

# Graphical Models

## Lecture 19: Structure Learning


Andrew McCallum  
mccallum@cs.umass.edu

Thanks to Noah Smith and Carlos Guestrin for some slide materials.

# Administration

- HW#4 posted Monday.  
Now on official course website.

# Learning Bayesian Networks

	Known structure	Unknown structure
Fully observed data		hard (today)
Missing data	hard (last class)	very hard

# Goal of Learning?

- **Density estimation:** return a model  $M$  that precisely captures  $P^*$
- **Prediction:** optimize quality of answers to specific queries
- **Knowledge discovery:** reveal facts about the domain.

# Learning GM Structure

## Three Approaches

- **Constraint-based approaches:**  
use statistical tests to determine all conditional independencies, then construct the PDAG (I-equivalence class).
- **Score-based approaches:**  
learning as model selection considering a hypothesis space of models, select according to score, e.g. data likelihood.
- **Bayesian model averaging:**  
generate an ensemble of possible structures.

# Learning GM Structure

## Three Approaches


- **Constraint-based approaches:**  
use statistical tests to determine all conditional independencies, then construct the PDAG (I-equivalence class).
- **Score-based approaches:**  
learning as model selection considering a hypothesis space of models, select according to score, e.g. data likelihood.
- **Bayesian model averaging:**  
generate an ensemble of possible structures.

# Identifying the Graph Skeleton

- Let  $d$  be the maximum number of parents in  $G^*$ .
- For each pair  $X_i, X_j$ :
  - $E_{i,j} = \text{true}$
  - For each  $\mathbf{W}$  such that  $\mathbf{W} \subseteq \mathbf{X} \setminus \{X_i, X_j\}$  and  $|\mathbf{W}| \leq d$ :
    - If  $X_i \perp X_j \mid \mathbf{W}$  then  $E_{i,j} = \text{false}$  (and store  $\mathbf{W}$  as  $\mathbf{W}_{i,j}$ )
  - If  $E_{i,j}$  then add  $X_i \iff X_j$  to the skeleton

Need independence test that is robust to limited training data.

Chapter 3: *Build-Minimal-I-Map* method assumes an ordering, still requires  $2^{i-1}$  subsets of  $X_1, \dots, X_{i-1}$ .



# Learning GM Structure

## Three Approaches

- **Constraint-based approaches:**  
use statistical tests to determine all conditional independencies, then construct the PDAG (I-equivalence class).
- **Score-based approaches:**  
learning as model selection considering a hypothesis space of models, select according to score, e.g. data likelihood.
- **Bayesian model averaging:**  
generate an ensemble of possible structures.



# Likelihood Score and BN Structures

$$\begin{aligned}\max_{\mathcal{G}, \boldsymbol{\theta}} \log P_{\mathcal{G}, \boldsymbol{\theta}}(\mathbf{X} = \mathbf{x}) &= \max_{\mathcal{G}} \max_{\boldsymbol{\theta}} \log P_{\mathcal{G}, \boldsymbol{\theta}}(\mathbf{X} = \mathbf{x}) \\ &= \max_{\mathcal{G}} \log P_{\mathcal{G}, \boldsymbol{\theta}_{\text{MLE}}(\mathcal{G})}(\mathbf{X} = \mathbf{x})\end{aligned}$$

- For every possible graph structure, consider it with its best possible parameters (MLE)
  - This is “optimistic” but still correct if our overall goal is to maximize likelihood!

# Deriving the Structure Score for $\mathcal{G}$

$$\log P_{\mathcal{G},\theta}(\mathbf{X} = \mathbf{x}) = \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \text{count}(x_i, \mathbf{u}; \mathbf{x}) \log \theta_{x_i|\mathbf{u}}$$

$$\begin{aligned}
 &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} m \hat{P}(x_i, \mathbf{u}) \log \hat{P}(x_i | \mathbf{u}) \\
 &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} m \hat{P}(x_i, \mathbf{u}) \log \left( \frac{\hat{P}(x_i, \mathbf{u})}{\hat{P}(\mathbf{u})} \frac{\hat{P}(x_i)}{\hat{P}(x_i)} \right) \\
 &= m \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \hat{P}(x_i, \mathbf{u}) \left( \log \left( \frac{\hat{P}(x_i, \mathbf{u})}{\hat{P}(\mathbf{u}) \hat{P}(x_i)} + \log \hat{P}(x_i) \right) \right) \\
 &= m \sum_i I_{P_{\mathcal{G},\theta}}(X_i; \text{Parents}(X_i)) + m \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \hat{P}(x_i, \mathbf{u}) \log \hat{P}(x_i) \\
 &= m \sum_i I_{P_{\mathcal{G},\theta}}(X_i; \text{Parents}(X_i)) + m \sum_i \sum_{x_i \in \text{Val}(X_i)} \hat{P}(x_i) \log \hat{P}(x_i) \\
 &= m \sum_i I_{P_{\mathcal{G},\theta}}(X_i; \text{Parents}(X_i)) - m \sum_i H_{P_{\mathcal{G},\theta}}(X_i)
 \end{aligned}$$

