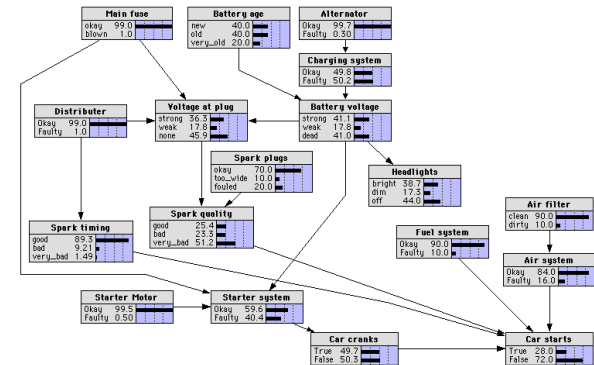


Cmpu466 / 551



# Learning Belief Net Structures

Readings: HTF ~Ch17

+ Bayesian Networks without the Tears (Charniak)

R Greiner

University of Alberta

Some material taken from C Guesterin (CMU), K Murphy (UBC)



# Outline

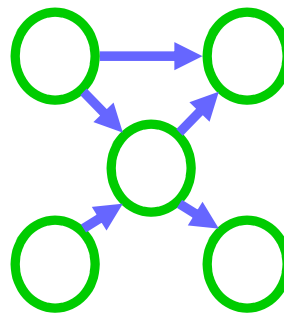
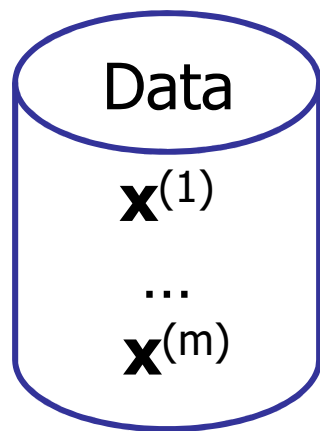
---

- Motivation
- What is a Belief Net?
- Learning a Belief Net
  - Goal?
  - Learning Parameters – Complete Data
  - Learning Parameters – Incomplete Data
  - Learning Structure
    - Learning best TREE Structure
- Dynamical Belief Nets ... HMMs

# Learning Belief Nets

Jump

		Structure	
		Known	Unknown
Data	Complete	✓ <b>Easy</b>	<b>NP-hard</b>
	Missing	✓ <b>Hard ... EM</b>	<b>Very hard!!</b>



structure

+ CPTs :  
 $\{ P(X_i | \mathbf{Pa}_{X_i}) \}$

parameters

# Learning the structure of a BN

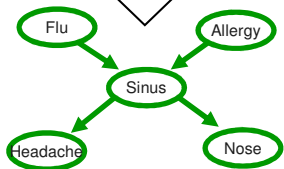
Data

$[x_1^{(1)}, \dots, x_n^{(1)}]$

...

$[x_1^{(m)}, \dots, x_n^{(m)}]$

Learn structure and parameters



## ■ Constraint-based approach

- BN encodes conditional independencies
- Test conditional independencies in data
- Find an I-map (?P-map?)
  - Only include link if dependency (given other links)

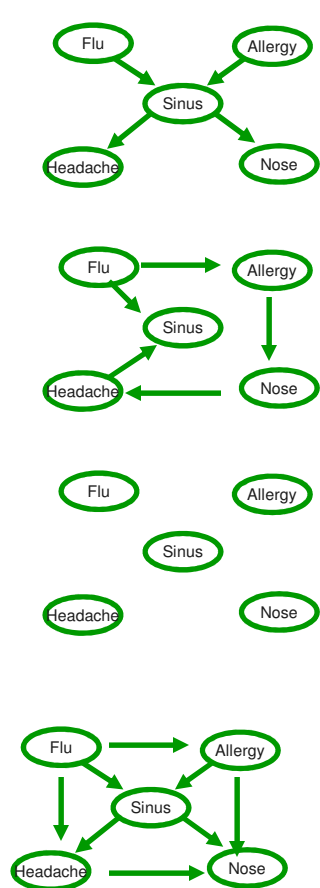
## ■ Score-based approach

- Finding structure + parameters = *density estimation*
- Evaluate model as we evaluated parameters
  - Maximum likelihood
  - Bayesian
  - etc.

