

Probabilistische Graphische Modelle und Wahrscheinlichkeitsverteilungen

$p(x_1, x_2) = p(x_1)p(x_2)$

- PGM G ist **Visualisierung** von **Wahrscheinlichkeitsverteilung** p
- Zufallsvariablen** x_i entsprechen **Knoten** V_i
- Kanten $E = (V_i, V_j)$ stellen Zusammenhänge zwischen x_i und x_j dar
- G ist ... von p , wenn ...
 - Dependency-Map (D-Map):**
 $x_i \perp x_j \Rightarrow$ es gibt keine Kante $E = (V_i, V_j)$ im Graph G
 - Independency-Map (I-Map):**
Es gibt keine Kante $E = (V_i, V_j) \Rightarrow x_i \perp x_j$
 - Perfect Map:**
 G ist D-Map und I-Map

\Rightarrow Wir sind in der Regel an I-Maps interessiert, wenn wir Modelle graphisch erstellen.

(Beierle/Kern-Isberner 2003, Koller/Friedman 2010)

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 8 / 65

Gerichtete Probabilistische Graphische Modelle: Faktorisierung und Visualisierung

- Wie kommen wir nun von einer Verteilung zu einem Graph – und umgekehrt?
- Idee: Faktorisiere Verteilung mit der **Kettenregel**: $p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | x_{i-1}, \dots, x_n)$
- $p(x_1, x_2, x_3, x_4, x_5) = p(x_1 | x_2, x_3, x_4, x_5) \cdot p(x_2 | x_3, x_4, x_5) \cdot p(x_3 | x_4, x_5) \cdot p(x_4 | x_5) \cdot p(x_5)$

- Diese Darstellung gilt für jede Verteilung und ist daher wenig hilfreich, ...
- ...aber wir können mit dem graphischen Modell **Unabhängigkeiten implizieren**.

\Rightarrow **Weniger Parameter, mit Vorannahmen angepasstes Modell**

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 9 / 65

Gerichtete Probabilistische Graphische Modelle: Unabhängigkeitsannahmen

$p(x_1, x_2, x_3, x_4, x_5) = p(x_1 | x_2) \cdot p(x_2 | x_3) \cdot p(x_3 | x_4, x_5) \cdot p(x_4 | x_5) \cdot p(x_5)$

- Erinnerung: Wir nutzen den graphischen Formalismus um Aussagen über Annahmen in der Wahrscheinlichkeitsverteilung zu machen.

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 10 / 65

Populäre Beispiel-PGMs

Naïve Bayes

$p(y, \vec{x}) = p(y) \prod_{i=1}^n p(x_i | y)$

(Domingos, 1997)

Hidden Markov Model

$p(\vec{y}, \vec{x}) = p(y_1) \prod_{i=1}^n p(y_i | y_{i-1}) p(x_i | y_i)$

(Rabiner, 1989)

Latent Dirichlet Allocation

$p(\vec{w}, \vec{z}, \vec{\theta}, \alpha, \beta) = p(\beta) \cdot \prod_{j=1}^M p(\theta_j | \alpha) \cdot \prod_{i=1}^N (p(z_{j,t} | \theta_j) p(w_{j,t} | z_{j,t}))$

(Blei, 2003)

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 11 / 65

Definition: Gerichtetes Graphisches Modell

Ein gerichtetes graphisches Modell, oder **Bayesisches Netzwerk**:

- ist ein gerichteter azyklischer Graph
 - Knoten** entsprechen **Variablen**
- Jeder Knoten hat eine **bedingte Wahrscheinlichkeitsverteilung**: $p(x_i | \text{Eltern}(x_i))$
- Das Netzwerk stellt dann eine gemeinsame Wahrscheinlichkeitsverteilung dar:

