

Wissensentdeckung in Datenbanken

Belief Propagation, Strukturlernen

Nico Piatkowski und Uwe Ligges

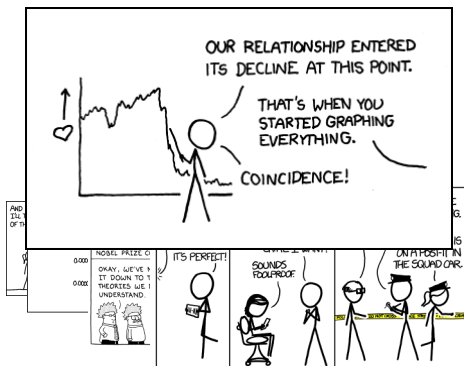
Informatik—Künstliche Intelligenz
Computergestützte Statistik
Technische Universität Dortmund

29.06.2017

Überblick

Was bisher geschah...

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



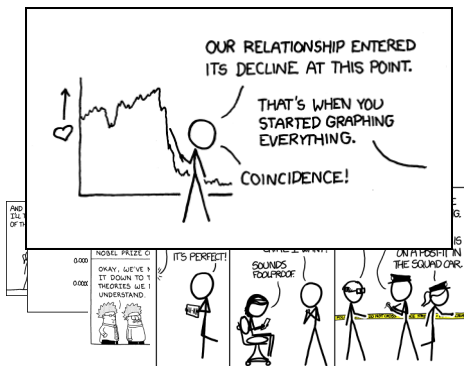
Heute

- Graphische Modelle—Inferenz und Strukturlernen

Überblick

Was bisher geschah...

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



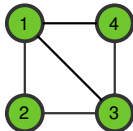
Heute

- Graphische Modelle—Inferenz und Strukturlernen

Mehr Eigenschaften von Graphen

Nachbarschaft: Für Knoten $v \in V$ in Graph $G = (V, E)$

$$\mathcal{N}(v) = \{u \in V \mid \{v, u\} \in E\}$$

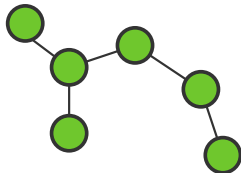


Pfad: Folge (v_1, v_2, \dots, v_m) von Knoten in der sich kein Knoten wiederholt

Kreis: Pfad (v_1, v_2, \dots, v_m) mit $\{v_1, v_m\} \in E$

Baum:

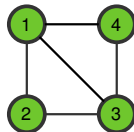
- Graph ohne Kreise (“kreisfrei”)
- Maximale Cliquengröße = 2



Mehr Eigenschaften von Graphen

Nachbarschaft: Für Knoten $v \in V$ in Graph $G = (V, E)$

$$\mathcal{N}(v) = \{u \in V \mid \{v, u\} \in E\}$$

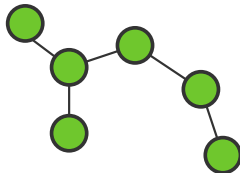


Pfad: Folge (v_1, v_2, \dots, v_m) von Knoten in der sich kein Knoten wiederholt

Kreis: Pfad (v_1, v_2, \dots, v_m) mit $\{v_1, v_m\} \in E$

Baum:

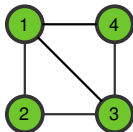
- Graph ohne Kreise (“kreisfrei”)
- Maximale Cliquengröße = 2



Mehr Eigenschaften von Graphen

Nachbarschaft: Für Knoten $v \in V$ in Graph $G = (V, E)$

$$\mathcal{N}(v) = \{u \in V \mid \{v, u\} \in E\}$$

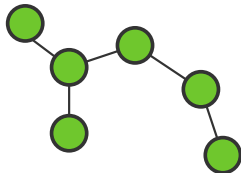


Pfad: Folge (v_1, v_2, \dots, v_m) von Knoten in der sich kein Knoten wiederholt

Kreis: Pfad (v_1, v_2, \dots, v_m) mit $\{v_1, v_m\} \in E$

Baum:

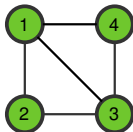
- Graph ohne Kreise (“kreisfrei”)
- Maximale Cliquengröße = 2



Mehr Eigenschaften von Graphen

Nachbarschaft: Für Knoten $v \in V$ in Graph $G = (V, E)$

$$\mathcal{N}(v) = \{u \in V \mid \{v, u\} \in E\}$$

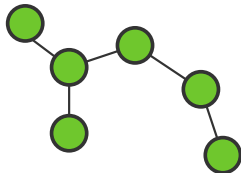


Pfad: Folge (v_1, v_2, \dots, v_m) von Knoten in der sich kein Knoten wiederholt

Kreis: Pfad (v_1, v_2, \dots, v_m) mit $\{v_1, v_m\} \in E$

Baum:

- Graph ohne Kreise (“kreisfrei”)
- Maximale Cliquengröße = 2



Gradient der Log-Likelihood von Exponentialfamilien

Gegeben Datensatz \mathcal{D} , Funktion ϕ (binär),

$$\tilde{\mathbb{E}}[\phi(\mathbf{X})] = (1/|\mathcal{D}|) \sum_{\mathbf{x} \in \mathcal{D}} \phi(\mathbf{x}), \mathbb{P}_{\beta}(\mathbf{x}) = \exp(\langle \beta, \phi(\mathbf{x}) \rangle - A(\beta))$$

Negative mittlere log-Likelihood:

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

Partielle Ableitung nach β_i :

$$\frac{\partial \ell(\beta; \mathcal{D})}{\partial \beta_i} = -\tilde{\mathbb{E}}[\phi(\mathbf{X})_i] + \frac{\partial}{\partial \beta_i} A(\beta)$$



Gradient der Log-Likelihood von Exponentialfamilien

Gegeben Datensatz \mathcal{D} , Funktion ϕ (binär),

$$\tilde{\mathbb{E}}[\phi(\mathbf{X})] = (1/|\mathcal{D}|) \sum_{\mathbf{x} \in \mathcal{D}} \phi(\mathbf{x}), \mathbb{P}_{\beta}(\mathbf{x}) = \exp(\langle \beta, \phi(\mathbf{x}) \rangle - A(\beta))$$

Negative mittlere log-Likelihood:

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

Partielle Ableitung nach β_i :

$$\frac{\partial \ell(\beta; \mathcal{D})}{\partial \beta_i} = -\tilde{\mathbb{E}}[\phi(\mathbf{X})_i] + \frac{\partial}{\partial \beta_i} A(\beta)$$

Gradient der Log-Likelihood von Exponentialfamilien

Gegeben Datensatz \mathcal{D} , Funktion ϕ (binär),

$$\tilde{\mathbb{E}}[\phi(\mathbf{X})] = (1/|\mathcal{D}|) \sum_{\mathbf{x} \in \mathcal{D}} \phi(\mathbf{x}), \mathbb{P}_{\beta}(\mathbf{x}) = \exp(\langle \beta, \phi(\mathbf{x}) \rangle - A(\beta))$$

Negative mittlere log-Likelihood:

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

Partielle Ableitung nach β_i :

$$\frac{\partial \ell(\beta; \mathcal{D})}{\partial \beta_i} = -\tilde{\mathbb{E}}[\phi(\mathbf{X})_i] + \frac{\partial}{\partial \beta_i} A(\beta)$$

Gradient der Log-Likelihood von Exponentialfamilien

Gegeben Datensatz \mathcal{D} , Funktion ϕ (binär),

$$\tilde{\mathbb{E}}[\phi(\mathbf{X})] = (1/|\mathcal{D}|) \sum_{\mathbf{x} \in \mathcal{D}} \phi(\mathbf{x}), \mathbb{P}_{\beta}(\mathbf{x}) = \exp(\langle \beta, \phi(\mathbf{x}) \rangle - A(\beta))$$

Negative mittlere log-Likelihood:

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

Partielle Ableitung nach β_i :

$$\frac{\partial \ell(\beta; \mathcal{D})}{\partial \beta_i} = -\tilde{\mathbb{E}}[\phi(\mathbf{X})_i] + \frac{\partial}{\partial \beta_i} A(\beta)$$

Gradient der Log-Likelihood von Exponentialfamilien

Gegeben Datensatz \mathcal{D} , Funktion ϕ (binär),

$$\tilde{\mathbb{E}}[\phi(\mathbf{X})] = (1/|\mathcal{D}|) \sum_{\mathbf{x} \in \mathcal{D}} \phi(\mathbf{x}), \mathbb{P}_{\boldsymbol{\beta}}(\mathbf{x}) = \exp(\langle \boldsymbol{\beta}, \phi(\mathbf{x}) \rangle - A(\boldsymbol{\beta}))$$

Negative mittlere log-Likelihood:

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

Partielle Ableitung nach β_i :

$$\frac{\partial \ell(\boldsymbol{\beta}; \mathcal{D})}{\partial \beta_i} = -\tilde{\mathbb{E}}[\phi(\mathbf{X})_i] + \frac{\partial}{\partial \beta_i} A(\boldsymbol{\beta})$$