Take Cmpu659...

# Score-based Approach

Possible DAG structures  
(gazillions)



**Data**

$[x_1^{(1)}, \dots, x_n^{(1)}]$   
...  
 $[x_1^{(m)}, \dots, x_n^{(m)}]$

Score of each Structure

-15,000

-10,000

-20,000

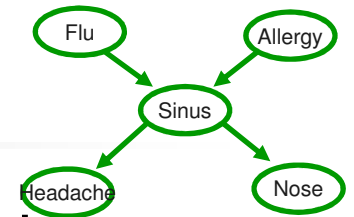
-10,500

Learn Parameters  
+  
Evaluate ...

# Just use MLE parameters

- $\max_{\mathcal{G}, \theta_{\mathcal{G}}} L( \langle \mathcal{G}, \theta_{\mathcal{G}} \rangle : \mathcal{D} ) =$   
 $\max_{\mathcal{G}} \max_{\theta_{\mathcal{G}}} L( \langle \mathcal{G}, \theta_{\mathcal{G}} \rangle : \mathcal{D} ) =$   
 $\max_{\mathcal{G}} L( \langle \mathcal{G}, \theta^*(\mathcal{G}) \rangle : \mathcal{D} )$ 
  - $\theta^*(\mathcal{G}) = \max_{\theta_{\mathcal{G}}} L( \langle \mathcal{G}, \theta_{\mathcal{G}} \rangle : \mathcal{D} )$
- So...  
seek the structure  $\mathcal{G}$  that achieves  
highest likelihood,  
given its MLE parameters  $\theta^*(\mathcal{G})$
- $\text{Score}(\mathcal{G}, \mathcal{D}) = \log L( \langle \mathcal{G}, \theta^*(\mathcal{G}) \rangle : \mathcal{D} )$

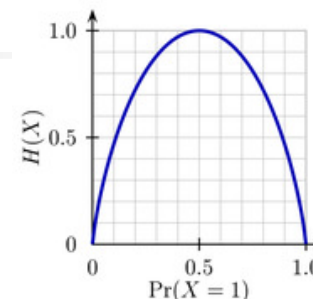
# Information-theoretic interpretation of maximum likelihood



- Given structure  $\mathcal{G}$ , parameters  $\theta_{\mathcal{G}}$ , log likelihood of data  $\mathcal{D}$ :

$$\begin{aligned}
 \log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G}) &= \sum_{j=1}^m \sum_{i=1}^n \log P \left( X_i = x_i^{(j)} \mid \mathbf{Pa}_{X_i} = \mathbf{x}^{(j)} \right) \\
 &= \sum_{i=1}^n \sum_{j=1}^m \log P \left( X_i = x_i^{(j)} \mid \mathbf{Pa}_{X_i} = \mathbf{x}^{(j)} \right) \\
 &= \sum_{i=1}^n \sum_{x_i, \mathbf{u}} \#(X_i = x_i, \mathbf{Pa}_{X_i} = \mathbf{u}) \log P(X_i = x_i \mid \mathbf{Pa}_{X_i} = \mathbf{u}) \\
 &= m \sum_{i=1}^n \sum_{x_i, \mathbf{u}} \frac{\#(X_i = x_i, \mathbf{Pa}_{X_i} = \mathbf{u})}{m} \log P(X_i = x_i \mid \mathbf{Pa}_{X_i} = \mathbf{u}) \\
 &= m \sum_{i=1}^n \sum_{x_i, \mathbf{u}} \hat{P}(X_i = x_i, \mathbf{Pa}_{X_i} = \mathbf{u}) \log P(X_i = x_i \mid \mathbf{Pa}_{X_i} = \mathbf{u})
 \end{aligned}$$

# Entropy & Conditional Entropy



## ■ Entropy of Distribution

- $H(X) = - \sum_i P(x_i) \log P(x_i)$
- “How ‘surprising’ variable is”
- Entropy = 0 when know everything... eg  $P(+x)=1.0$