$$p(x_1, \dots, x_n) = \prod_i p(x_i | \text{Eltern}(x_i))$$

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 12 / 65

Definition: Ungerichtete Probabilistische Graphische Modelle

- Ein **paarweises Markov-Netzwerk** ist ein **ungerichteter Graph** mit Knoten, welche Variablen x_1, \dots, x_n darstellen
 - Kante $E = (x_i, x_j)$ ist mit **Potentialfunktion** $\Psi(x_i, x_j)$ verknüpft
- Ein **Markov-Netzwerk** G mit **Cliquen-Faktorisierung** hat Potentialfunktionen welche (maximalen) Cliques entsprechen

$$p(\vec{x}) = \frac{1}{Z} \prod_{C \in \text{Cliques}(G)} \Psi_C(\vec{x}_C)$$

- Typische Formulierung von Potentialfunktionen:
 $\Psi_i(\vec{x}) = \exp(\sum_j \lambda_j f_j(\vec{x}))$

(Koller/Friedman 2010)

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 13 / 65

Faktor-Graph

Ein **Faktor-Graph** ist ein bipartiter Graph über Faktoren und Variablen

- Faktor Ψ_i berechnet ein Skalar über alle Variablen
- Seien \vec{x} Zufallsvariablen
- $\Psi_i(\vec{x}_i) = \exp\left(\sum_k \lambda_{ki} f_{ki}(\vec{x}_i)\right)$
- Wahrscheinlichkeitsverteilung:
 $p(\vec{x}) = \frac{1}{Z} \prod_i \Psi_i(\vec{x}_i)$

(Bishop, 2006; Kschischang, 2001)

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 14 / 65

Conditional Random Field als Faktorgraph

Conditional Random Field: Markovnetzwerk mit zusätzlichen, immer bekannten Parametern.

- Faktor Ψ_i berechnet ein Skalar
- Seien \vec{x} Eingabe- und \vec{y} Ausgabevariablen
- $\Psi_i(\vec{x}_i, \vec{y}_i) = \exp\left(\sum_k \lambda_{ki} f_{ki}(\vec{x}_i, \vec{y}_i)\right)$
- Wahrscheinlichkeitsverteilung:
 $p(\vec{y} | \vec{x}) = \frac{1}{Z(\vec{x})} \prod_i \Psi_i(\vec{x}_i, \vec{y}_i)$

$\vec{x} = (x_1, x_2)^T$
 $\vec{y} = (y_1, y_2, y_3, y_4)^T$

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 15 / 65

Einführung und Motivation Methodischer Überblick Verhältnis von PCM und KNN Integration von KNN in PCM Zusammenfassung/Ausblick

Lernen und Inferenz

Lernen

- Gerichtete Modelle:** Bestimmung der bedingten Wahrscheinlichkeiten durch Abzählen in Trainingsdaten \mathcal{D}
- Ungerichtete Modelle:** Iterative Maximierung von $\log \sum_{\vec{x} \in \mathcal{D}} p_{\vec{\lambda}}(\vec{x})$

Inferenzaufgaben

- Berechnen von Wahrscheinlichkeitsverteilungen einzelner Teilmengen von Knoten: **Sum-Product-Algorithm** (später, Spezialfall: Forward-Backward)
- Berechnen der Belegung, welche die Gesamtwahrscheinlichkeit maximiert: **Max-Product-Algorithm** (Spezialfall: Viterbi, kleine Variation zu Sum-Product)
- ⇒ **Belief-Propagation**

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 16 / 65

Einführung und Motivation Methodischer Überblick Verhältnis von PCM und KNN Integration von KNN in PCM Zusammenfassung/Ausblick

Beispielfaktorgraph

$\Psi(x=1, y=1) = 30$
 $\Psi(x=1, y=0) = 5$
 $\Psi(x=0, y=1) = 1$
 $\Psi(x=0, y=0) = 10$

$\Psi(x=1, z=1) = 100$
 $\Psi(x=1, z=0) = 1$
 $\Psi(x=0, z=1) = 1$
 $\Psi(x=0, z=0) = 100$

$\Psi(y=1, w=1) = 100$
 $\Psi(y=1, w=0) = 1$
 $\Psi(y=0, w=1) = 1$
 $\Psi(y=0, w=0) = 100$

