

Stat 521A
Lecture 24

Outline

- Scoring functions for DAGs with hidden vars (19.4.1)
- Structure search (19.4.2)
- Structural EM (19.4.3)
- Inventing hidden variables in DGMs (19.5)

Bayesian score

- Need a way to measure model quality; orthogonal to issue of how we search through space of models
- Bayesian score hard to compute since posterior is an exponential number of modes

$$\text{score}_{\mathcal{G}}(\mathcal{D}) = \log P(\mathcal{D} | \mathcal{G}) + \log P(\mathcal{G})$$

$$p(\mathcal{D} | G) = \int \prod_m p(\mathbf{o}[m] | \boldsymbol{\theta}, G) p(\boldsymbol{\theta} | G) d\boldsymbol{\theta}$$

$$p(\mathbf{o}[m] | \boldsymbol{\theta}, G) = \sum_{\mathbf{h}} p(\mathbf{o}[m], \mathbf{h} | \boldsymbol{\theta}, G)$$

- Approximations: asymptotic, variational, MCMC

Chib's candidate method

- Approximate $p(\mathcal{D}|\mathcal{G})$ using output of a standard MCMC run. For any θ (eg MAP) compute

$$P(\mathcal{D} | \mathcal{G}) = \frac{P(\mathcal{D} | \theta, \mathcal{G})P(\theta | \mathcal{G})}{P(\theta | \mathcal{D}, \mathcal{G})}.$$

- Requires that $p(\theta|\mathcal{D},\mathcal{G})$ cover chosen θ .
- This requires that MCMC mix over all posterior modes, even if symmetrical. If not, it will underestimate $p(\mathcal{D}|\mathcal{G})$. See rejected letter to editor by Radford Neal.*

* <http://www.cs.utoronto.ca/~radford/ftp/chib-letter.pdf>

RJMCMC

- Instead of doing discrete search, and integrating out params at each point, let us jointly sample in graph and param space
- Since the size of the cts space is changing, we need to use a change of measure when we move between dimensionalities
- This results in reversible jump MCMC
- Getting it working is delicate...

Laplace approximation

Box 19.F — Concept: Laplace Approximation. The Laplace approximation can be applied to any function of the form $f(\mathbf{w}) = e^{g(\mathbf{w})}$ for some vector \mathbf{w} . Our task is to compute the integral

$$F = \int f(\mathbf{w}) d\mathbf{w}$$

Using Taylor's expansion, we can expand an approximation of g around a point \mathbf{w}_0

$$g(\mathbf{w}) \approx g(\mathbf{w}_0) + \left[\frac{\partial g(\mathbf{w})}{\partial x_i} \right] \Big|_{\mathbf{w}=\mathbf{w}_0} (\mathbf{w} - \mathbf{w}_0) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_0)^T \left[\frac{\partial^2 g(\mathbf{w})}{\partial x_i \partial x_j} \right] \Big|_{\mathbf{w}=\mathbf{w}_0} (\mathbf{w} - \mathbf{w}_0),$$

where $\left[\frac{\partial g(\mathbf{w})}{\partial x_i} \right] \Big|_{\mathbf{w}=\mathbf{w}_0}$ denotes the vector of first derivatives and $\left[\frac{\partial^2 g(\mathbf{w})}{\partial x_i \partial x_j} \right] \Big|_{\mathbf{w}=\mathbf{w}_0}$ denotes the Hessian — the matrix of second derivatives.

If \mathbf{w}_0 is the maximum of $g(\mathbf{w})$, then the second term disappears. We now set

$$C = - \left[\frac{\partial^2 g(\mathbf{w})}{\partial x_i \partial x_j} \right] \Big|_{\mathbf{w}=\mathbf{w}_0}$$

to be the negative of the matrix of second derivatives of $g(\mathbf{w})$ at \mathbf{w}_0 . Since \mathbf{w}_0 is a maximum, this matrix is positive semi-definitive. Thus, we get the approximation

$$g(\mathbf{w}) \approx g(\mathbf{w}_0) - \frac{1}{2} (\mathbf{w} - \mathbf{w}_0)^T C (\mathbf{w} - \mathbf{w}_0).$$

