

Wissensentdeckung in Datenbanken

LBP, Strukturlernen

Nico Piatkowski und Uwe Ligges

Informatik—Künstliche Intelligenz
Computergestützte Statistik
Technische Universität Dortmund

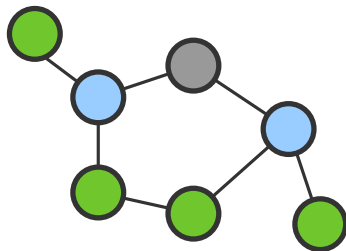
04.07.2017

Überblick

- Wiederholung: Belief Propagation
 - Berechnung der Randverteilung

- Strukturlernen—“Woher kommt der Graph?”
 - Was bedeutet der Graph?
 - Chow-Liu Bäumen
 - Regularisierung

- Merkmalsauswahl
 - Forward-Selection
 - Regularisierung



Rückblick (I)

- (1) Die Zielfunktion von probabilistischen Modellen ist die (mittlere) negative Log-Likelihood (02.05.)
- (2) Die Parameter β des graphischen Modells lernen wir mittels Gradientenabstieg (04.05. + Ü)
- (3) Der Gradient besteht aus partiellen Ableitungen (04.05.)
- (4) Im graphischen Modell entspricht eine partielle Ableitung bezüglich Parameter β_i dem Erwartungswert der i -ten Komponente der Zufallsvariable $\phi(\mathbf{X})$ (29.06.)

Rückblick (I)

- (1) Die Zielfunktion von probabilistischen Modellen ist die (mittlere) negative Log-Likelihood (02.05.)
- (2) Die Parameter β des graphischen Modells lernen wir mittels Gradientenabstieg (04.05. + Ü)
- (3) Der Gradient besteht aus partiellen Ableitungen (04.05.)
- (4) Im graphischen Modell entspricht eine partielle Ableitung bezüglich Parameter β_i dem Erwartungswert der i -ten Komponente der Zufallsvariable $\phi(\mathbf{X})$ (29.06.)

Rückblick (I)

- (1) Die Zielfunktion von probabilistischen Modellen ist die (mittlere) negative Log-Likelihood (02.05.)
- (2) Die Parameter β des graphischen Modells lernen wir mittels Gradientenabstieg (04.05. + Ü)
- (3) Der Gradient besteht aus partiellen Ableitungen (04.05.)
- (4) Im graphischen Modell entspricht eine partielle Ableitung bezüglich Parameter β_i dem Erwartungswert der i -ten Komponente der Zufallsvariable $\phi(\mathbf{X})$ (29.06.)

Rückblick (I)

- (1) Die Zielfunktion von probabilistischen Modellen ist die (mittlere) negative Log-Likelihood (02.05.)
- (2) Die Parameter β des graphischen Modells lernen wir mittels Gradientenabstieg (04.05. + Ü)
- (3) Der Gradient besteht aus partiellen Ableitungen (04.05.)
- (4) Im graphischen Modell entspricht eine partielle Ableitung bezüglich Parameter β_i dem Erwartungswert der i -ten Komponente der Zufallsvariable $\phi(\mathbf{X})$ (29.06.)

Rückblick (II)

(5) Bei diskreten Modellen ist $\phi(\mathbf{X})_i$ binär (22.06. ff.) und ihr Erwartungswert entspricht der Wahrscheinlichkeit $\mathbb{P}(\phi(\mathbf{X})_i = 1)$ (u.a. 29.06.)

(6) Jede Komponente von $\phi(\mathbf{X})$ entspricht einem Ereignis $\mathbf{X}_C = \mathbf{y}$, d.h. $\mathbb{P}(\phi(\mathbf{X})_i = 1) = \mathbb{P}(\mathbf{X}_C = \mathbf{y})$ mit $\mathbf{y} \in \mathcal{X}_C$

\Rightarrow Wir brauchen einen Algorithmus zur Berechnung von $\mathbb{P}(\mathbf{X}_C = \mathbf{y})!$

Rückblick (II)

(5) Bei diskreten Modellen ist $\phi(\mathbf{X})_i$ binär (22.06. ff.) und ihr Erwartungswert entspricht der Wahrscheinlichkeit $\mathbb{P}(\phi(\mathbf{X})_i = 1)$ (u.a. 29.06.)

(6) Jede Komponente von $\phi(\mathbf{X})$ entspricht einem Ereignis $\mathbf{X}_C = \mathbf{y}$, d.h. $\mathbb{P}(\phi(\mathbf{X})_i = 1) = \mathbb{P}(\mathbf{X}_C = \mathbf{y})$ mit $\mathbf{y} \in \mathcal{X}_C$

\Rightarrow Wir brauchen einen Algorithmus zur Berechnung von $\mathbb{P}(\mathbf{X}_C = \mathbf{y})!$

Rückblick (II)

(5) Bei diskreten Modellen ist $\phi(\mathbf{X})_i$ binär (22.06. ff.) und ihr Erwartungswert entspricht der Wahrscheinlichkeit $\mathbb{P}(\phi(\mathbf{X})_i = 1)$ (u.a. 29.06.)

(6) Jede Komponente von $\phi(\mathbf{X})$ entspricht einem Ereignis $\mathbf{X}_C = \mathbf{y}$, d.h. $\mathbb{P}(\phi(\mathbf{X})_i = 1) = \mathbb{P}(\mathbf{X}_C = \mathbf{y})$ mit $\mathbf{y} \in \mathcal{X}_C$

\Rightarrow Wir brauchen einen Algorithmus zur Berechnung von $\mathbb{P}(\mathbf{X}_C = \mathbf{y})!$



Loopy Belief Propagation

- Hier: Einschränkung auf Graphen mit $\max_{C \in \mathcal{C}(G)} |C| = 2$
- Dann: Menge der Cliques $\mathcal{C}(G) =$ Menge der Kanten E
- Der Gradient enthält dann $\mathbb{P}(X_{\{v,u\}} = \mathbf{y})$ für alle $\{v, u\} \in E$ und alle $\mathbf{y} \in \mathcal{X}_{\{v,u\}}$

Algorithmus:

