

# Advanced Topics in Machine Learning 5 Representation: Markov Networks

Prof. Dr. Steffen Staab Dr. Rafika Boutalbi Zihao Wang https://www.ipvs.uni-stuttgart.de/departments/ac/









#### **Learning Objectives**

Undirected Graphical Models / Markov (random) networks

- General Gibbs distribution
- Conditional Random Fields

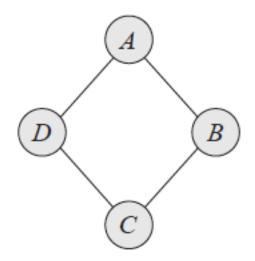
#### **Disclaimer**

Figures and examples not marked otherwise are taken from the book by Koller & Friedman

### 1 Undirected Graphical Models

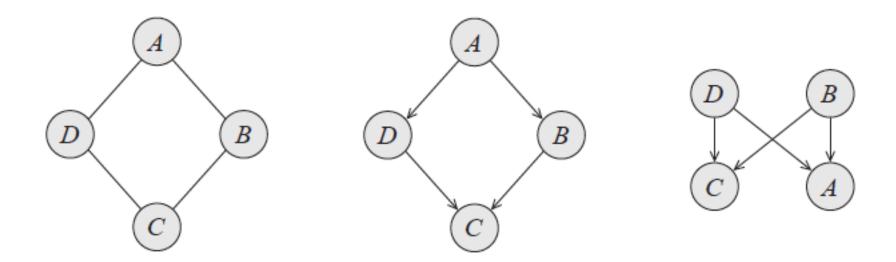
#### A Modeling Problem for Bayesian Network

- Professor was ambiguous, two interpretations: 0,1
- Four students Alice, Bob, Charlie, Debbie collaborating in pairs, but never Alice-Charlie, Bob-Debbie
- Target:  $(B \perp D \mid A, C), \ (A \perp C \mid B, D)$
- Question:
   How to model as Bayesian Network?



#### A Modeling Problem for Bayesian Network

 $\bullet$   $(B\perp D\,\big|\,A,C)$  and  $(A\perp C\,|\,B,D)$  cannot be achieved using Bayesian Network



#### **Undirected Graphical Model**

- Random variables:  $\mathcal{X} = \{X_1, ..., X_k\}$
- Assignment to  $\mathcal{X}$ :  $x_{\mathcal{X}}$ , or simply x
- Collection of subsets of  $\mathcal{X}$ :  $\mathcal{A} = \{A_1, ..., A_l | \forall i : A_i \subset \mathcal{X}\}$
- Assignment to  $A \in \mathcal{A}$ :  $x_A$
- Factors:  $\Phi = \{\phi_{A_1}, ..., \phi_{A_l}\}$

**Definition**: An undirected graphical model is a probability distribution that can be written using factors  $\phi_A$  with corresponding scopes A as

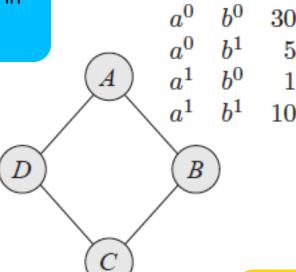
$$P(x) = \frac{1}{Z_{\Phi}} \prod_{A \in \mathcal{A}} \phi_A(x_A)$$

with normalization factor Z

$$Z_{\Phi} = \sum_{x \in \mathcal{X}} \prod_{A \in \mathcal{A}} \phi(x_A)$$

#### Simple Example: Pairwise Markov Network

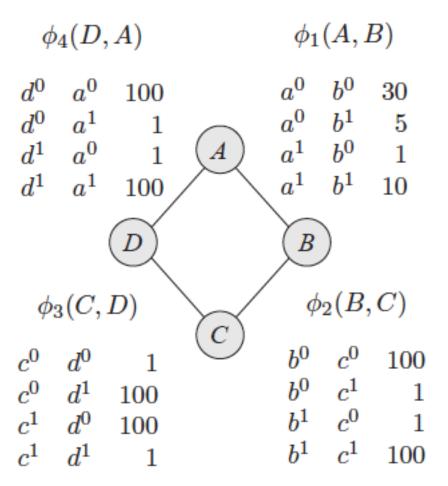
Factors have pairs of variables in their scope



 $\phi_1(A,B)$ 

No direction
No conditional probability tables

#### Simple Example: Pairwise Markov Network



affinity compatibility soft constraints

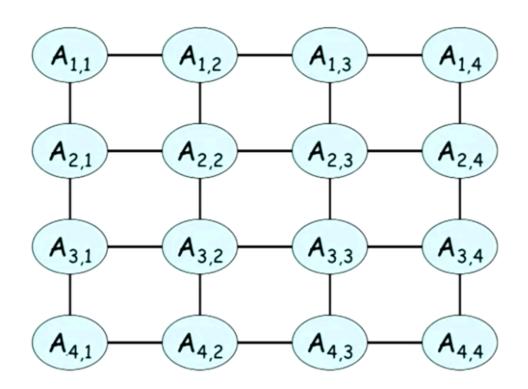
#### Simple Example: Pairwise Markov Network