$\Psi(z=1, w=1) = 1$
 $\Psi(z=1, w=0) = 100$
 $\Psi(z=0, w=1) = 100$
 $\Psi(z=0, w=0) = 1$

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 17 / 65

Einführung und Motivation Methodischer Überblick Verhältnis von PCM und KNN Integration von KNN in PCM Zusammenfassung/Ausblick

Beispielfaktorgraph

$\Psi(x=1, y=1) = 30$
 $\Psi(x=1, y=0) = 5$
 $\Psi(x=0, y=1) = 1$
 $\Psi(x=0, y=0) = 10$

$\Psi(x=1, z=1) = 100$
 $\Psi(x=1, z=0) = 1$
 $\Psi(x=0, z=1) = 1$
 $\Psi(x=0, z=0) = 100$

$\Psi(y=1, w=1) = 100$
 $\Psi(y=1, w=0) = 1$
 $\Psi(y=0, w=1) = 1$
 $\Psi(y=0, w=0) = 100$

$\Psi(z=1, w=1) = 1$
 $\Psi(z=1, w=0) = 100$
 $\Psi(z=0, w=1) = 100$
 $\Psi(z=0, w=0) = 1$

Variablenbelegung				Wert	Wahrscheinlichkeit
x	y	z	w		
0	0	0	0	100000	0.0159445
0	0	0	1	100000	0.0159445
0	0	1	0	100000	0.0159445
0	0	1	1	10	0.0000316
0	1	0	0	10000	0.0015945
0	1	0	1	100000	0.0159445
0	1	1	0	100	0.0000159
0	1	1	1	100	0.0000159
1	0	0	0	500	0.0000797
1	0	0	1	500	0.0000797
1	0	1	0	5000000	0.7972269
1	0	1	1	50000	0.0079723
1	1	0	0	30	0.0000046
1	1	0	1	300000	0.0478336
1	1	1	0	300000	0.0478336
1	1	1	1	300000	0.0478336

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 18 / 65

Einführung und Motivation Methodischer Überblick Verhältnis von PCM und KNN Integration von KNN in PCM Zusammenfassung/Ausblick

Beispiel-Faktorgraphen

Naïve Bayes Hidden Markov Model Linear-Chain Conditional Random Field

(Lafferty, 2001)

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 19 / 65

Einführung und Motivation Methodischer Überblick Verhältnis von PCM und KNN Integration von KNN in PCM Zusammenfassung/Ausblick

Inferenz auf einer Kette (I)

$p(\vec{x}) = \frac{1}{Z} \Psi_{1,2}(x_1, x_2) \Psi_{2,3}(x_2, x_3) \dots \Psi_{n-1,n}(x_{n-1}, x_n)$

- Annahme: Variablen können k diskrete Werte annehmen
- Jede Potentialfunktion hat also k^2 Parameter
- Anzahl aller Parameter für die gemeinsame Verteilung entsprechend der Faktoren: $(n-1)k^2$
- Inferenz der Randverteilung durch Aussummieren von x_i :

$$p(x_i) = \sum_{x_1} \dots \sum_{x_{i-1}} \sum_{x_{i+1}} \dots \sum_{x_n} p(\vec{x})$$

- k^n Werte für \vec{x} - schade.

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 20 / 65

Einführung und Motivation Methodischer Überblick Verhältnis von PCM und KNN Integration von KNN in PCM Zusammenfassung/Ausblick

Inferenz auf einer Kette (II)

$p(\vec{x}) = \frac{1}{Z} \Psi_{1,2}(x_1, x_2) \Psi_{2,3}(x_2, x_3) \dots \Psi_{n-1,n}(x_{n-1}, x_n)$

Idee: Umsortieren dieser Berechnung:

$$p(x_i) = \frac{1}{Z} \left[\sum_{x_{i-1}} \Psi_{i-1,i}(x_{i-1}, x_i) \dots \left[\sum_{x_2} \Psi_{2,3}(x_2, x_3) \left[\sum_{x_1} \Psi_{1,2}(x_1, x_2) \right] \right] \dots \right]$$

⇒ $O(nk^2)$ Berechnungen, wie wir gleich sehen werden.