- (1) Lege zwei Arrays m^{neu} und m^{alt} an (jeweils $\sum_{\{v,u\} \in E} (|\mathcal{X}|_v + |\mathcal{X}|_u)$ Einträge) und initialisiere mit 1
- (2) Berechne für alle Kanten $\{v, u\}$ in E , beide Kantenrichtungen $s \rightarrow t$ (d.h. $v \rightarrow u$ und $u \rightarrow v$) und jeden Zustand $x \in \mathcal{X}_t$:

$$m_{v \rightarrow u}^{\text{neu}}(x) = \sum_{y \in \mathcal{X}_v} \exp(\langle \beta_{uv}, \phi_{uv}(x, y) \rangle) \prod_{w \in \mathcal{N}(v) \setminus \{u\}} m_{w \rightarrow v}^{\text{alt}}(y)$$

- (3) Falls $\|m^{\text{neu}} - m^{\text{alt}}\| > \varepsilon$: $m^{\text{alt}} \leftarrow m^{\text{neu}}$ und GOTO (2)



Loopy Belief Propagation

- Hier: Einschränkung auf Graphen mit $\max_{C \in \mathcal{C}(G)} |C| = 2$
- Dann: Menge der Cliques $\mathcal{C}(G) =$ Menge der Kanten E
- Der Gradient enthält dann $\mathbb{P}(\mathbf{X}_{\{v,u\}} = \mathbf{y})$ für alle $\{v, u\} \in E$ und alle $\mathbf{y} \in \mathcal{X}_{\{v,u\}}$

Algorithmus:

- (1) Lege zwei Arrays m^{neu} und m^{alt} an (jeweils $\sum_{\{v,u\} \in E} (|\mathcal{X}|_v + |\mathcal{X}|_u)$ Einträge) und initialisiere mit 1
- (2) Berechne für alle Kanten $\{v, u\}$ in E , beide Kantenrichtungen $s \rightarrow t$ (d.h. $v \rightarrow u$ und $u \rightarrow v$) und jeden Zustand $x \in \mathcal{X}_t$:

$$m_{v \rightarrow u}^{\text{neu}}(x) = \sum_{y \in \mathcal{X}_v} \exp(\langle \beta_{uv}, \phi_{uv}(x, y) \rangle) \prod_{w \in \mathcal{N}(v) \setminus \{u\}} m_{w \rightarrow v}^{\text{alt}}(y)$$

- (3) Falls $\|m^{\text{neu}} - m^{\text{alt}}\| > \varepsilon$: $m^{\text{alt}} \leftarrow m^{\text{neu}}$ und GOTO (2)



Loopy Belief Propagation

- Hier: Einschränkung auf Graphen mit $\max_{C \in \mathcal{C}(G)} |C| = 2$
- Dann: Menge der Cliques $\mathcal{C}(G) =$ Menge der Kanten E
- Der Gradient enthält dann $\mathbb{P}(\mathbf{X}_{\{v,u\}} = \mathbf{y})$ für alle $\{v, u\} \in E$ und alle $\mathbf{y} \in \mathcal{X}_{\{v,u\}}$

Algorithmus:

- (1) Lege zwei Arrays m^{neu} und m^{alt} an (jeweils $\sum_{\{v,u\} \in E} (|\mathcal{X}|_v + |\mathcal{X}|_u)$ Einträge) und initialisiere mit 1
- (2) Berechne für alle Kanten $\{v, u\}$ in E , beide Kantenrichtungen $s \rightarrow t$ (d.h. $v \rightarrow u$ und $u \rightarrow v$) und jeden Zustand $x \in \mathcal{X}_t$:

$$m_{v \rightarrow u}^{\text{neu}}(x) = \sum_{y \in \mathcal{X}_v} \exp(\langle \beta_{uv}, \phi_{uv}(x, y) \rangle) \prod_{w \in \mathcal{N}(v) \setminus \{u\}} m_{w \rightarrow v}^{\text{alt}}(y)$$

- (3) Falls $\|m^{\text{neu}} - m^{\text{alt}}\| > \varepsilon$: $m^{\text{alt}} \leftarrow m^{\text{neu}}$ und GOTO (2)



Loopy Belief Propagation

- Hier: Einschränkung auf Graphen mit $\max_{C \in \mathcal{C}(G)} |C| = 2$
- Dann: Menge der Cliques $\mathcal{C}(G) =$ Menge der Kanten E
- Der Gradient enthält dann $\mathbb{P}(\mathbf{X}_{\{v,u\}} = \mathbf{y})$ für alle $\{v, u\} \in E$ und alle $\mathbf{y} \in \mathcal{X}_{\{v,u\}}$

Algorithmus:

- (1) Lege zwei Arrays m^{neu} und m^{alt} an (jeweils $\sum_{\{v,u\} \in E} (|\mathcal{X}_v| + |\mathcal{X}_u|)$ Einträge) und initialisiere mit 1
- (2) Berechne für alle Kanten $\{v, u\}$ in E , beide Kantenrichtungen $s \rightarrow t$ (d.h. $v \rightarrow u$ und $u \rightarrow v$) und jeden Zustand $x \in \mathcal{X}_t$:

$$m_{v \rightarrow u}^{\text{neu}}(x) = \sum_{y \in \mathcal{X}_v} \exp(\langle \beta_{uv}, \phi_{uv}(x, y) \rangle) \prod_{w \in \mathcal{N}(v) \setminus \{u\}} m_{w \rightarrow v}^{\text{alt}}(y)$$

- (3) Falls $\|m^{\text{neu}} - m^{\text{alt}}\| > \varepsilon$: $m^{\text{alt}} \leftarrow m^{\text{neu}}$ und GOTO (2)



Loopy Belief Propagation

- Hier: Einschränkung auf Graphen mit $\max_{C \in \mathcal{C}(G)} |C| = 2$
- Dann: Menge der Cliques $\mathcal{C}(G) =$ Menge der Kanten E
- Der Gradient enthält dann $\mathbb{P}(\mathbf{X}_{\{v,u\}} = \mathbf{y})$ für alle $\{v, u\} \in E$ und alle $\mathbf{y} \in \mathcal{X}_{\{v,u\}}$

Algorithmus:

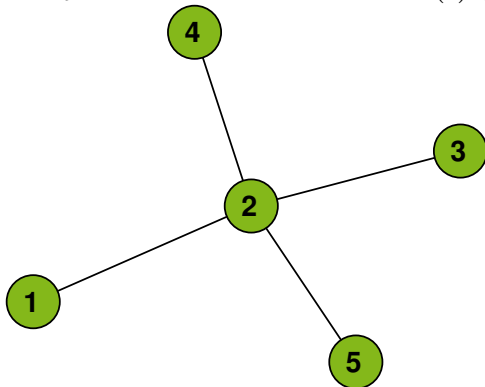
- (1) Lege zwei Arrays m^{neu} und m^{alt} an (jeweils $\sum_{\{v,u\} \in E} (|\mathcal{X}_v| + |\mathcal{X}_u|)$ Einträge) und initialisiere mit 1
- (2) Berechne für alle Kanten $\{v, u\}$ in E , beide Kantenrichtungen $s \rightarrow t$ (d.h. $v \rightarrow u$ und $u \rightarrow v$) und jeden Zustand $x \in \mathcal{X}_t$:

$$m_{v \rightarrow u}^{\text{neu}}(x) = \sum_{y \in \mathcal{X}_v} \exp(\langle \beta_{uv}, \phi_{uv}(x, y) \rangle) \prod_{w \in \mathcal{N}(v) \setminus \{u\}} m_{w \rightarrow v}^{\text{alt}}(y)$$

- (3) Falls $\|m^{\text{neu}} - m^{\text{alt}}\| > \varepsilon$: $m^{\text{alt}} \leftarrow m^{\text{neu}}$ und GOTO (2)

Loopy Belief Propagation: Visualisierung (1. Iteration)

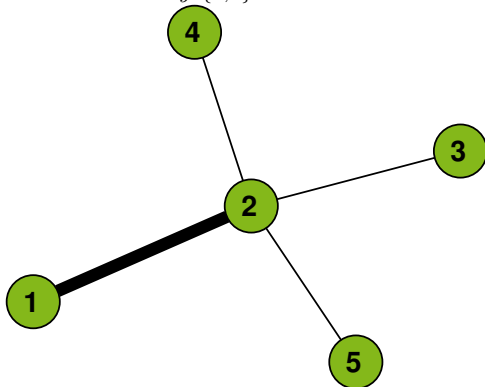
$$m_{v \rightarrow u}^{\text{neu}}(x) = \sum_{y \in \mathcal{X}_v} \exp(\langle \beta_{uv}, \phi_{uv}(x, y) \rangle) \prod_{w \in \mathcal{N}(v) \setminus \{u\}} m_{w \rightarrow v}^{\text{alt}}(y)$$



$$\mathcal{X}_v = \{0, 1\}, v \in \{1, 2, \dots, 5\}$$

Loopy Belief Propagation: Visualisierung (1. Iteration)

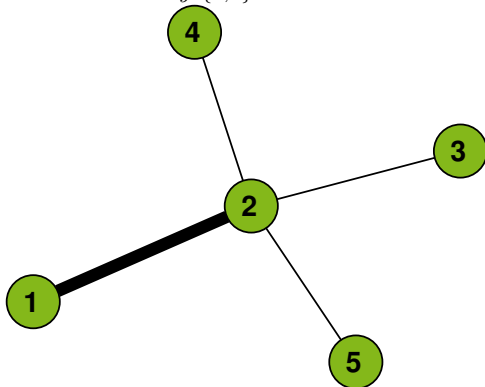
$$m_{1 \rightarrow 2}^{\text{neu}}(0) = \sum_{y \in \{0,1\}} \exp(\langle \beta_{12}, \phi_{12}(y, 0) \rangle)$$



$$\mathcal{X}_v = \{0, 1\}, v \in \{1, 2, \dots, 5\}$$

Loopy Belief Propagation: Visualisierung (1. Iteration)

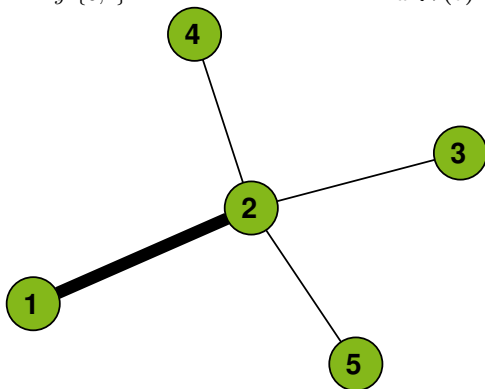
$$m_{1 \rightarrow 2}^{\text{neu}}(1) = \sum_{y \in \{0,1\}} \exp(\langle \beta_{12}, \phi_{12}(y, 1) \rangle)$$



$$\mathcal{X}_v = \{0, 1\}, v \in \{1, 2, \dots, 5\}$$

Loopy Belief Propagation: Visualisierung (1. Iteration)

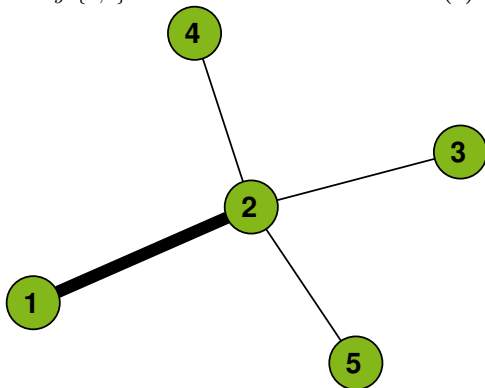
$$m_{2 \rightarrow 1}^{\text{neu}}(0) = \sum_{y \in \{0,1\}} \exp(\langle \beta_{12}, \phi_{12}(0, y) \rangle) \prod_{w \in \mathcal{N}(v) \setminus \{1\}} m_{w \rightarrow 2}^{\text{alt}}(y)$$



$$\mathcal{X}_v = \{0, 1\}, v \in \{1, 2, \dots, 5\}$$

Loopy Belief Propagation: Visualisierung (1. Iteration)

