

$$P(\vec{X}) = P(X_1, X_2, \dots, X_d) \quad \vec{X} \in \mathbb{R}^d$$

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

Sampling

$P_\theta(X)$  known  $\Rightarrow$  Samples:  $X^1, X^2, \dots, X^m \in \mathbb{R}^d$

Learning

Data  $X^1, X^2, \dots, X^m$  known  $\Rightarrow$  find  $P_\theta(X)$

- find  $P_\theta(X)$ :
  - find  $\theta$ : parameter learning
  - find structure + parameters
    - BN: directed edges
    - MRF: undirected edges

parameter learning (find  $\theta$ )

find  $\theta$  that maximizes the probability of occurrence of  $\{X^1, X^2, \dots, X^m\} = D$

$$\theta^* = \operatorname{argmax}_\theta \Pr(X^1, X^2, \dots, X^m)$$

$X^1, X^2, \dots, X^m$  are i.i.d samples of  $P(X)$

$$\Pr(X^1, X^2, \dots, X^m) = \Pr(X^1) \Pr(X^2) \dots \Pr(X^m)$$

independent

$$= P_\theta(X^1) P_\theta(X^2) \dots P_\theta(X^m)$$

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

$$= \prod_{i=1}^m P_\theta(X^i) = \mathcal{L}(\theta; D)$$

parameter  $\theta$       data  $D$

$$\mathcal{L}(\theta; D) = \mathcal{L}(\theta) \quad \text{likelihood function}$$

$X^1, X^2, \dots, X^m$   
H T T H T H

$$\theta = \Pr(X=H)$$

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

$$\begin{aligned} \mathcal{L}(\theta) &= \prod_{i=1}^m \Pr(X=X^i) \\ &= \theta^{\#(X^i=H)} (1-\theta)^{\#(X^i=T)} \\ &= \theta^{n_H} (1-\theta)^{n_T} \end{aligned}$$

$$n_H + n_T = m$$

Assume  $n_H, n_T > 0$

$$\text{maximize } \mathcal{L}(\theta) = \max_{\theta} \theta^{n_H} (1-\theta)^{n_T}$$

$$\begin{aligned} &= \frac{d}{d\theta} \mathcal{L}(\theta) = n_H \theta^{n_H-1} (1-\theta)^{n_T} - n_T \theta^{n_H} (1-\theta)^{n_T-1} \\ &= \theta^{n_H-1} (1-\theta)^{n_T-1} [n_H (1-\theta) - n_T \theta] \end{aligned}$$

$\theta \in [0, 1] \Rightarrow$  check for  $\theta=1, 0$

$\theta=0$   
 $\theta=1$

$$n_H (1-\theta) - n_T \theta = 0 \Rightarrow n_H = \theta (n_H + n_T)$$

$$\Rightarrow \theta = \frac{n_H}{n_H + n_T} = \frac{n_H}{m}$$

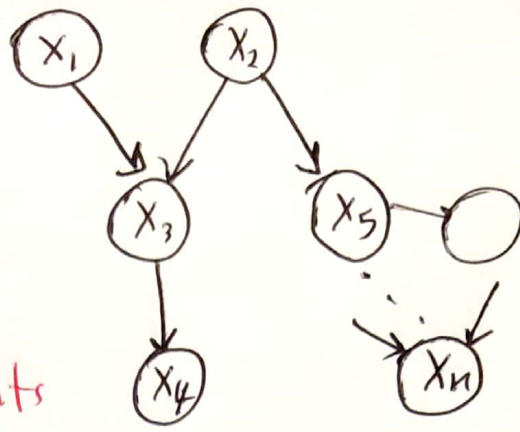
$$\theta=0, \theta=1, \theta = \frac{n_H}{m}$$

$$\mathcal{L}(\theta) = \theta^{n_H} (1-\theta)^{n_T} = \begin{matrix} \theta=0 \Rightarrow 0 \\ \theta=1 \Rightarrow 0 \end{matrix}$$

$$\theta = \frac{n_H}{m} \Rightarrow \left(\frac{n_H}{m}\right)^{n_H} \left(1 - \frac{n_H}{m}\right)^{n_T} > 0$$

if  $n_H, n_T > 0$

$$\Rightarrow \theta^* = \frac{n_H}{m}$$



$$P_{\theta}(X) = P_{\theta}(X_1, X_2, \dots, X_n)$$

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

parents of  $X_i$

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

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

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

1: Each CPD has its own parameters

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

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

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

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

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

local log-likelihood

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

gradient w.r.t.  $\theta_j$

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

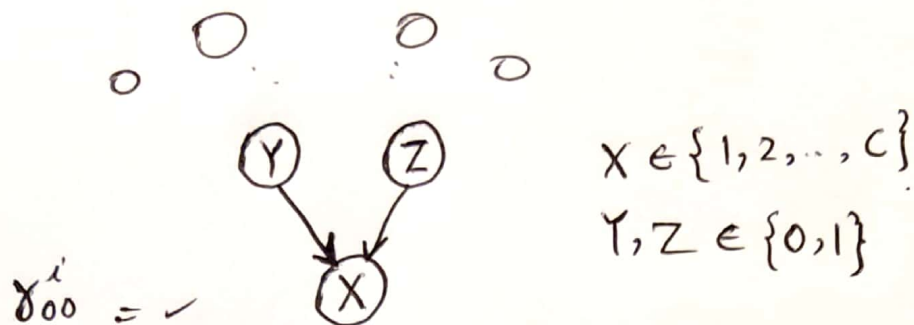
$$= \sum_{k=1}^m \frac{\frac{\partial}{\partial \theta_j} P_{\theta_j}(X_i^k | X_{P_i}^k)}{P_{\theta_j}(X_i^k | X_{P_i}^k)}$$

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



no shared parameters among CPD, +

table representation



$\delta_{00}^i = \checkmark$

$\delta_{01}^i = \Pr(X=i | Y=0, Z=1)$        $\sum_i \delta_{10}^i = 1$

$\delta_{10}^i = \Pr(X=i | Y=1, Z=0)$

$\delta_{11}^i = \Pr(X=i | Y=1, Z=1)$

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

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

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

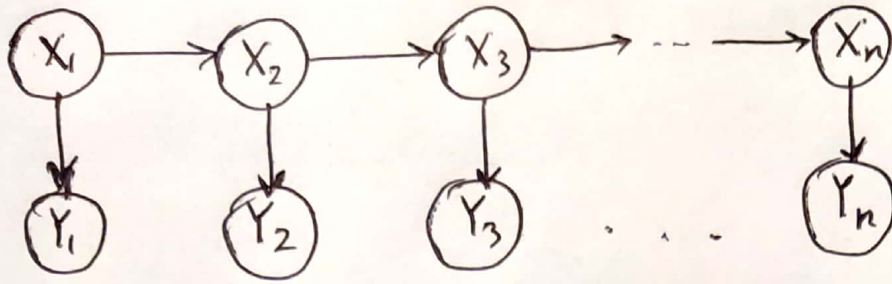
$\sum_x \gamma_{yz}^x = 1$

sufficient statistics

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

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

# Shared Parameters



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

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

$$\begin{aligned} \ell(\theta) = & \sum_{k=1}^m \log P_{\alpha}(X_i^k) + \sum_{k=1}^m \sum_{i=2}^n \log P_{\beta}(X_i^k | X_{i-1}^k) \\ & + \sum_{k=1}^m \sum_{i=1}^n \log P_{\gamma}(Y_i^k | X_i^k) \end{aligned}$$

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

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

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

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

ML solution

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

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