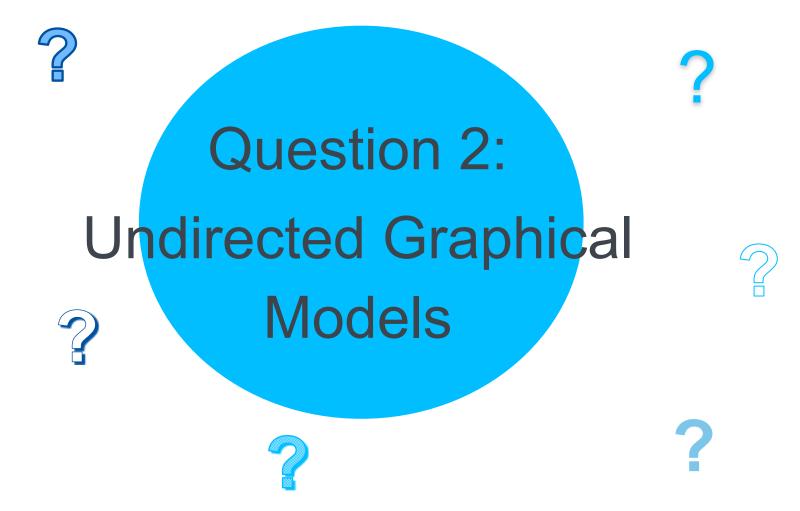
$$\begin{split} \tilde{P}(A,B,C,D) &= \phi_1(A,B)\phi_2(B,C)\phi_3(C,D)\phi_4(D,A) \quad \phi_4(D,A) \\ P(A,B,C,D) &= \frac{1}{Z}\tilde{P}(A,B,C,D) \\ & \begin{vmatrix} Assignment & Unnormalized & Normalized \\ a^0 & b^0 & c^0 & d^1 & 300,000 & 0.04 \\ a^0 & b^0 & c^1 & d^0 & 300,000 & 0.04 \\ a^0 & b^0 & c^1 & d^0 & 300,000 & 0.04 \\ a^0 & b^0 & c^1 & d^1 & 300,000 & 0.04 \\ a^0 & b^1 & c^1 & d^0 & 300,000 & 0.04 \\ a^0 & b^1 & c^1 & d^1 & 500 & 6.9 \cdot 10^{-5} \\ a^0 & b^1 & c^1 & d^1 & 500 & 6.9 \cdot 10^{-5} \\ a^1 & b^0 & c^0 & d^1 & 1,000,000 & 0.14 \\ a^1 & b^0 & c^1 & d^1 & 100,000 & 0.14 \\ a^1 & b^1 & c^0 & d^1 & 100,000 & 0.014 \\ a^1 & b^1 & c^0 & d^1 & 100,000 & 0.014 \\ a^1 & b^1 & c^1 & d^0 & 100,000 & 0.014 \\ a^1 & b^1 & c^1 & d^1 & 100,000 & 0.014 \\ a^1$$

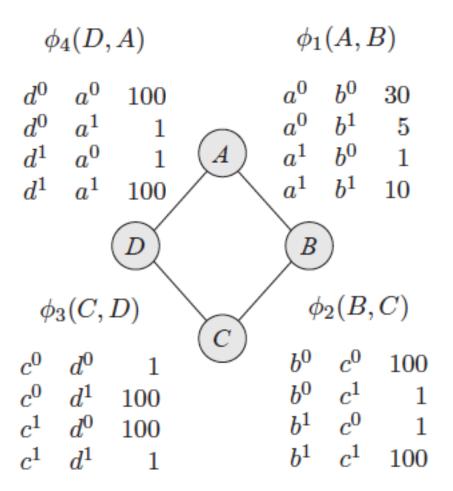
#### Probability diverges from individual factors

						$\phi_4(D,A)$				$\phi_1(A,B)$			
Assignment				Unnormalized	Normalized								
$a^0$	$b^0$	$c^0$	$d^0$	300,000	0.04		$d^0$	$a^0$	100		$a^0$	$b^0$	30
$a^0$	$b^0$	$c^0$	$d^1$	300,000	0.04				100				
$a^0$	$b^0$	$c^1$	$d^0$	300,000	0.04		$d^0$	$a^1$	1		$a^0$	$b^1$	5
$a^0$	$b^0$	$c^1$	$d^1$	30	$4.1 \cdot 10^{-6}$		$d^1$	$a^0$	1	(A)	$a^1$	$b^0$	1
$a^0$	$b^1$	$c^0$	$d^0$	500	$6.9 \cdot 10^{-5}$	$A = a^{0}, =$	0.691		100	$\sim$	$a^1$	$b^1$	10
$a^0$	$b^1$	$c^0$	$d^1$	500	$6.9 \cdot 10^{-5}$	$B = b^1$	$a^{\perp}$	$a^1$	100	/ \	$a^{-}$	0	10
$a^0$	$b^1$	$c^1$	$d^0$	5,000,000	0.69				$\sim$		$\overline{}$		
$a^0$	$b^1$	$c^1$	$d^1$	500	$6.9\cdot10^{-5}$				(D)		(B)	)	
$a^{\scriptscriptstyle 1}$	$b^0$	$c^0$	$d^0$	100	$1.4\cdot 10^{-5}$								
$a^1$	$b^0$	$c^0$	$d^1$	1,000,000	0.14				\	\ /			
$a^1$	$b^0$	$c^1$	$d^0$	100	$1.4\cdot 10^{-5}$		$\phi$	$_{3}(C,$	D)	\_/	$\phi$	$\rho_2(B,$	C
$a^1$	$b^0$	$c^1$	$d^1$	100	$1.4 \cdot 10^{-5}$		7.	,	,	$\begin{pmatrix} C \end{pmatrix}$		_ ( ,	,
$a^1$	$b^1$	$c^0$	$d^0$	10	$1.4 \cdot 10^{-6}$		0	70			$b^0$	$c^0$	100
$a^1$	$b^1$	$c^0$	$d^1$	100,000	0.014		$c^0$	$d^0$	1	_			100
$a^1$	$b^1$	$c^1$	$d^0$	100,000	0.014		$c^0$	$d^1$	100		$b^0$	$c^1$	1
$a^1$	$b^1$	$c^1$	$d^1$	100,000	0.014		$c^1$				$b^1$	$c^0$	1
,		-	. '	•		•		$d^0$	100				1
							$c^1$	$d^1$	1		$b^1$	$c^1$	100