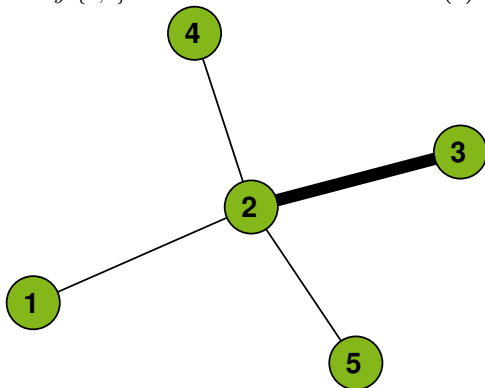
$$m_{2 \rightarrow 1}^{\text{neu}}(1) = \sum_{y \in \{0,1\}} \exp(\langle \beta_{12}, \phi_{12}(1, y) \rangle) \prod_{w \in \mathcal{N}(v) \setminus \{1\}} m_{w \rightarrow 2}^{\text{alt}}(y)$$



$$\mathcal{X}_v = \{0, 1\}, v \in \{1, 2, \dots, 5\}$$

Loopy Belief Propagation: Visualisierung (1. Iteration)

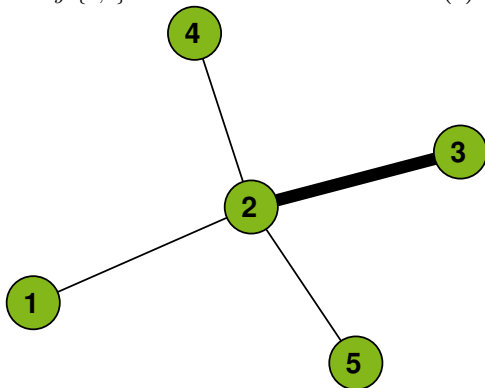
$$m_{2 \rightarrow 3}^{\text{neu}}(0) = \sum_{y \in \{0,1\}} \exp(\langle \beta_{23}, \phi_{23}(y, 0) \rangle) \prod_{w \in \mathcal{N}(v) \setminus \{3\}} m_{w \rightarrow 2}^{\text{alt}}(y)$$



$$\mathcal{X}_v = \{0, 1\}, v \in \{1, 2, \dots, 5\}$$

Loopy Belief Propagation: Visualisierung (1. Iteration)

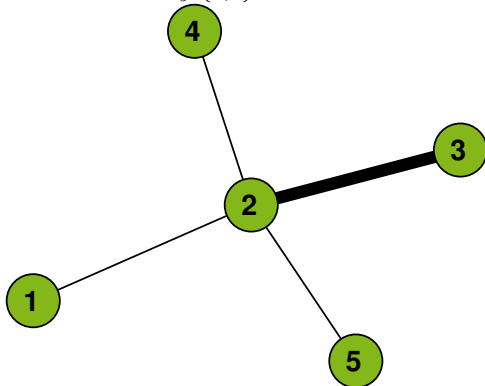
$$m_{2 \rightarrow 3}^{\text{neu}}(1) = \sum_{y \in \{0,1\}} \exp(\langle \beta_{23}, \phi_{23}(y, 1) \rangle) \prod_{w \in \mathcal{N}(v) \setminus \{3\}} m_{w \rightarrow 2}^{\text{alt}}(y)$$



$$\mathcal{X}_v = \{0, 1\}, v \in \{1, 2, \dots, 5\}$$

Loopy Belief Propagation: Visualisierung (1. Iteration)

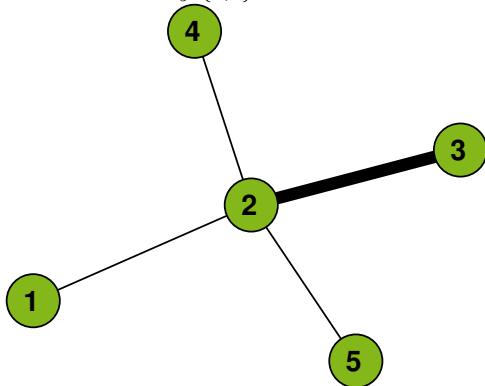
$$m_{3 \rightarrow 2}^{\text{neu}}(0) = \sum_{y \in \{0,1\}} \exp(\langle \beta_{23}, \phi_{23}(0, y) \rangle)$$



$$\mathcal{X}_v = \{0, 1\}, v \in \{1, 2, \dots, 5\}$$

Loopy Belief Propagation: Visualisierung (1. Iteration)

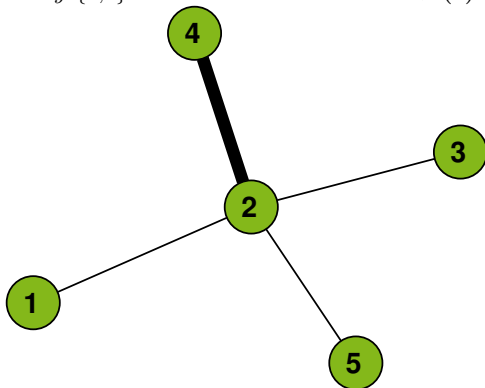
$$m_{3 \rightarrow 2}^{\text{neu}}(1) = \sum_{y \in \{0,1\}} \exp(\langle \beta_{23}, \phi_{23}(1, y) \rangle)$$



$$\mathcal{X}_v = \{0, 1\}, v \in \{1, 2, \dots, 5\}$$

Loopy Belief Propagation: Visualisierung (1. Iteration)

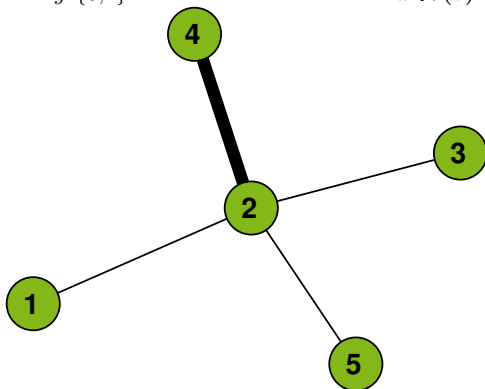
$$m_{2 \rightarrow 4}^{\text{neu}}(0) = \sum_{y \in \{0,1\}} \exp(\langle \beta_{24}, \phi_{24}(y, 0) \rangle) \prod_{w \in \mathcal{N}(v) \setminus \{4\}} m_{w \rightarrow 2}^{\text{alt}}(y)$$



$$\mathcal{X}_v = \{0, 1\}, v \in \{1, 2, \dots, 5\}$$

Loopy Belief Propagation: Visualisierung (1. Iteration)