# Deriving the Structure Score for $\mathcal{G}$

$\theta$  is the MLE!

$$\begin{aligned} \log P_{\mathcal{G},\theta}(\mathbf{X} = \mathbf{x}) &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \text{count}(x_i, \mathbf{u}; \mathbf{x}) \log \theta_{x_i|\mathbf{u}} \\ &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} m \hat{P}(x_i, \mathbf{u}) \log \hat{P}(x_i | \mathbf{u}) \end{aligned}$$

$$\begin{aligned} &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \theta_{x_i|\mathbf{u}} \\ &= m \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \frac{\text{count}(x_i, \mathbf{u})}{\text{count}(\mathbf{u})} \\ &= m \sum_i I_{P_{\mathcal{G},\theta}}(X_i; \text{Parents}(X_i)) \frac{\text{count}(x_i, \mathbf{u})}{m} = \hat{P}(x_i, \mathbf{u}) \\ &= m \sum_i I_{P_{\mathcal{G},\theta}}(X_i; \text{Parents}(X_i)) + m \sum_i \sum_{x_i \in \text{Val}(X_i)} \hat{P}(x_i) \log \hat{P}(x_i) \\ &= m \sum_i I_{P_{\mathcal{G},\theta}}(X_i; \text{Parents}(X_i)) - m \sum_i H_{P_{\mathcal{G},\theta}}(X_i) \end{aligned}$$

# Deriving the Structure Score for $\mathcal{G}$

$$\begin{aligned}
 \log P_{\mathcal{G},\theta}(\mathbf{X} = \mathbf{x}) &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \text{count}(x_i, \mathbf{u}; \mathbf{x}) \log \theta_{x_i|\mathbf{u}} \\
 &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} m \hat{P}(x_i, \mathbf{u}) \log \hat{P}(x_i | \mathbf{u}) \\
 &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} m \hat{P}(x_i, \mathbf{u}) \log \left( \frac{\hat{P}(x_i, \mathbf{u})}{\hat{P}(\mathbf{u})} \frac{\hat{P}(x_i)}{\hat{P}(x_i)} \right)
 \end{aligned}$$

$$\begin{aligned}
 &= m \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \hat{P}(x_i, \mathbf{u}) \left( \log \left( \frac{\hat{P}(x_i, \mathbf{u})}{\hat{P}(\mathbf{u}) \hat{P}(x_i)} + \log \hat{P}(x_i) \right) \right) \\
 &= m \sum_i I_{P_{\mathcal{G},\theta}}(X_i; \text{Parents}(X_i)) + m \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \hat{P}(x_i, \mathbf{u}) \log \hat{P}(x_i) \\
 &= m \sum_i I_{P_{\mathcal{G},\theta}}(X_i; \text{Parents}(X_i)) + m \sum_i \sum_{x_i \in \text{Val}(X_i)} \hat{P}(x_i) \log \hat{P}(x_i) \\
 &= m \sum_i I_{P_{\mathcal{G},\theta}}(X_i; \text{Parents}(X_i)) - m \sum_i H_{P_{\mathcal{G},\theta}}(X_i)
 \end{aligned}$$

# Deriving the Structure Score for $G$

$$\begin{aligned}
 \log P_{G,\theta}(\mathbf{X} = \mathbf{x}) &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \text{count}(x_i, \mathbf{u}; \mathbf{x}) \log \theta_{x_i|\mathbf{u}} \\
 &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} m \hat{P}(x_i, \mathbf{u}) \log \hat{P}(x_i | \mathbf{u}) \\
 &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} m \hat{P}(x_i, \mathbf{u}) \log \left( \frac{\hat{P}(x_i, \mathbf{u})}{\hat{P}(\mathbf{u})} \frac{\hat{P}(x_i)}{\hat{P}(x_i)} \right) \\
 &= m \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \hat{P}(x_i, \mathbf{u}) \left( \log \left( \frac{\hat{P}(x_i, \mathbf{u})}{\hat{P}(\mathbf{u}) \hat{P}(x_i)} + \log \hat{P}(x_i) \right) \right) \\
 &= m \sum_i I_{P_{G,\theta}}(X_i; \text{Parents}(X_i)) + m \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \hat{P}(x_i, \mathbf{u}) \log \hat{P}(x_i) \\
 &= m \sum_i I_{P_{G,\theta}}(X_i; \text{Parents}(X_i)) + m \sum_i \sum_{x_i \in \text{Val}(X_i)} \hat{P}(x_i) \log \hat{P}(x_i) \\
 &= m \sum_i I_{P_{G,\theta}}(X_i; \text{Parents}(X_i)) - m \sum_i H_{P_{G,\theta}}(X_i)
 \end{aligned}$$