Plugging this approximation into the definition of $f(x)$, we can write

$$\int f(\mathbf{w}) d\mathbf{w} \approx f(\mathbf{w}_0) \int e^{-\frac{1}{2} (\mathbf{w} - \mathbf{w}_0)^T C (\mathbf{w} - \mathbf{w}_0)} d\mathbf{w}.$$

The integral is identical to the integral of an unnormalized Gaussian distribution with covariance matrix $\Sigma = C^{-1}$. We can therefore solve this integral analytically and obtain:

$$\int f(\mathbf{w}) d\mathbf{w} \approx f(\mathbf{w}_0) |C|^{-\frac{1}{2}} (2\pi)^{\frac{1}{2} \dim(C)}$$

where $\dim(C)$ is the dimension of the matrix C .

Laplace approximation cont'd

- Let $g(w) = \log p(D, w | G)$.
- Laplace approximation to $p(D, G)$ is

$$\text{score}_{\text{Laplace}}(\mathcal{G} : \mathcal{D}) = \log P(\mathcal{G}) + \log P(\mathcal{D} | \tilde{\theta}_{\mathcal{G}}, \mathcal{G}) + \frac{\dim(C)}{2} \log 2\pi - \frac{1}{2} \log |C|,$$

C is negative Hessian: requires inference on x_i, x_j, u_i, u_j

$$-\frac{\partial^2 \log P(\mathcal{D} | \theta, \mathcal{G})}{\partial \theta_{x_i | u_i} \partial \theta_{x_j | u_j}} \Big|_{\tilde{\theta}_{\mathcal{G}}} = -\sum_m \frac{\partial^2 \log P(o[m] | \theta, \mathcal{G})}{\partial \theta_{x_i | u_i} \partial \theta_{x_j | u_j}} \Big|_{\tilde{\theta}_{\mathcal{G}}},$$

BIC score

- BIC is the limit of Laplace as $M \rightarrow \infty$.

$$C = \sum_{m=1}^M C_m \quad C = M \frac{1}{M} \sum_{m=1}^M C_m.$$

$$\det(C) = M^{\dim(C)} \det\left(\frac{1}{M} \sum_{m=1}^M C_m\right) \approx M^{\dim(C)} \det(E_{P^*}[C_o]).$$

$$\log \det(C) \approx \dim(C) \log M + \log \det(E_{P^*}[C_o]).$$

Theorem 19.4.1: *As $M \rightarrow \infty$, we have that:*

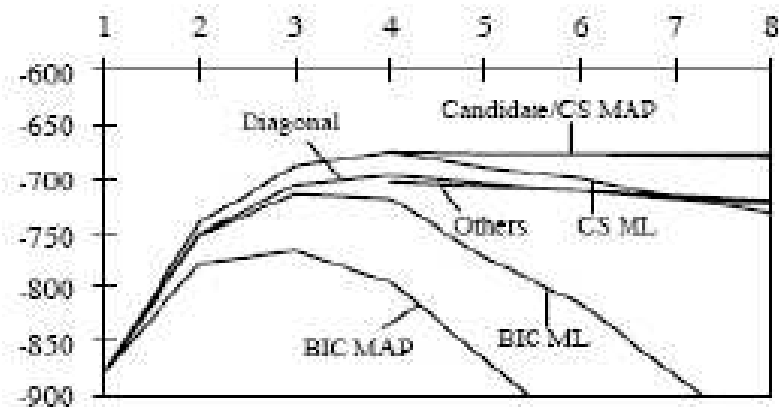
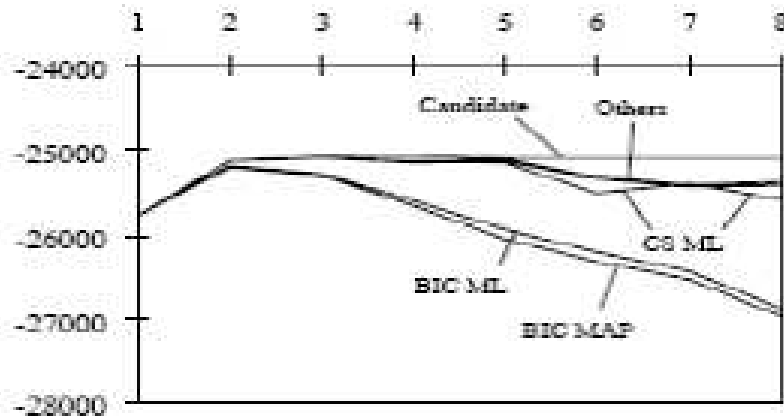
$$\text{score}_{\text{Laplace}}(\mathcal{G} ; \mathcal{D}) = \text{score}_{\text{BIC}}(\mathcal{G} ; \mathcal{D}) + O(1)$$

