu}} \rangle + A(\boldsymbol{\beta})
 \end{aligned}$$

Partielle Ableitung:

$$\frac{\partial \ell(\boldsymbol{\beta}, \mathcal{D})}{\partial \beta_i} = -\tilde{\mu}_i + \frac{\partial}{\partial \beta_i} A(\boldsymbol{\beta}) = \hat{\mu}_i - \tilde{\mu}_i$$



Parameterlernen durch Gradientenabstieg

Gegeben Datensatz \mathcal{D} , Funktion ϕ (binär) Verlustfunktion:
Negative mittlere log-Likelihood

$$\begin{aligned}\ell(\boldsymbol{\beta}, \mathcal{D}) &= -\frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \log \mathbb{P}_{\boldsymbol{\beta}}(\mathbf{x}) \\ &= -\frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \langle \boldsymbol{\beta}, \phi(\mathbf{x}) \rangle + A(\boldsymbol{\beta}) \\ &= -\langle \boldsymbol{\beta}, \tilde{\boldsymbol{\mu}} \rangle + A(\boldsymbol{\beta})\end{aligned}$$

Partielle Ableitung:

$$\frac{\partial \ell(\boldsymbol{\beta}, \mathcal{D})}{\partial \beta_i} = -\tilde{\mu}_i + \frac{\partial}{\partial \beta_i} A(\boldsymbol{\beta}) = \hat{\mu}_i - \tilde{\mu}_i$$

Marginalisierung

Wenn ϕ binär, dann ist $\mu_i = \mathbb{E}[\phi_i(\mathbf{X})]$ die Wahrscheinlichkeit für $\phi_i(\mathbf{X}) = 1$

Annahme: Paarweises Modell \equiv Nur die Kantengewichte sind relevant

$$\mathbb{P}(\mathbf{X}_v = x) = \sum_{\mathbf{y} \in \mathcal{X}_{V \setminus \{v\}}} \mathbb{P}(\mathbf{y}, x)$$

Ausnutzen der Faktorisierung sowie der Distributivität:

$$\mathbb{P}(\mathbf{X}_v = x) = \frac{1}{Z} \sum_{\mathbf{y} \in \mathcal{X}_{V \setminus \{v\}}} \prod_{C \in \mathcal{C}(G)} \exp(\langle \beta_C, \phi_C(\mathbf{x}_C) \rangle)$$

wobei $\mathbf{x} = (\mathbf{y}, x)$



Marginalisierung

Wenn ϕ binär, dann ist $\mu_i = \mathbb{E}[\phi_i(\mathbf{X})]$ die Wahrscheinlichkeit für $\phi_i(\mathbf{X}) = 1$

Annahme: Paarweises Modell \equiv Nur die Kantengewichte sind relevant

$$\mathbb{P}(\mathbf{X}_v = x) = \sum_{\mathbf{y} \in \mathcal{X}_{V \setminus \{v\}}} \mathbb{P}(\mathbf{y}, x)$$

Ausnutzen der Faktorisierung sowie der Distributivität:

$$\mathbb{P}(\mathbf{X}_v = x) = \frac{1}{Z} \sum_{\mathbf{y} \in \mathcal{X}_{V \setminus \{v\}}} \prod_{C \in \mathcal{C}(G)} \exp(\langle \beta_C, \phi_C(\mathbf{x}_C) \rangle)$$

wobei $\mathbf{x} = (\mathbf{y}, x)$

Marginalisierung

Wenn ϕ binär, dann ist $\mu_i = \mathbb{E}[\phi_i(\mathbf{X})]$ die Wahrscheinlichkeit für $\phi_i(\mathbf{X}) = 1$

Annahme: Paarweises Modell \equiv Nur die Kantengewichte sind relevant

$$\mathbb{P}(\mathbf{X}_v = x) = \sum_{\mathbf{y} \in \mathcal{X}_{V \setminus \{v\}}} \mathbb{P}(\mathbf{y}, x)$$

Ausnutzen der Faktorisierung sowie der Distributivität:

$$\mathbb{P}(\mathbf{X}_v = x) = \frac{1}{Z} \sum_{\mathbf{y} \in \mathcal{X}_{V \setminus \{v\}}} \prod_{C \in \mathcal{C}(G)} \exp(\langle \boldsymbol{\beta}_C, \phi_C(\mathbf{x}_C) \rangle)$$

wobei $\mathbf{x} = (\mathbf{y}, x)$



Gibbs Sampling

- Idee: Erzeuge neue Stichprobe gemäß \mathbb{P}_β und berechne $\hat{\mu}_i$ durch “abzählen”
- Aber: Wie erzeugt man neue Samples aus $\mathbb{P}_\beta(\mathbf{X})$?

⇒ Ausnutzung bedingter Unabhängigkeiten!



Gibbs Sampling

- Idee: Erzeuge neue Stichprobe gemäß \mathbb{P}_β und berechne $\hat{\mu}_i$ durch “abzählen”
- Aber: Wie erzeugt man neue Samples aus $\mathbb{P}_\beta(\mathbf{X})$?

⇒ Ausnutzung bedingter Unabhängigkeiten!