# Deriving the Structure Score for $G$

$$\begin{aligned}
 \log P_{G,\theta}(\mathbf{X} = \mathbf{x}) &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \text{count}(x_i, \mathbf{u}; \mathbf{x}) \log \theta_{x_i|\mathbf{u}} \\
 &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} m \hat{P}(x_i, \mathbf{u}) \log \hat{P}(x_i | \mathbf{u}) \\
 &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} m \hat{P}(x_i, \mathbf{u}) \log \left( \frac{\hat{P}(x_i, \mathbf{u})}{\hat{P}(\mathbf{u})} \frac{\hat{P}(x_i)}{\hat{P}(x_i)} \right) \\
 &= m \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \hat{P}(x_i, \mathbf{u}) \left( \log \left( \frac{\hat{P}(x_i, \mathbf{u})}{\hat{P}(\mathbf{u}) \hat{P}(x_i)} + \log \hat{P}(x_i) \right) \right) \\
 &= m \sum_i I_{P_{G,\theta}}(X_i; \text{Parents}(X_i)) + m \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \hat{P}(x_i, \mathbf{u}) \log \hat{P}(x_i)
 \end{aligned}$$

$$= m \sum_i I_{P_{G,\theta}}(X_i; \text{Parents}(X_i)) + m \sum_i \sum_{x_i \in \text{Val}(X_i)} \hat{P}(x_i) \log \hat{P}(x_i)$$

$$= m \sum_i I_{P_{G,\theta}}(X_i; \text{Parents}(X_i)) - m \sum_i H_{P_{G,\theta}}(X_i)$$

# Deriving the Structure Score for $G$

$$\begin{aligned}
 \log P_{G,\theta}(\mathbf{X} = \mathbf{x}) &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \text{count}(x_i, \mathbf{u}; \mathbf{x}) \log \theta_{x_i|\mathbf{u}} \\
 &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} m \hat{P}(x_i, \mathbf{u}) \log \hat{P}(x_i | \mathbf{u}) \\
 &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} m \hat{P}(x_i, \mathbf{u}) \log \left( \frac{\hat{P}(x_i, \mathbf{u})}{\hat{P}(\mathbf{u})} \frac{\hat{P}(x_i)}{\hat{P}(x_i)} \right) \\
 &= m \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \hat{P}(x_i, \mathbf{u}) \left( \log \left( \frac{\hat{P}(x_i, \mathbf{u})}{\hat{P}(\mathbf{u}) \hat{P}(x_i)} + \log \hat{P}(x_i) \right) \right) \\
 &= m \sum_i I_{P_{G,\theta}}(X_i; \text{Parents}(X_i)) + m \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \hat{P}(x_i, \mathbf{u}) \log \hat{P}(x_i) \\
 &= m \sum_i I_{P_{G,\theta}}(X_i; \text{Parents}(X_i)) + m \sum_i \sum_{x_i \in \text{Val}(X_i)} \hat{P}(x_i) \log \hat{P}(x_i) \\
 &= m \sum_i I_{P_{G,\theta}}(X_i; \text{Parents}(X_i)) - m \sum_i H_{P_{G,\theta}}(X_i)
 \end{aligned}$$

# Deriving the Structure Score for $G$

$$\begin{aligned}
 \log P_{G,\theta}(\mathbf{X} = \mathbf{x}) &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \text{count}(x_i, \mathbf{u}; \mathbf{x}) \log \theta_{x_i|\mathbf{u}} \\
 &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} m \hat{P}(x_i, \mathbf{u}) \log \hat{P}(x_i | \mathbf{u}) \\
 &= \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} m \hat{P}(x_i, \mathbf{u}) \log \left( \frac{\hat{P}(x_i, \mathbf{u})}{\hat{P}(\mathbf{u})} \frac{\hat{P}(x_i)}{\hat{P}(x_i)} \right) \\
 &= m \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \hat{P}(x_i, \mathbf{u}) \left( \log \left( \frac{\hat{P}(x_i, \mathbf{u})}{\hat{P}(\mathbf{u}) \hat{P}(x_i)} + \log \hat{P}(x_i) \right) \right) \\
 &= m \sum_i I_{P_{G,\theta}}(X_i; \text{Parents}(X_i)) + m \sum_i \sum_{x_i \in \text{Val}(X_i)} \sum_{\mathbf{u} \in \text{Val}(\text{Parents}(X_i))} \hat{P}(x_i, \mathbf{u}) \log \hat{P}(x_i) \\
 &= m \sum_i I_{P_{G,\theta}}(X_i; \text{Parents}(X_i)) + m \sum_i \sum_{x_i \in \text{Val}(X_i)} \hat{P}(x_i) \log \hat{P}(x_i) \\
 &= m \sum_i I_{P_{G,\theta}}(X_i; \text{Parents}(X_i)) - m \sum_i H_{P_{G,\theta}}(X_i)
 \end{aligned}$$





# Decomposition

$$\log P_{\mathcal{G},\theta}(\mathbf{X} = \mathbf{x}) = m \sum_i I_{P_{\mathcal{G},\theta}}(X_i; \text{Parents}(X_i)) - m \sum_i H_{P_{\mathcal{G},\theta}}(X_i)$$

- Structure's likelihood score *decomposes* by family.
  - Good for efficiency!

$$\text{familyscore}(X_i, \text{Parents}(X_i)) = m I_{P_{\mathcal{G},\theta}}(X_i; \text{Parents}(X_i)) - m H_{P_{\mathcal{G},\theta}}(X_i)$$

# Bad News

- Property of mutual information:
  - Unless conditional independence holds *exactly* in the data, more connections are always better!

$$I(\mathbf{X}; \mathbf{Y}) \geq I(\mathbf{X}; \mathbf{Y} \cup \mathbf{Z})$$

- For structures, MLE will overfit.
  - This will happen even for non-table CPDs.
  - Need additional mechanism to disallow overly complicated structures, e.g. fixed indegree.