where $\text{score}_{\text{BIC}}(\mathcal{G} ; \mathcal{D})$ is the BIC score

$$\text{score}_{\text{BIC}}(\mathcal{G} ; \mathcal{D}) = \log P(\mathcal{D} | \tilde{\theta}_{\mathcal{G}}, \mathcal{G}) - \frac{\log M}{2} \text{Dim}[\mathcal{G}] + \log P(\mathcal{G}) + \log P(\tilde{\theta}_{\mathcal{G}} | \mathcal{G}).$$

Cheeseman-Stutz approximation

- CS approx to $\log p(D|G)$ is more accurate than BIC, yet faster than Laplace
- Matt Beal's thesis proves CS is a lower bound
- Example: we plot $\log p(D|K)$ vs K for a mixture of Bernoullis for different methods; 'candidate' is a 'gold standard' MCMC method



CS approx

- Idea 1: If D^* is complete, $p(D^*|G)$ just relies on sufficient statistics, so use ESS instead

$$P(D_{\mathcal{G}, \hat{\theta}_{\mathcal{G}}}^* | \mathcal{G}) = \int p(D_{\mathcal{G}, \hat{\theta}_{\mathcal{G}}}^* | \theta, \mathcal{G}) P(\theta | \mathcal{G}) d\theta$$

- Unfortunately this does not work well, since it sums over 1 (imputed) dataset whereas $p(D|G)$ sums over an exponential number

$$P(\mathcal{D} | \mathcal{G}) = \int \sum_{\mathcal{H}} p(\mathcal{D}, \mathcal{H} | \theta, \mathcal{G}) P(\theta | \mathcal{G}) d\theta = \sum_{\mathcal{H}} \int p(\mathcal{D}, \mathcal{H} | \theta, \mathcal{G}) P(\theta | \mathcal{G}) d\theta.$$

- Idea 2: add an approximate correction term

$$\log P(\mathcal{D} | \mathcal{G}) = \log P(D_{\mathcal{G}, \hat{\theta}_{\mathcal{G}}}^* | \mathcal{G}) + \underbrace{\log P(\mathcal{D} | \mathcal{G}) - \log P(D_{\mathcal{G}, \hat{\theta}_{\mathcal{G}}}^* | \mathcal{G})}_{\text{approximate correction term}}$$

Approximate with BIC

CS approx

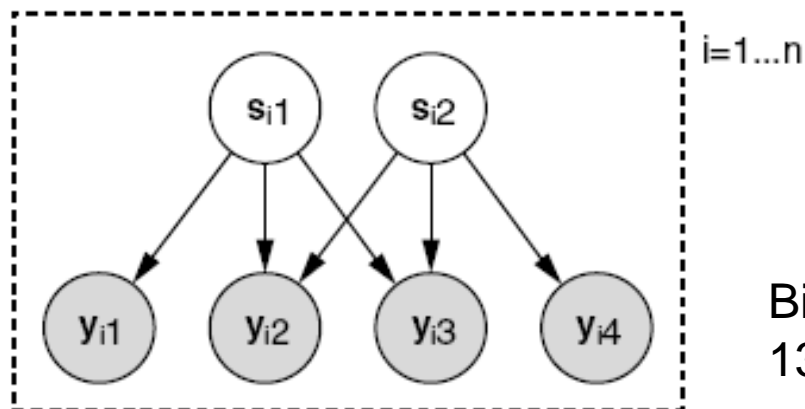
$$\begin{aligned}\log P(\mathcal{D} | \mathcal{G}) - \log P(\mathcal{D}_{\mathcal{G}, \tilde{\theta}_{\mathcal{G}}}^* | \mathcal{G}) &\approx \left[\log P(\mathcal{D} | \tilde{\theta}_{\mathcal{G}}, \mathcal{G}) - \frac{1}{2} \text{Dim}[\tilde{\theta}_{\mathcal{G}}] \log M \right] \\ &\quad - \left[\log P(\mathcal{D}_{\mathcal{G}, \tilde{\theta}_{\mathcal{G}}}^* | \tilde{\theta}_{\mathcal{G}}, \mathcal{G}) - \frac{1}{2} \text{Dim}[\tilde{\theta}_{\mathcal{G}}] \log M \right] \\ &= \log P(\mathcal{D} | \tilde{\theta}_{\mathcal{G}}, \mathcal{G}) - \log P(\mathcal{D}_{\mathcal{G}, \tilde{\theta}_{\mathcal{G}}}^* | \tilde{\theta}_{\mathcal{G}}, \mathcal{G}).\end{aligned}$$

