# Probabilistic Graphical Models

# Lectures 22

## Introduction to Learning
## Learning Bayesian Networks

# Sampling vs. Learning

$$P_\theta(X) = P(X; \theta) = P(X \mid \theta)$$

Sampling $\quad P_\theta(X) \overset{\text{given}}{=} \checkmark \implies$ generate data $X^1, X^2, \ldots, X^m$

$$\text{i.i.d samples}$$

Learning $\quad$ data $X^1, X^2, \ldots, X^m \overset{\text{given}}{=} \checkmark \implies$ find $P_\theta(X)$

find $P_\theta(X)$ $\begin{cases} \text{find parameters } \theta \\ \text{find structure (connections in BN/MRF)} \\ \qquad\qquad\qquad\qquad\qquad\qquad \text{graph} \end{cases}$

# Likelihood function

given data $X^1, X^2, \ldots, X^m$ what is the probability of

occurrence of $X^1, \underline{\quad} X^m$?

$$Pr(X^1, X^2, \ldots, X^m) = \prod_{i=1}^{m} Pr(X^i) = \prod_{i=1}^{m} P_\theta(X^i)$$

independet
samples

all are
samples
of $P_\theta(X)$

$$\ell(\theta) = \prod_{i=1}^{m} P_\theta(X^i) \quad \text{likelihood}$$

# Maximum Likelihood Solution

$i=1$

One Solution: choose $\theta$ that maximized the likelihood

$$\theta^* = \text{argmax}_\theta \, \ell(\theta) = \text{argmax}_\theta \prod_{i=1}^{m} P_\theta(X^i)$$

maximum-likelihood solution

4

# Maximum Likelihood Solution



One Solution: choose $\theta$ that maximized the likelihood

$$\theta^* = \text{argmax}_\theta \; \ell(\theta) = \text{argmax}_\theta \; \prod_{i=1}^{m} P_\theta(X^i)$$

maximum-likelihood solution

# Example: tossing a coin

(unfair) coin $\qquad$ $P_\theta(X)$ $\quad X \in \{H, T\}$

$$\theta = Pr(X = H) = P_\theta(H) \qquad Pr(X = T) = 1 - \theta = P_\theta(T)$$

$$X^1, X^2, \ldots, X^m$$

$$H \quad H \quad T \quad H \quad T \quad T \ldots$$

# Example: tossing a coin

$$\ell(\theta) = \prod_{i=1}^{m} P_\theta(X^i)$$

$$P_\theta(X^i) = \begin{cases} \theta & X^i = H \\ 1-\theta & X^i = T \end{cases}$$

$$= \prod_{i=1}^{m} \left( \theta \, \mathbb{1}(X^i = H) + (1-\theta) \, \mathbb{1}(X^i = T) \right)$$

$n_H =$ no. of $H$

$n_T =$ no. of $T$

$m = n_H + n_T$

$$\implies \ell(\theta) = \theta^{n_H} (1-\theta)^{n_T}$$

# Example: tossing a coin

$$\ell\left(\theta = \frac{n_H}{m}\right) = \left(\frac{n_H}{m}\right)^{n_H} \left(1 - \frac{n_H}{m}\right)^{n_T}$$

$$= \left(\frac{n_H}{m}\right)^{n_H} \left(\frac{n_T}{m}\right)^{n_T} \quad (> 0 \quad \text{if } n_H, n_T > 0)$$

$$\ell(\theta = 0) = \ell(\theta = 1) = 0 \quad (\text{if } n_H, n_T > 0)$$

$$\theta^* = \frac{n_H}{n_H + n_T} = \frac{n_H}{m}, \quad \text{maximum-likelihood solution}$$

$$\text{log-likelihood} = \ell\ell(\theta)$$

# Example: tossing a coin - log-likelihood

$$\theta^* = \underset{\theta}{\arg\max} \; \ell(\theta) = \underset{\theta}{\arg\max} \; \log \ell(\theta) = \underset{\theta}{\arg\max} \; \log \prod_{i=1}^{m} P_\theta(x^i)$$

log-likelihood $= \ell\ell(\theta)$

$n_H + n_T$

$$= \underset{\theta}{\arg\max} \; \sum_{i=1}^{m} \log P_\theta(x^i)$$

$$= \underset{\theta}{\arg\max} \; n_H \log \theta + n_T \log(1-\theta)$$

$$\max \ell\ell(\theta) \Rightarrow \frac{d}{d\theta} \ell\ell(\theta) = \frac{d}{d\theta} n_H \log\theta + n_T \log(1-\theta) \Rightarrow \frac{n_H}{\theta} - \frac{n_T}{1-\theta}$$

$$= \frac{n_H(1-\theta) - n_T \theta}{\theta(1-\theta)} = 0 \Rightarrow \boxed{\theta^* = \frac{n_H}{n_H + n_T}}$$