## ■ Conditional Entropy $H(X | \mathbf{U})$

- $H(X|\mathbf{U}) = - \sum_{\mathbf{u}} P(\mathbf{u}) \sum_i P(x_i|\mathbf{u}) \log P(x_i|\mathbf{u})$
- How much uncertainty is left in  $X$ , after observing  $\mathbf{U}$

$$H(X_i | \mathbf{Pa}_{X_i}) = - \sum_{x_i, \mathbf{u}} \hat{P}(X_i = x_i, \mathbf{Pa}_{X_i} = \mathbf{u}) \log P\left(X_i = x_i^{(j)} \mid \mathbf{Pa}_{X_i} = \mathbf{u}\right)$$





## Information-theoretic interpretation of maximum likelihood ... 2

- Given structure  $\mathcal{G}$ , parameters  $\theta_{\mathcal{G}}$ , log likelihood of data  $\mathcal{D}$  is...

$$\begin{aligned} \uparrow \log \hat{P}(\mathcal{D} \mid \theta, \mathcal{G}) &= m \sum_i \sum_{x_i, \mathbf{u}} \hat{P}(x_i, \mathbf{Pa}_{x_i, \mathcal{G}} = \mathbf{u}) \log \hat{P}(x_i \mid \mathbf{Pa}_{x_i, \mathcal{G}} = \mathbf{u}) \\ &= m \sum_i -\hat{H}(X_i \mid \mathbf{Pa}_{x_i, \mathcal{G}}) \\ &= -m \sum_i \boxed{\hat{H}(X_i \mid \mathbf{Pa}_{x_i, \mathcal{G}})} \quad \downarrow \end{aligned}$$

So  $\log P(\mathcal{D} \mid \theta, \mathcal{G})$  is LARGEST

when each  $H(X_i \mid \mathbf{Pa}_{x_i, \mathcal{G}})$  is SMALL...

...ie, when parents of  $X_i$  are very INFORMATIVE about  $X_i$  !



# Easier Form

---

- Note

$$H_p(X | \mathbf{U}) = \sum_{\mathbf{x}, \mathbf{u}} P(\mathbf{x}, \mathbf{u}) \log P(\mathbf{x} | \mathbf{u})$$

is not symmetric in  $X, \mathbf{U}$

- Better to use symmetric

$$I_P(X, \mathbf{U}) = \sum_{\mathbf{x}, \mathbf{u}} P(\mathbf{x}, \mathbf{u}) \log \frac{P(\mathbf{x}, \mathbf{u})}{P(\mathbf{x}) P(\mathbf{u})}$$



# (Conditional) Mutual Information

- **Mutual information:**  $I_P(X, U) = \sum_{x,u} P(x, u) \log \frac{P(x, u)}{P(x) P(u)}$

- **Mutual information and independence:**

- $X$  and  $U$  independent if and only if  $I(X, U) = 0$
- $X \perp U \iff P(x, u) = P(x) P(u) \iff \log[ P(x, u) / P(x) P(u) ] = 0$

- **Conditional mutual information:**

$$I_P(X, Y | Z) = E_Z[ I(X, Y | Z = z) ] = \sum_z \sum_{x,y} P(x, y | z) \log \frac{P(x, y | z)}{P(x | z) P(y | z)}$$

- $X \perp Y | Z \iff P(X, Y | Z) = P(X | Z) P(Y | Z) \iff I(X, Y | Z) = 0$

- Using the data  $D$

- Empirical distribution:  $\hat{P}(x, y) = \frac{\text{Count}_D(X=x, Y=y)}{|D|}$
- Mutual information:  $I_{\hat{P}}(X, Y) = \sum_{x,y} \hat{P}(x, y) \log \frac{\hat{P}(x, y)}{\hat{P}(x) \hat{P}(y)}$



# Mutual Information

- **Mutual information:**  $I_P(X, U) = \sum_{x,u} P(x, u) \log \frac{P(x,u)}{P(x) P(u)}$

- Mutual information and independence:

- $X$  and  $U$  independent if and only if  $I(X, U) = 0$

- $X \perp U \Leftrightarrow \forall x, u \ P(x, u) = P(x) P(u) \Leftrightarrow \forall x, u \ \log \left[ \frac{P(x,u)}{P(x)P(u)} \right] = 0$