$$\begin{aligned}\log P(\mathcal{D} | \mathcal{G}) &= \log P(\mathcal{D}_{\mathcal{G}, \tilde{\theta}_{\mathcal{G}}}^* | \mathcal{G}) + \log P(\mathcal{D} | \mathcal{G}) - \log P(\mathcal{D}_{\mathcal{G}, \tilde{\theta}_{\mathcal{G}}}^* | \mathcal{G}) \\ &\approx \log P(\mathcal{D}_{\mathcal{G}, \tilde{\theta}_{\mathcal{G}}}^* | \mathcal{G}) + \log P(\mathcal{D} | \tilde{\theta}_{\mathcal{G}}, \mathcal{G}) - \log P(\mathcal{D}_{\mathcal{G}, \tilde{\theta}_{\mathcal{G}}}^* | \tilde{\theta}_{\mathcal{G}}, \mathcal{G}).\end{aligned}$$

$$\text{score}_{CS}(\mathcal{G} : \mathcal{D}) = \log P(\mathcal{G}) + \log P(\mathcal{D}_{\mathcal{G}, \tilde{\theta}_{\mathcal{G}}}^* | \mathcal{G}) + \log P(\mathcal{D} | \tilde{\theta}_{\mathcal{G}}, \mathcal{G}) - \log P(\mathcal{D}_{\mathcal{G}, \tilde{\theta}_{\mathcal{G}}}^* | \tilde{\theta}_{\mathcal{G}}, \mathcal{G})$$

Variational lower bound

EM for MAP estimation	Variational Bayesian EM
<p>Goal: maximise $p(\theta y, m)$ w.r.t. θ</p> <p>E Step: compute $q_x^{(t+1)}(x) = p(x y, \theta^{(t)})$</p> <p>M Step: $\theta^{(t+1)} = \arg \max_{\theta} \int dx q_x^{(t+1)}(x) \ln p(x, y, \theta)$</p>	<p>Goal: lower bound $p(y m)$</p> <p>VBE Step: compute $q_x^{(t+1)}(x) = p(x y, \bar{\phi}^{(t)})$</p> <p>VBM Step: $q_{\theta}^{(t+1)}(\theta) \propto \exp \int dx q_x^{(t+1)}(x) \ln p(x, y, \theta)$</p>