# One Solution: **Chow-Liu**

- Assume that each node can only have at most one parent.\*
- We now have a much simpler decision to make; consider  $I(X_i ; X_j)$  to be the score of putting an edge between  $X_i$  and  $X_j$
- Find the maximum-scoring spanning tree.
  - Number of trees?  $2^{\Theta(n \log n)}$

\*Interesting assumption. Can we do better?

# Chow-Liu: Attractions

- Maximum-scoring spanning tree gives us a skeleton.
- Two trees with the same skeleton will have the same ...
  - conditional independence assertions
  - mutual information score
- No need to worry about V-structures here!

# Taking Stock

- MLE is easy to generalize for Bayesian networks with a fixed structure.
  - Assumes parameters are disjoint for each CPD.
- Maximum likelihood structure learning selects a trivial, fully-connected Bayesian network.
  - Greedy assumption about parameters (MLE).
  - Heuristic solution: give each random variable one parent.

# Learning GM Structure

## Three Approaches

- **Constraint-based approaches:**  
use statistical tests to determine all conditional independencies, then construct the PDAG (I-equivalence class).
- **Score-based approaches:**  
learning as model selection considering a hypothesis space of models, select according to score, e.g. data *maximum a posteriori* likelihood (with prior on model complexity)
- **Bayesian model averaging:**  
generate an ensemble of possible structures.

# Bayesian Mantra

- If you are uncertain about something, put a probability distribution over it!
- In parameter learning, this applies to the parameters.
- In structure learning, this applies to the parameters *and* the structure.

# Most Probable vs. “Average” Parameters

- Before we saw the Bayesian approach, we could calculate the MLE and the likelihood score.

$$(\arg) \max_{\theta} P(\mathbf{X} = \mathbf{x} \mid \theta)$$

- Bayesian setting:

$$\arg \max_{\theta} P(\mathbf{X} = \mathbf{x} \mid \theta) P(\theta)$$

$$P(\mathbf{X} = \mathbf{x}) = \int P(\mathbf{X} = \mathbf{x} \mid \theta) P(\theta) d\theta$$



# Notation Issue

- This notation is really hiding some things:
  - Bayesian network *structure*
  - Prior over  $\boldsymbol{\theta}$

$$\arg \max_{\boldsymbol{\theta}} P(\mathbf{X} = \mathbf{x} \mid \boldsymbol{\theta}) P(\boldsymbol{\theta})$$

$$P(\mathbf{X} = \mathbf{x}) = \int P(\mathbf{X} = \mathbf{x} \mid \boldsymbol{\theta}) P(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

# Notation Issue

Better:

$$P(\mathbf{X} = \mathbf{x} \mid \mathcal{G}, \boldsymbol{\alpha}) = \int P(\mathbf{X} = \mathbf{x} \mid \boldsymbol{\theta}, \mathcal{G}) P(\boldsymbol{\theta} \mid \mathcal{G}, \boldsymbol{\alpha}) d\boldsymbol{\theta}$$

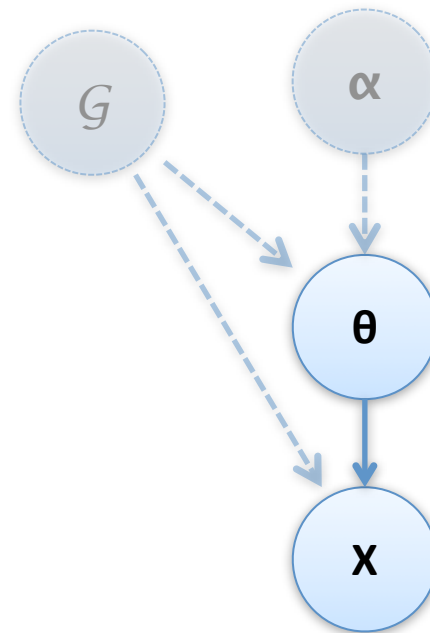
- This notation is really hiding some things:
  - Bayesian network *structure* ( $\mathcal{G}$ )
  - Prior over  $\boldsymbol{\theta}$  ( $\boldsymbol{\alpha}$ )

$$\arg \max_{\boldsymbol{\theta}} P(\mathbf{X} = \mathbf{x} \mid \boldsymbol{\theta}) P(\boldsymbol{\theta})$$