### **Application for Pairwise Markov Network: Modeling Pixels**







Consider the pairwise factor  $\phi_1(A,B)$ . That potential is proportional to:

- $\circ$  The marginal probability P(A,B)
- $\circ$  The conditional probability  $P(A \mid B)$
- $\bigcirc$  The conditional probability P(A, B | C, D)
- None of the above

#### 2 Gibbs Distribution

#### Are pairwise Markov Networks fully expressive?

Consider a fully connected pairwise Markov network over  $X_1, ... X_n$  where each  $X_i$  has d values.

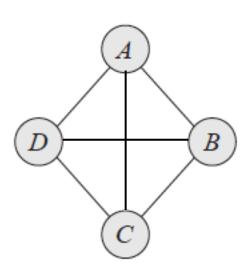
How many parameters does the network have?

How many edges: 
$$\frac{n}{2}(n-1)$$

How many value combinations on each edge:  $d^2$ 

$$\frac{n}{2}(n-1)d^2 \ll d^n$$

Not every distribution can be represented as a pairwise Markov network



#### **Undirected Graphical Model**

- Random variables:  $\mathcal{X} = \{X_1, ..., X_k\}$
- Assignment to  $\mathcal{X}$ :  $x_{\mathcal{X}}$ , or simply x
- Collection of subsets of  $\mathcal{X}$ :  $\mathcal{A} = \{A_1, ..., A_l \mid \forall i : A_i \subset \mathcal{X}\}$
- Assignment to  $A \in \mathcal{A}$ :  $x_A$
- Factors:  $\Phi = \{\phi_{A_1}, ..., \phi_{A_l}\}$

**Definition**: An undirected graphical model is a probability distribution that can be written using factors  $\phi_A$  with corresponding scopes A as

$$P(x) = \frac{1}{Z_{\Phi}} \prod_{A \in \mathcal{A}} \phi_A(x_A)$$

with normalization factor Z

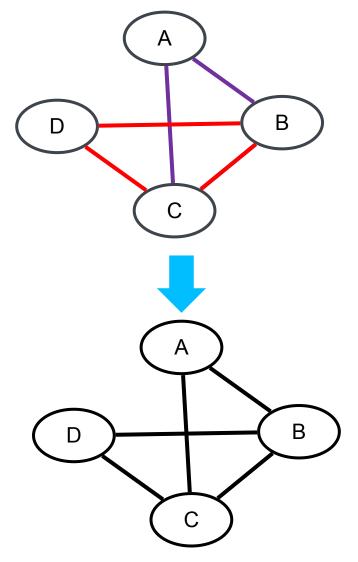
$$Z_{\Phi} = \sum_{x \in \mathcal{X}} \prod_{A \in \mathcal{A}} \phi(x_A)$$

Gibbs Distribution  $\widetilde{P}_{\Phi}(\mathcal{X}) = \prod_{A \in \mathcal{A}} \phi_A(A)$   $Z_{\Phi} = \sum_{\mathcal{X}} \widetilde{P}_{\Phi}(\mathcal{X})$   $P_{\Phi}(\mathcal{X}) = \frac{1}{Z_{\Phi}} \widetilde{P}_{\Phi}(\mathcal{X})$ 

#### **Induced Markov Network**

 $\phi_1(A, B, C), \phi_2(B, C, D)$ 

Induced Markov network  $H_{\Phi}$  has an edge  $X_i - X_j$  when there exists  $\phi \in \Phi$  such that  $\{X_i, X_i\} \subseteq \operatorname{scope}(\phi)$ 



#### **Factorization**

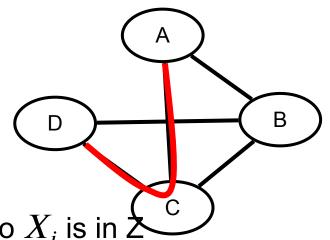
 $P(\mathcal{X})$  factorizes over H if there exists  $\Phi$  such that

$$P(\mathcal{X}) = P_{\Phi}(\mathcal{X})$$
 and  $H$  is the induced graph for  $\Phi$ 