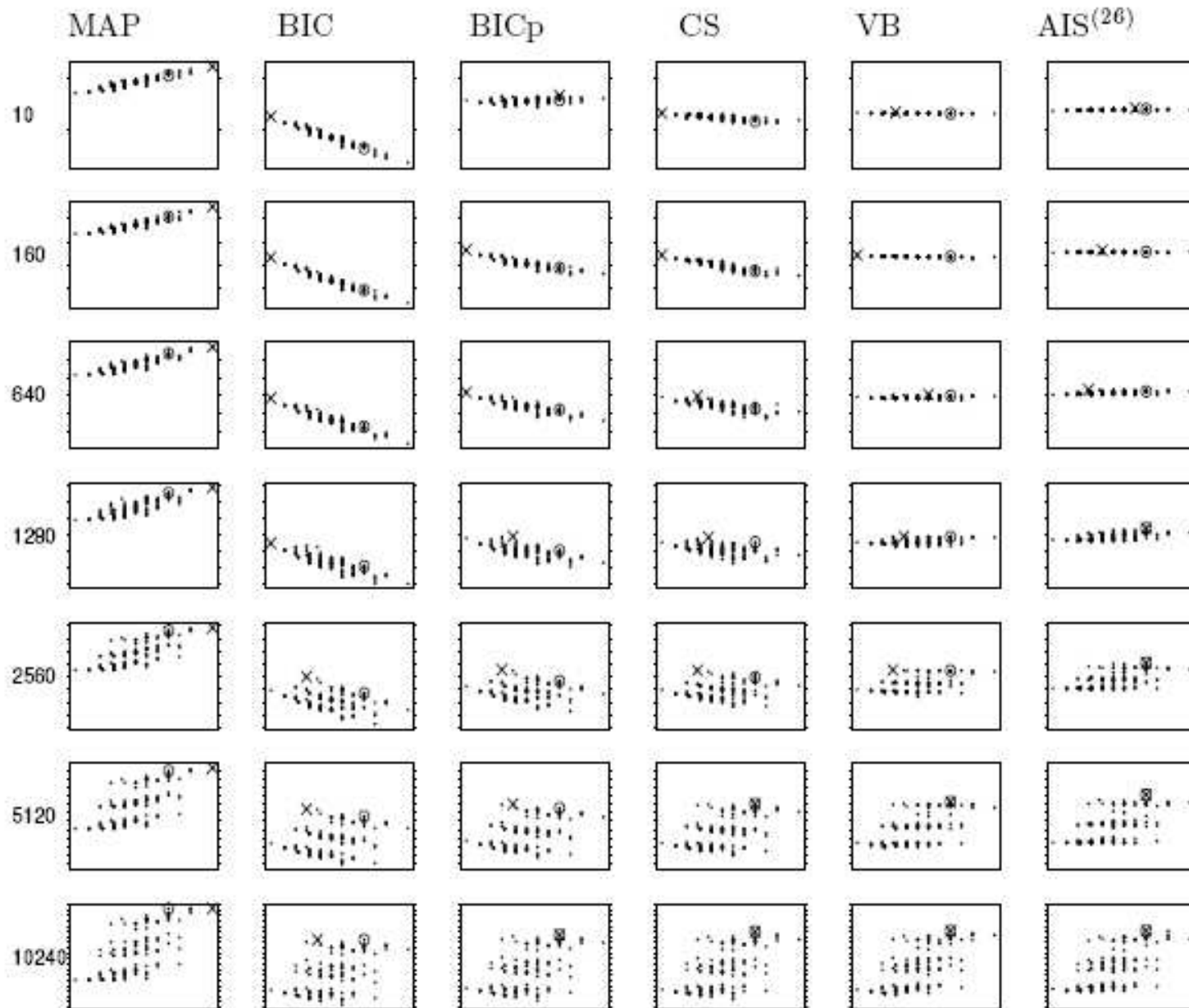
VB provably tighter lower Bound than CS

Binary hidden nodes, 5-ary obs nodes
136 distinct DAGs

Beal, M.J. and Ghahramani, Z.

Variational Bayesian Learning of Directed Graphical Models with Hidden Variables
[Bayesian Analysis](#) 1(4), 2006.

Log p(D|G) vs dof(G)

