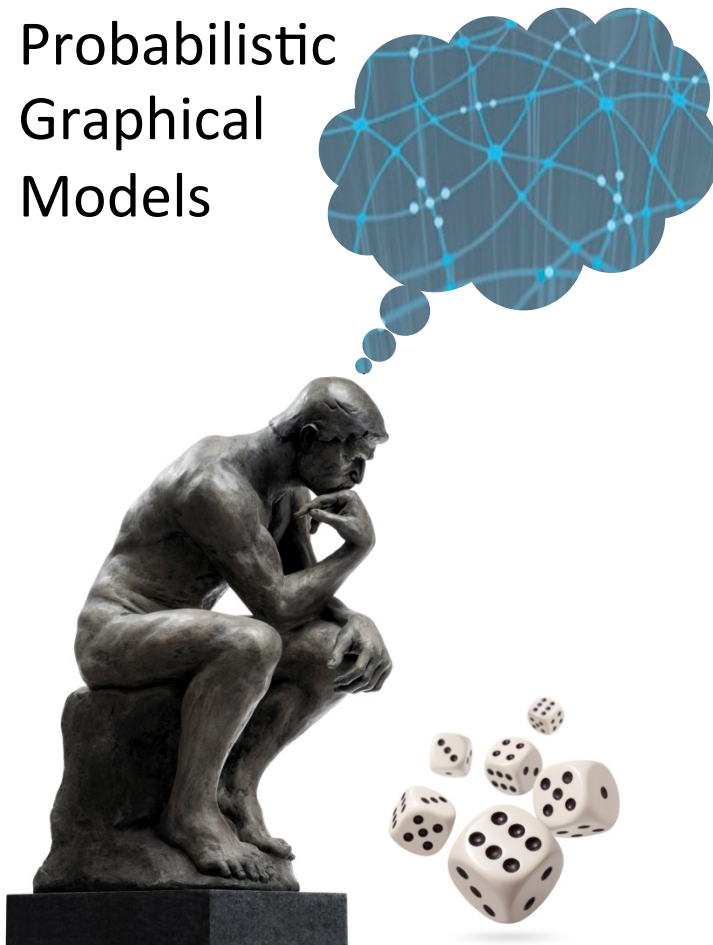


Probabilistic  
Graphical  
Models



Learning

---

BN Structure

---

# Bayesian Score

# Bayesian Score

Marginal likelihood

Prior over structures

$$\rightarrow \underbrace{P(\mathcal{G} \mid \mathcal{D})}_{\text{independent of } \mathcal{G}} = \frac{\underbrace{P(\mathcal{D} \mid \mathcal{G})}_{\text{independent of } \mathcal{G}} \underbrace{P(\mathcal{G})}_{\text{independent of } \mathcal{G}}}{\underbrace{P(\mathcal{D})}_{\text{independent of } \mathcal{G}}}$$

independent  
of  $\mathcal{G}$

Marginal probability of Data

$$\underline{\text{score}_B(\mathcal{G} : \mathcal{D})} = \underline{\log P(\mathcal{D} \mid \mathcal{G})} + \underline{\log P(\mathcal{G})}$$

# Marginal Likelihood of Data Given $\mathcal{G}$

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

Likelihood

$P(\mathcal{D} | \mathcal{G}, \theta_{\mathcal{G}})$   
Prior over parameters

$$\underline{P(\mathcal{D} | \mathcal{G})} = \int \underline{P(\mathcal{D} | \mathcal{G}, \theta_{\mathcal{G}})} \underline{P(\theta_{\mathcal{G}} | \mathcal{G})} d\theta_{\mathcal{G}}$$

# Marginal Likelihood Intuition

$$\underline{P(\mathcal{D} \mid \mathcal{G})} = \underline{P(x[1], \dots, x[M] \mid \mathcal{G})}$$

$$P(x[1] \mid \mathcal{G}) \leftarrow$$

$$P(x[2] \mid x[1], \mathcal{G}) \leftarrow$$

...

$$P(\underline{x[M]} \mid \underline{x[1], \dots, x[M-1]}, \mathcal{G}) \leftarrow$$