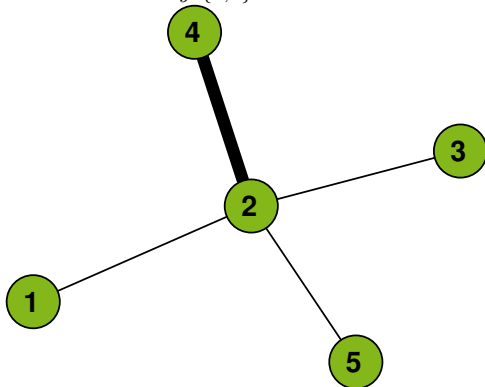
$$m_{2 \rightarrow 4}^{\text{neu}}(1) = \sum_{y \in \{0,1\}} \exp(\langle \beta_{24}, \phi_{24}(y, 1) \rangle) \prod_{w \in \mathcal{N}(v) \setminus \{4\}} m_{w \rightarrow 2}^{\text{alt}}(y)$$



$$\mathcal{X}_v = \{0, 1\}, v \in \{1, 2, \dots, 5\}$$

Loopy Belief Propagation: Visualisierung (1. Iteration)

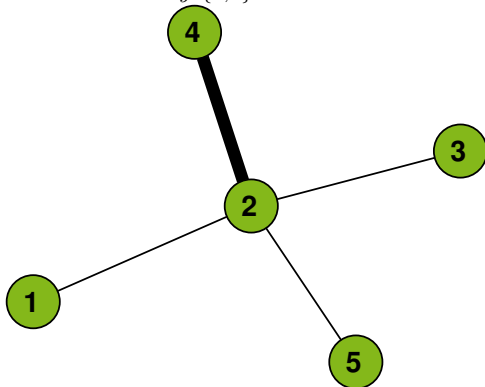
$$m_{4 \rightarrow 2}^{\text{neu}}(0) = \sum_{y \in \{0,1\}} \exp(\langle \beta_{24}, \phi_{24}(0, y) \rangle)$$



$$\mathcal{X}_v = \{0, 1\}, v \in \{1, 2, \dots, 5\}$$

Loopy Belief Propagation: Visualisierung (1. Iteration)

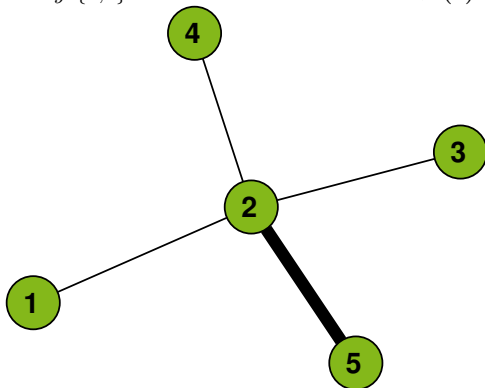
$$m_{4 \rightarrow 2}^{\text{neu}}(1) = \sum_{y \in \{0,1\}} \exp(\langle \beta_{24}, \phi_{24}(1, y) \rangle)$$



$$\mathcal{X}_v = \{0, 1\}, v \in \{1, 2, \dots, 5\}$$

Loopy Belief Propagation: Visualisierung (1. Iteration)

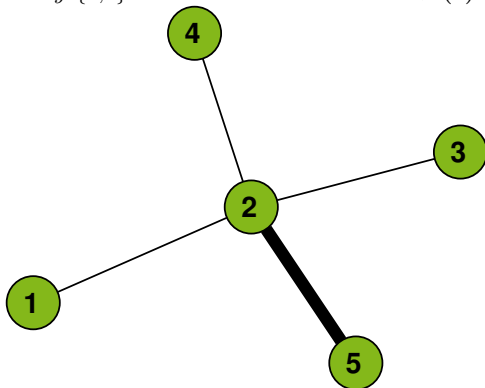
$$m_{2 \rightarrow 5}^{\text{neu}}(0) = \sum_{y \in \{0,1\}} \exp(\langle \beta_{25}, \phi_{25}(y, 0) \rangle) \prod_{w \in \mathcal{N}(v) \setminus \{5\}} m_{w \rightarrow 2}^{\text{alt}}(y)$$



$$\mathcal{X}_v = \{0, 1\}, v \in \{1, 2, \dots, 5\}$$

Loopy Belief Propagation: Visualisierung (1. Iteration)

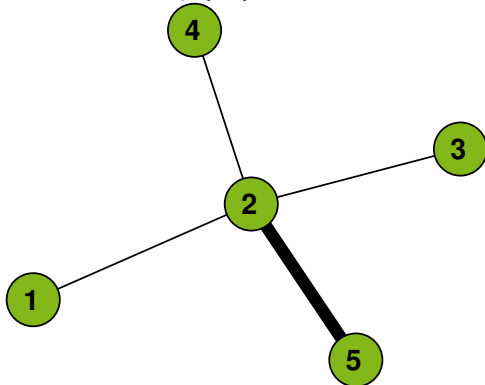
$$m_{2 \rightarrow 5}^{\text{neu}}(1) = \sum_{y \in \{0,1\}} \exp(\langle \beta_{25}, \phi_{25}(y, 1) \rangle) \prod_{w \in \mathcal{N}(v) \setminus \{5\}} m_{w \rightarrow 2}^{\text{alt}}(y)$$



$$\mathcal{X}_v = \{0, 1\}, v \in \{1, 2, \dots, 5\}$$

Loopy Belief Propagation: Visualisierung (1. Iteration)

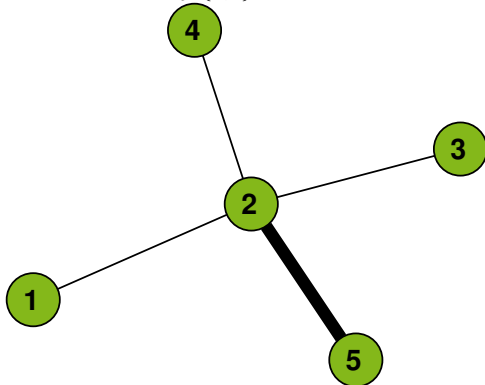
$$m_{5 \rightarrow 2}^{\text{neu}}(0) = \sum_{y \in \{0,1\}} \exp(\langle \beta_{25}, \phi_{25}(0, y) \rangle)$$



$$\mathcal{X}_v = \{0, 1\}, v \in \{1, 2, \dots, 5\}$$

Loopy Belief Propagation: Visualisierung (1. Iteration)

$$m_{5 \rightarrow 2}^{\text{neu}}(1) = \sum_{y \in \{0,1\}} \exp(\langle \beta_{25}, \phi_{25}(1, y) \rangle)$$



$$\mathcal{X}_v = \{0, 1\}, v \in \{1, 2, \dots, 5\}$$



Loopy Belief Propagation: “Konvergenz”