Gradient der Log-Likelihood von Exponentialfamilien

Gegeben Datensatz \mathcal{D} , Funktion ϕ (binär),

$$\tilde{\mathbb{E}}[\phi(\mathbf{X})] = (1/|\mathcal{D}|) \sum_{\mathbf{x} \in \mathcal{D}} \phi(\mathbf{x}), \mathbb{P}_{\beta}(\mathbf{x}) = \exp(\langle \beta, \phi(\mathbf{x}) \rangle - A(\beta))$$

Negative mittlere log-Likelihood:

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

Partielle Ableitung nach β_i :

$$\frac{\partial \ell(\beta; \mathcal{D})}{\partial \beta_i} = -\tilde{\mathbb{E}}[\phi(\mathbf{X})_i] + \frac{\partial}{\partial \beta_i} A(\beta)$$



Ableitung der Normalisierung

Erinnerung: Form von $A(\beta)$ ($= \log Z(\beta)$) ist abhängig von \mathcal{X} (diskret?, reel?, endlich?). Hier:

$$\begin{aligned}\frac{\partial}{\partial \beta_i} A(\beta) &= \frac{\partial}{\partial \beta_i} \log \sum_{x \in \mathcal{X}} \exp(\langle \beta, \phi(x) \rangle) \\ &= \frac{1}{Z(\beta)} \frac{\partial}{\partial \beta_i} \sum_{x \in \mathcal{X}} \exp(\langle \beta, \phi(x) \rangle) \\ &= \frac{1}{Z(\beta)} \sum_{x \in \mathcal{X}} \exp(\langle \beta, \phi(x) \rangle) \frac{\partial}{\partial \beta_i} \langle \beta, \phi(x) \rangle \\ &= \sum_{x \in \mathcal{X}} \exp(\langle \beta, \phi(x) \rangle - A(\beta)) \phi(x)_i \\ &= \mathbb{E}_\beta[\phi(\mathbf{X})_i]\end{aligned}$$



Ableitung der Normalisierung

Erinnerung: Form von $A(\beta)$ ($= \log Z(\beta)$) ist abhängig von \mathcal{X} (diskret?, reel?, endlich?). Hier:

$$\begin{aligned}\frac{\partial}{\partial \beta_i} A(\beta) &= \frac{\partial}{\partial \beta_i} \log \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle) \\ &= \frac{1}{Z(\beta)} \frac{\partial}{\partial \beta_i} \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle) \\ &= \frac{1}{Z(\beta)} \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle) \frac{\partial}{\partial \beta_i} \langle \beta, \phi(\mathbf{x}) \rangle \\ &= \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle - A(\beta)) \phi(\mathbf{x})_i \\ &= \mathbb{E}_{\beta}[\phi(\mathbf{X})_i]\end{aligned}$$



Ableitung der Normalisierung

Erinnerung: Form von $A(\beta)$ ($= \log Z(\beta)$) ist abhängig von \mathcal{X} (diskret?, reel?, endlich?). Hier:

$$\begin{aligned}\frac{\partial}{\partial \beta_i} A(\beta) &= \frac{\partial}{\partial \beta_i} \log \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle) \\ &= \frac{1}{Z(\beta)} \frac{\partial}{\partial \beta_i} \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle) \\ &= \frac{1}{Z(\beta)} \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle) \frac{\partial}{\partial \beta_i} \langle \beta, \phi(\mathbf{x}) \rangle \\ &= \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle - A(\beta)) \phi(\mathbf{x})_i \\ &= \mathbb{E}_{\beta}[\phi(\mathbf{X})_i]\end{aligned}$$



Ableitung der Normalisierung

Erinnerung: Form von $A(\beta)$ ($= \log Z(\beta)$) ist abhängig von \mathcal{X} (diskret?, reel?, endlich?). Hier:

$$\begin{aligned}\frac{\partial}{\partial \beta_i} A(\beta) &= \frac{\partial}{\partial \beta_i} \log \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle) \\ &= \frac{1}{Z(\beta)} \frac{\partial}{\partial \beta_i} \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle) \\ &= \frac{1}{Z(\beta)} \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle) \frac{\partial}{\partial \beta_i} \langle \beta, \phi(\mathbf{x}) \rangle \\ &= \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle - A(\beta)) \phi(\mathbf{x})_i \\ &= \mathbb{E}_{\beta}[\phi(\mathbf{X})_i]\end{aligned}$$