#### Flow of Influence & Active Trails

$$\phi_1(A, B, C), \phi_2(B, C, D)$$

Influence can flow along any trail regardless of the form of the factors



A trail  $X_1 - \ldots - X_n$  is active given Z if no  $X_i$  is in  $\Sigma$ 

#### Examples:

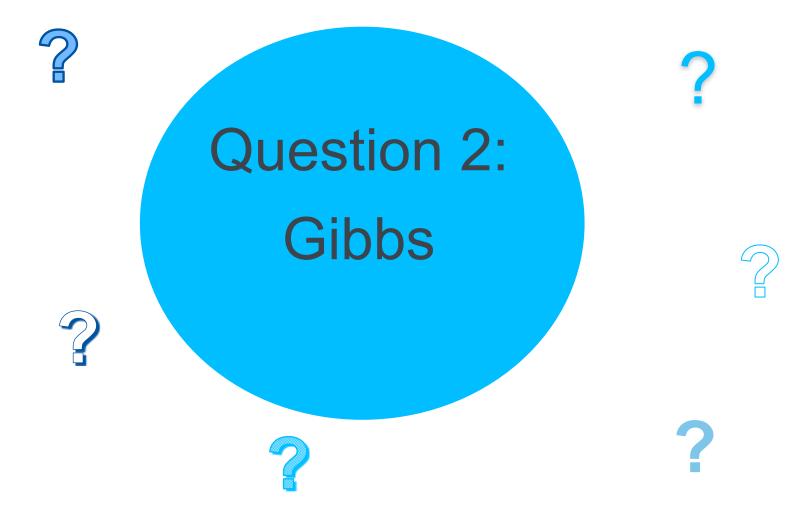
$$A - C - D$$
 is active if  $Z = \{\}$ 

$$A-C-D$$
 is inactive if  $Z=\{C\}$ 

There is no active trail from A to D if  $Z = \{B, C\}$ 

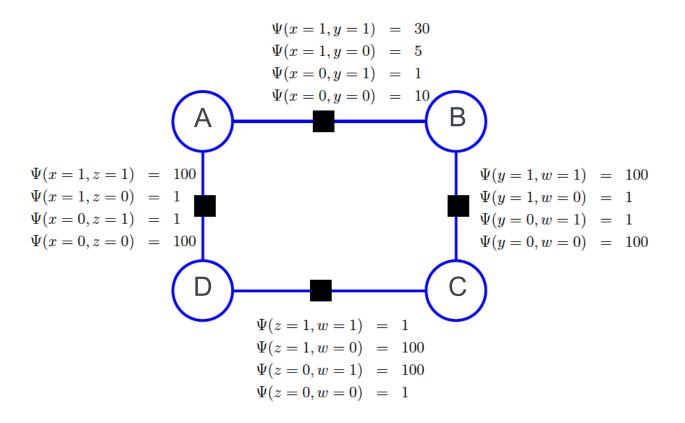
#### **Summary**

- Gibbs distribution represents distribution as a product of factors
- Induced Markov network connects every pair of nodes that are in the same factor
- Markov network structure does not fully specify the factorization of  $P(\mathcal{X})$
- Active trails depend only on graph structure



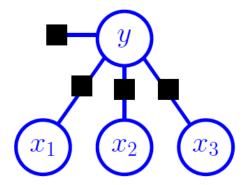
### **Intermezzo:** Factor Graphs

#### **Example factor graph**

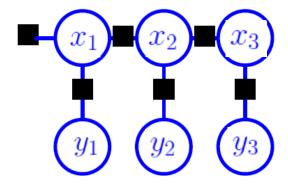


#### Further example factor graphs

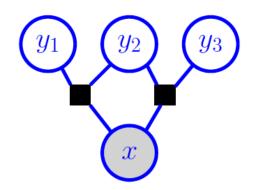
Naïve Bayes



Hidden Markov Model

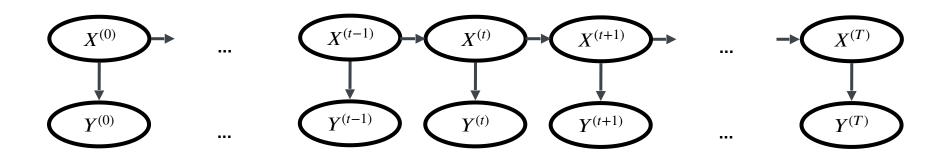


Linear-Chain Conditional Random Field



## 3 Conditional Random Field (CRF)

#### Problem 1: Re-consider Hidden Markov Model



#### Example application:

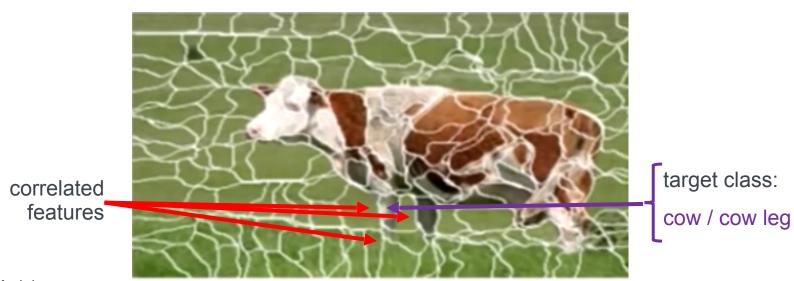
"The horse chased by the barn fell."