- $I_P(X, U) = \sum_{x,u} P(x, u) \log \frac{P(x,u)}{P(x) P(u)} = \sum_{x,u} P(x, u) \log \frac{P(x | u)}{P(x)}$

$$= \sum_{x,u} P(x, u) \log P(x | u) - \sum_{x,u} P(x, u) \log P(x)$$

$$= \sum_{x,u} P(x, u) \log P(x | u) - \sum_x P(x) \log P(x) \quad \cancel{\sum_u P(u|x)}$$

$$= -H(X | U) + H(X)$$



# Score for Belief Network

- $\mathcal{J}(X, \mathbf{U}) = H(X) - H(X \mid \mathbf{U})$

$$\Rightarrow H(X \mid \mathbf{Pa}_{X, \mathcal{G}}) = H(X) - \mathcal{J}(X, \mathbf{Pa}_{X, \mathcal{G}})$$

Doesn't involve the structure,  $\mathcal{G}$  !

- Log data likelihood

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

- So use score:  $\sum_i \mathcal{J}(X_i, \mathbf{Pa}_{X_i, \mathcal{G}})$



# Best Tree Structure

$\log P(D | \theta, G)$  is monotonic with  $\sum_i I(x_i, \text{Pa}_{X_i, G})$

- Identify tree with set  $\mathcal{F} = \{ \text{Pa}(X) \}$ 
  - each  $\text{Pa}(X)$  is either  $\{\}$ , or another variable
- Optimal tree, given data, is
$$\text{argmax}_{\mathcal{F}} \sum_i I(X_i, \text{Pa}(X_i))$$
- So ... want parents  $\mathcal{F}$  s.t.
  - tree structure
  - maximizes  $\sum_i I(X_i, \text{Pa}(X_i))$

# Chow-Liu Tree Learning Alg

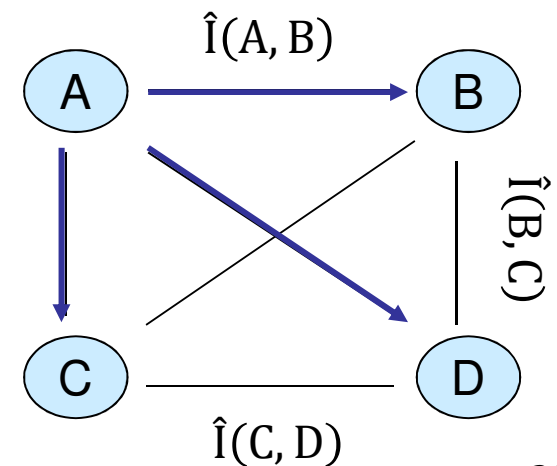
- For each pair of variables  $X_i, X_j$ 
  - Compute empirical distribution:

$$\hat{P}(x_i, x_j) = \frac{\text{count}(x_i, x_j)}{m}$$

- 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  $X_1, \dots, X_n$
  - Edge  $(i, j)$  gets weight  $\hat{I}(X_i, X_j)$
- Find Maximal Spanning Tree
- Pick a node for root, dangle...



