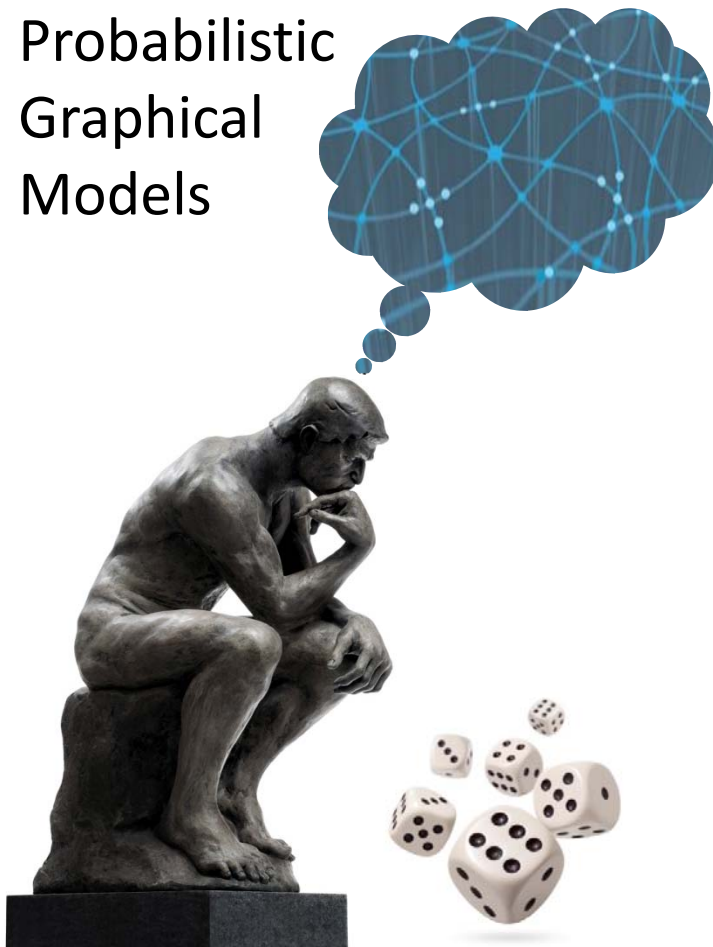Probabilistic Graphical Models

Learning

Parameter Estimation
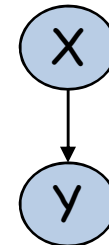
# Max-Likelihood for BNs

# MLE for Bayesian Networks

- Parameters: $\rightarrow \theta_{x^0}, \theta_{x^1}$
  $$\theta_{y^0|x^0}, \theta_{y^1|x^0}, \theta_{y^0|x^1}, \theta_{y^1|x^1}$$

- Data instances: <x[m],y[m]>

| X | |
|---|---|
| x⁰ | x¹ |
| 0.7 | 0.3 |

P(x)

X

Y

| X | Y | |
|---|---|---|
| | y⁰ | y¹ |
| x⁰ | 0.95 | 0.05 |
| x¹ | 0.2 | 0.8 |

P(Y|x)

Daphne Koller

# MLE for Bayesian Networks

$\theta_X$

- ## Parameters:

$$\{\theta_x : x \in Val(X)\}$$

$$\{\theta_{y|x} : x \in Val(X), y \in Val(Y)\}$$

$X$

$\theta_{Y|X}$

$Y$

Data $d$

$$L(\Theta : D) = \prod_{m=1}^{M} P(x[m], y[m] : \theta)$$

chain rule for BNs

$$= \prod_{m=1}^{M} P(x[m] : \theta) P(y[m] \mid x[m] : \theta)$$

$$= \left( \prod_{m=1}^{M} P(x[m] : \theta) \right) \left( \prod_{m=1}^{M} P(y[m] \mid x[m] : \theta) \right)$$

product of two local likelihood

$$= \left( \prod_{m=1}^{M} P(x[m] : \theta_X) \right) \left( \prod_{m=1}^{M} P(y[m] \mid x[m] : \theta_{Y|X}) \right)$$

Daphne Koller

# MLE for Bayesian Networks

- Likelihood for Bayesian network

$$L(\Theta : D) \quad = \prod_m P(x[m] : \Theta)$$

parents of $X_i$

chain rule

$$= \prod_m \prod_i P(x_i[m] \mid U_i[m] : \Theta_i)$$

$$= \prod_i \prod_m P(x_i[m] \mid U_i[m] : \Theta_i)$$

local likelihood $\Longrightarrow$ $$= \prod_i L_i(\Theta_i : D) \quad L_i(\theta_{x_i} : D)$$

$\Rightarrow$ if $\theta_{X_i \mid U_i}$ are disjoint, then MLE can be computed by maximizing each local likelihood separately