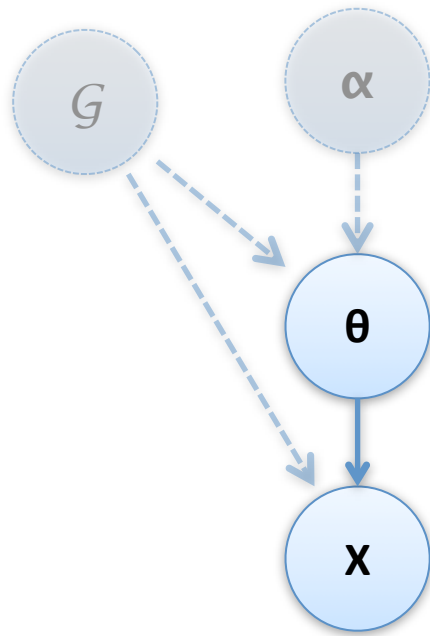
$$P(\mathbf{X} = \mathbf{x}) = \int P(\mathbf{X} = \mathbf{x} \mid \boldsymbol{\theta}) P(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

# Meta Bayesian Network

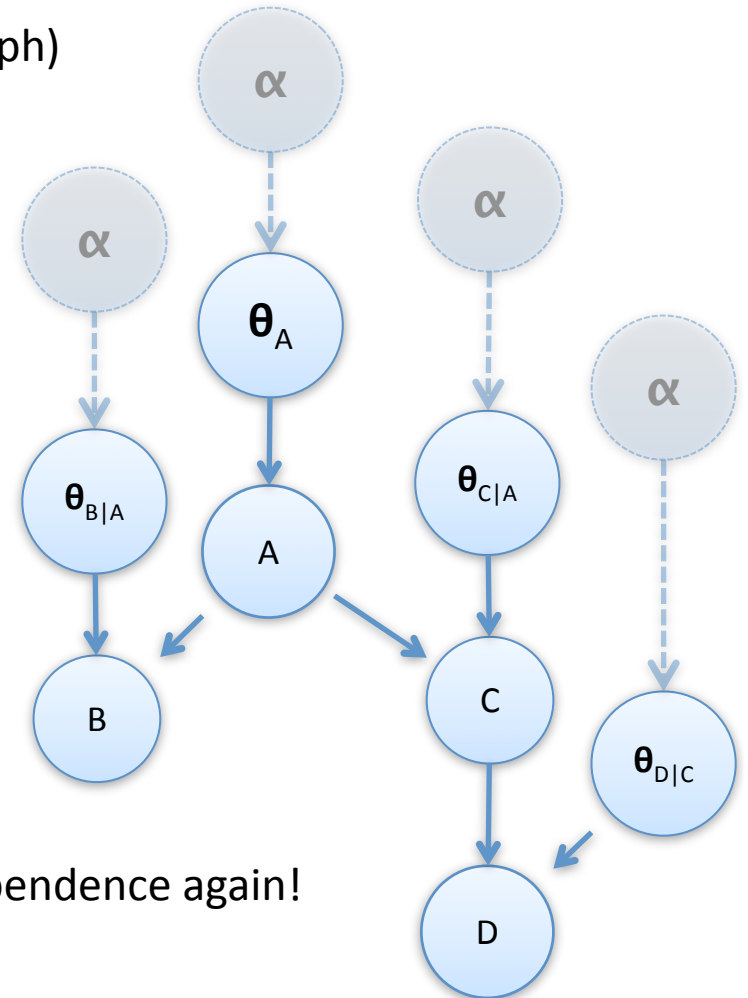
$$P(\mathbf{X} = \mathbf{x} \mid \mathcal{G}, \boldsymbol{\alpha}) = \int P(\mathbf{X} = \mathbf{x} \mid \boldsymbol{\theta}, \mathcal{G}) P(\boldsymbol{\theta} \mid \mathcal{G}, \boldsymbol{\alpha}) d\boldsymbol{\theta}$$



# Example

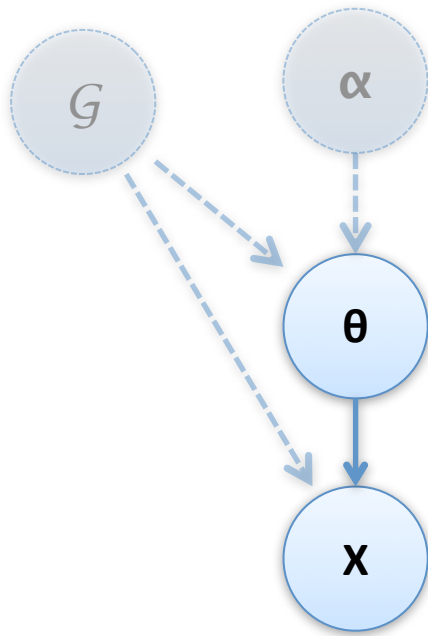


(for one graph)