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 21 / 65

Einführung und Motivation Methodischer Überblick Verhältnis von PCM und KNN Integration von KNN in PCM Zusammenfassung/Ausblick

Inferenz auf einer Kette (III)

$$p(x_i) = \frac{1}{Z} \left[\sum_{x_{i-1}} \Psi_{i-1,i}(x_{i-1}, x_i) \dots \left[\sum_{x_2} \Psi_{2,3}(x_2, x_3) \left[\sum_{x_1} \Psi_{1,2}(x_1, x_2) \right] \right] \dots \right]$$

$\mu_{\alpha}(x_i) = \sum_{x_{i-1}} \Psi_{i-1,i}(x_{i-1}, x_i) \left[\sum_{x_{i-2}} \dots \right] = \sum_{x_{i-1}} \Psi_{i-1,i}(x_{i-1}, x_i) \mu_{\alpha}(x_{i-1})$

$\mu_{\beta}(x_i) = \sum_{x_{i+1}} \Psi_{i,i+1}(x_i, x_{i+1}) \left[\sum_{x_{i+2}} \dots \right] = \sum_{x_{i+1}} \Psi_{i,i+1}(x_i, x_{i+1}) \mu_{\beta}(x_{i+1})$

- Jede Nachricht ist eine Menge von k Werten
- Eine Wahrscheinlichkeit für jeden Wert

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 22 / 65

Einführung und Motivation Methodischer Überblick Verhältnis von PCM und KNN Integration von KNN in PCM Zusammenfassung/Ausblick

Inferenz auf einer Kette (IV)

$\mu_{\alpha}(x_i) = \sum_{x_{i-1}} \Psi_{i-1,i}(x_{i-1}, x_i) \left[\sum_{x_{i-2}} \dots \right] = \sum_{x_{i-1}} \Psi_{i-1,i}(x_{i-1}, x_i) \mu_{\alpha}(x_{i-1})$

$\mu_{\beta}(x_i) = \sum_{x_{i+1}} \Psi_{i,i+1}(x_i, x_{i+1}) \left[\sum_{x_{i+2}} \dots \right] = \sum_{x_{i+1}} \Psi_{i,i+1}(x_i, x_{i+1}) \mu_{\beta}(x_{i+1})$

Rekursive Evaluation von Nachrichten:

- Erste Evaluation: $\mu_{\alpha}(x_2) = \sum_{x_1} \Psi_{1,2}(x_1, x_2)$

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 23 / 65

Stand der Dinge

- Probabilistische Graphische Modelle stellen **Unabhängigkeitsannahmen über Wahrscheinlichkeitsverteilungen** dar
- **Inferenz ist effizient**, wenn Graph ein **Baum** ist
- Bei **beliebigen Graphen** ist Inferenz **NP schwer**: **Approximationsalgorithmen** und **Samplingverfahren** existieren

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 24 / 65

Outline

- 1 Einführung und Motivation
- 2 Methodischer Überblick
 - Probabilistische Graphische Modelle (PGM)
 - Künstliche Neuronale Netze (KNN)
- 3 Verhältnis von PGM und KNN
 - Propagieralgorithmen
 - Verwandtschaft der Formulierungen
 - Hidden CRF
- 4 Integration von KNN in PGM
 - Fallstudien
 - Werkzeuge
- 5 Zusammenfassung und Ausblick

Probabilistische Graphische Modelle und Neuronale Netze

Beide Modelle werden oft als Netzwerk von Knoten mit Kanten dargestellt

Directed Probabilistic Graphical Model

Feed-forward Neural Network

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 26 / 65

Ein sehr einfaches künstliches neuronales "Netzwerk"

Eingabe x_1, x_2, x_3 → $z = \sum_i w_i x_i$ → Ausgabe $y = f(z)$

$$f\left(\sum_i w_i x_i\right) = y$$

- $z = \sum_i w_i x_i$: interner Zustand
- $f(\cdot)$: Aktivierungsfunktion (z.B. linear, tanh, sigmoid, ReLU)