past participle

Long range dependency

 If a word is unknown ("Emthe"), its identity is not a useful feature, neighboring words are needed, context is needed

#### **Problem 2: Correlated Features**



- Naive Models:
  - redundant features swamp decision
  - label bias: prior training occludes observations (Sutton & McCallum 2007)
- Model correlations:
  - · too many correlations to model
    - · intractable if too many dependencies

#### Generative vs. discriminative classification

- Generative approach
  - Assume functional form for P(Y), P(X|Y)
  - Estimate parameters of P(Y),  $P(X \mid Y)$  from training data
  - Use Bayes rule to calculate P(Y|X)
- Bayesian Networks
  - Naive Bayes
  - HMM
  - Markov random fields
- Model  $P(\mathcal{Y}, \mathcal{X})$ 
  - then query this model

- Discriminative approaches
  - Assume some functional form for P(Y | X)
  - Estimate parameters of P(Y|X) directly from training data
- Approaches
  - Logistic regression
  - Support vector machines
  - (most) neural networks
  - Conditional random fields
- Model  $P(\mathcal{Y} | \mathcal{X})$ 
  - then query this model

Ng, Andrew Y., and Michael I. Jordan. "On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes." *Advances in neural information processing systems*. 2002.

#### **Core Ideas for Conditional Random Field**

#### 1. Distinguish two sets of random variables:

- $\mathcal{X} = \{X_1, ..., X_k\}$  are always observed
- $\mathcal{Y} = \{Y_1, ..., Y_l\}$  are the target variables

#### 2. Model $P(\mathcal{Y} \mid \mathcal{X})$

- Model factors
  - with scopes  $A \subset (\mathcal{Y} \cup \mathcal{X})$
  - no scope A such that  $A \subseteq \mathcal{X}$
- . Normalize over  $\sum_{\mathbf{y}} P(\mathbf{y}, \mathbf{\mathcal{X}})$

A conditional model specifies the probabilities of possible **label sequences** given an **observation sequence**.

Different from Gibbs distribution

#### Core Ideas for Conditional Random Field: DO NOTs

- Do not model  $P(\mathcal{Y}, \mathcal{X})$  at all
- Do not train by maximizing joint likelihood:

$$\underset{x,y \in Data}{\operatorname{argmax}} P(x, y \mid \theta)$$

ullet Do not model correlations in  ${\mathcal X}$ 

Do not normalize over P(9/ 9')

Avoid sources of intractability (as far as possible)

#### Task specific prediction

input variables  ${\mathcal X}$ 

target variable

y

Image segmentation

pixel values, processed features (SIFT, SURF, CNN) class for every pixel (cow,...)

**Text** 

words in sentence

labels of words (person, location,...)

#### **Observation Sequence**

- Sequence / Matrix / other structure
- Raw observation:
  - "The horse chased by the barn fell."
- Extended observation:
  - "The/det horse/n chased/pp by/prep the/det barn/n fell/v."
- Transformed observation:
  - Sequence of BERT embedding vectors for each of the words in this sentence:

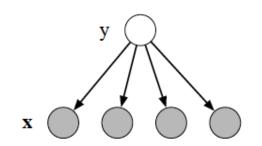
$$(\overrightarrow{w}_1 \overrightarrow{w}_2 ... \overrightarrow{w}_7)$$

feature 'engineering

#### From Naive Bayes to Naive Markov model (I)

#### Remember Naive Bayes:

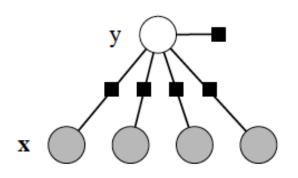
$$P(y, x_1, ...x_k) = P(y, \mathbf{x}) = P(y) \prod_{i=1}^k P(x_i | y)$$



#### Reformulate as factor graph:

$$\phi_Y(y) = P(y), \phi_i(y, x_i) = P(x_i | y)$$

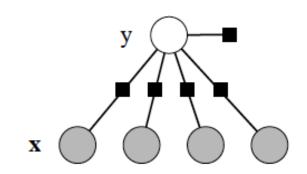
Naive Markov model is a special case of CRF



#### From Naive Bayes to Naive Markov model (II)

#### Reformulate as factor graph:

$$\phi_0(y) = P(y),$$
  
$$\phi_i(y, x_i) = P(x_i | y)$$



#### Model P(y | x) as in logistic regression:

$$\log P(y, \mathbf{x}) = \lambda_{y,0} + \sum_{i=1}^{k} \lambda_{y,i} x_i,$$

$$P(y, \mathbf{x}) = e^{\left(\lambda_{y,0} + \lambda_{y} \mathbf{x}\right)}$$

$$P(x) = Z(x) = \sum_{y} e^{\left(\lambda_{y,0} + \lambda_{y} x\right)}$$

$$P(y \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} e^{\left(\lambda_{y,0} + \lambda_y \mathbf{x}\right)}$$

Different from Gibbs distribution

#### From Naive Bayes to Naive Markov model (III)

Model P(y | x) as in logistic regression:

$$P(y \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} e^{(\lambda_{y,0} + \lambda_y \mathbf{x})}$$

Random fields came from 70ies physics.
They invented names and notations that you may encounter nowadays such as this second notation.