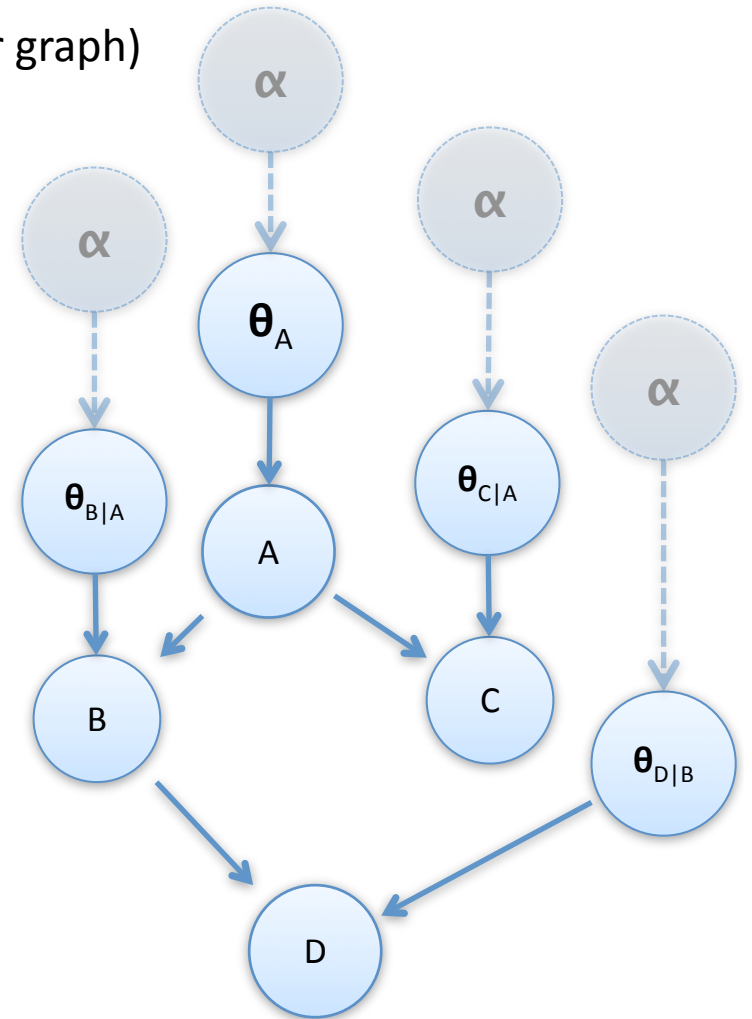


Note that it's helpful to assume global parameter independence again!

# Example

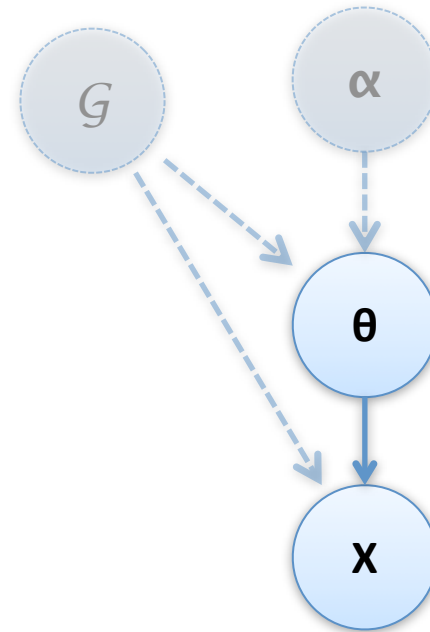


(for another graph)



# Next Steps

- Put a prior on  $\mathcal{G}$ ?
- Use full distribution over  $\theta$  to select  $\mathcal{G}$ ?



# Bayesian Score for Structures

$$\begin{aligned}P(\mathcal{G} \mid \mathbf{x}) &= \frac{P(\mathbf{x} \mid \mathcal{G})P(\mathcal{G})}{P(\mathbf{x})} \\ &\propto P(\mathbf{x} \mid \mathcal{G})P(\mathcal{G}) \\ &= \left( \int P(\mathbf{x} \mid \mathcal{G}, \boldsymbol{\theta})P(\boldsymbol{\theta} \mid \mathcal{G}) d\boldsymbol{\theta} \right) P(\mathcal{G})\end{aligned}$$

$$\log P(\mathcal{G} \mid \mathbf{x}) \propto \boxed{\log \int P(\mathbf{x} \mid \mathcal{G}, \boldsymbol{\theta})P(\boldsymbol{\theta} \mid \mathcal{G}) d\boldsymbol{\theta}} + \boxed{\log P(\mathcal{G})}$$

this part is harder

this part is  
easier

# Bayesian Score

$$\log \int P(\mathbf{x} \mid \mathcal{G}, \boldsymbol{\theta}) P(\boldsymbol{\theta} \mid \mathcal{G}) d\boldsymbol{\theta} + \log P(\mathcal{G})$$

- Assume global parameter independence:

$$\sum_i \log \int \prod_m P(x_i^{(m)} \mid \text{Par}_{\mathcal{G}}(x_i^{(m)}), \boldsymbol{\theta}_{X_i \mid \text{Par}_{\mathcal{G}}(X_i)}) P(\boldsymbol{\theta}_{X_i \mid \text{Par}_{\mathcal{G}}(X_i)} \mid \mathcal{G}) d\boldsymbol{\theta}_{X_i \mid \text{Par}_{\mathcal{G}}(X_i)}$$

– Note decomposition!