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 27 / 65

ODER mit einem Neuron

OR

x_1	x_2	$\sum w_i x_i$	Sigmoid
0	0	-5	0.007
0	1	5	0.993
1	0	5	0.993
1	1	15	0.999

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 28 / 65

UND mit einem Neuron

AND

x_1	x_2	$\sum w_i x_i$	Sigmoid
0	0	-15	0.000
0	1	-5	0.007
1	0	-5	0.007
1	1	5	0.993

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 29 / 65

ENTWEDER-ODER als Netzwerk von UND, NUND, und ODER

XOR

x_1	x_2	$\sum x_i w_{1,i}$	y_1	$\sum x_i w_{2,i}$	y_2	$\sum x_i w_{3,i}$	y_3
0	0	-10	≈ 0	30	≈ 1	≈ -10	≈ 0
0	1	10	≈ 1	10	≈ 1	≈ 10	≈ 1
1	0	10	≈ 1	10	≈ 1	≈ 10	≈ 1
1	1	30	≈ 1	-10	≈ 0	≈ -10	≈ 0

• XOR = AND(OR(x_1, x_2), NAND(x_1, x_2))

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 30 / 65

Semantik von vorwärtsgerichteten Neuronalen Netzen

Schicht 1 2 3 4

- Feed-forward Neuronale Netze stellen eine Funktionskomposition dar
- $\vec{y} = \vec{f}_4(\vec{f}_3(\vec{f}_2(\vec{f}_1(\vec{x}, W_1), W_2), W_3), W_4)$
- Vorwärtspropagierung durch das Netz: Ebenenweise Berechnung der jeweiligen Ausgabewerte

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 31 / 65

Methodischer Überblick

Inferenz und Lernen

- Inferenz:** **Forward Propagation**, berechne Ergebnisse Ebene für Ebene
- Lernen:** Gewichtsanzpassung durch Minimierung eines Fehlers der Vorhersage auf Trainingsdaten:

$$E(\mathbf{W}, \mathcal{D}) = \sum_{\tilde{x}, \tilde{y} \in \mathcal{D}} \frac{1}{2} \|\tilde{y}(\tilde{x}, \mathbf{W}) - \tilde{y}^*\|^2$$
- Minimierung durch Gradientenabstieg, Berechnung des Gradienten:
Backpropagation

(Bishop 2006)

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 32 / 65

Methodischer Überblick

Backpropagation

- $w_{a,b}^t = w_{a,b}^{t-1} + \alpha \cdot \delta_b \cdot y_a$
- $w_{a,b}$: Gewicht von Neuron a zu Neuron b
- a : Neuron in Schicht i
- b : Neuron in Schicht $i + 1$
- y_a : Ausgabe von Neuron a
- Lokale Gradienten δ_b :
 - $\delta_b = \begin{cases} (y_b^* - y_b) \cdot f'(z_b) & \text{wenn } b \text{ in Ebene } n \\ (\sum_{c \in \text{Succ}(b)} \delta_c \cdot w_{b,c}) \cdot f'(z_b) & \text{sonst.} \end{cases}$
 - z_b : interner Zustand von b
 - Rekursive Berechnung

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 33 / 65

Methodischer Überblick

Beispiele für Architekturen Neuronaler Netze

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 34 / 65

Methodischer Überblick

Stand der Dinge

- KNN** stellen **Funktionskompositionen** dar
- Berechnung der Ausgabe:** Vorwärtspropagierung
- Lernen:** Gradientenabstieg, Backpropagation
- Komplexe Zusammenhänge durch latente Variablen lernbar

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 35 / 65

Methodischer Überblick

Outline

- 1 Einführung und Motivation
- 2 Methodischer Überblick
 - Probabilistische Graphische Modelle (PGM)
 - Künstliche Neuronale Netze (KNN)
- 3 Verhältnis von PGM und KNN
 - Propagierungsalgorithmen
 - Verwandtschaft der Formulierungen Hidden CRF
- 4 Integration von KNN in PGM
 - Fallstudien
 - Werkzeuge
- 5 Zusammenfassung und Ausblick

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 37 / 65

Methodischer Überblick

Vergleich der Propagierungsalgorithmen

	Modell-Paradigma	
Aufgabe	PGM	FF-KNN
Inferenz	Belief Propagation (Max-Product, Max-Sum)	Vorwärtspropagierung
Lernen	Zählen, Gradientenaufstieg (Normalisierungsberechnung: Max-Product)	Gradientenabstieg (Gradientenbestimmung: Backpropagation)

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 37 / 65

Methodischer Überblick

Ist Belief Propagation das Gleiche wie Backpropagation?

