

**Universität Stuttgart**

IPVS – Institute for Parallel and Distributed Systems

Analytic Computing

# Advanced Topics in Machine Learning

## 8 Learning - Part 1

Prof. Dr. Steffen Staab

Dr. Rafika Boutalbi

Zihao Wang

<https://www.ipvs.uni-stuttgart.de/departments/ac/>



# Learning Objectives

Learning Tasks: Queries, Prediction, Knowledge discovery

Sufficient statistics: multinomial, Gaussian

Data Fragmentation and Overfitting

Maximum likelihood for complete observations of BNs

Maximum likelihood for complete observations of MNs

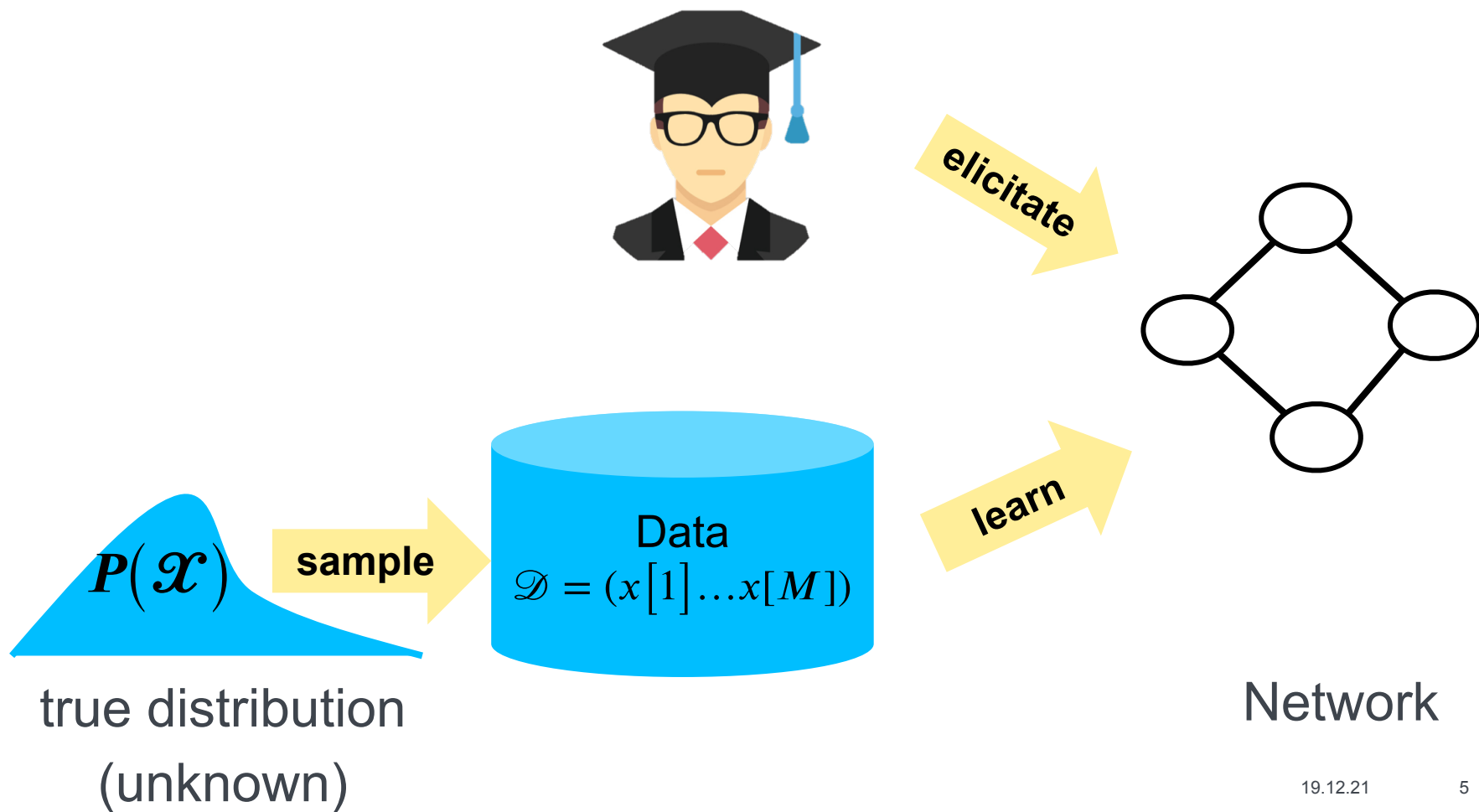
Log-linear models, feature functions

# Disclaimer

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

# 1 Learning Tasks in PGMs

# Learning PGMs



# Why PGM Learning?

- Predictions of structured objects
  - sequences, graphs, trees
  - exploit correlations between several predicted variables
- Incorporate prior knowledge into model
  - Expert knowledge!
- Knowledge discovery
  - Explainability!
- One model, many tasks!
  - From cause to conclusion
  - from conclusion to plausible cause

# Classification of challenges

- **Network structure known**
  - Induce only the factors
- **Network structure unknown**
  - Induce structure
  - Induce factors
- **All** random variables **completely observed**
- **All** random variables **observed**, but not completely
- **Some** random variables **completely observed**
  - others: **Latent Variables**
- **Some** random variables **observed**
  - others: **Latent Variables**

# PGM Learning Task 1: Queries

Goal: Answer general probabilistic queries about new instances

- Targeted quality: Generalization
  - Metric: Test set likelihood:  $P(\mathcal{D}' ; \mathcal{M})$



