

2/6/14

1

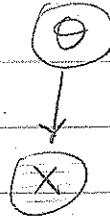
BN: PARAMETER SHARING & STRUCTURE LEARNING

① MetaBN (or Bayesian view of BN)

→ explicitly draw parameters as variables in the BN

e.g #1 $X \sim \text{Cat}(\vec{\theta})$
 $x \in \{1, \dots, k\}$

$$P(\vec{\theta}) = \text{Dir}(\vec{\alpha})$$



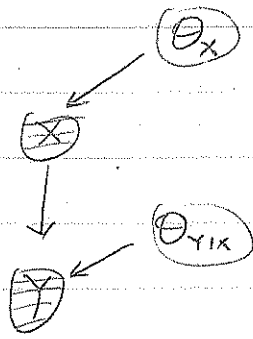
$$P(X|\vec{\theta}) = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_k \end{bmatrix}$$

$$\text{s.t. } \left. \begin{array}{l} \sum_{i=1}^k \theta_i = 1 \\ \theta_i \geq 0 \end{array} \right\} \text{ Simplex}$$

Here $\vec{\alpha}$ is a hyper parameter = fixed.

If we start modeling uncertainty over $\vec{\alpha}$ to start estimating $\vec{\alpha}$ from data, then we can put $\vec{\alpha}$ in the BN.

eg #2

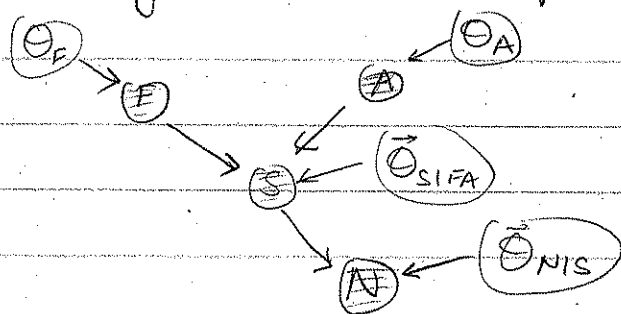


① Global Parameter Independence

- very common assumption
- all CPT parameters are indep.

$$P(\theta_x, \theta_{Y|X}) = P(\theta_x) P(\theta_{Y|X})$$

Revisiting an old example:



Assumption: $P(\vec{\Theta}) = P(\Theta_F) P(\Theta_A) P(\Theta_{SIFA}) P(\Theta_{NIS})$

$\Rightarrow \log P(\vec{\Theta}) = \log P(\Theta_F) + \log P(\Theta_A) + \log P(\Theta_{SIFA}) + \log P(\Theta_{NIS})$

fully factorized

$\log P(D|\Theta)$

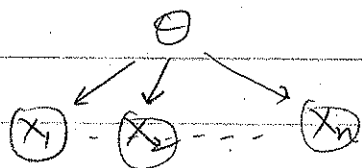
factorizes according to G

These two together lead to $\vec{\Theta}_{MAP}$ becoming separate problems at each variable.

② Parameter sharing

\rightarrow what if x_1, \dots, x_n are outcome of n tosses of SAME coin?

Appropriate BN:



③ Plate Notation



Variable inside plate is repeated n times
 Variable outside plate is not (It's shared)

→ There can be nested plates \equiv hierarchical Bayesian models



④ Structure Learning

Setup:

→ Random Vars: $\vec{X} = \{x_1, \dots, x_n\}$ (all categorical for now)

→ Fully observed dataset of samples from P^θ

$$\begin{bmatrix} x_1^{(1)} & \dots & x_n^{(1)} \\ \vdots & & \vdots \\ x_1^{(M)} & \dots & x_n^{(M)} \end{bmatrix}_{M \times n}$$

→ Goal: Estimate G (& associated θ CPTs)

→ Formulation: Score-based approach

$$\hat{G} = \operatorname{argmax}_{G \in \left\{ \begin{array}{l} \text{DAGs} \\ \text{on } n \\ \text{nodes} \end{array} \right\}} \text{Score}(G; D)$$

problem: # (undirected) graphs on n nodes = $2^{\binom{n}{2}} \approx 2^{O(n^2)}$
 # " trees " " " = $n^{n-2} \approx 2^{O(n \log n)}$
 # " spanning " " " = $n^{n-2} \approx 2^{O(n \log n)}$

(4.1) What's a good score?

