

Readings:

K&F: 17.3, 17.4, 17.5.1, 8.1, 12.1

## Structure Learning (The Good), The Bad, The Ugly)

A little inference too...

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

October 8<sup>th</sup>, 2008

10-708 – Carlos Guestrin 2006-2008

1

## Decomposable score

- Log data likelihood

$$\log \hat{P}(\mathcal{D} \mid \theta, \mathcal{G}) = m \sum_i \hat{I}(X_i, \mathbf{Pa}_{X_i}) - m \sum_i \hat{H}(X_i)$$

- Decomposable score:

- Decomposes over families in BN (node and its parents)
- Will lead to significant computational efficiency!!!
- $\text{Score}(G : D) = \sum_i \text{FamScore}(X_i \mid \mathbf{Pa}_{X_i} : D)$

for MLE  $\text{FamScore}(X_i \mid \mathbf{Pa}_{X_i} : D) = m \hat{I}(X_i \mid \mathbf{Pa}_{X_i}) - m \hat{H}(X_i)$

10-708 – Carlos Guestrin 2006-2008

2

# Chow-Liu tree learning algorithm 1

- For each pair of variables  $X_i, X_j$

- Compute empirical distribution:

$$\hat{P}(x_i, x_j) \stackrel{\text{MLE}}{=} \frac{\text{Count}(x_i, x_j)}{m}$$

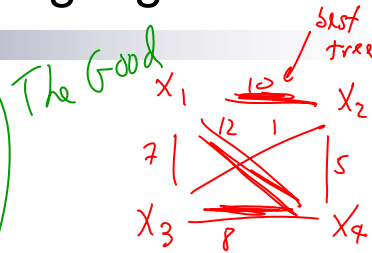
- Compute mutual information:

$$\hat{I}(X_i, X_j) = \sum_{x_i, x_j} \hat{P}(x_i, x_j) \log \frac{\hat{P}(x_i, x_j)}{\hat{P}(x_i) \hat{P}(x_j)}$$

- Define a graph

- Nodes  $X_1, \dots, X_n$   $w_{ij}$
  - Edge  $(i, j)$  gets weight  $\hat{I}(X_i, X_j)$

find maximum spanning tree



max  $\uparrow$  score(tree)  
 $= \sum_{i,j} I(X_i, X_j)$   
 $= \sum_{i,j} w_{ij}$

best tree BN

## Maximum likelihood score overfits!

$$\uparrow \log \hat{P}(\mathcal{D} | \theta, \mathcal{G}) = m \sum_i \hat{I}(X_i, \text{Pa}_{X_i, \mathcal{G}}) - m \sum_i \hat{H}(X_i)$$

- Information never hurts:

$$\uparrow \hat{I}(X_i, \text{Pa}_{X_i}) = \hat{H}(X_i) - \hat{H}(X_i | \text{Pa}_{X_i})$$

$$H(A|B) \leq H(A|C) \quad C \subseteq B$$

the more parents  
the higher  
 $\hat{I}(X_i, \text{Pa}_{X_i})$

- Adding a parent always increases score!!!

MLE  $\Rightarrow$  complete Graph

# Bayesian score

## ■ Prior distributions:

- Over structures ✓
- Over parameters of a structure ✓

## ■ Posterior over structures given data:

note LD  $P(D|G, \theta_G)$

prior over graphs, eg  $P(G) \propto e^{-c(\text{number of edges})}$

prior over CPT parameters

$$P(G|D) = \frac{P(D|G)P(G)}{P(D)}$$

$$= \frac{\int_{\theta_G} P(D|G, \theta_G) P(\theta_G|G) P(G) d\theta_G}{P(D)}$$

prior over graphs ✓

posterior

$$\log P(G|D) = \log P(G) + \log \int_{\theta_G} P(D|G, \theta_G) P(\theta_G|G) d\theta_G$$