# Bayesian Score

$$\log \int P(\mathbf{x} \mid \mathcal{G}, \boldsymbol{\theta}) P(\boldsymbol{\theta} \mid \mathcal{G}) d\boldsymbol{\theta} + \log P(\mathcal{G})$$

- Assume global and local parameter independence:

$$\sum_i \sum_{\mathbf{u} \in \text{Val}(\text{Par}_{\mathcal{G}}(X_i))} \log \int \prod_m P(x_i^{(m)} \mid \mathbf{u}, \boldsymbol{\theta}_{X_i | \text{Par}_{\mathcal{G}}(X_i) = \mathbf{u}}) P(\boldsymbol{\theta}_{X_i | \text{Par}_{\mathcal{G}}(X_i) = \mathbf{u}} \mid \mathcal{G}) d\boldsymbol{\theta}_{X_i | \text{Par}_{\mathcal{G}}(X_i) = \mathbf{u}}$$

– Note decomposition!

# Bayesian Score

- Want: decomposable score for structures!
- Global and local parameter independence are part of what we need.
- Also need:
  - Parameter modularity: if  $X$  has the same parents in  $\mathcal{G}$  and  $\mathcal{G}'$ ,  $\theta_{X|\text{Parents}(X)}$  has the same prior.
  - Structure modularity:  $P(\mathcal{G})$  decomposes into families.

$$\log \int P(\mathbf{x} | \mathcal{G}, \boldsymbol{\theta}) P(\boldsymbol{\theta} | \mathcal{G}) d\boldsymbol{\theta} + \log P(\mathcal{G})$$

# Bayesian Score

- Want: decomposable score for structures!
- Global and local parameter independence are part of what we need.
- Also need:
  - Parameter modularity: if  $X$  has the same parents in  $\mathcal{G}$  and  $\mathcal{G}'$ ,  $\theta_{X|\text{Parents}(X)}$  has the same prior.
  - Structure modularity:  $P(\mathcal{G})$  decomposes into families.

$$\log \int P(\mathbf{x} | \mathcal{G}, \boldsymbol{\theta}) P(\boldsymbol{\theta} | \mathcal{G}) d\boldsymbol{\theta} + \log P(\mathcal{G})$$

# Parameter Priors, Revisited

- We can't design separate priors over  $\theta$  for every different structure.
  - Parameter modularity helps some, but there are still issues.
- Assume discrete.
- One idea (**K2 prior**) lets every multinomial have a symmetric Dirichlet prior with  $\alpha$  imaginary counts for each event.

# K2 Prior

- More parents for  $X_i$  implies more distributions over  $|\text{Val}(X_i)|$  outcomes (condition on more parent configurations).
- More effective imaginary counts! Let  $k = |\text{Val}(X_i)|$  for all  $X_i$ .
  - Zero parents:  $k\alpha$  counts
  - One parent:  $k^2\alpha$  counts
  - Two parents:  $k^3\alpha$  counts
- Different graph structures are getting different priors with different effective sample sizes.

# Bayesian Dirichlet Equivalent (BDe) Prior

- Basic idea: use a joint probability distribution  $P'$  over the space  $(X_i, \text{Parents}(X_i))$ , and then add imaginary counts proportional to  $P'$ .
- Fix the total to  $\alpha$ , so that

$$\alpha_{x|\mathbf{u}} = \alpha P'(x, \mathbf{u})$$

# Bayesian Score

- Want: decomposable score for structures!
- Global and local parameter independence are part of what we need.
- Also need:
  - Parameter modularity: if  $X$  has the same parents in  $\mathcal{G}$  and  $\mathcal{G}'$ ,  $\theta_{X|\text{Parents}(X)}$  has the same prior.
  - Structure modularity:  $P(\mathcal{G})$  decomposes into families.

$$\log \int P(\mathbf{x} | \mathcal{G}, \boldsymbol{\theta}) P(\boldsymbol{\theta} | \mathcal{G}) d\boldsymbol{\theta} + \log P(\mathcal{G})$$

# $P(\mathcal{G})$

- Simple idea:  $P(\mathcal{G}) \propto c^{\#\text{edges}(\mathcal{G})}$ 
  - Structure modularity!



## Aside: Approximating the Bayesian Score

- General pattern:
  - Bayesian approach prefers *simpler* structures, but willing to allow a more complex structure if there's enough data.
- For Dirichlet priors over parameters the **Bayesian Information Criterion (BIC)** score approximates  $\log P(\mathcal{G} \mid \mathbf{x})$ :

$$m \sum_i I_{P_{\mathcal{G}, \theta_{\text{MLE}}}}(X_i, \text{Par}_{\mathcal{G}}(X_i)) - m \sum_i H_{P_{\mathcal{G}, \theta_{\text{MLE}}}}(X_i) - \frac{\log m}{2} \text{Dim}(\mathcal{G})$$

# BIC

- $\text{Dim}(\mathcal{G})$  = model dimension (number of independent parameters)
- Penalty:  $(\log m)/2$  per parameter
- What happens as we see more data?

$$m \sum_i I_{P_{\mathcal{G}}, \theta_{\text{MLE}}} (X_i, \text{Par}_{\mathcal{G}}(X_i)) - m \sum_i H_{P_{\mathcal{G}}, \theta_{\text{MLE}}} (X_i) - \frac{\log m}{2} \text{Dim}(\mathcal{G})$$

# Structure Score & Search

- Now we have a way to score a candidate structure.