Make  $\lambda$  independent of y by using feature functions:

• for bias: 
$$\lambda_{y,0}$$
 write  $f_{y'}(y, \boldsymbol{x}) = \mathbf{1}_{\{y=y'\}}$ 

• for weights:  $\lambda_{y,i}$  write  $f_{y',i}(y, \boldsymbol{x}) = \mathbf{1}_{\{y=y'\}} x_i$ 

• re-index 
$$f$$
 
$$P(y \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} e^{\left(\sum_{i=0}^{k} \lambda_i f_i(y, \mathbf{x})\right)}$$

Note that  $\mathbf{1}_{\{y=y'\}}$  is a function of y, which is 1 if the condition is true

Difference:
Now weights are
shared between
classes!
(before: different
per class)

# Sigmoid / Softmax

For 
$$Y = \{0,1\}$$

**Unnormalized:** 

$$\widetilde{P}(y=0 \mid \mathbf{x}) = e^{(\lambda_{0,0} + \lambda_0 \mathbf{x})} = e^{z_0}$$

$$\widetilde{P}(y=1 \mid \mathbf{x}) = e^{(\lambda_{1,0} + \lambda_1 \mathbf{x})} = e^{z_1}$$

Therefore:

$$P(y = 0 \mid x) = \frac{e^{z_0}}{e^{z_0} + e^{z_1}}$$

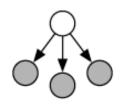
Sigmoid

For 
$$Y = \{1...n\}$$

Generalized:

$$P(y=1 \mid x) = \frac{e^{z_1}}{\sum e^{z_i}}$$

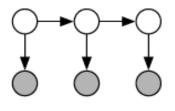
Softmax



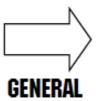
Naive Bayes



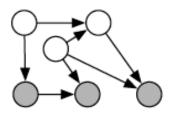
SEQUENCE



**HMMs** 

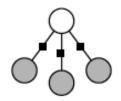


GENERAL Graphs



Generative directed models



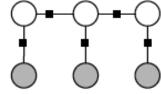


Logistic Regression



SEQUENCE

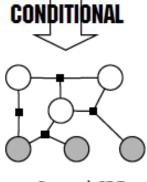




Linear-chain CRFs



**GRAPHS** 

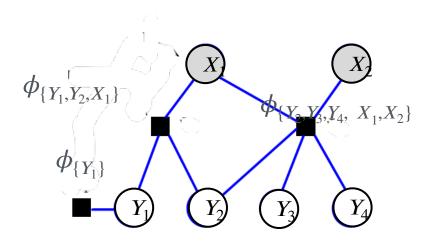


General CRFs

### **Conditional Random Field as Factor Graph**

A factor graph for a CRF is a bi-partite graph G = (V, F, E) of

- variable nodes  $V = \mathcal{X} \cup \mathcal{Y}$ ,
- factor nodes  $F = \{\phi_A\}$ , with  $\operatorname{scope}(\phi_A) = A \subset V$
- edges  $E, (\phi_A, v) \in E, \text{ iff } v \in A$



### **Conditional Random Field as Factor Graph**

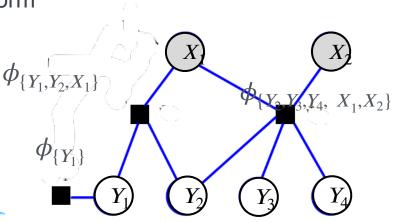
A factor graph for a CRF is a bi-partite graph G = (V, F, E) of

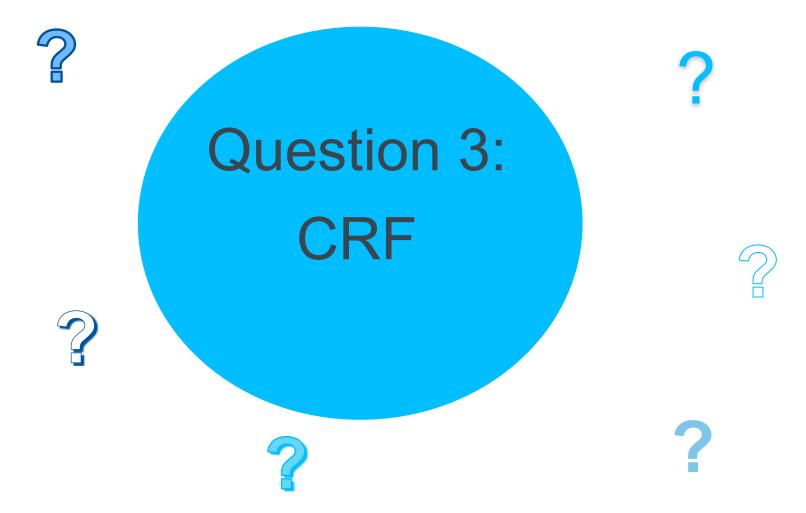
- variable nodes  $V = \mathcal{X} \cup \mathcal{Y}$ ,
  - Input variables:  $\mathcal{X} = \{X_1, X_2\}$ , Target variables:  $\mathcal{Y} = \{Y_1, Y_2, Y_3, Y_4\}$
- factor nodes  $F = \{\phi_A\}$ , with  $\operatorname{scope}(\phi_A) = A \subset V$ 
  - Given feature functions:  $\{f_{A,k}\}$
  - Local function / factor  $\phi_A(x_A, y_A)$  has the form