# Sufficient statistics

$$\text{data } x^1, x^2, \ldots, x^m \Rightarrow \ell(\theta) = \theta^{n_H}(1-\theta)^{n_T}$$

$$\ell\ell(\theta) = n_H \log \theta + n_T \log (1-\theta)$$

$$n_H = \sum 1(x^i = H)$$
$$n_T = \sum_{i=1}^{m} 1(x^i = T)$$

for calculating the likelihood
we only need $n_H$, $n_T$ (not the
complete data)

$n_H$, $n_T$ : sufficient statistics

Example 2: normal distribution

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

27

10

# Example 2: Normal Distribution

Example 2: normal distribution

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

data $x^1, x^2, \ldots, x^m$    find maximum-likelihood solutions for $\mu, \sigma$

# Example 2: Normal Distribution

$$ll(\theta) = \sum_{i=1}^{m} \log P_{\theta}(x^i) = \sum_{i=1}^{m} \left( -\log \sqrt{2\pi} - \log \sigma - \frac{1}{2} \frac{(x^i - \mu)^2}{\sigma^2} \right)$$

$$ll(\theta) = C - m \log \sigma - \frac{1}{2} \frac{\sum_{i=1}^{m} \left[ (x^i)^2 + \mu^2 - 2\mu x^i \right]}{\sigma^2}$$

$$= C - m \log \sigma - \frac{1}{2\sigma^2} \left( \sum_{i=1}^{m} (x^i)^2 - 2\mu \sum_{i=1}^{m} x^i + m\mu^2 \right)$$

sufficient statistics $\left( \sum (x^i)^2, \sum x^i, m \right)$

# Example 2: Normal Distribution

$$\frac{\partial}{\partial \mu} \, \ell\ell(\theta) = \frac{\partial}{\partial \mu} \ell\ell(\mu, \sigma) = -\frac{1}{2\sigma^2}\left(-2\sum_{i=1}^{m} x^i + 2m\mu\right) = \emptyset$$

$$\Rightarrow 2m\mu = 2\sum_{i=1}^{m} x^i \Rightarrow \boxed{\mu^* = \frac{1}{m}\sum_{i=1}^{m} x^i}$$

$$\frac{\partial}{\partial \sigma} \ell\ell(\theta) = \frac{-m}{\sigma} - \frac{1}{2} \times (-2)\frac{\sum (x^i - \mu^*)^2}{\sigma^3} = 0$$

$$m\sigma^2 = \sum_{i=1}^{m}(x^i - \mu^*)^2$$

$$\Rightarrow \boxed{\sigma^2 = \frac{1}{m}\sum (x^i - \mu^*)^2} \qquad \sigma = \sqrt{\sigma^2}$$

# Example 3: Tossing a dice

$$P_\theta(X) \qquad X \in \{1, 2, \cdots, q\}$$

$$X \in \{1, 2, 3, 4, 5, 6\} \qquad \theta = (\theta_1, \theta_2, \cdots, \theta_q)$$

$$\theta_1 = Pr(X=1)$$
$$\theta_2 = Pr(X=2)$$
$$\vdots$$
$$\theta_q = Pr(X=q) = 1 - \theta_1 - \theta_2 - \cdots - \theta_{q-1}$$

$$\sum_{j=}^{q} \theta_j = 1 \qquad\qquad data = X^1, X^2, \cdots, X^m$$

$$\ell(\theta) = \ell(\theta_1 - \theta_q) = \prod_{i=1}^{m} P_\theta(X) = \theta_1^{h_1}, \theta_2^{h_2} \cdots \theta_q^{h_q}$$