Nachdem m^{neu} basierend auf m^{alt} berechnet wurde, überprüfen wir die Konvergenz mittels

$$\|m^{\text{neu}} - m^{\text{alt}}\| > \varepsilon$$

Die obige Bedingung kann für jedes $\varepsilon \geq 0$ erfüllt werden, falls

- der Graph G ein Baum ist

Andernfalls muss ε “groß genug” (anwendungsspezifisch) gewählt werden.



Loopy Belief Propagation: “Konvergenz”

Nachdem m^{neu} basierend auf m^{neu} berechnet wurde, überprüfen wir die Konvergenz mittels

$$\|m^{\text{neu}} - m^{\text{alt}}\| > \varepsilon$$

Die obige Bedingung kann für jedes $\varepsilon \geq 0$ erfüllt werden, falls

- der Graph G ein Baum ist

Andernfalls muss ε “groß genug” (anwendungsspezifisch) gewählt werden.

Strukturlernen

“Woher kommt der Graph?”

- **Gegeben:** Menge an Variablen V , Datensatz $\mathcal{D} \rightarrow \tilde{\mathbb{E}}[\phi(\mathbf{X})] = \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \phi(\mathbf{x})$
- **Idee:** Ausnutzen bedingter Unabhängigkeiten zwischen den Variablen

Erinnerung:

$$\mathbf{X}_v \perp\!\!\!\perp \mathbf{X}_{V \setminus (\mathcal{N}(v) \cup \{v\})} \mid \mathbf{X}_{\mathcal{N}(v)}, \forall v \in V$$

$$\mathbf{X}_U \perp\!\!\!\perp \mathbf{X}_W \mid \mathbf{X}_S$$

$\forall U, W \subset V$, falls jeder Pfad von U nach W durch S führt.



Strukturlernen

“Woher kommt der Graph?”

- **Gegeben:** Menge an Variablen V , Datensatz $\mathcal{D} \rightarrow \tilde{\mathbb{E}}[\phi(\mathbf{X})] = \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \phi(\mathbf{x})$
- **Idee:** Ausnutzen bedingter Unabhängigkeiten zwischen den Variablen

Erinnerung:

$$\mathbf{X}_v \perp\!\!\!\perp \mathbf{X}_{V \setminus (\mathcal{N}(v) \cup \{v\})} \mid \mathbf{X}_{\mathcal{N}(v)}, \forall v \in V$$

$$\mathbf{X}_U \perp\!\!\!\perp \mathbf{X}_W \mid \mathbf{X}_S$$

$\forall U, W \subset V$, falls jeder Pfad von U nach W durch S führt.



“Woher kommt der Graph?”

- **Gegeben:** Menge an Variablen V , Datensatz $\mathcal{D} \rightarrow \tilde{\mathbb{E}}[\phi(\mathbf{X})] = \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \phi(\mathbf{x})$
- **Idee:** Ausnutzen bedingter Unabhängigkeiten zwischen den Variablen

Erinnerung:

$$\mathbf{X}_v \perp\!\!\!\perp \mathbf{X}_{V \setminus (\mathcal{N}(v) \cup \{v\})} \mid \mathbf{X}_{\mathcal{N}(v)}, \forall v \in V$$

$$\mathbf{X}_U \perp\!\!\!\perp \mathbf{X}_W \mid \mathbf{X}_S$$

$\forall U, W \subset V$, falls jeder Pfad von U nach W durch S führt.



Option1: Chow-Liu Bäume

Idee: Finde besten **Baum** T , so dass Abstand (Kullback-Leiber Divergenz) zwischen echter Funktion \mathbb{P}^* und \mathbb{P}_T minimal ist.

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

Kann umgeformt werden zu

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



Option1: Chow-Liu Bäume

Idee: Finde besten **Baum** T , so dass Abstand (Kullback-Leiber Divergenz) zwischen echter Funktion \mathbb{P}^* und \mathbb{P}_T minimal ist.

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

Kann umgeformt werden zu

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



Option1: Chow-Liu Bäume

Idee: Finde besten **Baum** T , so dass Abstand (Kullback-Leiber Divergenz) zwischen echter Funktion \mathbb{P}^* und \mathbb{P}_T minimal ist.

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

Kann umgeformt werden zu

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



Mutual Information

Erinnerung: Zwei Variablen \mathbf{X}_v und \mathbf{X}_u sind unabhängig, falls $\mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y) = \mathbb{P}(\mathbf{X}_v = x)\mathbb{P}(\mathbf{X}_u = y)$ für alle x, y . Das heißt falls

$$\frac{\mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y)}{\mathbb{P}(\mathbf{X}_v = x)\mathbb{P}(\mathbf{X}_u = y)} = 1, \forall x \in \mathcal{X}_v, y \in \mathcal{X}_y$$

Das gilt möglicherweise nicht für alle x, y , **aber:**
Allgemeines Maß für Unabhängigkeit durch Gewichtung jedes Paares x, y mittels $\mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y)$

$$\sum_{x \in \mathcal{X}_v} \sum_{y \in \mathcal{X}_u} \mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y) \frac{\mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y)}{\mathbb{P}(\mathbf{X}_v = x)\mathbb{P}(\mathbf{X}_u = y)} = I(\mathbf{X}_v, \mathbf{X}_u)$$

(I heißt “Mutual Information”)



Mutual Information

Erinnerung: Zwei Variablen \mathbf{X}_v und \mathbf{X}_u sind unabhängig, falls $\mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y) = \mathbb{P}(\mathbf{X}_v = x)\mathbb{P}(\mathbf{X}_u = y)$ für alle x, y . Das heißt falls

$$\frac{\mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y)}{\mathbb{P}(\mathbf{X}_v = x)\mathbb{P}(\mathbf{X}_u = y)} = 1, \forall x \in \mathcal{X}_v, y \in \mathcal{X}_y$$