$\hat{\theta}_{\mathcal{G}}$  depends  
on all of  $\mathcal{D}$

# Marginal Likelihood: BayesNets

$$\underbrace{P(\mathcal{D} \mid \mathcal{G})}_{\text{variables}} = \prod_i \left( \prod_{\mathbf{u}_i \in \text{Val}(\text{Pa}_{X_i}^{\mathcal{G}})} \underbrace{\left[ \frac{\Gamma(\alpha_{X_i|\mathbf{u}_i})}{\Gamma(\alpha_{X_i|\mathbf{u}_i} + M[\mathbf{u}_i])} \right]}_{\substack{\text{prior} \\ \text{sum stats}}} \prod_{x_i^j \in \text{Val}(X_i)} \underbrace{\left[ \frac{\Gamma(\alpha_{x_i^j|\mathbf{u}_i} + M[x_i^j, \mathbf{u}_i])}{\Gamma(\alpha_{x_i^j|\mathbf{u}_i})} \right]}_{\substack{\text{sum stats} \\ \text{prior}}} \right)$$

$$\Gamma(x) = \int_0^{\infty} t^{x-1} e^{-t} dt$$

$$\Gamma(x) = x \cdot \Gamma(x-1)$$

# Marginal Likelihood Decomposition

$$P(\mathcal{D} \mid \mathcal{G}) = \prod_i \left( \prod_{\mathbf{u}_i \in \text{Val}(\mathbf{Pa}_{X_i}^{\mathcal{G}})} \frac{\Gamma(\alpha_{X_i|\mathbf{u}_i})}{\Gamma(\alpha_{X_i|\mathbf{u}_i} + M[\mathbf{u}_i])} \prod_{x_i^j \in \text{Val}(X_i)} \left[ \frac{\Gamma(\alpha_{x_i^j|\mathbf{u}_i} + M[x_i^j, \mathbf{u}_i])}{\Gamma(\alpha_{x_i^j|\mathbf{u}_i})} \right] \right),$$

$$\log P(\mathcal{D} \mid \mathcal{G}) = \sum_i \text{FamScore}_B(X_i \mid \mathbf{Pa}_{X_i}^{\mathcal{G}} : \mathcal{D})$$

# Structure Priors

$$\text{score}_B(\mathcal{G} : \mathcal{D}) = \log P(\mathcal{D} \mid \mathcal{G}) + \log P(\mathcal{G})$$

- Structure prior  $P(\mathcal{G})$ 
  - Uniform prior:  $P(\mathcal{G}) \propto \text{constant}$
  - Prior penalizing # of edges:  $P(\mathcal{G}) \propto c^{|\mathcal{G}|}$  ( $0 < c < 1$ )
  - Prior penalizing # of parameters
- Normalizing constant across networks is  $P(\mathcal{G})$  similar and can thus be ignored

# Parameter Priors

- Parameter prior  $P(\theta|G)$  is usually the BDe prior
  - $\alpha$ : equivalent sample size
  - $B_0$ : network representing prior probability of events
  - Set  $\alpha(x_i, pa_i^G) = \alpha P(x_i, pa_i^G | B_0)$ 
    - Note:  $pa_i^G$  are not the same as parents of  $X_i$  in  $B_0$
- A single network provides priors for all candidate networks
- Unique prior with the property that I-equivalent networks have the same Bayesian score



# BDe and BIC

- As  $M \rightarrow \infty$ , a network  $G$  with Dirichlet priors satisfies

$$\log P(\mathcal{D} \mid \mathcal{G}) = \underbrace{\ell(\hat{\theta}_{\mathcal{G}} : \mathcal{D})}_{\substack{\text{likelihood} \\ \text{score} \\ \text{of } \mathcal{D} \text{ given MLE } \hat{\theta}_{\mathcal{G}}}} - \underbrace{\frac{\log M}{2} \text{Dim}[\mathcal{G}]}_{\substack{\text{\#instances} \\ \text{\#independent} \\ \text{parameters}}} + \underbrace{O(1)}_{\substack{\text{constant relative} \\ \text{to } m - \text{ doesn't} \\ \text{grow with } m}}$$

$$\text{score}_{BIC}(\mathcal{G} : \mathcal{D}) = \ell(\hat{\theta}_{\mathcal{G}} : \mathcal{D}) - \frac{\log M}{2} \text{Dim}[\mathcal{G}]$$

as  $m \rightarrow \infty$  score is consistent

# Summary

- Bayesian score averages over parameters to avoid overfitting
- Most often instantiated as BDe
  - BDe requires assessing prior network
  - Can naturally incorporate prior knowledge
  - I-equivalent networks have same score
- Bayesian score
  - Asymptotically equivalent to BIC (as  $m \rightarrow \infty$ )
  - Asymptotically consistent learns correct network as  $m \rightarrow \infty$
  - But for small  $M$ , BIC tends to underfit