# PGM Learning Task 2: Prediction

Goal: Specific prediction task on new instances


- Predict target variables  $\mathbf{Y}$  from observed variables  $\mathbf{X}$ 
  - Typically MAP Assignment
    - examples: image segmentation, speech recognition
- Targeted quality: correct assignment
  - conditional likelihood  $\prod_m P(\mathbf{y}[m] \mid \mathbf{x}[m] ; \mathcal{M})$
  - model evaluated on „true“ assignment over test data (“gold standard”)

# PGM Learning Task 3: Knowledge Discovery

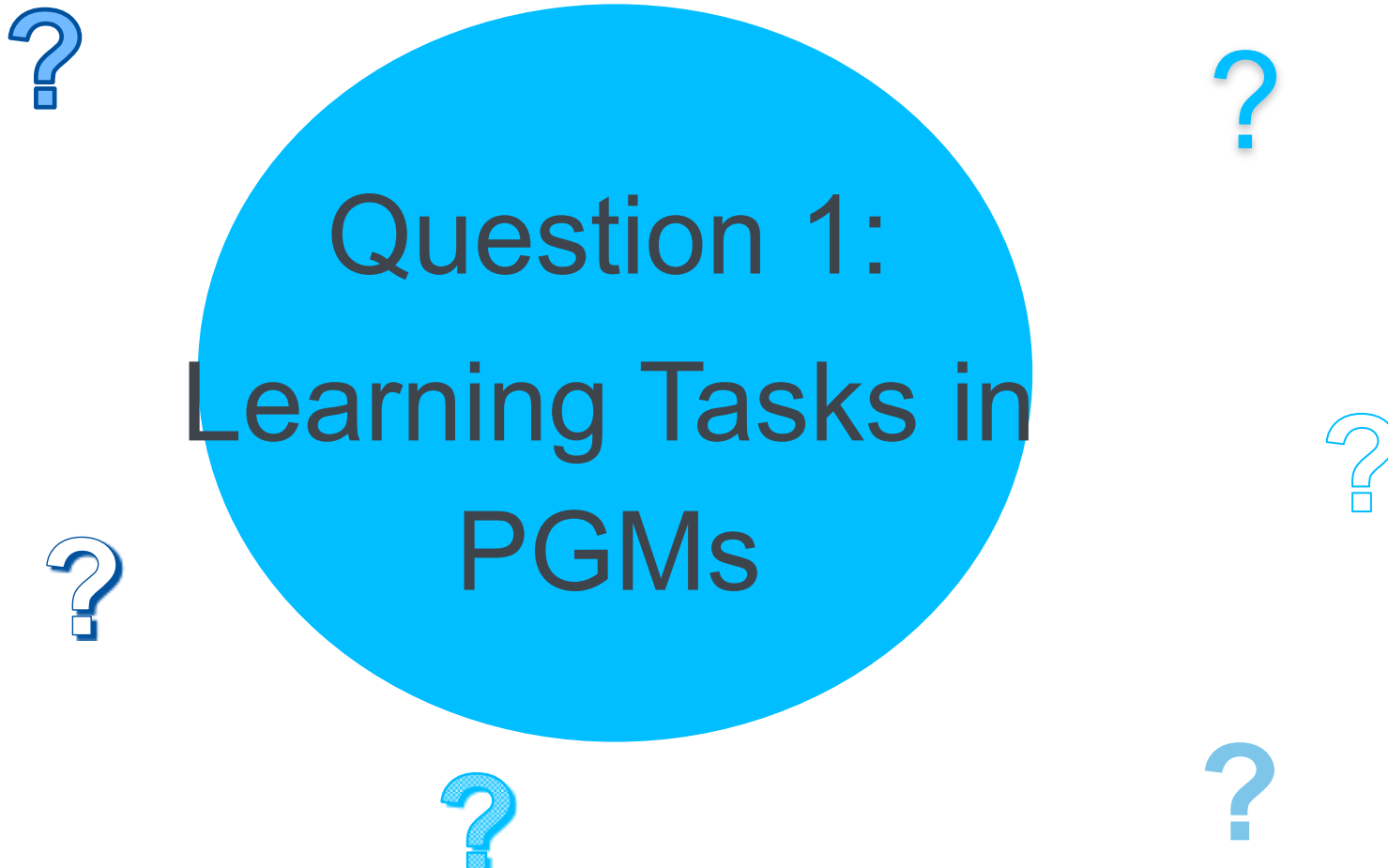
Goal: Knowledge discovery of  $\mathcal{M}^*$

- Distinguish direct vs. indirect dependencies
- Possibly directionality of edges (causality!)
- Presence and location of hidden variables
  - e.g. confounding variables / factors
- Trained using likelihood
  - does not reflect structural accuracy
- Evaluate by comparing to prior knowledge

# Overfitting

- Same problems as discussed in ML lecture
- Trade-off
  - more complex model better fits the training data
  - more complex model tends to overfit more
- Two kinds of overfitting
  - parameter overfitting
  - structure overfitting

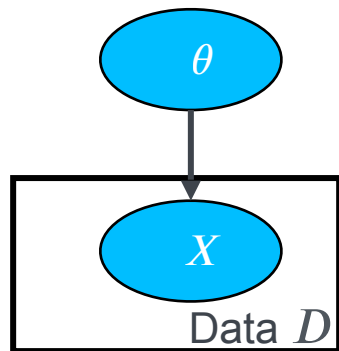
regularization / penalizing model complexity



# Question 1: Learning Tasks in PGMs

## **2 Maximum Likelihood for Bayesian Networks with Complete Observation of All Variables**

# Remember: MLE for Binary Coin Toss



- $\mathcal{X} = (X_1 \dots X_m)$  with  $\text{val}(X_i) = \{0, 1\}$ ,  $X_i$  are iid
- Which  $\theta$  best explains  $\mathcal{D} = \{x[1] \dots x[M]\}$ ?
  - $k$  times head up,  $M - k$  times head down