Nein.

Unterschiede:

- Belief Propagation:** Inferenz entsprechend Struktur des graphischen Modells
- Backpropagation:** Berechnung von Gradienten entsprechend eines Fehlers in der Ausgabe

Aber:

- Es existieren Arbeiten, die die Verfahren aufeinander abbilden. (Eisner, 2016; Dauwels, 2005; Dauwels, 2006)
- Diese Arbeiten führen nicht zu einer Gleichstellung von PGMs und KNNs.

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 38 / 65

Methodischer Überblick

Outline

- 1 Einführung und Motivation
- 2 Methodischer Überblick
 - Probabilistische Graphische Modelle (PGM)
 - Künstliche Neuronale Netze (KNN)
- 3 Verhältnis von PGM und KNN
 - Propagierungsalgorithmen
 - Verwandtschaft der Formulierungen Hidden CRF
- 4 Integration von KNN in PGM
 - Fallstudien
 - Werkzeuge
- 5 Zusammenfassung und Ausblick

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 38 / 65

Sind KNN und PGM verwandt? Beispiel: XOR

- XOR mit KNN:
- Faktor-Graph-Idee 1:

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 40 / 65

Sind KNN und PGM verwandt? Beispiel: XOR

- $\Psi_i(\cdot) = \exp(\sum_j \lambda_{ij} f_{ij}(\cdot))$
 - $f_{i1}(x_a, x_b) = \{1, \text{wenn } x_a = 1 \wedge x_b = 1 \wedge y = 1; 0 \text{ sonst}\}$
 - $f_{i2}(x_a, x_b) = \{1, \text{wenn } x_a = 1 \wedge x_b = 0 \wedge y = 1; 0 \text{ sonst}\}$
 - $f_{i3}(x_a, x_b) = \{1, \text{wenn } x_a = 0 \wedge x_b = 1 \wedge y = 1; 0 \text{ sonst}\}$
 - $f_{i4}(x_a, x_b) = \{1, \text{wenn } x_a = 0 \wedge x_b = 0 \wedge y = 1; 0 \text{ sonst}\}$
- $\Psi_{OR}: \lambda_{i1} = 1, \lambda_{i2} = 1, \lambda_{i3} = 1, \lambda_{i4} = 0$
- $\Psi_{NAND}: \lambda_{i1} = 0, \lambda_{i2} = 1, \lambda_{i3} = 1, \lambda_{i4} = 0$
- $\Psi_{AND}: \lambda_{i1} = 1, \lambda_{i2} = 0, \lambda_{i3} = 0, \lambda_{i4} = 0$

• Vergleichbare Darstellung möglich, aber Inferenz und Lernen dennoch vollkommen unterschiedlich.

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 41 / 65

Sind KNN und PGM verwandt? Beispiel: XOR

Idee 2:

- $\Psi_{XOR}(\cdot) = \exp(\sum_j \lambda_{ij} f_{ij}(\cdot))$
 - $f_{i1}(x_a, x_b) = \{1, \text{wenn } x_a = 1 \wedge x_b = 1 \wedge y = 1; 0 \text{ sonst}\}$
 - $f_{i2}(x_a, x_b) = \{1, \text{wenn } x_a = 1 \wedge x_b = 0 \wedge y = 1; 0 \text{ sonst}\}$
 - $f_{i3}(x_a, x_b) = \{1, \text{wenn } x_a = 0 \wedge x_b = 1 \wedge y = 1; 0 \text{ sonst}\}$
 - $f_{i4}(x_a, x_b) = \{1, \text{wenn } x_a = 0 \wedge x_b = 0 \wedge y = 1; 0 \text{ sonst}\}$
- $\lambda_1 = 0, \lambda_2 = 1, \lambda_3 = 1, \lambda_4 = 0$