14

# Example 3: Tossing a dice

$$l(\theta) = l(\theta_1 - \theta_q) = \prod_{i=1}^{m} P_\theta(X) = \theta_1^{n_1} \cdot \theta_2^{n_2} \cdots \theta_q^{n_q}$$

suff. stat: $n_j = \sum_{i=1}^{m} 1(X^{a_i} = j) = \# X^{a_i} = j$

$$ll(\theta) = \sum_{i=1}^{m} n_i \log \theta_i$$

# Example 3: Tossing a dice

$$\theta^* = \underset{\theta_1, \theta_2 - \theta_q}{\arg\max} \sum_{j=1}^{q} n_j \log \theta_j \quad \text{subject to} \quad \sum_{j=1}^{q} \theta_j = 1$$

$$\begin{cases} \text{lagrange multipliers} \\ \quad\quad\quad \text{or} \\ \text{or } \theta_q = 1 - \sum_{j=1}^{q-1} \theta_j \end{cases}$$

$$\implies \theta_j^* = \frac{n_j}{\sum_{k=1}^{q} n_k} = \frac{n_j}{m}$$

16

# PGM problems

$$\text{PGM} \qquad P_\theta(X) = P_\theta(X_1, X_2, \cdots, X_n) \quad n:\text{large}$$

$$\text{data:} \quad X^1 = (X_1', X_2', X_3', \cdots, X_n')$$

$$X^2 = (X_1', X_2', \cdots, X_n')$$

$$\vdots$$

$$X^m = (X_1', X_2', \cdots, X_n')$$

# Bayes Nets Parameter Learning

pgm 22

$$P_\theta(X) = P_\theta(X_1, X_2, \ldots, X_n)$$

$$= \prod_{i=1}^{n} P_\theta(X_i \mid X_{P_i})$$

<span style="color:red">→ parents of $X_i$</span>

Data: $X^1, X^2, \ldots, X^m = (X_1^1, X_2^1, \ldots, X_n^1), (X_1^2, X_2^2, \ldots, X_n^2), \ldots, (X_1^m, X_2^m, \ldots, X_n^m)$

$$ll(\theta) = \sum_{k=1}^{m} \log P_\theta(X^k) = \sum_{k=1}^{m} \log \prod_{i=1}^{n} P_\theta(X_i^k \mid X_{P_i}^k)$$

$$= \sum_{k=1}^{m} \sum_{i=1}^{n} \log P_\theta(X_i^k \mid X_{P_i}^k)$$

18

# No shared parameters

1: Each CPD has its own parameters

$$\theta = (\theta_1, \theta_2, \dots, \theta_n)$$

$$P_\theta(X_i \mid X_{P_i}) = P_{\theta_i}(X_i \mid X_{P_i})$$

$$\Rightarrow \ell\ell(\theta) = \sum_{k=1}^{m} \sum_{i=1}^{n} \log P_\theta(X_i^k \mid X_{P_i}^k)$$

$$= \sum_{k=1}^{m} \sum_{i=1}^{n} \log P_{\theta_i}(X_i^k \mid X_{P_i}^k)$$

$$= \sum_{i=1}^{n} \underbrace{\sum_{k=1}^{m} \log P_{\theta_i}(X_i^k \mid X_{P_i}^k)}_{\text{local log-likelihood}}$$

19

# No shared parameters

$$\Rightarrow \ell\ell(\theta) = \sum_{k=1}^{m} \sum_{i=1}^{n} \log P_\theta \left( X_i^k \mid X_{P_i}^k \right)$$

$$= \sum_{k=1}^{m} \sum_{i=1}^{n} \log P_{\theta_i} \left( X_i^k \mid X_{P_i}^k \right)$$

$$= \sum_{i=1}^{n} \underbrace{\sum_{k=1}^{m} \log P_{\theta_i} \left( X_i^k \mid X_{P_i}^k \right)}_{\text{local log-likelihood}}$$

$$\frac{\partial}{\partial \theta_j} \ell\ell(\theta) = \frac{\partial}{\partial \theta_j} \ell\ell(\theta_1, \theta_2, \dots, \theta_n)$$

gradient
w.r.t.
$\theta_j$

$$= \sum_{k=1}^{m} \frac{\partial}{\partial \theta_j} \log P_{\theta_j} \left( X_i^k \mid X_{P_i}^k \right)$$

$$= \sum_{k=1}^{m} \frac{\frac{\partial}{\partial \theta_j} P_{\theta_j} \left( X_i^k \mid X_{P_i}^k \right)}{P_{\theta_j} \left( X_i^k \mid X_{P_i}^k \right)}$$