- Likelihood:

$$L(\mathcal{D}; \theta) = \prod_m P(x[m]; \theta) = \theta^k (1 - \theta)^{M-k}$$

- Log-likelihood:

$$\log L(\mathcal{D}; \theta) = k \log \theta + (M - k) \log(1 - \theta)$$

- Maximize by setting derivative to 0

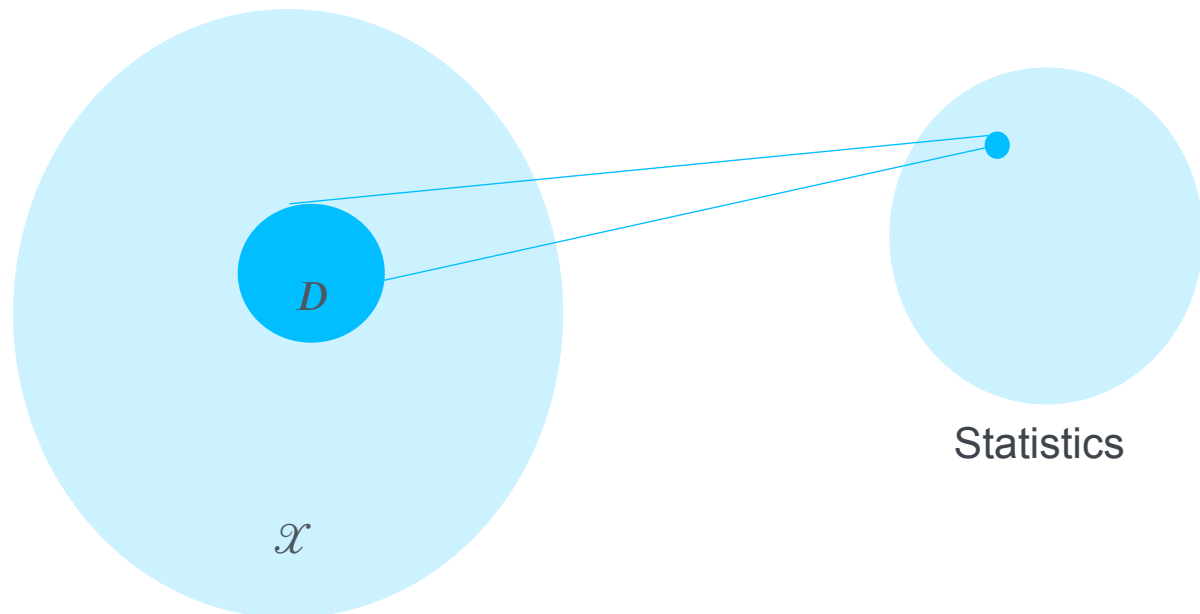
$$\frac{d \log L(\mathcal{D}; \theta)}{d\theta} = \frac{k}{\theta} - \frac{M - k}{1 - \theta} \stackrel{!}{=} 0$$

$$k(1 - \theta) = (M - k)\theta$$

$$k - k\theta = M\theta - k\theta$$

$$\theta = \frac{k}{M}$$

# Sufficient statistics



$s(x) = \begin{cases} (1 \ 0), & \text{for head up} \\ (0 \ 1), & \text{for head down} \end{cases}$   
is sufficient statistics for  
binary coin toss

A function  $s(\mathcal{D})$  is a sufficient statistics from instances to a vector in  $\mathbb{R}^k$  if for any two datasets  $\mathcal{D}$  and  $\mathcal{D}'$  and any  $\theta \in \Theta$  we have

$$\sum_{x \in \mathcal{D}} s(x) = \sum_{x \in \mathcal{D}'} s(x) \implies L(\mathcal{D}; \theta) = L(\mathcal{D}'; \theta)$$

# Sufficient Statistic for Multinomial

- Example: Bag of Words

$$\mathcal{D} = \{\text{dog, cat, dog, bee, lion, dog, dog, cat}\}$$

- For a dataset  $\mathcal{D}$  over variable  $X$  with  $|\text{val}(X)| = K$ , the sufficient statistics are counts  $\langle k_1, \dots, k_K \rangle$  where  $k_i$  is the number of times that  $x[\dots] = x^i$  in  $\mathcal{D}$

- $s(x^i) = (0 \dots 0 \ 1 \ 0 \dots 0)$ , with 1 in  $i$ -th place

- $$L(\mathcal{D}; \theta) = \prod_i \theta_i^{k_i}$$

- MLE: 
$$\hat{\theta}_i = \frac{k_i}{M}$$



# Sufficient Statistic for Gaussian

- Gaussian distribution:

$$P(X) \sim N(\mu, \sigma^2)$$

$$P(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} =$$

$$= \frac{1}{\sqrt{2\pi}\sigma} e^{\left(-x^2 \frac{1}{2\sigma^2} + x \frac{\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2}\right)}$$

