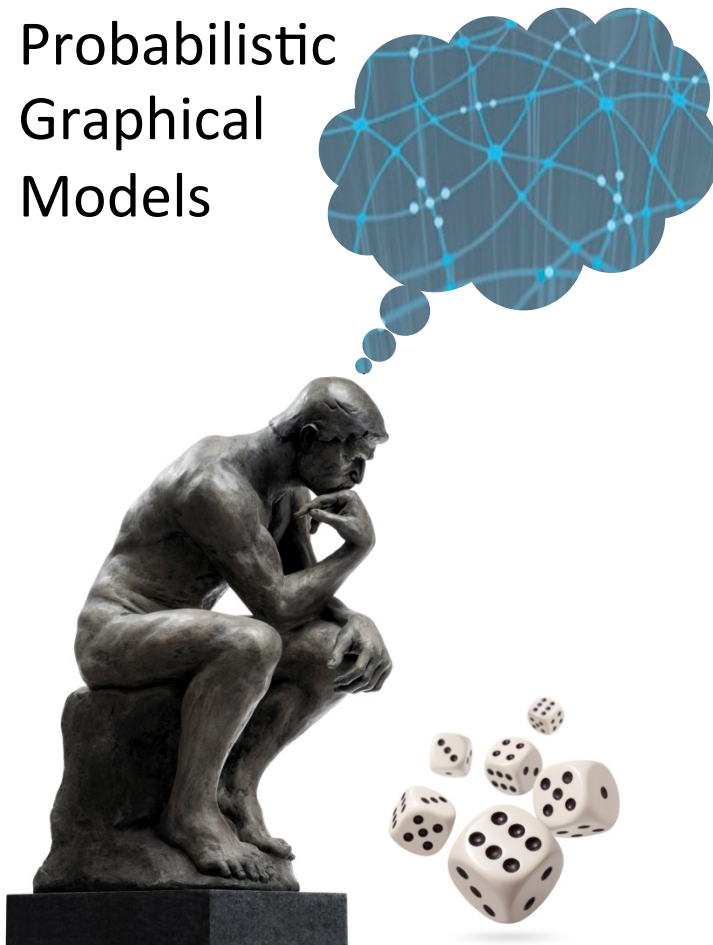


Probabilistic
Graphical
Models



Learning

BN Structures

Likelihood
Structure
Score

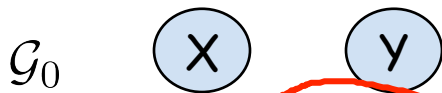
Likelihood Score

- Find (G, θ) that maximize the likelihood

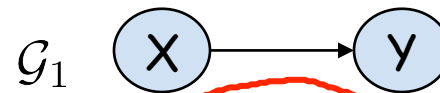
$$\underline{\text{score}_L(\mathcal{G} : \mathcal{D})} = \underline{\ell((\hat{\theta}, \mathcal{G}) : \mathcal{D})}$$

$\hat{\theta}$ = MLE of params. given \mathcal{G} and \mathcal{D}

Example



$$\text{score}_L(\mathcal{G}_0 : \mathcal{D}) = \sum_m (\log \hat{\theta}_{x[m]} + \log \hat{\theta}_{y[m]})$$



$$\text{score}_L(\mathcal{G}_1 : \mathcal{D}) = \sum_m (\log \hat{\theta}_{x[m]} + \log \hat{\theta}_{y[m]|x[m]})$$

$$\text{score}_L(\mathcal{G}_1 : \mathcal{D}) - \text{score}_L(\mathcal{G}_0 : \mathcal{D}) = \sum_m (\log \hat{\theta}_{y[m]|x[m]} - \log \hat{\theta}_{y[m]})$$

$$= \sum_{x,y} M[x,y] \log \hat{\theta}_{y|x} - \sum_y M[y] \log \hat{\theta}_y$$

\hat{P} = empirical distribution

$$M[x,y] = M \hat{P}[x,y]$$

$$= M \sum_{x,y} \hat{P}(x,y) \log \hat{P}(y|x) - M \sum_y \hat{P}(y) \log \hat{P}(y)$$

$$\sum_x \hat{P}(x,y) = \hat{P}(y)$$

$$= M \left(\sum_{x,y} \hat{P}(x,y) \log \hat{P}(y|x) - \sum_{x,y} \hat{P}(x,y) \log \hat{P}(y) \right)$$

$$= M \left(\sum_{x,y} \hat{P}(x,y) \log \frac{\hat{P}(x,y)}{\hat{P}(x)\hat{P}(y)} \right) = M \cdot \mathbf{I}_{\hat{P}}(X; Y)$$

mutual information

General Decomposition

- The Likelihood score decomposes as:

$$\text{score}_L(\mathcal{G} : \mathcal{D}) = M \sum_{i=1}^n \underbrace{\mathbf{I}_{\hat{P}}(X_i; \text{Pa}_{X_i}^{\mathcal{G}})}_{\text{mutual information}} - M \sum_i \underbrace{H_{\hat{P}}(X_i)}_{\text{independent of } \mathcal{G}}$$

$$\mathbf{I}_P(\mathbf{X}; \mathbf{Y}) = \sum_{\mathbf{x}, \mathbf{y}} P(\mathbf{x}, \mathbf{y}) \log \frac{P(\mathbf{x}, \mathbf{y})}{P(\mathbf{x})P(\mathbf{y})}$$

Score is higher if X_i is correlated with parents

$$H_P(\mathbf{X}) = - \sum_{\mathbf{x}} P(\mathbf{x}) \log P(\mathbf{x})$$

Limitations of Likelihood Score



$$\underline{\text{score}_L(\mathcal{G}_1 : \mathcal{D})} - \underline{\text{score}_L(\mathcal{G}_0 : \mathcal{D})} = \underline{MI_{\hat{P}}(X; Y)}$$

- Mutual information is always ≥ 0
- Equals 0 iff X, Y are independent
 - In empirical distribution \hat{P} $I_{\hat{P}}(X, Y) > 0$ almost always
- Adding edges can't hurt, and almost always helps
- Score maximized for fully connected network

Avoiding Overfitting

- Restricting the hypothesis space
 - restrict # of parents or # of parameters
- Scores that penalize complexity:
 - Explicitly ←
 - Bayesian score averages over all possible parameter values

Summary

- Likelihood score computes log-likelihood of D relative to G , using MLE parameters ℓ_G
 - Parameters optimized for D
- Nice information-theoretic interpretation in terms of (in)dependencies in G
- Guaranteed to overfit the training data (if we don't impose constraints)