each $\theta_j$
can be
found independently
of other $\theta_i$-s

20

# No shared parameters - table representation

$$\sum_{i=1}^{n} \sum_{k=1}^{m} \log P_{\theta_i}(X_i^k \mid X_{P_i}^k)$$

$\underbrace{\qquad\qquad\qquad\qquad}_{\text{local log-likelihood}}$

<u>table representation</u>

Example $P(X_i \mid Y_i)$ $X_i, X_i \in \{0, 1\}$

| $X_i$ | $Y_i$ | $P(X_i \mid Y_i)$ |
|-------|-------|-------------------|
| 1 | 0 | $\theta_0$ |
| 0 | 0 | $1-\theta_0$ |
| 1 | 1 | $\theta_1$ |
| 0 | 1 | $1-\theta_1$ |

→ local likelihood

$$\sum_{k=1}^{m} \log P(X_i^k \mid Y_i^k)$$

$$\frac{\partial}{\partial \theta_0} \sum_{k=1}^{m} \log P(X_i^k \mid Y_i^k)$$

$$\frac{\partial}{\partial \theta_0} \sum_{\substack{k=1 \\ Y_i^k=0}}^{m} \log P(X_i^k \mid 0) + \sum_{\substack{k=1 \\ Y_i^k=1}}^{m} \log P(X_i^k \mid 1)$$

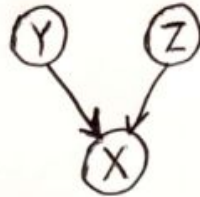a function of $\theta_0$ \qquad a function of $\theta_1$

$$= \frac{\partial}{\partial \theta_0} \left( n_{10} \log \theta_0 + n_{00} \log(1-\theta_0) \right) \Rightarrow \theta_0^* = \frac{n_{10}}{n_{10} + n_{11}}$$

$$n_{10} = \#\left( X_i^k = 1, Y_i^k = 0 \right)$$

$$n_{00} = \#\left( X_i^k = 0, Y_i^k = 0 \right)$$

21

# No shared parameters - table representation



$$X \in \{1, 2, \ldots, C\}$$
$$Y, Z \in \{0, 1\}$$

$$\gamma_{00}^i = \checkmark$$
$$\gamma_{01}^i = Pr(X = i \mid Y = 0, Z = 1)$$
$$\gamma_{10}^i = Pr(X = i \mid Y = 1, Z = 0)$$
$$\gamma_{11}^i = Pr(X = i \mid Y = 1, Z = 1)$$

$$\sum_i \gamma_{10}^i = 1$$

local likelihood
$$\sum_{k=1}^{m} \log P_\gamma(X^k \mid Y^k, Z^k)$$

$$\sum_{k=1}^{m} \log \gamma_{Y^k, Z^k}^{X^k}$$

# No shared parameters - table representation

local likelihood $\sum_{k=1}^{m} \log P_\gamma(X^k \mid \Upsilon, Z)$

$$\sum_{k=1}^{m} \log \; \gamma_{\Upsilon, Z^k}^{X^k}$$

$$\sum_{y=0}^{1} \sum_{z=0}^{1} \sum_{x=1}^{c} \left( \log \gamma_{yz}^{x} \right) \cdot \left[ \#\left( X^k = x, \Upsilon^k = y, Z^k = z \right) \right]$$

sufficient statistics

$$\boxed{\sum_{x} \gamma_{yz}^{x} = 1}$$

$$\gamma_{yz}^{x} = \frac{\#(X^k = x, \Upsilon^k = y, Z^k = z)}{\#(\Upsilon^k = y, Z^k = z)}$$

23

# No shared parameters - table representation

$$\sum_{y=0}^{1} \sum_{z=0}^{1} \sum_{x=1}^{c} \left( \log \gamma_{yz}^{x} \right) \cdot \left[ \# \left\{ \left( X^{k}=x, Y^{k}=y, Z^{k}=z \right) \right. \right]$$

sufficient statistics

$$\boxed{\sum_{x} \gamma_{yz}^{x} = 1}$$

$$\gamma_{yz}^{x} = \frac{\#\left( X^{k}=x, Y^{k}=y, Z^{k}=z \right)}{\#\left( Y^{k}=y, Z^{k}=z \right)}$$

$$= \frac{\sum_{k=1}^{m} 1\left( X^{k}=x, Y^{k}=y, Z^{k}=z \right)}{\sum_{k=1}^{m} 1\left( Y^{k}=y, Z^{k}=z \right)}$$