- Sufficient statistics for Gaussian:

$$s(x) = \langle 1, x, x^2 \rangle$$

Remember: Gaussian MLE

$$\hat{\mu} = \frac{1}{M} \sum_{m=1}^M x[m]$$

$$\hat{\sigma}^2 = \frac{1}{M} \sum_{m=1}^M \left( x[m] - \hat{\mu} \right)^2$$

# MLE for Bayesian Networks

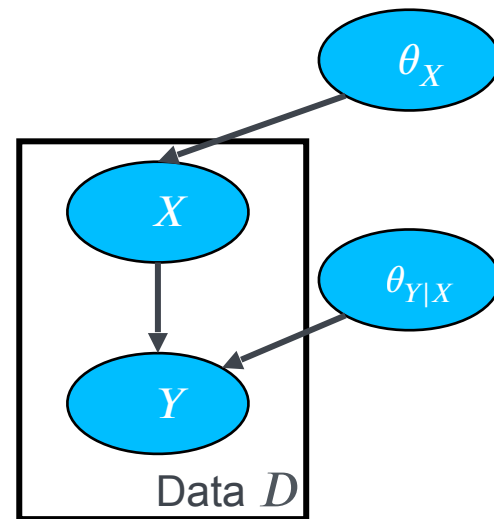
Parameters

$$\{\theta_x : x \in \text{val}(X)\}$$

$$\{\theta_{y|x} : (y, x) \in \text{val}(Y) \times \text{val}(X)\}$$

$$\begin{aligned} L(\mathcal{D}; \Theta) &= \prod_{m=1}^M P(x[m], y[m]; \Theta) = \prod_{m=1}^M P(x[m]; \Theta) P(y[m] | x[m]; \Theta) = \\ &= \left( \prod_{m=1}^M P(x[m]; \Theta) \right) \left( \prod_{m=1}^M P(y[m] | x[m]; \Theta) \right) \\ &= \left( \prod_{m=1}^M P(x[m]; \theta_X) \right) \left( \prod_{m=1}^M P(y[m] | x[m]; \theta_{Y|X}) \right) \end{aligned}$$

$\underbrace{\hspace{10em}}$   
optimized like a multinomial      also optimized like a multinomial



local likelihoods

# Likelihood for Bayesian Network in General

$$\begin{aligned}
 L(\mathcal{D}; \Theta) &= \prod_{m=1}^M P(x[m]; \Theta) = \\
 &= \prod_{m=1}^M \prod_i P(x_i[m] \mid \text{pa}_{X_i}[m]; \Theta) = \\
 &= \prod_i \prod_{m=1}^M \underbrace{P(x_i[m] \mid \text{pa}_{X_i}[m]; \Theta)}_{\substack{\text{value assignments to parents of } X_i \\ \text{in data item } x[m]}} = \\
 &= \prod_i L_i(\mathcal{D}; \Theta_i)
 \end{aligned}$$

product of local  
likelihoods

If  $\theta_{X_i \mid \text{pa}_{X_i}}$  are disjoint,

then MLE can be computed by maximizing each local likelihood separately

with

$$L_i(\mathcal{D}; \Theta_i) = P(x_i[m] \mid \text{pa}_{X_i}[m]; \theta_{X_i \mid \text{pa}_{X_i}})$$

# Likelihood for Table CPDs

Given:

- Data  $\mathcal{D} = \left\{ \mathbf{x}[1] \dots \mathbf{x}[M] \right\} \sim \mathbf{P}(\mathcal{X})$ ,  $\mathcal{X} = (X_1 \times \dots \times X_l)$
- structure of tables representing  $P(X_i | \text{pa}_{X_i})$ , entries unknown
- $k(\mathbf{u}, x)$  counts how often value combination  $\mathbf{u}, x$  is observed in  $D$

Output:  $\hat{\theta}_{X_i | \text{pa}_{X_i}}$  entries for all  $P(X_i | \text{pa}_{X_i})$

Approach: local likelihood function

$$\begin{aligned} L_{X_i}(\mathcal{D}; \theta_{X_i | \text{pa}_{X_i}}) &= \prod_{m=1}^M \theta_{x_i[m] | \text{pa}_{X_i}[m]} = \\ &= \prod_{\mathbf{u} \in \text{val}(\text{pa}_{X_i})} \left[ \prod_{x \in \text{val}(X_i)} \left( \theta_{x | \mathbf{u}} \right)^{k(\mathbf{u}, x)} \right] \end{aligned}$$

Multinomial MLE:

$$\hat{\theta}_{x | \mathbf{u}} = \frac{k(\mathbf{u}, x)}{k(\mathbf{u})}$$

# MLE for Linear Gaussian Bayesian Network

$$P(X_i | \mathbf{u}) = \mathcal{N}(\beta_0 + \beta_1 u_1 + \dots + \beta_l u_l; \sigma^2)$$

$$\ell_X(\mathcal{D}; \theta_{X|\text{pa}_X}) = \log L_X(\mathcal{D}; \theta_{X|\text{pa}_X}) =$$

$$= \sum_{m=1}^M \left[ -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\beta_0 + \beta_1 u_1[m] + \dots + \beta_l u_l[m])^2 \right]$$

Closed Form  
Solution!

Computing the partial derivatives  $\frac{\partial}{\partial \beta_i}$  and setting the result to 0

leads to a set of linear equations that can be solved.