Ableitung der Normalisierung

Erinnerung: Form von $A(\beta)$ ($= \log Z(\beta)$) ist abhängig von \mathcal{X} (diskret?, reel?, endlich?). Hier:

$$\begin{aligned}\frac{\partial}{\partial \beta_i} A(\beta) &= \frac{\partial}{\partial \beta_i} \log \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle) \\ &= \frac{1}{Z(\beta)} \frac{\partial}{\partial \beta_i} \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle) \\ &= \frac{1}{Z(\beta)} \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle) \frac{\partial}{\partial \beta_i} \langle \beta, \phi(\mathbf{x}) \rangle \\ &= \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle - A(\beta)) \phi(\mathbf{x})_i \\ &= \mathbb{E}_{\beta}[\phi(\mathbf{X})_i]\end{aligned}$$



Ableitung der Normalisierung

Erinnerung: Form von $A(\beta)$ ($= \log Z(\beta)$) ist abhängig von \mathcal{X} (diskret?, reel?, endlich?). Hier:

$$\begin{aligned}\frac{\partial}{\partial \beta_i} A(\beta) &= \frac{\partial}{\partial \beta_i} \log \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle) \\ &= \frac{1}{Z(\beta)} \frac{\partial}{\partial \beta_i} \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle) \\ &= \frac{1}{Z(\beta)} \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle) \frac{\partial}{\partial \beta_i} \langle \beta, \phi(\mathbf{x}) \rangle \\ &= \sum_{\mathbf{x} \in \mathcal{X}} \exp(\langle \beta, \phi(\mathbf{x}) \rangle - A(\beta)) \phi(\mathbf{x})_i \\ &= \mathbb{E}_{\beta}[\phi(\mathbf{X})_i]\end{aligned}$$

Gradient der Log-Likelihood von Exponentialfamilien

Also gilt für die **partielle Ableitung** nach β_i :

$$\begin{aligned}
 \frac{\partial \ell(\beta; \mathcal{D})}{\partial \beta_i} &= -\tilde{\mathbb{E}}[\phi(\mathbf{X})_i] + \frac{\partial}{\partial \beta_i} A(\beta) \\
 &= \mathbb{E}_\beta[\phi(\mathbf{X})_i] - \tilde{\mathbb{E}}[\phi(\mathbf{X})_i]
 \end{aligned}$$

Ableitung ist beschränkt $\frac{\partial \ell(\beta; \mathcal{D})}{\partial \beta_i} \in [-1; 1] \Rightarrow \ell$ ist Lipschitz stetig
 Man kann zeigen: $\nabla \ell(\beta; \mathcal{D})$ auch Lipschitz stetig ($L = 2|\mathcal{C}(G)|$)

Jetzt: Wie berechnet man $\mathbb{E}_\beta[\phi(\mathbf{X})_i]$?

Jedes i entspricht einem Paar von Clique und Zustand, d.h.

$\phi(\mathbf{X})_i = \phi(\mathbf{X})_{C=y}$ für ein $C \in \mathcal{C}(g)$ und $y \in \mathcal{X}_C$

Gradient der Log-Likelihood von Exponentialfamilien

Also gilt für die **partielle Ableitung** nach β_i :

$$\begin{aligned}
 \frac{\partial \ell(\beta; \mathcal{D})}{\partial \beta_i} &= -\tilde{\mathbb{E}}[\phi(\mathbf{X})_i] + \frac{\partial}{\partial \beta_i} A(\beta) \\
 &= \mathbb{E}_\beta[\phi(\mathbf{X})_i] - \tilde{\mathbb{E}}[\phi(\mathbf{X})_i]
 \end{aligned}$$

Ableitung ist beschränkt $\frac{\partial \ell(\beta; \mathcal{D})}{\partial \beta_i} \in [-1; 1] \Rightarrow \ell$ ist Lipschitz stetig
 Man kann zeigen: $\nabla \ell(\beta; \mathcal{D})$ auch Lipschitz stetig ($L = 2|\mathcal{C}(G)|$)

Jetzt: Wie berechnet man $\mathbb{E}_\beta[\phi(\mathbf{X})_i]$?

Jedes i entspricht einem Paar von Clique und Zustand, d.h.

$\phi(\mathbf{X})_i = \phi(\mathbf{X})_{C=y}$ für ein $C \in \mathcal{C}(g)$ und $y \in \mathcal{X}_C$

Gradient der Log-Likelihood von Exponentialfamilien

Also gilt für die **partielle Ableitung** nach β_i :

$$\begin{aligned}
 \frac{\partial \ell(\beta; \mathcal{D})}{\partial \beta_i} &= -\tilde{\mathbb{E}}[\phi(\mathbf{X})_i] + \frac{\partial}{\partial \beta_i} A(\beta) \\
 &= \mathbb{E}_\beta[\phi(\mathbf{X})_i] - \tilde{\mathbb{E}}[\phi(\mathbf{X})_i]
 \end{aligned}$$

Ableitung ist beschränkt $\frac{\partial \ell(\beta; \mathcal{D})}{\partial \beta_i} \in [-1; 1] \Rightarrow \ell$ ist Lipschitz stetig
 Man kann zeigen: $\nabla \ell(\beta; \mathcal{D})$ auch Lipschitz stetig ($L = 2|\mathcal{C}(G)|$)

Jetzt: Wie berechnet man $\mathbb{E}_\beta[\phi(\mathbf{X})_i]$?