# MLE for Table CPDs

$$\prod_{m=1}^{M} P(x[m] \mid \boldsymbol{u}[m] : \theta) = \prod_{m=1}^{M} P(x[m] \mid \boldsymbol{u}[m] : \theta_{X|U})$$

$$= \prod_{x,\boldsymbol{u}} \left( \prod_{m:x[m]=x, \boldsymbol{u}[m]=\boldsymbol{u}} P(x[m] \mid \boldsymbol{u}[m] : \theta_{X|U}) \right)$$

$$P(x[m]=x \mid \boldsymbol{u}[m]=\boldsymbol{u} : \theta_{x|u}) = \theta_{x|u}$$

$$= \prod_{x,\boldsymbol{u}} \left( \prod_{m:x[m]=x, \boldsymbol{u}[m]=\boldsymbol{u}} \theta_{x|u} \right)$$

fraction of $X=x$ among cases where $\bar{u}=\bar{u}$

$$= \prod_{x,\boldsymbol{u}} \theta_{x|u}^{M[x,\boldsymbol{u}]}$$

$$P(x|u)$$

$$\theta_{x|u} = \frac{M[x,\boldsymbol{u}]}{\sum_{x'} M[x',\boldsymbol{u}]} = \frac{M[x,\boldsymbol{u}]}{M[\boldsymbol{u}]}$$

Daphne Koller

# Shared Parameters

$\theta_{S'|S}$

$S^{(0)} \rightarrow S^{(1)} \rightarrow S^{(2)} \rightarrow S^{(3)}$

$$L(\theta : S^{(0:T)}) = \prod_{t=1}^{T} P(S^{(t)} \mid S^{(t-1)} : \theta)$$

$$= \prod_{i,j} \prod_{t:S^{(t)}=s^i, S^{(t+1)}=s^j} P(S^{(t+1)} \mid S^{(t)} : \theta_{S'|S})$$

$s^i \rightarrow s^j$

$$= \prod_{i,j} \prod_{t:S^{(t)}=s^i, S^{(t+1)}=s^j} \theta_{s^i \rightarrow s^j}$$

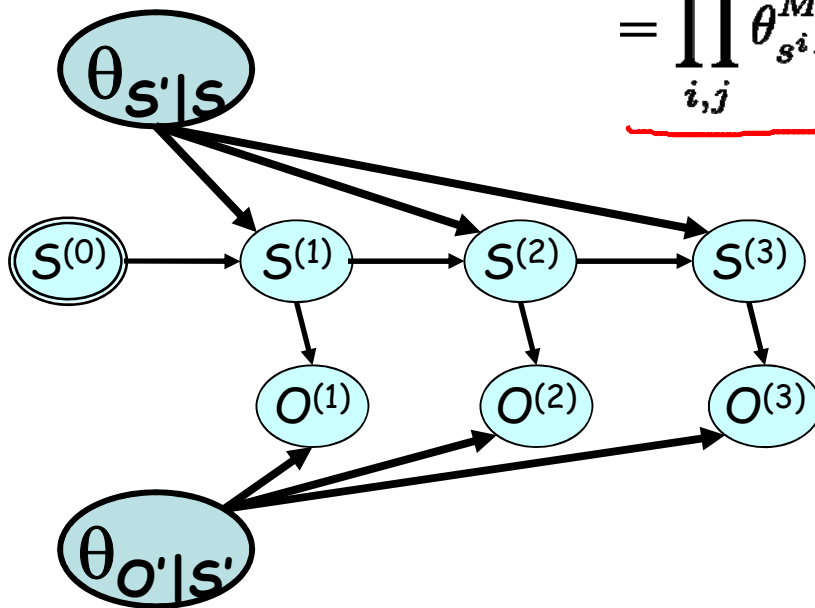$$= \prod_{i,j} \theta_{s^i \rightarrow s^j}^{M[s^i \rightarrow s^j]}$$

$$\hat{\theta}_{s^i \rightarrow s^j} = \frac{M[s^i \rightarrow s^j]}{M[s^i]}$$

$$M[s^i \rightarrow s^j] = |\{t \ : \ S^{(t)} = s^i, S^{(t+1)} = s^j\}|$$

Daphne Koller

# Shared Parameters

$$L(\Theta : S^{(0:T)}, O^{(0:T)}) = \prod_{t=1}^{T} P(S^{(t)} \mid S^{(t-1)} : \theta_{S'|S}) \prod_{t=1}^{T} P(O^{(t)} \mid S^{(t)} : \theta_{O'|S'})$$

$$= \prod_{i,j} \theta_{s^i \rightarrow s^j}^{M[s^i \rightarrow s^j]} \prod_{i,k} \theta_{o^k|s^i}^{M[o^k, s^i]}$$



$$M[s^i \rightarrow s^j] = |\{t \;:\; S^{(t)} = s^i, S^{(t+1)} = s^j\}|$$

$$M[o^k, s^i] = |\{t \;:\; S^{(t)} = s^i, O^{(t)} = o^k\}|$$

Daphne Koller

# Summary

- For BN with disjoint sets of parameters in CPDs, likelihood decomposes as product of local likelihood functions, one per variable

- For table CPDs, local likelihood further decomposes as product of likelihood for multinomials, one for each parent combination

- For networks with shared CPDs, sufficient statistics accumulate over all uses of CPD

Daphne Koller

# Fragmentation & Overfitting

$$\theta_{x|u} = \frac{M[x,u]}{\sum_{x'} M[x',u]} = \frac{M[x,u]}{M[u]}$$

- \# of "buckets" increases exponentially with |U|
- For large |U|, most "buckets" will have very few instances
  - ⇨ **very poor parameter estimates** ⇦
- **With limited data, we often get better generalization with simpler structures**
  even when wrong

Daphne Koller