# Data Fragmentation and Overfitting

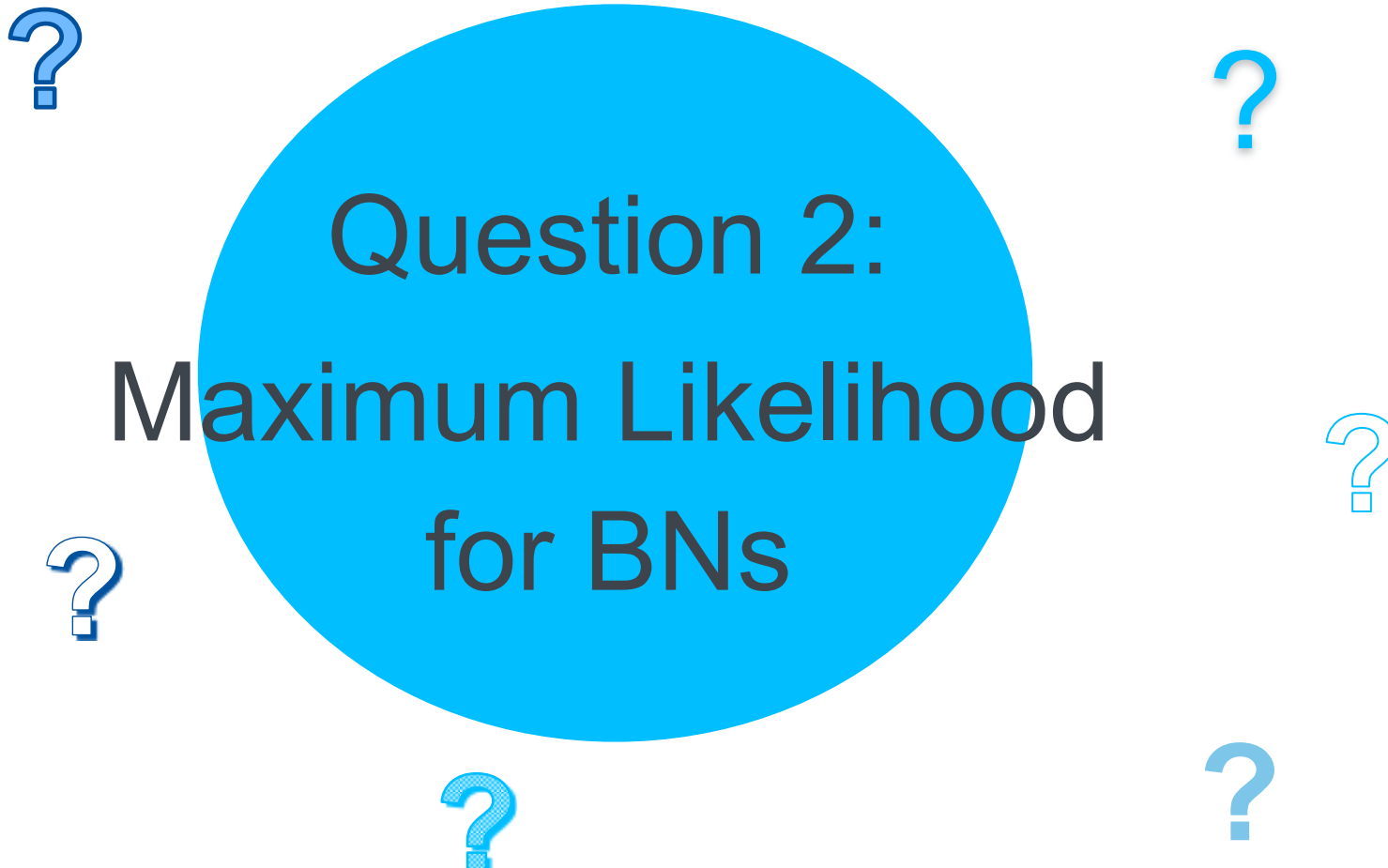
- Data  $\mathcal{D} = \{ \mathbf{x}[1] \dots \mathbf{x}[M] \}$
- $\hat{\theta}_{x|u}$  is estimated based on  $\frac{M}{k(u)}$  data points
- Possible value assignments  $\mathbf{u}$  to  $\text{pa}_{X_i}$  grows
  - exponentially with number of parents
  - larger set of values  $\longrightarrow$  larger basis for the exponent

Data Fragmentation leads to Overfitting.  
Simpler network structures may prevent fragmentation and overfitting.

Example:  
Naive Bayes classifier  
avoids complex structure

# Summary

- For Bayesian Networks with disjoint sets of parameters in CPD, likelihood decomposes as product of local likelihood functions
- For table CPDs, local likelihood further decomposes as product of likelihood for multinomials – one for each parent combination
- For networks with shared CPDs (e.g. HMMs), statistics accumulate over all uses of CPDs