How to find the best structure?

- Search!

# Structure Search

- Think of every Bayesian network structure as a *state*.
- Define a search operator that lets us move among the states.
  - Tradeoff: interconnectivity vs. speed
  - Typical: add, delete, reverse edges
- Initial state?
- Search procedure?

# Massive Literature on Search

- Computational cost
  - Key: score decomposition is crucial
- Local maxima
- Representation: search DAGs or l-equivalence classes of PDAGs?
  
- Start with K&F 18.4.3-4 (pp. 811-824).

# Learning GM Structure

## Three Approaches

- **Constraint-based approaches:**  
use statistical tests to determine all conditional independencies, then construct the PDAG (I-equivalence class).
- **Score-based approaches:**  
learning as model selection considering a hypothesis space of models, select according to score, e.g. data likelihood.
- **Bayesian model averaging:**  
generate an ensemble of possible structures.

# Average Across Several Structures

- Several structures may have similar scores. Suggests several may be close to “true” structure. Shouldn’t just pick one.

$$P(x) = \int \int P(x|\theta, \mathcal{G})P(\theta|\mathcal{G})P(\mathcal{G})d\theta d\mathcal{G}$$

- Can approximate this integrate by **sampling** (Markov-chain Monte Carlo) over structures.
  - More on this next week: *Dirichlet processes*.
  - Often used inappropriately(?) for structure search.





# Learning the Structure of a Markov Network

# Constraint-Based Structure Learning

- As in the Bayesian network case, we can use independence tests.
- Markov network independence assertions are not easy to localize.
  - Computational efficiency, also may not have sufficient data to perform tests involving many variables.
- If we assume the “true” network has bounded degree, we can use ...

# Identifying the Skeleton

- Let  $d$  be the maximum number of parents in  $G^*$ .
- For each pair  $X_i, X_j$ :
  - $E_{i,j} = \text{true}$
  - For each  $\mathbf{W}$  such that  $\mathbf{W} \subseteq \mathbf{X} \setminus \{X_i, X_j\}$  and  $|\mathbf{W}| \leq d$ :
    - If  $X_i \perp X_j \mid \mathbf{W}$  then  $E_{i,j} = \text{false}$  (and store  $\mathbf{W}$  as  $\mathbf{W}_{i,j}$ )
  - If  $E_{i,j}$  then add  $X_i \iff X_j$  to the skeleton

Earlier this Lecture

$$\sum_{k=0}^d \binom{n-2}{k}$$

# Score-Based Structure Learning

- Different levels of granularity:
  - Markov network (complexity measured in clique sizes; optimize tree-width?) – coarse
  - Factor graph (complexity measured in factor sizes)
  - Log-linear features (complexity measured in number of features included) – fine
- Coarser focus lets us think about the network and efficient inference.
- Finer focus lets us choose a parameterization that avoids overfitting. *This will be our focus.*

# Likelihood Score

- A log-linear model structure  $\mathcal{M}$  corresponds to a choice of features.

$$\max_{\mathcal{M}} \left( \max_{\mathbf{w}} \sum_t \log P(\mathbf{x}^{(t)} \mid \mathcal{M}, \mathbf{w}) \right)$$

- The likelihood score will always prefer more complex models, just like with Bayesian networks.
  - Only makes sense with strict constraints on the expressive power of the model: e.g., structure of the Markov network or number of features.

# Alternatives

- Bayesian scores, e.g., BIC:

$$\max_{\mathcal{M}} \left( \max_{\mathbf{w}} \sum_t \log P(\mathbf{x}^{(t)} | \mathcal{M}, \mathbf{w}) - \frac{\dim(\mathcal{M})}{2} \log T \right)$$

$\dim(\mathcal{M})$  is the number of degrees of freedom in the model; number of non-redundant features.

- Maximum *a posteriori*: use MAP estimation instead of MLE in the inner loop:

$$\max_{\mathcal{M}} \left( \max_{\mathbf{w}} \sum_t \log P(\mathbf{x}^{(t)} | \mathcal{M}, \mathbf{w}) + \log P(\mathbf{w} | \mathcal{M}) \right)$$

# Priors and Structural Sparsity

- We saw how the *Laplacian* prior can achieve feature-level sparsity.
- We can achieve *group* sparsity via block  $L_1$  regularization:
  - For each factor, collect the features that have scope over the factor's random variables.
  - Call these **groups**.
  - Penalize the  $L_1$  norm of group  $L_2$  norms:

$$- \sum_i \left| \sqrt{\sum_{j: f_j \in \text{Group}(\phi_i)} w_j^2} \right|$$

- As the overall magnitude of the group goes down, the weights are pushed more strongly to zero.

# Summary: Finding MN Structure

- Search!
  - Greedy, or not.
  - Putting parameter learning in the inner loop can be expensive; warm starts can help.
- Some tricks can be used to approximate search steps (see K&F 20.7.4).
- One option is to simply use  $L_1$  regularization and variations on it to thin out features (and therefore factors, maybe).