Das gilt möglicherweise nicht für alle x, y , **aber:**

Allgemeines Maß für Unabhängigkeit durch Gewichtung jedes Paares x, y mittels $\mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y)$

$$\sum_{x \in \mathcal{X}_v} \sum_{y \in \mathcal{X}_u} \mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y) \frac{\mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y)}{\mathbb{P}(\mathbf{X}_v = x)\mathbb{P}(\mathbf{X}_u = y)} = I(\mathbf{X}_v, \mathbf{X}_u)$$

(I heißt “Mutual Information”)



Mutual Information

Erinnerung: Zwei Variablen \mathbf{X}_v und \mathbf{X}_u sind unabhängig, falls $\mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y) = \mathbb{P}(\mathbf{X}_v = x)\mathbb{P}(\mathbf{X}_u = y)$ für alle x, y . Das heißt falls

$$\frac{\mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y)}{\mathbb{P}(\mathbf{X}_v = x)\mathbb{P}(\mathbf{X}_u = y)} = 1, \forall x \in \mathcal{X}_v, y \in \mathcal{X}_y$$

Das gilt möglicherweise nicht für alle x, y , **aber:**
Allgemeines Maß für Unabhängigkeit durch Gewichtung jedes Paares x, y mittels $\mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y)$

$$\sum_{x \in \mathcal{X}_v} \sum_{y \in \mathcal{X}_u} \mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y) \frac{\mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y)}{\mathbb{P}(\mathbf{X}_v = x)\mathbb{P}(\mathbf{X}_u = y)} = I(\mathbf{X}_v, \mathbf{X}_u)$$

(I heißt “Mutual Information”)



Mutual Information

Erinnerung: Zwei Variablen \mathbf{X}_v und \mathbf{X}_u sind unabhängig, falls $\mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y) = \mathbb{P}(\mathbf{X}_v = x)\mathbb{P}(\mathbf{X}_u = y)$ für alle x, y . Das heißt falls

$$\frac{\mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y)}{\mathbb{P}(\mathbf{X}_v = x)\mathbb{P}(\mathbf{X}_u = y)} = 1, \forall x \in \mathcal{X}_v, y \in \mathcal{X}_y$$

Das gilt möglicherweise nicht für alle x, y , **aber:**
Allgemeines Maß für Unabhängigkeit durch Gewichtung jedes Paares x, y mittels $\mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y)$

$$\sum_{x \in \mathcal{X}_v} \sum_{y \in \mathcal{X}_u} \mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y) \frac{\mathbb{P}(\mathbf{X}_v = x, \mathbf{X}_u = y)}{\mathbb{P}(\mathbf{X}_v = x)\mathbb{P}(\mathbf{X}_u = y)} = I(\mathbf{X}_v, \mathbf{X}_u)$$

(I heißt “Mutual Information”)



Option1: Chow-Liu Bäume

Idee: Finde besten **Baum** T , so dass Abstand (Kullback-Leiber Divergenz) zwischen echter Funktion \mathbb{P}^* und \mathbb{P}^T minimal ist.

$$\min_T \text{KL}(\mathbb{P}^*, \mathbb{P}_T) = \max_T \sum_{vu \in E(T)} I(\mathbf{X}_v, \mathbf{X}_u)$$

Berechnung beispielsweise mittels Algorithmus von Kruskal
(Maximaler Spannbaum)



Option2: Graphen mittels Regularisierung

Baobachtung: Falls Parametervektor $\beta_{\{v,u\}}$ einer Kante = 0, so hat diese Kante keinen Einfluss auf $\mathbb{P}(\mathbf{X} = \mathbf{x})!$

$$\mathbb{P}(\mathbf{X} = \mathbf{x}) = \frac{1}{Z(\boldsymbol{\beta})} \exp(\langle \boldsymbol{\beta}, \phi(\mathbf{x}) \rangle) = \frac{\prod_{C \in \mathcal{C}(G)} \exp(\langle \boldsymbol{\beta}_C, \phi_C(\mathbf{x}_C) \rangle)}{\sum_{\mathbf{x}' \in \mathcal{X}} \prod_{C \in \mathcal{C}(G)} \exp(\langle \boldsymbol{\beta}_C, \phi_C(\mathbf{x}'_C) \rangle)}$$

Falls alle Cliquenfaktoren in paarweise Funktionen zerfallen:

$$\mathbb{P}(\mathbf{X} = \mathbf{x}) = \frac{\prod_{\{v,u\} \in E} \exp(\langle \boldsymbol{\beta}_{\{v,u\}}, \phi_{\{v,u\}}(\mathbf{x}_{\{v,u\}}) \rangle)}{\sum_{\mathbf{x}' \in \mathcal{X}} \prod_{\{v,u\} \in E} \exp(\langle \boldsymbol{\beta}_{\{v,u\}}, \phi_{\{v,u\}}(\mathbf{x}'_{\{v,u\}}) \rangle)}$$

Falls $\boldsymbol{\beta}_{\{v,u\}} = \mathbf{0}$ für eine Kante $\{v, u\}$, so ist $\exp(\langle \boldsymbol{\beta}_{\{v,u\}}, \phi_{\{v,u\}}(\mathbf{x}_{\{v,u\}}) \rangle) = 1$.