# shared parameters



$$P_\theta(X_1 - X_n, Y_1 - Y_n) = \underset{\alpha}{P}(X_1) \prod_{i=2}^{n} \underset{\beta}{P}(X_i | X_{i-1}) \prod_{i=1}^{n} \underset{\gamma}{P}(Y_i | X_i)$$

$$\theta = (\alpha, \beta, \gamma)$$

$$ll(\theta) = \sum_{k=1}^{m} \log \underset{\alpha}{P}(X_1^k) + \sum_{k=1}^{m} \sum_{i=2}^{n} \log \underset{\beta}{P}(X_i^k | X_{i-1}^k)$$

$$+ \sum_{k=1}^{m} \sum_{i=1}^{n} \log \underset{\gamma}{P}(Y_i^k | X_i^k)$$

$$\frac{\partial ll(\theta)}{\partial \beta} = \sum_{k=1}^{m} \sum_{i=2}^{n} \frac{\partial}{\partial \beta} \log P_\beta(X_i^k | X_{i-1}^k)$$

25

table representation $\quad X_i \in \{0,1\}$

$\beta_0 = \beta_{00}, \beta_{10}, \beta_{11}, \beta_{01}$

$\beta_{01} = P(X_i = 0 \mid X_{i-1} = 1)$ independent of $i$

ML solution

$$\beta_{01}^* = \frac{\#(X_i^k = 0, X_{i-1}^k = 1)}{\#(X_{i-1}^k = 1)}$$

$$= \frac{\sum_{i=2}^{n} \sum_{k=1}^{m} \mathbb{1}(X_i^k = 0, X_{i-1}^k = 1)}{\sum_{i=2}^{n} \sum_{k=1}^{m} \mathbb{1}(X_{i-1}^k = 1)}$$
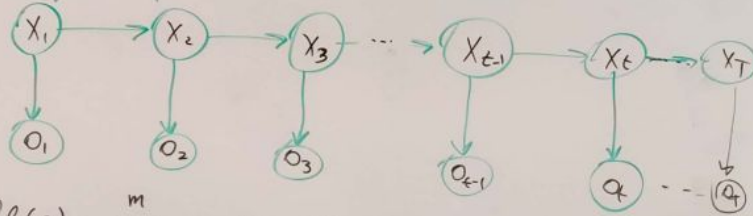
# shared parameters - table representation



27

# shared parameters - table representation



$$P(X_t = k \mid X_{t-1} = \ell) = f_\theta(k, \ell)$$

$$P(X_t \mid X_{t-1}) = f_\theta(X_t, X_{t-1}) =$$

$$P(O_t \mid X_t) = g_\gamma(O_t, X_t)$$

table representation:

$$P(X_t = j \mid X_{t-1} = \ell) = \lambda_{j, \ell}$$
$$P(O_t = j \mid X_t = \ell) = \gamma_{j, \ell}$$

$$\frac{\partial \ell\ell(\lambda, \gamma)}{\partial \lambda_{\ell, n}} = \sum_{k=1}^{m} \sum_{i=2}^{T} \log P(X_i^k \mid X_{i-1}^k)$$

Data  $\quad X_1^1, X_2^1, \dots, X_T^1, O_1^1, O_2^1, \dots, O_T^1$

$\quad X_1^2, X_2^2, \quad , X_T^2, O_1^2, O_2^2, \dots, O_T^2$

$\quad \vdots$

$\quad X_1^m, X_2^m, \dots X_T^m, O_1^m, \dots O_T^m$

$$X_t \in \{1, 2, \dots, q\}$$

$$= \sum_{j=1}^{q} \sum_{\ell=1}^{q} \sum_{\substack{i=2 \\ X_i^k = j \\ X_{i-1}^k = \ell}}^{T} \log P(X_i^k \mid X_{i-1}^k)$$

$$\sum_{j=1}^{q} \sum_{\ell=1}^{q} \sum_{X_i^k = j, X_{i-1}^k = \ell} \log P(j \mid \ell)$$

$$\sum_{j=1}^{q} \sum_{\ell=1}^{q} \sum_{\substack{X_i^k = j \\ X_{i-1}^k = \ell}} \log \lambda_{j\ell}$$

$$\lambda_{j\ell}^* = \frac{\#(X_i = j, X_{i-1} = \ell)}{\#(X_{i-1} = \ell)}$$

28