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 42 / 65

Restricted Boltzmann-Machine

- Lernt Zusammenhänge zwischen sichtbaren Variablen \vec{v} mit Hilfe der latenten Variablen \vec{h}
- $p(\vec{v}, \vec{h}) = \frac{1}{Z} \cdot \exp(-E(\vec{v}, \vec{h}))$
- $E(\vec{v}, \vec{h}) = -\vec{b}^T \vec{v} - \vec{b}^T \vec{h} - \vec{v}^T W \vec{h}$
- $p(\vec{h} | \vec{v}) = \prod_i p(h_i | \vec{v})$ und $p(\vec{v} | \vec{h}) = \prod_i p(v_i | \vec{h})$
- Inferenz/Training mit Block-Gibbs-Sampling
- Baustein zur Erstellung tiefer Netzwerke (Goodfellow, 2016)

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 43 / 65

Beispiele eines Deep Belief Networks

⇒ PGM
 ⇒ Struktur entsprechend tiefem neuronalen Netz
 ⇒ Semantik versteckter Variablen nicht a priori definiert

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 44 / 65

Probabilistisches Modell mit latenten Variablen: Hidden CRF

- Eingabevariablen x (immer beobachtbar)
- Ausgabevariablen z (zum Trainingszeitpunkt beobachtbar)
- Ausgabevariablen y (niemals beobachtbar)
- ⇒ Vergleichbare Situation zu Feed-Forward Neural Networks
 - Semantik der latenten Variablen ist aber vordefiniert.
- Training: Gradientenstieg (aber nicht konvex) oder Expectation Maximization (Quattoni, 2007; Tackstrom, 2011)

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 45 / 65

Stand der Dinge

- PGMs und KNNs sind nicht das Gleiche, haben aber ähnliche Komponenten
- KNNs nutzen **latente Variablen** um komplexe Zusammenhänge zu lernen
 - PGMs können das ebenfalls: Hidden CRF
- PGMs sind Formalismen um Annahmen über Variablenzusammenhänge zu modellieren
 - Eher unüblich in KNNs
- Gewichte in (manchen) neuronalen Netzen können als Faktoren in PGMs gesehen werden: Repräsentationslernen mit Hilfe von RBM

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 46 / 65

Outline

- 1 Einführung und Motivation
- 2 Methodischer Überblick
 - Probabilistische Graphische Modelle (PGM)
 - Künstliche Neuronale Netze (KNN)
- 3 Verhältnis von PGM und KNN
 - Propagierungsalgorithmen
 - Verwandtschaft der Formulierungen
 - Hidden CRF
- 4 Integration von KNN in PGM
 - Fallstudien
 - Werkzeuge
- 5 Zusammenfassung und Ausblick

Einbindung von KNN-Faktoren in PGMs?

Idee mit Beispiel:

- Modellierung der Zusammenhänge entsprechend Expertenwissen für Variablen mit bekannter Semantik
- Einbindung von neuronalen Faktoren für automatisches Lernen komplexer Zusammenhänge
- $\Psi_1(x_1^1, x_1^2, x_2^2) = \exp(\sum_i \lambda_i f(x_1^1, x_1^2, x_2^2))$
- $\Psi_4(x_1^3, x_2^3, x_3^3) = \text{RNN}(\cdot)$
- $\Psi_5(\bar{x}) = \text{FFNN}(\cdot)$
- $p(\bar{x}) = \frac{1}{Z} \cdot \prod_i \Psi_i(\cdot)$
 - Welche Realisierungen gibt es von solchen Modellen?
 - Ich stelle nun einige Fallstudien vor.

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 48 / 65

LSTM-CRF

Idee: Modelliere Eigenschaften der Ausgabevariablen mit einem PGM, nutze flexibles RNN um Dateneigenschaften zu lernen.