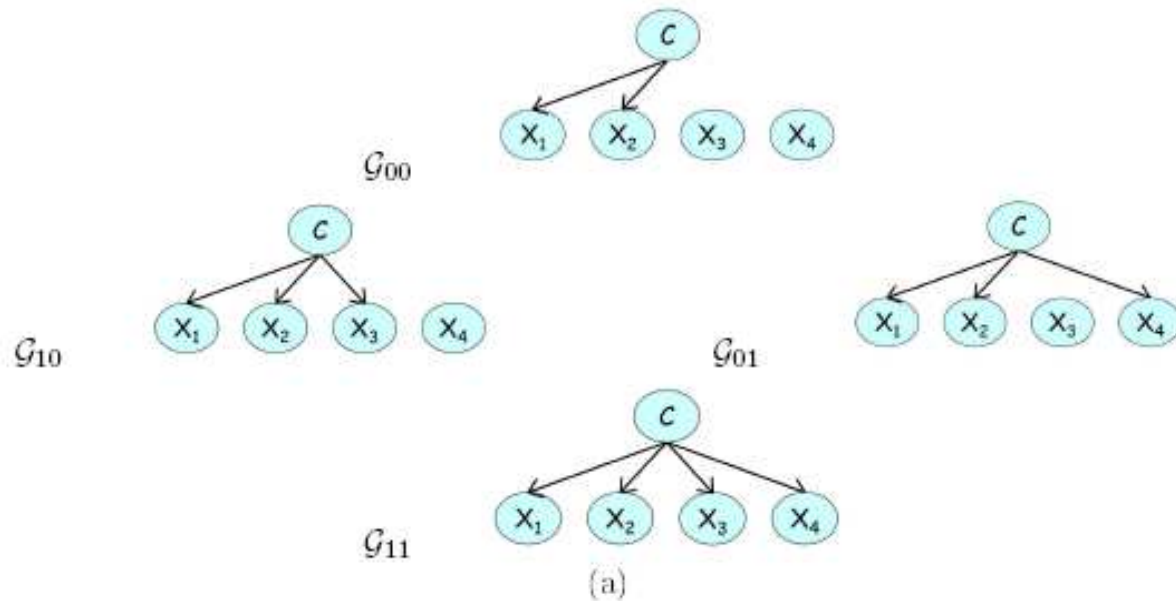




Structure search

- $P(D|G)$ does not factorize across families, unlike the fully observed case
- Cannot find (easily) optimal tree or optimal DAG given ordering.
- For local search, evaluating score of neighbors is expensive – score does not decompose, so need to find MAP estimate for each graph just to compute its BIC score

Illustration of non-decomposability



{1,2} and 3: weak corr
3 and 4: strong corr

Network	ΔLL	ΔCS
\mathcal{G}_{10} (add $C \rightarrow X_3$)	+3	-0.4
\mathcal{G}_{01} (add $C \rightarrow X_4$)	+10.6	+7.2
\mathcal{G}_{11} (add $C \rightarrow X_3, C \rightarrow X_4$)	+24.1	+17.4

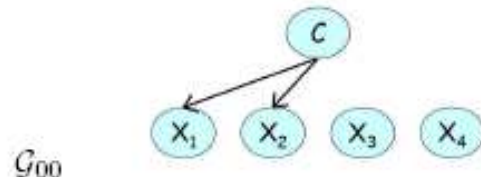
Structural EM

- Given current graph G_t and MAP params $\theta(t)$, compute ESS for all possible families (potentially in a lazy fashion – may need out-of-clq queries)
- Evaluate BIC score for $G(t+1)$ using ESS| $G(t)$
- Thm: increasing expected BIC score increases true BIC score

Theorem 19.4.3: *Let \mathcal{G}_0 be a graph structure and $\tilde{\theta}_0$ be the MAP parameters for \mathcal{G}_0 given a dataset \mathcal{D} . Then for any graph structure \mathcal{G} :*

$$\text{score}_{BIC}(\mathcal{G} : \mathcal{D}_{\mathcal{G}_0, \tilde{\theta}_0}^*) - \text{score}_{BIC}(\mathcal{G}_0 : \mathcal{D}_{\mathcal{G}_0, \tilde{\theta}_0}^*) \leq \text{score}_{BIC}(\mathcal{G} : \mathcal{D}) - \text{score}_{BIC}(\mathcal{G}_0 : \mathcal{D}).$$

Sparse mixture model



- Run parameter estimation (such as EM or gradient ascent) to learn parameters $\tilde{\theta}_t$ for \mathcal{G}_t .
- Construct a new structure \mathcal{G}_{t+1} so that \mathcal{G}_{t+1} contains the edge $C \rightarrow X_i$ if

$$\text{FamScore}(X_i, \{C\} : \mathcal{D}_{\mathcal{G}_t, \tilde{\theta}_t}^*) > \text{FamScore}(X_i, \emptyset : \mathcal{D}_{\mathcal{G}_t, \tilde{\theta}_t}^*).$$