+ constant  $\leftarrow \log P(D)$

# Bayesian learning for multinomial

## ■ What if you have a k sided coin???

## ■ Likelihood function if multinomial:

- $P(D|\theta_1, \dots, \theta_k) = \theta_1^{m_1} \theta_2^{m_2} \dots \theta_k^{m_k}$
  - $\sum_i \theta_i = 1, \theta_i \geq 0$
- $m_i \leftarrow \# \text{ observations of class, value or side } i$
- equivalent sample size  $\sum_i m_i$

## ■ Conjugate prior for multinomial is Dirichlet:

- $\theta \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k) \sim \prod_i \theta_i^{\alpha_i - 1}$
- $\alpha_i \geq 0$

## ■ Observe $m$ data points, $m_i$ from assignment $i$ , posterior:

$$P(\theta_1, \dots, \theta_k | m_1, \dots, m_k) \propto P(m_1, \dots, m_k | \theta_1, \dots, \theta_k) P(\theta)$$

$$\equiv \text{Dirichlet}(d_1 + m_1, d_2 + m_2, \dots, d_k + m_k)$$

## ■ Prediction:

$$E[\theta_i] = \frac{m_i + \alpha_i}{\sum_j (m_j + \alpha_j)}$$

# Global parameter independence, d-separation and local prediction

Independencies in **meta BN**:  $\log P(\theta | D) = \sum_i \log P(\theta_{X_i} | \text{Pa}_{X_i}; D)$

add prior vars to the BN

$P(\theta) = P(\theta_F) P(\theta_A) P(\theta_{SIFA}) P(\theta_{NIS}) P(\theta_{HIS})$

**Proposition:** For fully observable data  $D$ , if prior satisfies global parameter independence, then

$P(\theta | D) = \prod_i P(\theta_{X_i} | \text{Pa}_{X_i} | D)$

params indep. given data

## Priors for BN CPTs

(more when we talk about structure learning)

- Consider each CPT:  $P(X | U=u)$
- Conjugate prior:
  - Dirichlet( $\alpha_{X=1|U=u}, \dots, \alpha_{X=k|U=u}$ )  $\equiv$  Dirichlet( $\text{Count}'(X=1, U=u), \dots, \text{Count}'(X=k, U=u)$ )
- More intuitive:
  - "prior data set"  $D'$  with  $m'$  equivalent sample size
  - "prior counts":  $\text{Count}'(X=x, U=u)$  or  $m' \cdot P'(X=x, U=u)$
  - prediction:
 
$$E[\theta_{X=x|U=u}] = \frac{\text{Count}(X=x, U=u) + \text{Count}'(X=x, U=u)}{\text{Count}(U=u) + \text{Count}'(U=u)}$$

[illegible]

- Bayesian parameter learning: ✓
  - motivation for Bayesian approach
  - Bayesian prediction
  - conjugate priors, equivalent sample size
  - Bayesian learning  $\Rightarrow$  smoothing
- Bayesian learning for BN parameters ✓
  - Global parameter independence
  - Decomposition of prediction according to CPTs
  - Decomposition within a CPT

- Bayesian parameter learning:
  - motivation for Bayesian approach
  - Bayesian prediction
  - conjugate priors, equivalent sample size
  - Bayesian learning  $\Rightarrow$  smoothing
- Bayesian learning for BN parameters
  - Global parameter independence
  - Decomposition of prediction according to CPTs
  - Decomposition within a CPT

# Bayesian score and model complexity

*prefer simpler graphs over more complex*

$\log P(D | G) = \log \int_{\theta_G} P(D | G, \theta_G) P(\theta_G | G) d\theta_G$

True model:

```

graph TD
    X((X)) --> Y((Y))
  
```

Structure 1: X and Y independent

$$\log P(D | G) = \log \int_{\theta_x} \int_{\theta_y} P(D_x, D_y | G_1, \theta_x, \theta_y) P(\theta_x) P(\theta_y) d\theta_x d\theta_y$$

$$= \log \int_{\theta_x} P(D_x | G_1, \theta_x) P(\theta_x) d\theta_x + \log \int_{\theta_y} P(D_y | G_1, \theta_y) P(\theta_y) d\theta_y$$