Jedes i entspricht einem Paar von Clique und Zustand, d.h.

$\phi(\mathbf{X})_i = \phi(\mathbf{X})_{C=y}$ für ein $C \in \mathcal{C}(g)$ und $\mathbf{y} \in \mathcal{X}_C$



Marginalisierung

Wenn Z binäre Zufallsvariable ($\mathcal{Z} = \{0, 1\}$), dann
 $\mathbb{E}[Z] = \mathbb{P}(Z = 1)$

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

Allgemein:

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

Ausnutzen der Faktorisierung sowie der Distributivität:

$$\mathbb{P}_\beta(\mathbf{X}_C = \mathbf{y}) = \frac{1}{Z(\beta)} \sum_{\mathbf{x} \in \mathcal{X}_{V \setminus C}} \prod_{U \in \mathcal{C}(G)} \exp(\langle \beta_U, \phi_U(z_U) \rangle)$$

wobei $\mathbf{z} = (\mathbf{y}, \mathbf{x})$ und z_U sind die Werte der Knoten in $U \subseteq V$

Marginalisierung

Wenn Z binäre Zufallsvariable ($\mathcal{Z} = \{0, 1\}$), dann
 $\mathbb{E}[Z] = \mathbb{P}(Z = 1)$

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

Allgemein:

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

Ausnutzen der Faktorisierung sowie der Distributivität:

$$\mathbb{P}_\beta(\mathbf{X}_C = \mathbf{y}) = \frac{1}{Z(\beta)} \sum_{\mathbf{x} \in \mathcal{X}_{V \setminus C}} \prod_{U \in \mathcal{C}(G)} \exp(\langle \beta_U, \phi_U(z_U) \rangle)$$

wobei $z = (\mathbf{y}, \mathbf{x})$ und z_U sind die Werte der Knoten in $U \subseteq V$



Marginalisierung

Wenn Z binäre Zufallsvariable ($\mathcal{Z} = \{0, 1\}$), dann
 $\mathbb{E}[Z] = \mathbb{P}(Z = 1)$

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

Allgemein:

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

Ausnutzen der Faktorisierung sowie der Distributivität:

$$\mathbb{P}_{\beta}(\mathbf{X}_C = \mathbf{y}) = \frac{1}{Z(\beta)} \sum_{\mathbf{x} \in \mathcal{X}_{V \setminus C}} \prod_{U \in \mathcal{C}(G)} \exp(\langle \beta_U, \phi_U(\mathbf{z}_U) \rangle)$$

wobei $\mathbf{z} = (\mathbf{y}, \mathbf{x})$ und \mathbf{z}_U sind die Werte der Knoten in $U \subseteq V$

Marginalisierung in Bäumen

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

Wenn G ein Baum ist mit $V = \{1, 2, \dots, n\}$, dann ist

$$\mathbb{P}_{\beta}(\mathbf{x}) = \frac{1}{Z(\beta)} \prod_{uv \in E} \exp(\langle \beta_{uv}, \phi_{uv}(\mathbf{x}_{uv}) \rangle)$$

und

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

→ Distributivität ausnutzen!

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

- Problem: Belief Propagation nur exakt wenn G ein Baum ist!
- 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!

Beobachtung: Wenn ganze Nachbarschaft eines Knotens v beobachtet ist, kann $\mathbb{P}_v(x \mid \mathbf{x}_{\mathcal{N}(v)})$ einfach berechnet werden!

Gibbs Sampling

- Problem: Belief Propagation nur exakt wenn G ein Baum ist!
- 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!

Beobachtung: Wenn ganze Nachbarschaft eines Knotens v beobachtet ist, kann $\mathbb{P}_v(x \mid \mathbf{x}_{\mathcal{N}(v)})$ einfach berechnet werden!



Gibbs Sampling: Algorithmus

- 1 Erzeuge x zufällig Gleichverteilt (das entspricht **NICHT** \mathbb{P}_β !)
- 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}_β !

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}_β !)
- 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}_β !

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”



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 β_{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:

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