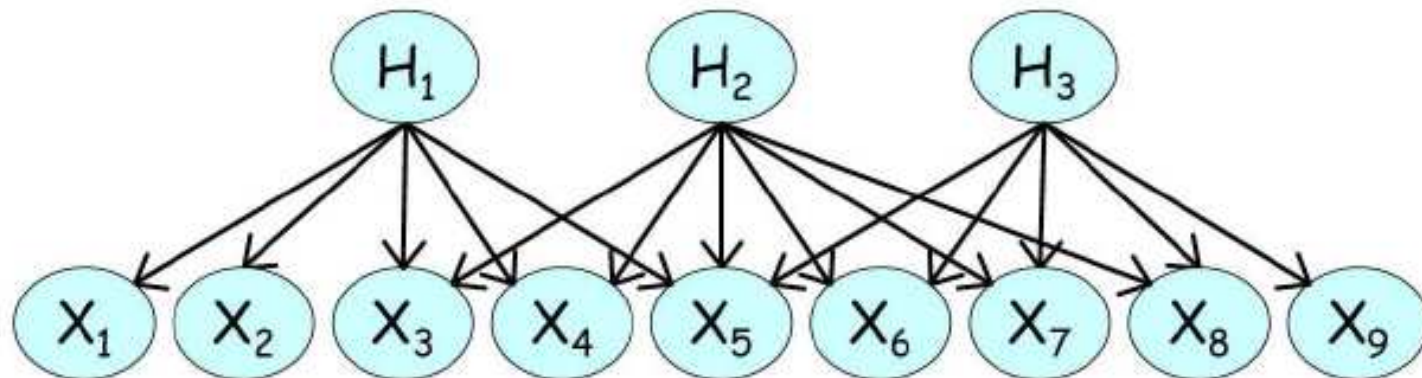
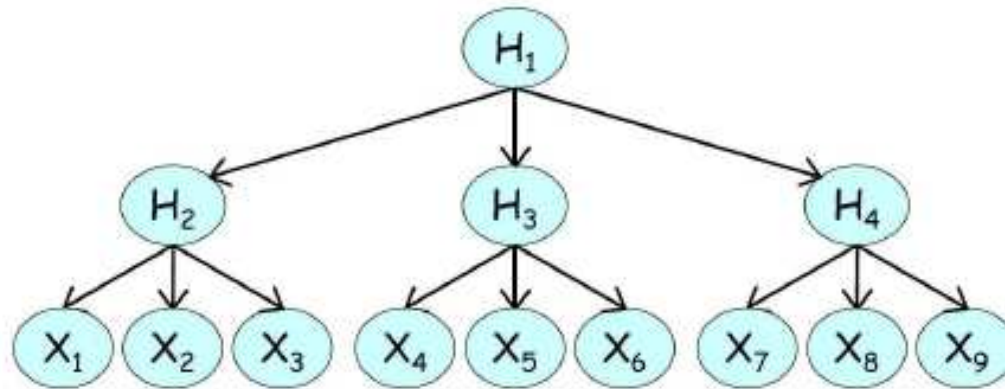
$$\begin{aligned} \bar{M}_{\mathcal{D}_{\mathcal{G}_t, \tilde{\theta}_t}^*} [x_i, c] &= \sum_m P(C[m] = c, X_i[m] = x_i \mid \mathbf{o}[m], \mathcal{G}_t, \tilde{\theta}_t) \\ &= \sum_{m, X_i[m]=x_i} P(C[m] = c \mid \mathbf{o}[m], \mathcal{G}_t, \tilde{\theta}_t). \end{aligned}$$

Initialization: if start from no children, will never add any! So start from all Children or random subset.