$$\phi_A(x_A, y_A) = e^{(\sum_k \theta_{A,k} f_{A,k}(x_{A,k}, y_{A,k}))}$$

• edges  $E, (\phi_A, v) \in E, \text{ iff } v \in A$ 

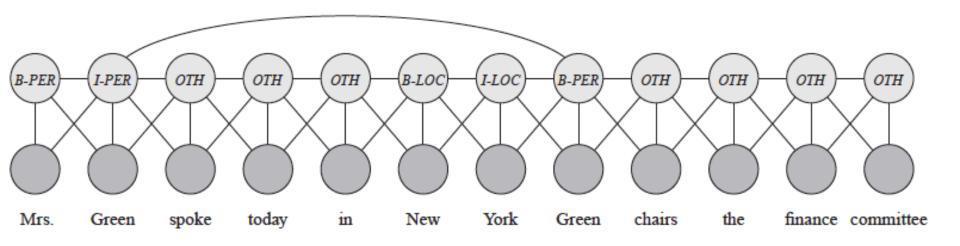
$$P(y \mid x) = \frac{1}{Z(x)} \prod_{A} \phi_{A}(x, y)$$





# **4 CRF Applications**

### **CRFs for Language**



### KEY

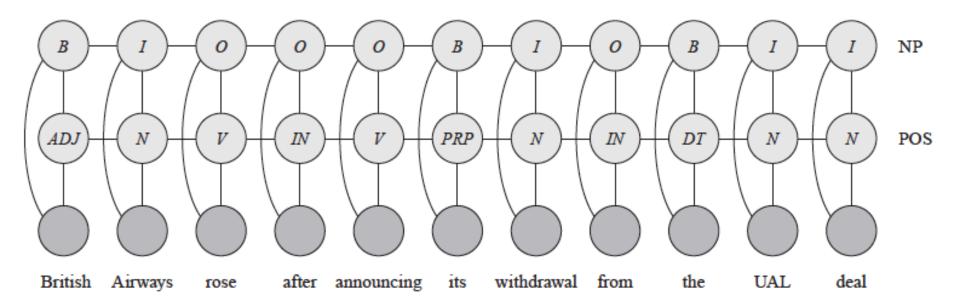
Begin person name I-PER

Within person name

Begin location name

Within location name *I-LOC* 

OTHNot an entitiy



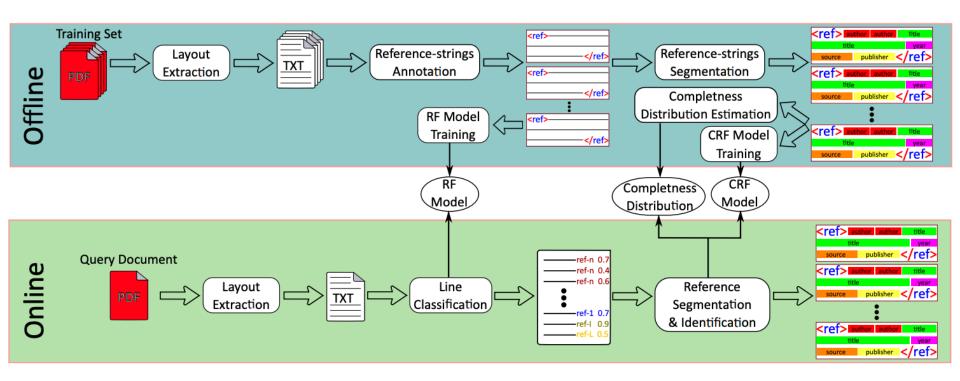
### KEY

В	Begin noun phrase	V	Verb
I	Within noun phrase	$I\!N$	Preposition
0	Not a noun phrase	PRP	Possesive pronoun
N	Noun	DT	Determiner (e.g., a, an, the)
ADJ	Adjective		

A .

### **Example: Excite**

https://excite.informatik.uni-stuttgart.de/



Z. Boukhers, S. Ambhore and S. Staab, "An End-to-End Approach for Extracting and Segmenting High-Variance References from PDF Documents," *2019 ACM/IEEE Joint Conference on Digital Libraries (JCDL)*, Champaign, IL, USA, 2019, pp. 186-195, doi: 10.1109/JCDL.2019.00035.

### Try out yourself!





#### Updated: EXparser v1.0.1 is released and running from 22/02/2019 (test it below!)

EXCITE (Extraction of Citations from PDF Documents) is a toolchain of citation extraction software and particular focus on the Germanlanguage social sciences and this is a public service for the project. In the background of this page we are using CERMINE for extracting content from PDF files and Exparser for reference string extraction and segmentation.