# Chow-Liu Tree Learning Alg ... 2

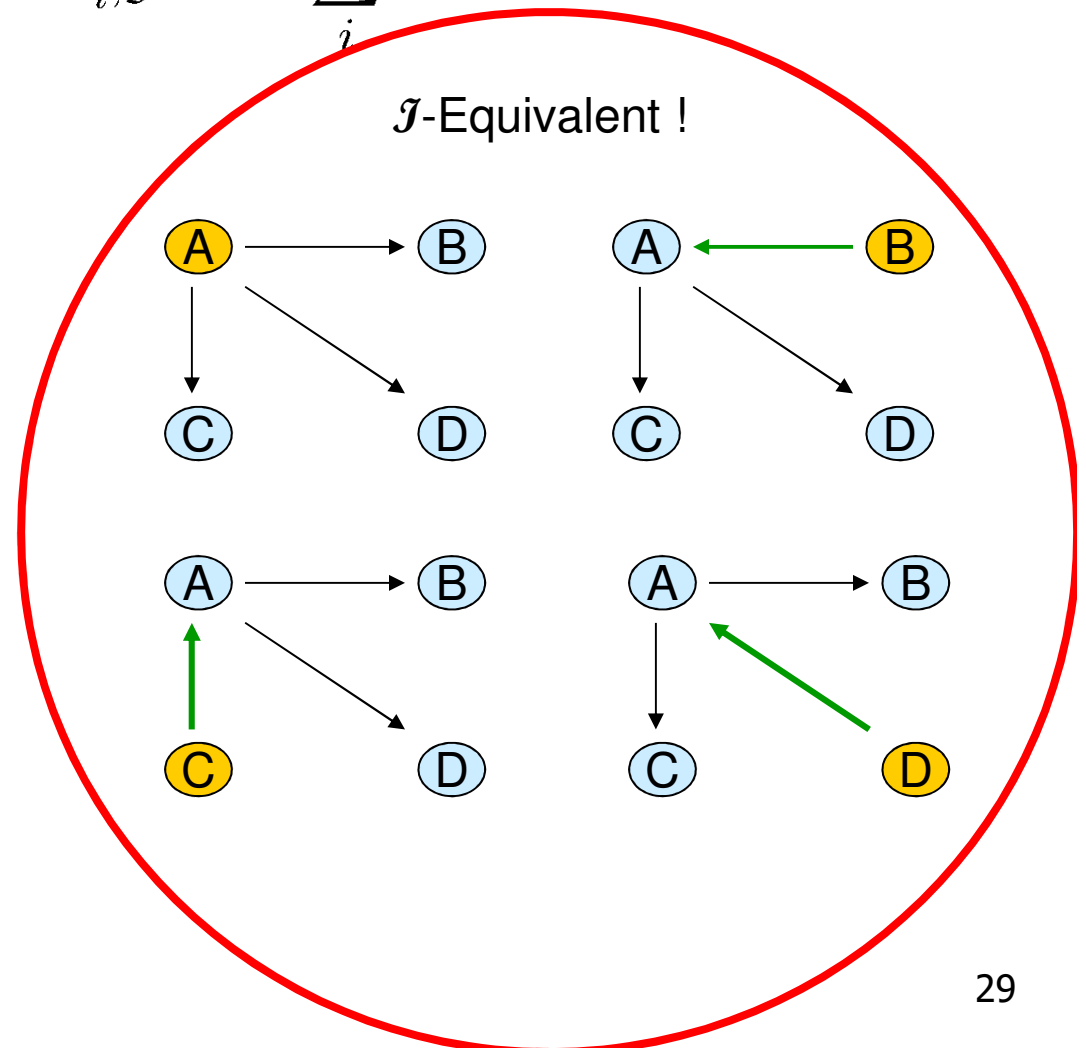
$$\log \hat{P}(\mathcal{D} \mid \theta, \mathcal{G}) = m \sum_i \hat{I}(x_i, \text{Pa}_{x_i, \mathcal{G}}) - m \sum_i \hat{H}(X_i)$$

## ■ Optimal tree BN

- ...
- Compute maximum weight spanning tree
- Directions in BN:
  - pick any node as root, ...doesn't matter which!
  - breadth-first-search defines directions

## ■ Score Equivalence:

If  $\mathcal{G}$  and  $\mathcal{G}'$  are  $\mathcal{I}$ -equiv, then scores are same







# Chow-Liu (CL) Results

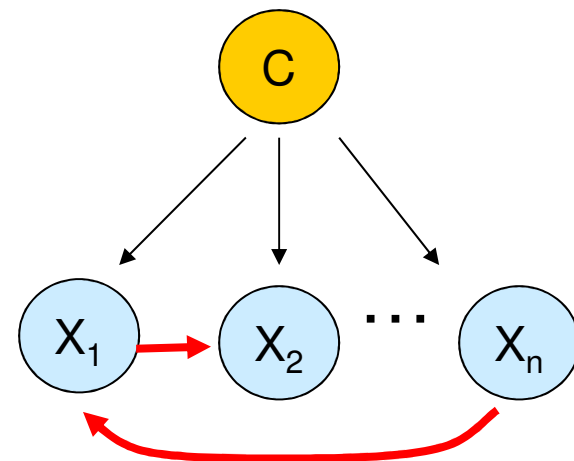
---

- If distribution  $P$  is tree-structured, CL finds CORRECT one
- If distribution  $P$  is NOT tree-structured, CL finds tree structured  $Q$  that has min'l KL-divergence:  $\operatorname{argmin}_Q \text{KL}(P; Q)$
- Even though  $2^{\theta(n \log n)}$  trees, CL finds BEST one in poly time  $O(n^2 [m + \log n])$ 
  - number of variables
  - number of instances