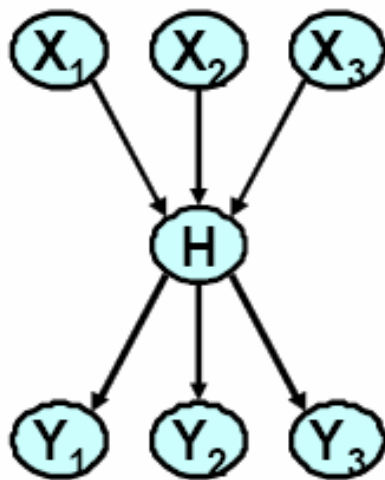
Inventing hidden variables

- Can add hidden variables in 'canonical' places

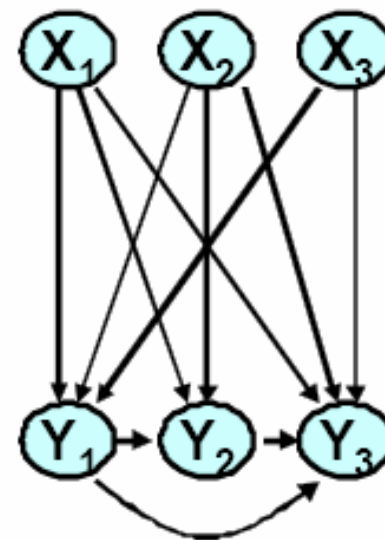


Structural signatures

- Can learn structure with no hidden vars, then look for ‘semi-cliques’.
- Unfortunately original model discourages nodes with high fan-in.



17 parameters



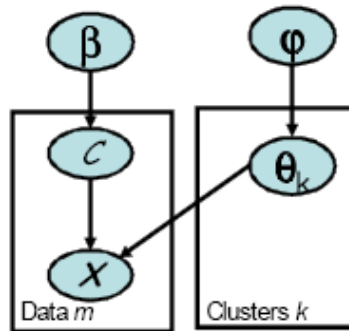
59 parameters

Can also look for signatures in the data - eg FCI* algorithm

Cardinality of hidden nodes

- Need to choose number of states.
- Can use an “infinite” number using Dirichlet processes.
- Let us first consider DP mixture models.

Marginalizing out θ



Collapsed Gibbs sampling
Cf DP mixtures

$$P(\lambda | c) = \text{Dirichlet}(\alpha_0/K + |I_1(c)|, \dots, \alpha_0/K + |I_K(c)|).$$

$$P(\theta_k | c, \mathcal{D}, \phi) = Q(\theta_k | \mathcal{D}_{I_k(c)}, \phi) \propto P(\theta_k | \phi) \prod_{m \in I_k(c)} P(x[m] | \theta_k),$$

$$P(C[m'] = k | c_{-m'}, \mathcal{D}, \phi) \propto P(C[m'] = k | \lambda, c_{-m'}) P(x[m'] | C[m'] = k, x[I_k(c_{-m'})], \phi).$$

$$P(C[m'] = k | c_{-m'}, \mathcal{D}, \phi) \propto (|I_k(c_{-m'})| + \alpha_0/K) Q(X | \mathcal{D}_{I_k(c_{-m'})}, \phi).$$

O(M K) per iter

DP mixture model (p865)

- Identity of clusters does not matter. Let $\sigma = \{I_1, \dots, I_L\}$ be a partition, $I_c =$ cases in cluster c . For case m' , either join existing cluster or create new one $O(ML)$ per iter

$$P(I \leftarrow I \cup \{m'\} \mid \sigma_{-m'}, \mathcal{D}, \phi) \propto \left(|I| + \frac{\alpha_0}{K}\right) Q(x[m'] \mid \mathcal{D}_I, \phi)$$

$$P(\sigma \leftarrow \sigma \cup \{\{m'\}\} \mid \sigma_{-m'}, \mathcal{D}, \phi) \propto (K - L) \frac{\alpha_0}{K} Q(x[m'] \mid \phi),$$

- Now let $K \rightarrow \infty$.

$$P(I \leftarrow I \cup \{m'\} \mid \sigma_{-m'}, \mathcal{D}, \phi) \propto |I| \cdot Q(x[m'] \mid \mathcal{D}_I, \phi)$$

$$P(\sigma \leftarrow \sigma \cup \{\{m'\}\} \mid \sigma_{-m'}, \mathcal{D}, \phi) \propto \alpha_0 \cdot Q(x[m'] \mid \phi).$$

- More likely to join a cluster if it is already crowded.
- Chinese Restaurant process.