Score doesn't depend on alpha

Structure 2:  $X \rightarrow Y$

$$\log P(D | G) = \log \int_{\theta_x} \int_{\theta_y} P(D_x, D_y | G_2, \theta_x, \theta_y) P(\theta_x) P(\theta_y | X) d\theta_x d\theta_y$$

$$= \log \int_{\theta_x} P(D_x | G_2, \theta_x) P(\theta_x) d\theta_x + \int_{\theta_y | X=t} P(D_y | G_2, \theta_y | X=t) P(\theta_y | X=t) d\theta_y + \int_{\theta_y | X=f} P(D_y | G_2, \theta_y | X=f) P(\theta_y | X=f) d\theta_y$$

$P(Y=t|X=t) = 0.5 + \alpha$   
 $P(Y=t|X=f) = 0.5 - \alpha$

$P(Y=t) = P(Y=t|X=t)(0.5 + \alpha) + P(Y=t|X=f)(0.5 - \alpha)$

$P(Y=f) = P(Y=f|X=t)(0.5 + \alpha) + P(Y=f|X=f)(0.5 - \alpha)$

$= 0.5$  (doesn't depend on alpha)

□ Data points split between  $P(Y=t|X=t)$  and  $P(Y=t|X=f)$

□ For fixed  $M$ , only worth it for large  $\alpha$

- Because posterior over parameter will be more diffuse with less data

*more data posterior*

## Bayesian, a decomposable score

$\log P(D | G) = \log \int_{\theta_G} P(D | G, \theta_G) P(\theta_G | G) d\theta_G$

As with last lecture, assume:

- Parameter independence  $P(\theta) = \prod P(\theta_{x_i} | \text{pa}_{x_i})$

Also, prior satisfies **parameter modularity**:

- If  $X_i$  has same parents in  $G$  and  $G'$ , then parameters have same prior

Finally, structure prior  $P(G)$  satisfies **structure modularity**

- Product of terms over families  $x_i, \text{pa}_{x_i}$
- E.g.,  $P(G) \propto c^{|G|}$  (number of edges  $|G|$ )

Bayesian score decomposes along families!

$$\text{Score}_{\text{Bayes}}(G; D) = \sum_i \text{Score}_{\text{Fam}}(x_i; \text{pa}_{x_i}; D)$$

$$P(D|G) = \int_{\theta} P(D|G, \theta) P(\theta|G) d\theta$$

## BIC approximation of Bayesian score

- Bayesian has difficult integrals
- For Dirichlet prior, can use simple Bayes information criterion (BIC) approximation
  - In the limit, we can forget prior!
  - Theorem:** for Dirichlet prior, and a BN with Dim(G) independent parameters, as  $m \rightarrow \infty$ :
    - changes at rate  $\frac{\log m}{2}$

$$\log P(D | G) = \log P(D | G, \theta_G) - \frac{\log m}{2} \text{Dim}(G) + O(1)$$

as  $m \rightarrow \infty$  data likelihood regularizer

$$m \left[ \sum_i (\hat{I}(x_i, Pa_{x_i}) - \hat{H}(x_i)) \right]$$

Dim(G)  
 (X) (Y)  
 2 params

X → Y  
 3 params

## BIC approximation, a decomposable score

all vars have  
 k assignments  
 $\text{Dim}(CPT) = (K-1)K | Pa_{x_i}$

BIC:  $\text{Score}_{\text{BIC}}(G : D) = \log P(D | G, \theta_G) - \frac{\log m}{2} \text{Dim}(G)$

$$\text{Dim}(G) = \sum_i \text{Dim}(CPT(x_i | Pa_{x_i}))$$

- Using information theoretic formulation:

$$\text{Score}_{\text{BIC}}(G : D) = m \sum_i \hat{I}(X_i, Pa_{X_i, G}) - m \sum_i \hat{H}(X_i) - \frac{\log m}{2} \sum_i \text{Dim}(P(X_i | Pa_{X_i, G}))$$

$$= \sum_i (m \hat{I}(x_i, Pa_{x_i, G}) - m \hat{H}(x_i) - \frac{\log m}{2} \text{Dim}(P(x_i | Pa_{x_i, G})))$$

$$\text{FamScore}_{\text{BIC}}(x_i | Pa_{x_i}, D)$$

# Consistency of BIC and Bayesian scores

Consistency is limiting behavior, says nothing about finite sample size!!!