# Using Chow-Liu to Improve NB

- Naïve Bayes model

- $X_i \perp X_j \mid C$
- Ignores correlation between features
- What if  $X_1 = X_2$  ? **Double count...**



- Avoid by conditioning features on one another

- Tree Augmented Naïve bayes (TAN)

[Friedman et al. '97]

- Same as Chow-Liu, but score edges with:

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

All but ONE feature have 2 parents: C,  $X_i$



# Can we extend Chow-Liu ?

---

- (Approximately learning)  
models with tree-width up to  $k$ 
  - [Narasimhan & Bilmes '04]
  - But,  $O(n^{k+1})$ ...
    - and more subtleties



# More Elaborate Models

---

- Space is sooo large
- Heuristic search
  - Pick a decomposable score
  - Pick a reasonable start, check neighbourhood
    - Add arc
    - Delete arc
    - Reverse arc
  - All easy checks, compare 1 or 2 CPDs, (cache  $\Delta$ 's)
  - Move to better score
    - If no better: terminate

# Likelihood Overfits To Data



- Consider

- $\log P_{G1}(\mathbf{D}) = I(X;Y) - H(X) - H(Y)$
- $\log P_{G2}(\mathbf{D}) = \quad \quad \quad - H(X) - H(Y)$

$\Rightarrow$  difference is mutual information

- Note:

$$\log P_{G1}(\mathbf{D}) - \log P_{G2}(\mathbf{D}) = I(X;Y) \geq 0$$

$\Rightarrow$  G1 always preferred than G2

# Maximum likelihood score overfits!

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

- Adding a parent never decreases score!!!

- *Facts:*  $H(X \mid \text{Pa}_{X, \mathcal{G}}) = H(X) - I(X, \text{Pa}_{X, \mathcal{G}})$

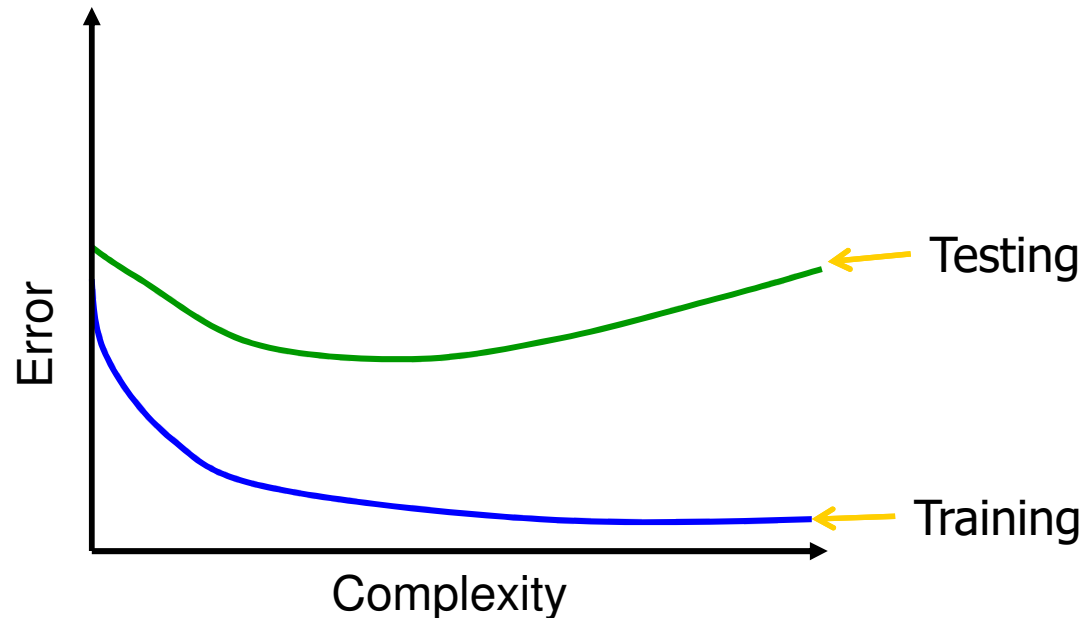
$$H(X \mid A) \geq H(X \mid A \cup Y)$$

- $I(X_i, \text{Pa}_{X_i, \mathcal{G}} \cup Y) = H(X_i) - H(X_i \mid \text{Pa}_{X_i, \mathcal{G}} \cup Y)$   
 $\geq H(X_i) - H(X_i \mid \text{Pa}_{X_i, \mathcal{G}})$   
 $= I(X_i, \text{Pa}_{X_i, \mathcal{G}})$

- So score increases as we add edges!
  - Best is COMPLETE Graph
  - Maximum likelihood score overfits !