# Question 2: Maximum Likelihood for BNs



# **3 Maximum Likelihood for Log-linear Models**

# Log-likelihood for Markov Networks



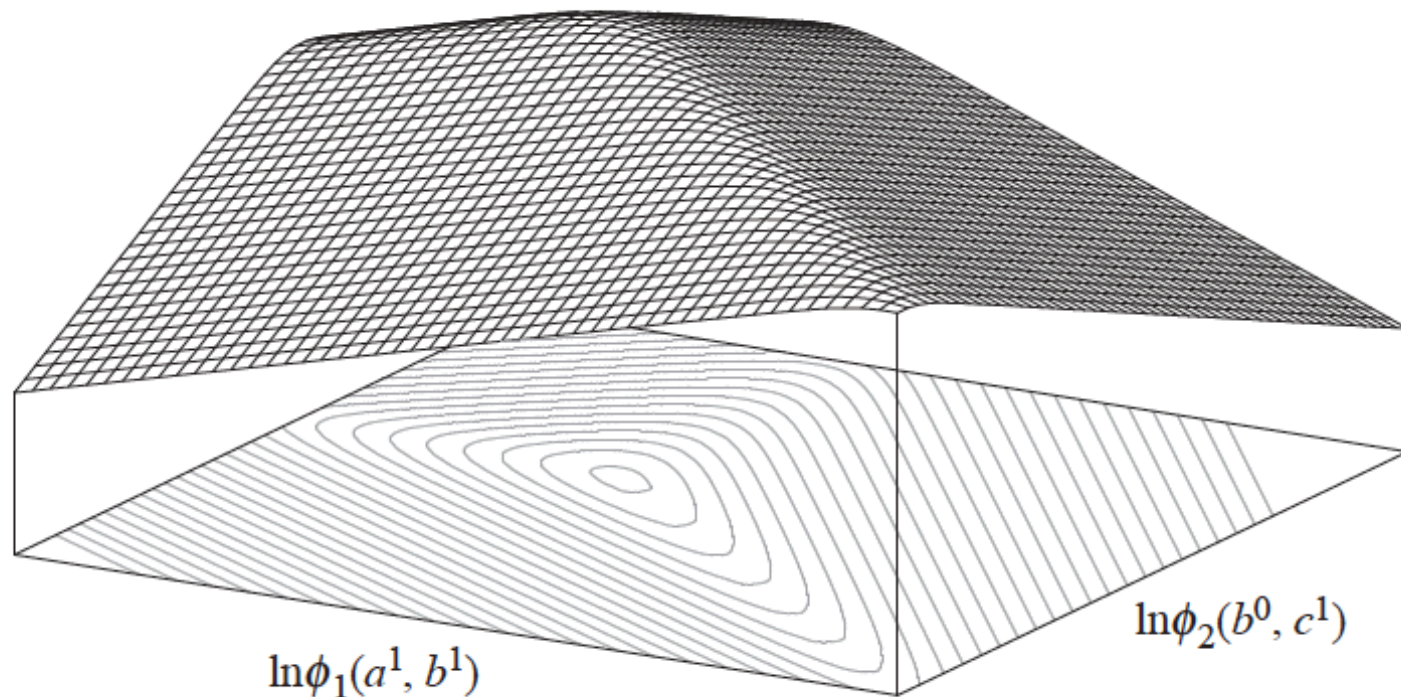
$$P(a, b, c; \Theta) = \frac{1}{Z(\Theta)} \phi_1(a, b) \phi_2(b, c)$$

log-likelihood

$$\begin{aligned} \ell(\mathcal{D}; \theta) &= \\ &= \sum_{m=1}^M (\ln \phi_1(a[m], b[m]) + \ln \phi_2(b[m], c[m]) - \ln Z(\Theta)) = \\ &= \sum_{a,b} k(A = a, B = b) \ln \phi_1(a, b) + \\ &+ \sum_{b,c} k(B = b, C = c) \ln \phi_2(b, c) - M \ln Z(\Theta) \\ Z(\Theta) &= \sum_{a,b,c} \phi_1(a, b) \phi_2(b, c) \end{aligned}$$

Partition function couples parameters

- No decomposition of likelihood
- No closed form solution



**Figure 20.1** Log-likelihood surface for the Markov network  $A-B-C$ , as a function of  $\ln \phi_1(a^1, b^1)$  ( $x$ -axis) and  $\ln \phi_2(b^0, c^1)$  ( $y$ -axis); all other parameters in both potentials are set to 1. Surface is viewed from the  $(+\infty, +\infty)$  point toward the  $(-, -)$  quadrant. The data set  $\mathcal{D}$  has  $M = 100$  instances, for which  $M[a^1, b^1] = 40$  and  $M[b^0, c^1] = 40$ . (The other sufficient statistics are irrelevant, since all of the other log-parameters are 0.)

# Log-linear models

A distribution  $P(\mathcal{X})$  is a log-linear model over a Markov Network  $\mathcal{H}$  if it is associated with:

- a set of features  $\mathcal{F} = \{f_1(D_1), \dots, f_l(D_l)\}$ ,
  - where each  $D_i$  is a complete subgraph in  $\mathcal{H}$
  - where each feature  $f_i: \text{val}(D_i) \rightarrow \mathbb{R}$
- a set of weights  $w_1, \dots, w_l$

such that

$$P(X_1, \dots, X_n) = \frac{1}{Z} e^{\left(-\sum_{i=1}^l w_i f_i(D_i)\right)}$$

**The log-linear model expresses the factor product as the sum of features in an exponential.  
For some distributions this representation is more compact.**

# Log-likelihood for log-linear model

$$\ell(\mathcal{D}; \theta) = \sum_{i=1}^l \theta_i \left( \sum_{m=1}^M f_i(x[m]) \right) - M \ln Z(\theta)$$

log-sum-exp:

$$\ln Z(\theta) = \ln \sum_{x \in \mathcal{X}} e^{\left( \sum_i \theta_i f_i(x) \right)}$$

exponentially large space

# Log-Partition Function

Theorem:

$$\frac{\partial}{\partial \theta_i} \ln Z(\theta) = \mathbb{E}_\theta[f_i] = \sum_{\mathbf{x} \in \mathcal{X}} P_\theta(\mathbf{x}) f_i(\mathbf{x})$$

$$\frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln Z(\theta) = \text{Cov}_\theta[f_i; f_j]$$

The Hessian of the log-sum-exp is the covariance matrix,  
which is always positive semi-definite and  
therefore weakly convex

$$\frac{\partial}{\partial \theta_i} \ln Z(\theta) = \mathbb{E}_\theta[f_i]$$

Proof Part 1:

$$\begin{aligned} \frac{\partial}{\partial \theta_i} \ln Z(\theta) &= \frac{1}{Z(\theta)} \sum_{\mathbf{x} \in \mathcal{X}} \frac{\partial}{\partial \theta_i} e^{\left(\sum_j \theta_j f_j(\mathbf{x})\right)} = \\ &= \frac{1}{Z(\theta)} \sum_{\mathbf{x} \in \mathcal{X}} f_i(\mathbf{x}) e^{\left(\sum_j \theta_j f_j(\mathbf{x})\right)} = \sum_{\mathbf{x} \in \mathcal{X}} \frac{1}{Z(\theta)} e^{\left(\sum_j \theta_j f_j(\mathbf{x})\right)} f_i(\mathbf{x}) = \\ &= \sum_{\mathbf{x} \in \mathcal{X}} P_\theta(\mathbf{x}) f_i(\mathbf{x}) = \mathbb{E}_\theta[f_i] \end{aligned}$$