Option2: Graphen mittels Regularisierung

Baobachtung: Falls Parametervektor $\beta_{\{v,u\}}$ einer Kante = 0, so hat diese Kante keinen Einfluss auf $\mathbb{P}(\mathbf{X} = \mathbf{x})!$

$$\mathbb{P}(\mathbf{X} = \mathbf{x}) = \frac{1}{Z(\boldsymbol{\beta})} \exp(\langle \boldsymbol{\beta}, \phi(\mathbf{x}) \rangle) = \frac{\prod_{C \in \mathcal{C}(G)} \exp(\langle \boldsymbol{\beta}_C, \phi_C(\mathbf{x}_C) \rangle)}{\sum_{\mathbf{x}' \in \mathcal{X}} \prod_{C \in \mathcal{C}(G)} \exp(\langle \boldsymbol{\beta}_C, \phi_C(\mathbf{x}'_C) \rangle)}$$

Falls alle Cliquenfaktoren in paarweise Funktionen zerfallen:

$$\mathbb{P}(\mathbf{X} = \mathbf{x}) = \frac{\prod_{\{v,u\} \in E} \exp(\langle \boldsymbol{\beta}_{\{v,u\}}, \phi_{\{v,u\}}(\mathbf{x}_{\{v,u\}}) \rangle)}{\sum_{\mathbf{x}' \in \mathcal{X}} \prod_{\{v,u\} \in E} \exp(\langle \boldsymbol{\beta}_{\{v,u\}}, \phi_{\{v,u\}}(\mathbf{x}'_{\{v,u\}}) \rangle)}$$

Falls $\beta_{\{v,u\}} = \mathbf{0}$ für eine Kante $\{v, u\}$, so ist $\exp(\langle \boldsymbol{\beta}_{\{v,u\}}, \phi_{\{v,u\}}(\mathbf{x}_{\{v,u\}}) \rangle) = 1$.

Option2: Graphen mittels Regularisierung

Baobachtung: Falls Parametervektor $\beta_{\{v,u\}}$ einer Kante = 0, so hat diese Kante keinen Einfluss auf $\mathbb{P}(\mathbf{X} = \mathbf{x})!$

$$\mathbb{P}(\mathbf{X} = \mathbf{x}) = \frac{1}{Z(\beta)} \exp(\langle \beta, \phi(\mathbf{x}) \rangle) = \frac{\prod_{C \in \mathcal{C}(G)} \exp(\langle \beta_C, \phi_C(\mathbf{x}_C) \rangle)}{\sum_{\mathbf{x}' \in \mathcal{X}} \prod_{C \in \mathcal{C}(G)} \exp(\langle \beta_C, \phi_C(\mathbf{x}'_C) \rangle)}$$

Falls alle Cliquenfaktoren in paarweise Funktionen zerfallen:

$$\mathbb{P}(\mathbf{X} = \mathbf{x}) = \frac{\prod_{\{v,u\} \in E} \exp(\langle \beta_{\{v,u\}}, \phi_{\{v,u\}}(\mathbf{x}_{\{v,u\}}) \rangle)}{\sum_{\mathbf{x}' \in \mathcal{X}} \prod_{\{v,u\} \in E} \exp(\langle \beta_{\{v,u\}}, \phi_{\{v,u\}}(\mathbf{x}'_{\{v,u\}}) \rangle)}$$

Falls $\beta_{\{v,u\}} = \mathbf{0}$ für eine Kante $\{v, u\}$, so ist $\exp(\langle \beta_{\{v,u\}}, \phi_{\{v,u\}}(\mathbf{x}_{\{v,u\}}) \rangle) = 1$.



Option2: Graphen mittels Regularisierung

Baobachtung: Falls Parametervektor $\beta_{\{v,u\}}$ einer Kante = 0, so hat diese Kante keinen Einfluss auf $\mathbb{P}(\mathbf{X} = \mathbf{x})!$

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

Links- und rechtsseitige Grenzwerte stimmen am Nullpunkt (0) nicht überein!

$\|\cdot\|_1$ nicht differenzierbar(!) \rightarrow Gradient kann nicht berechnet werden!!

Aber: Proximaler Gradientenabstieg möglich.



Proximaler Gradientenabstieg

Falls Funktion f nicht differenzierbar: Für die Minimierung von

$$F(\beta; \mathcal{D}) = \ell(\beta^t; \mathcal{D}) + f(\beta)$$

Nutzen wir anstatt

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

jetzt

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

Funktion F wird aufgeteilt in den differenzierbaren Teil ℓ und f
Intuition an der Tafel...



Proximaler Gradientenabstieg

Falls Funktion f nicht differenzierbar: Für die Minimierung von

$$F(\beta; \mathcal{D}) = \ell(\beta^t; \mathcal{D}) + f(\beta)$$

Nutzen wir anstatt

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

jetzt

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

Funktion F wird aufgeteilt in den differenzierbaren Teil ℓ und f
Intuition an der Tafel...



Proximaler Gradientenabstieg für l_1 -Regularisierung

Formal:

$$\text{prox}_{\lambda\|\cdot\|_1}(\gamma_i) = \min_a \left\{ \|a\|_1 + \frac{1}{2\lambda} \|a - \gamma_i\|_2^2 \right\}$$

Lösung:

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

Algorithmus zur Bestimmung von G' :

- (1) Wähle vollständigen Graphen G' (mit allen möglichen Kanten)
- (2) Löse $\min_{\beta} \ell(\beta; \mathcal{D}) + \lambda \|\beta\|_1$ mittels proximalem Gradientenabstieg
- (3) Entferne alle Kanten $\{v, u\}$ mit $\beta_{\{v,u\}} = 0$ aus G'



Proximaler Gradientenabstieg für l_1 -Regularisierung

Formal:

$$\text{prox}_{\lambda\|\cdot\|_1}(\gamma_i) = \min_a \left\{ \|a\|_1 + \frac{1}{2\lambda} \|a - \gamma_i\|_2^2 \right\}$$

Lösung:

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

Algorithmus zur Bestimmung von G :

- (1) Wähle vollständigen Graphen G' (mit allen möglichen Kanten)
- (2) Löse $\min_{\beta} \ell(\beta; \mathcal{D}) + \lambda \|\beta\|_1$ mittels proximalem Gradientenabstieg
- (3) Entferne alle Kanten $\{v, u\}$ mit $\beta_{\{v,u\}} = 0$ aus G'



Proximaler Gradientenabstieg für l_1 -Regularisierung

Formal:

$$\text{prox}_{\lambda \|\cdot\|_1}(\gamma_i) = \min_a \left\{ \|a\|_1 + \frac{1}{2\lambda} \|a - \gamma_i\|_2^2 \right\}$$

Lösung:

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

Algorithmus zur Bestimmung von G' :

- (1) Wähle vollständigen Graphen G' (mit allen möglichen Kanten)
- (2) Löse $\min_{\beta} \ell(\beta; \mathcal{D}) + \lambda \|\beta\|_1$ mittels proximalem Gradientenabstieg
- (3) Entferne alle Kanten $\{v, u\}$ mit $\beta_{\{v,u\}} = 0$ aus G'