# Likelihood Overfits To Data

- If additional arcs always favoured  
Then prefer “Complete Model”,  $K_N$
- With fixed data:
  - more complex models will overfit





# Bayesian Score

- Prior distributions:

- Over structures
- Over parameters of a structure

Goal: Prefer simpler structures... regularization ...

- Posterior over structures given data:

- $P(\mathcal{G}|\mathcal{D}) \propto P(\mathcal{D}|\mathcal{G}) \times P(\mathcal{G})$

↑  
Posterior

↑  
Likelihood

↑  
Prior over Graphs

↓  
Prior over Parameters

- $P(\mathcal{D}|\mathcal{G}) = \int_{\Theta} P(\mathcal{D} | \mathcal{G}, \Theta) P(\Theta|\mathcal{G}) d\Theta$

$$\log P(\mathcal{G} | D) \approx \log P(\mathcal{G}) + \log \int_{\theta_{\mathcal{G}}} P(D | \mathcal{G}, \theta_{\mathcal{G}}) P(\theta_{\mathcal{G}}|\mathcal{G}) d\theta_{\mathcal{G}}$$



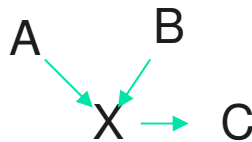
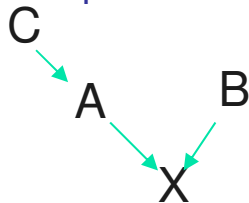
# Towards a decomposable Bayesian score

$$\log P(\mathcal{G} \mid D) \approx \log P(\mathcal{G}) + \log \int_{\theta_{\mathcal{G}}} P(D \mid \mathcal{G}, \theta_{\mathcal{G}}) P(\theta_{\mathcal{G}} \mid \mathcal{G}) d\theta_{\mathcal{G}}$$

- **Local and global parameter independence**  $\theta_{Y|+X} \perp \theta_X$

- Prior satisfies **parameter modularity**:

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



$\Theta(X; A, B)$  same in both structures

- Structure prior  $P(\mathcal{G})$  satisfies **structure modularity**

- Product of terms over families
  - Eg,  $P(\mathcal{G}) / c^{|\mathcal{G}|}$   $|\mathcal{G}| = \# \text{edges}; \quad c < 1$

- ... then: Bayesian score decomposes along families!

- $\log P(\mathcal{G} \mid \mathcal{D}) = \sum_X \text{ScoreFam}(X \mid \text{Pa}_X : \mathcal{D})$

Skip



# Priors for General Graphs

- For finite datasets, prior is important!
- Prior over structure satisfying prior modularity
  - Eg,  $P(\mathcal{G}) / c^{|\mathcal{G}|}$      $|\mathcal{G}| = \# \text{edges}$ ;     $c < 1$
- What is good prior over *all* parameters?
  - *K2 prior*: fix  $\alpha \in \mathbb{R}^+$ , set  $\theta_{X_i | \text{Pa}(X_i)} \sim \text{Dirichlet}(\alpha, \dots, \alpha)$
  - Effective sample size, wrt  $X_i$ ?
    - If 0 parents:  $k \times \alpha$
    - If 1 binary parent:  $2 \times k \times \alpha$
    - If  $d$   $k$ -ary parents:  $k^d \times k \times \alpha$
  - So  $X_i$  "effective sample size" depends on #parental assignments
    - More parents  $\Rightarrow$  strong prior... doesn't make sense!
  - K2 is "inconsistent"



# Summary wrt Learning BN Structure

---

- Decomposable scores
  - Data likelihood
  - Information theoretic interpretation
  - Bayesian
  -
- Priors
  - Structure and parameter assumptions
- Best tree (Chow-Liu)
- Best TAN
- 
- 
- 
- 
- Bayesian model averaging