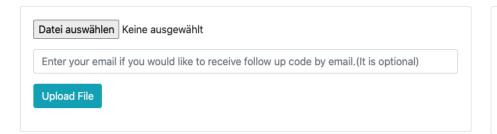
#### How to start the process:

- First: Choose a pdf file by click on "Choose File" button. (size of the file should be less than 20 MB)
- Second: Click on "Upload file" button to start the process. Please enter your email address if you would like to receive follow up code by email. (It is optional)
- tip: After uploading a file a code will be displayed on the screen. This code is necessary for follow up the result of reference extraction. If you would like to load last follow up code related to last uploaded file in your browser.

#### How to check the result:

- First: Enter the follow-up code in the appropriate box in the right-hand of page then click on "Display References" button.
- Second: The result will be displayed on screen if it available.
- tip: Extracting References process will take a little time, at least 30 seconds, and it completely depends on the size of the file.
- tip: Click on "Load last follow up code" If you would like load last follow up code related to last uploaded file in your browser.

This is a first version of the our public web service. How we can improve it? Please let us know if you have any feedback for us by clicking on this link.



If you already have a code, Enter code to load data.							
Enter code to display References.							
Loa	d Last follow up code	Display References					

### **Current pattern in 2020**

Einführung und Motivation I

Methodischer Überblick

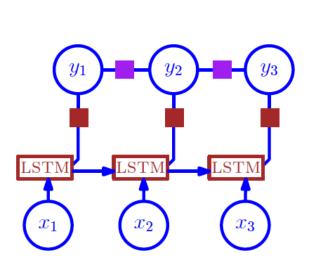
Verhältnis von PGM und KNN o 00000000000

Integration von KNN in PGM

Zusammenfassung/Ausblick

### LSTM-CRF

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



- Standard linear-chain CRF:  $p(\vec{y} \mid \vec{x}) = \frac{1}{Z(\vec{x})} \prod_{i} \exp(\sum_{i} \lambda_{j} f_{j}(y_{i-1}, y_{i}, \vec{x}, i))$
- Andere Formulierung:

$$p(\vec{y} \mid \vec{x}) = \frac{1}{Z(\vec{x})} \cdot \prod_{i} \exp(\sum_{j} \lambda_{j} f(y_{i-1}, y_{i})) \cdot \prod_{i} \exp(\sum_{j} \lambda_{j} f(\vec{x}, y_{i}, i))$$

• Ersetze datenbezogenes Log-lineares Modell durch LSTM:

$$p(\vec{y} \mid \vec{x}) = \frac{1}{Z(\vec{x})} \cdot \prod_{i} \exp(\sum_{j} \lambda_{j} f(y_{i-1}, y_{i})) \cdot \prod_{i} LSTM(\vec{x}_{i}, y_{i})$$

Dies ist ein etabliertes Modell.

# LSTM-CRF Anwendungsbeispiele (1)

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

$\vec{x}$ =	the	Severe	acute	respiratory	syndrome	coronavirus	2	
$\vec{y}$ =	0	В			1		Е	

CRF-Schicht zeigt Fehlerreduktion von etwa 10 %

Lineare CRF-Schicht, z.B. PyTorch oder AllenNLP

## LSTM-CRF Anwendungsbeispiele (2)

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



Original image (hover to highlight segmented parts)



Semantic segmentation

Verbesserungen insbesondere bei filigranen Strukturen

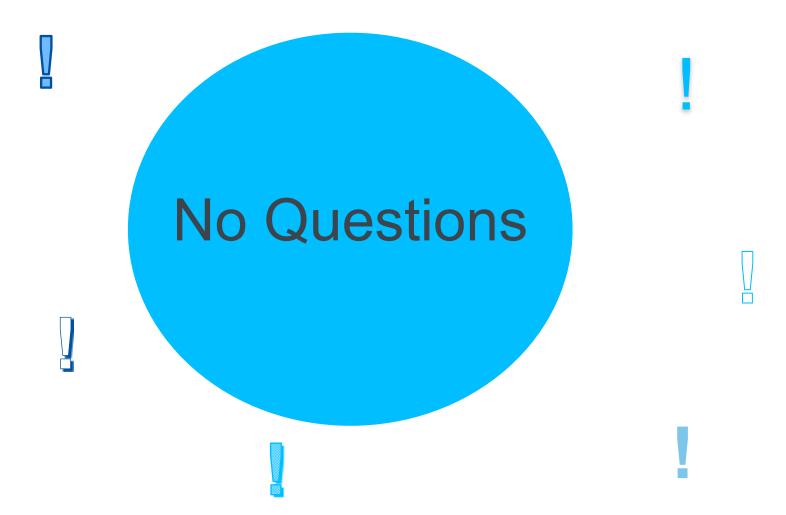
 $(\mathsf{http://www.robots.ox.ac.uk/}{\sim} \mathsf{szheng/crfasrnndemo/})$ 

Institut für Maschinelle Sprachverarbeitung, Universität Stuttgart

Roman Klinger

17. Juli 2020

51 / 65





### Thank you!



### **Steffen Staab**

E-Mail Steffen.staab@ipvs.uni-stuttgart.de
Telefon +49 (0) 711 685e be defined
www.ipvs.uni-stuttgart.de/departments/ac/

Universität Stuttgart Analytic Computing, IPVS Universitätsstraße 32, 50569 Stuttgart