- Standard linear-chain CRF: $p(\bar{y} | \bar{x}) = \frac{1}{Z(\bar{x})} \prod_i \exp(\sum_j \lambda_j f_j(y_{i-1}, y_i, \bar{x}, i))$
- Andere Formulierung: $p(\bar{y} | \bar{x}) = \frac{1}{Z(\bar{x})} \cdot \prod_i \exp(\sum_j \lambda_j f(y_{i-1}, y_i)) \cdot \prod_i \exp(\sum_j \lambda_j f(\bar{x}, y_i, i))$
- Ersetze datenbezogenes Log-lineares Modell durch LSTM: $p(\bar{y} | \bar{x}) = \frac{1}{Z(\bar{x})} \cdot \prod_i \exp(\sum_j \lambda_j f(y_{i-1}, y_i)) \cdot \prod_i \text{LSTM}(\bar{x}_i, y_i)$

Dies ist ein etabliertes Modell.

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 49 / 65

LSTM-CRF Anwendungsbeispiele (1)

Labeling von Wörtern in Text mit einer IOBE-Sequenz (Huang et al., 2015)

\bar{x}	=	the	Severe	acute	respiratory	syndrome	coronavirus	2	...
\bar{y}	=	O	B	I	I	I	I	E	...

CRF-Schicht zeigt Fehlerreduktion von etwa 10%

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 50 / 65

LSTM-CRF Anwendungsbeispiele (2)

2D-Variante: Bildsegmentierung (Zheng et al., 2015)

Original image (lower to highlight segmented parts) Semantic segmentation

Verbesserungen insbesondere bei filigranen Strukturen (<http://www.robots.ox.ac.uk/~szheng/crfasrnn demo/>)

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 51 / 65

LSTM-CRF Anwendungsbeispiele (3)

Manoeverentscheidungen bei autonomem Fahren (Wang et al., 2018)

	Lane Keeping		Change prep		Left change		Right change		F1
	P	R	P	R	P	R	P	R	
LSTM	.99	.80	.39	.97	.67	.97	.61	.97	.75
LSTM+CRF	.99	.78	.42	.99	.80	.98	.79	.99	.81

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 52 / 65

Mehrere in Relation stehende Sequenzen

Faktorielle Conditional Random Fields (Sutton, 2007)

- Entspricht Kombination zweier linearer Ketten, welche mit einer Eingabe konditioniert werden
- Aber: Vorhersage ist nicht unabhängig, daher messen Faktoren Kompatibilität zwischen Ausgabesequenzen
- Loopy Graph!

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 53 / 65

Mehrere in Relation stehende Sequenzen: Morphologisches Tagging

Neural Factor Graph Models for Cross-lingual Morphological Tagging (Malaviya et al., 2018)

POS: Det, Definite: Ind, Gender: Masc, Number: Sing, PronType: Art
 POS: Noun, Gender: Masc, Number: Sing
 POS: Adj, Gender: Masc, Number: Sing
 ES: PT, Un, renacimiento, revivimiento, refrescante, refrescante

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 54 / 65

Mehrere in Relation stehende Sequenzen: Morphologisches Tagging, Ergebnisse

Language	Model	tgt_size=100			tgt_size=1000		
		Accuracy	F1-Micro	F1-Macro	Accuracy	F1-Micro	F1-Micro
SV	Baseline	15.11	8.36	10.37	68.64	76.36	76.50
	Ours	29.47	54.09	54.36	71.32	84.42	84.46
BG	Baseline	29.05	14.32	29.62	59.20	67.22	67.12
	Ours	27.81	40.97	42.43	39.25	60.23	60.84
HU	Baseline	21.97	13.30	16.67	50.75	58.68	62.79
	Ours	33.32	54.88	54.69	45.90	74.05	73.38
PT	Baseline	18.91	7.10	10.33	74.22	81.62	81.87
	Ours	58.82	73.67	74.07	76.26	87.13	87.22

Table 3: Token-wise accuracy and F1 scores on mono-lingual experiments

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart Roman Klinger 17. Juli 2020 55 / 65