Proof Part 2 (Koller & Friedman, pp. 948)

## Optimizing $\ell(\mathcal{D}; \theta)$

Consider: 
$$\ell(\mathcal{D}; \theta) = \sum_{i=1}^l \theta_i \left( \sum_{m=1}^M f_i(x[m]) \right) - M \ln Z(\theta)$$

The first term is linear in  $\theta$ .

The second term is concave in  $\theta$  (“- convex”).

$\implies$  The overall sum is concave

$\implies$  local optimum is global optimum

$\implies$  easy to optimize using gradient ascent



# Maximum Likelihood Estimation

$$\frac{1}{M} \ell(\mathcal{D}; \theta) = \sum_{i=1}^l \theta_i \left( \frac{1}{M} \sum_{m=1}^M f_i(x[m]) \right) - \ln Z(\theta)$$

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\mathcal{D}; \theta) = \mathbb{E}_{\mathcal{D}}[f_i(\mathcal{X})] - \mathbb{E}_{\theta}[f_i]$$

Expectation of  $f_i$  in  $\mathcal{D}$       Expectation of  $f_i$  in  $P_{\theta}(\mathcal{X})$

Theorem:  $\hat{\theta}$  is the MLE if and only if  $\mathbb{E}_{\mathcal{D}}[f_i(\mathcal{X})] = \mathbb{E}_{\hat{\theta}}[f_i]$

# Computation of Gradient Ascent

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\mathcal{D}; \theta) = \mathbb{E}_{\mathcal{D}}[f_i(\mathcal{X})] - \mathbb{E}_{\theta}[f_i]$$

- Use gradient ascent
  - e.g. a quasi-Newton method like L-BFGS (
    - avoids expensive computation of the Hessian)
  - e.g. stochastic gradient descent with momentum (see ML-Chapter 9)
- Needed for gradient: expected feature counts
  - in data ✓
  - **relative to current model**
    - ⇒ expensive inference step at each gradient step 🙄

# Summary

- Partition function couples parameters in likelihood
- No closed form solution but convex optimization
  - solved using gradient ascent
- Gradient computation requires inference at each gradient step to compute expected feature counts
- Features are always within cluster in cluster graph
  - one calibration suffices for all feature expectations



# Question 3: Maximum Likelihood for Log-Linear Models



# 4 Maximum Likelihood for Conditional Random Fields

# Estimation for CRFs

$$\mathcal{D} = \{\mathbf{x}[m], \mathbf{y}[m]\}_{m=1}^M$$

$$P_{\theta}(\mathbf{Y} | \mathbf{x}) = \frac{1}{Z_{\mathbf{x}}(\theta)} \tilde{P}_{\theta}(\mathbf{x}, \mathbf{Y}) \quad Z_{\mathbf{x}}(\theta) = \sum_{\mathbf{Y}} \tilde{P}_{\theta}(\mathbf{x}, \mathbf{Y})$$

$$\ell_{\mathbf{Y} | \mathbf{x}}(\mathcal{D}; \theta) = \sum_{m=1}^M \ln P_{\theta}(\mathbf{y}[m] | \mathbf{x}[m]; \theta)$$

Considering single data point:

$$\ell_{\mathbf{Y} | \mathbf{x}}(\theta; \mathbf{x}[m], \mathbf{y}[m]) = \left( \sum_i \theta_i f_i(\mathbf{x}[m], \mathbf{y}[m]) \right) - \ln Z_{\mathbf{x}[m]}(\theta)$$

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell_{\mathbf{Y} | \mathbf{x}}(\mathcal{D}; \theta) = \frac{1}{M} \sum_{m=1}^M \left( f_i(\mathbf{x}[m], \mathbf{y}[m]) - \mathbb{E}_{\theta}[f_i | \mathbf{x}[m]] \right)$$

Basically the same formulas as before – with conditions  $\mathbf{x}[m]$  thrown in

# Comparing Computations

## MRF

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\mathcal{D}; \theta) = \mathbb{E}_{\mathcal{D}}[f_i(\mathcal{X})] - \mathbb{E}_{\theta}[f_i]$$

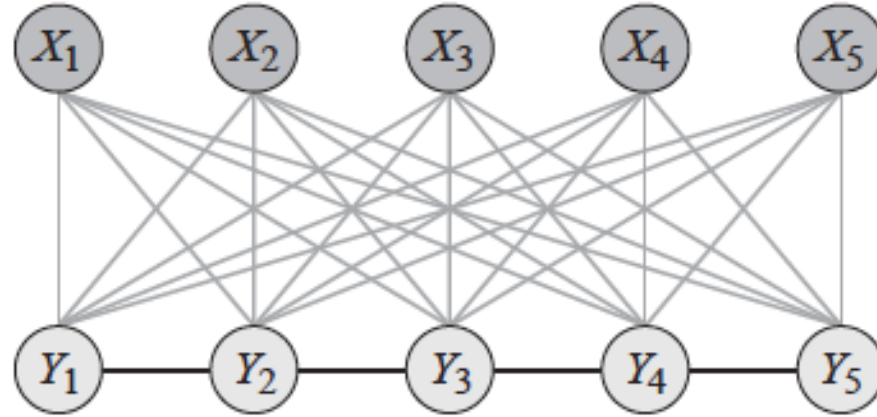
- Requires inference at each gradient step