- A scoring function is **consistent** if, for true model  $G^*$ , as  $m \rightarrow \infty$ , with probability 1
  - $G^*$  maximizes the score
  - All structures **not I-equivalent** to  $G^*$  have strictly lower score
- **Theorem:** BIC score is consistent
- **Corollary:** the Bayesian score is consistent
- What about maximum likelihood score? **NO**

MLE: not consistent

$$\text{Score}_{\text{MLE}}(\text{Complete Graph}) = \text{Score}_{\text{MLE}}(G^*) \quad \left| \quad \text{BIC:} \right. \\ \text{Penalty}(\text{Complete Graph}) > \text{Penalty}(G^*)$$

# Priors for general graphs

- For finite datasets, prior is important!
- Prior over structure satisfying prior modularity
- What about prior over parameters, how do we represent it?
  - K2 prior: fix an  $\alpha$ ,  $P(\theta_{X|PaX}) = \text{Dirichlet}(\alpha, \dots, \alpha)$  *for each*
  - K2 is "inconsistent"

all  
vars  
have  $K$  values

$P(\alpha_{X,i})$	"equivalent sample size"
0	$K\alpha$
1	for each assign. prob. $P(\theta_{X PaX} = u) \text{Dir}(\alpha, \alpha, \dots)$ $K^2\alpha$
2	$K^3\alpha$



## BDe prior

- Remember that Dirichlet parameters analogous to "fictitious samples"
- Pick a fictitious sample size  $m'$
- For each possible family, define a prior distribution  $P'(X_i, \mathbf{Pa}_{X_i})$ 
  - Represent with a BN  $(x_i) \quad \text{---}$
  - Usually independent (product of marginals)  $P'(x_i, \mathbf{Pa}_{x_i}) = P'(x_i) \prod_{j \in \mathbf{Pa}_{x_i}} P'(x_j)$
- BDe prior:** *most common uniform*  
 $P(\theta_{x_i} | \mathbf{Pa}_{x_i} = u) \sim \text{Dirichlet}(m' P'(x_i=1, \mathbf{Pa}_{x_i}=u), \dots, m' P'(x_i=k, \mathbf{Pa}_{x_i}=u))$
- Has "consistency property":

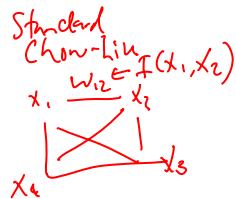
$\text{Score}_{\text{BDe}}(G; D)$  is consistent

## Score equivalence

- If G and G' are I-equivalent then they have same score
- Theorem 1:** Maximum likelihood score and BIC score satisfy score equivalence
- Theorem 2:**
  - If P(G) assigns same prior to I-equivalent structures (e.g., edge counting)
  - and parameter prior is dirichlet
  - then **Bayesian score satisfies score equivalence** if and only if prior over parameters represented as a BDe prior!!!!

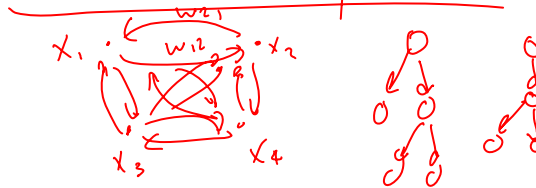
# Chow-Liu for Bayesian score

- Edge weight  $w_{X_j \rightarrow X_i}$  is advantage of adding  $X_j$  as parent for  $X_i$



Bayesian Chow-Liu

BIC  
 $w_{12} = w_{21}$   
 $D_{im}(P(X_1, X_2)) =$   
 $D_{im}(P(X_2 | X_1))$



- Now have a directed graph, need directed spanning forest
  - Note that adding an edge can hurt Bayesian score – choose forest not tree
  - Maximum spanning forest algorithm works