## CRF

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell_{Y|\mathbf{X}}(\mathcal{D}; \theta) = \frac{1}{M} \sum_{m=1}^M \left( f_i(\mathbf{x}[m], \mathbf{y}[m]) - \mathbb{E}_{\theta}[f_i | \mathbf{x}[m]] \right)$$

- Requires inference **for each  $\mathbf{x}[m]$**  at each gradient step

# What is better?

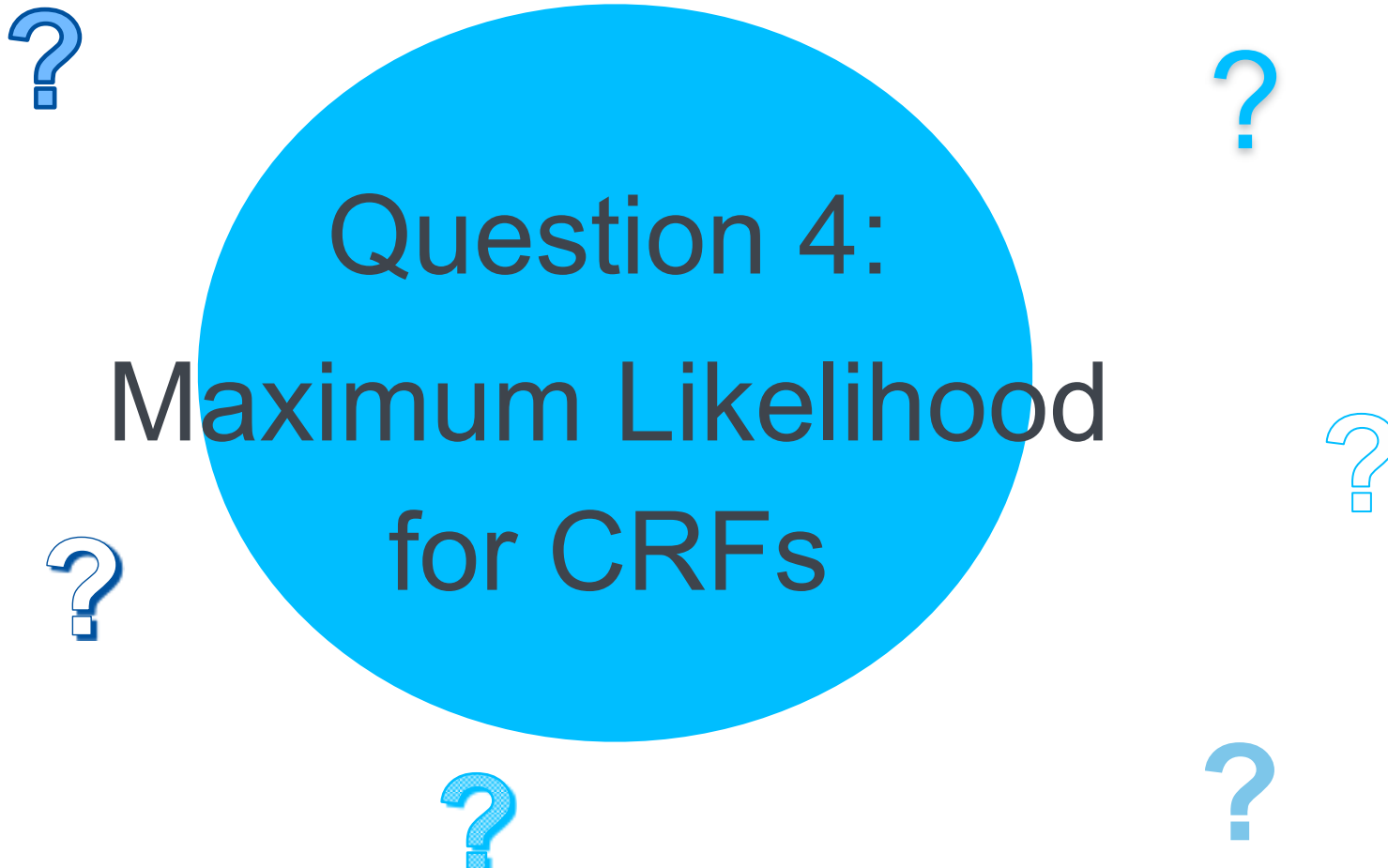


**Figure 20.2** A highly connected CRF that allows simple inference when conditioned: The edges that disappear in the reduced Markov network after conditioning on  $X$  are marked in gray; the remaining edges form a simple linear chain.

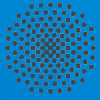


# Summary

- CRF learning very similar to MRF learning
  - likelihood function is concave
  - optimized using gradient ascent
- Gradient computation requires inference:
  - one per gradient step, data instance
    - cf: once per gradient step for MRFs
- But conditional model is often much simpler, so inference cost for CRF may even be lower



# Question 4: Maximum Likelihood for CRFs



Universität Stuttgart  
IPVS

# Thank you!



**Steffen Staab**

E-Mail [Steffen.staab@ipvs.uni-stuttgart.de](mailto:Steffen.staab@ipvs.uni-stuttgart.de)

Telefon +49 (0) 711 685-~~56~~ be defined

[www.ipvs.uni-stuttgart.de/departments/ac/](http://www.ipvs.uni-stuttgart.de/departments/ac/)

Universität Stuttgart

Analytic Computing, IPVS

Universitätsstraße 32, 50569 Stuttgart