How about the log-likelihood?

$$\text{Specifically } \text{Score}(G; D) = \log P(D | G, \hat{\theta}_{MLE})$$

A graph is "good" if it maximizes the chances of observing the dataset I observed (under MLE parameters)
(Very imp; Later we'll change this)

(4.2) e.g #1

G : (X) (Y)

D :

	X	Y
1	0	0
2	0	1
3	0	0
4	1	1
5	1	0
6	0	0

$M \times 2$

$$\text{Score}(G) = \sum_{j=1}^M \log P(X=x^{(j)}, Y=y^{(j)} | \hat{\theta}_{MLE}^G, G)$$

$$= \sum_{j=1}^M \left[\log P(X=x^{(j)} | \hat{\theta}_{MLE}^G, G) + \log P(Y=y^{(j)} | \hat{\theta}_{MLE}^G, G) \right]$$

$$= \sum_{x \in \{0,1\}} \text{Count}(x) \log P(x | \hat{\theta}_{MLE}^G) + (\quad)$$

$$= \sum_x \text{Count}(x) \log \hat{P}(x) + (\quad)$$

empirical distribution estimated from D
Why? $\because \hat{\theta}_{MLE}$!

$$= \sum_x M \cdot \hat{P}(x) \log \hat{P}(x) + (\cdot)$$

(3)

$$\Rightarrow \text{Score}(G) = M \sum_x \hat{P}(x) \log \hat{P}(x) + M \sum_y \hat{P}(y) \log \hat{P}(y)$$

$$= M \cdot \left[-H_{\hat{P}}(X) - H_{\hat{P}}(Y) \right]$$

Entropy of
 $\hat{P}(x)$

(4.3) e.g #2

 $G: (X) \rightarrow (Y)$

$$\text{Score}(G) = \sum_{j=1}^M \log P(X=x^{(j)}, Y=y^{(j)} | \hat{\Theta}_{\text{MLE}}, G)$$

$$= \sum_{j=1}^M \left[\log P(X=x^{(j)} | \hat{\Theta}_{\text{MLE}}, G) + \log P(Y=y^{(j)} | X=x^{(j)}, \hat{\Theta}_{\text{MLE}}, G) \right]$$

$$= -M \cdot H_{\hat{P}}(X) + \sum_x \sum_y \text{count}(x,y) \log P(Y=y | X=x, \hat{\Theta}_{\text{MLE}}, G)$$

$$= -M \cdot H_{\hat{P}}(X) + \sum_x \sum_y M \cdot \hat{P}(x,y) \log \hat{P}(y|x)$$

$$= -M H_{\hat{P}}(X) + M \sum_x \sum_y \hat{P}(x,y) \log \frac{\hat{P}(x,y)}{\hat{P}(x)} \cdot \frac{\hat{P}(y)}{\hat{P}(y)}$$

$$= -M \cdot H_{\hat{P}}(X) + M \sum_x \sum_y \hat{P}(x,y) \log \frac{\hat{P}(x,y)}{\hat{P}(x)\hat{P}(y)} + M \sum_x \sum_y \hat{P}(x,y) \log \hat{P}(y)$$

$$= M \cdot \left[I_{\hat{P}}(X, Y) - H_{\hat{P}}(X) - H_{\hat{P}}(Y) \right]$$

Mutual Information

≡ KL-divergence between $\underbrace{\hat{P}(X,Y)}_{\text{Joint}}$ & $\underbrace{\hat{P}(X)\hat{P}(Y)}_{\text{Factorized}}$

$$\equiv H(X) - H(X|Y)$$

$$\equiv H(Y) - H(Y|X)$$

In general

$$\text{Score}(G; D) = M \cdot \left[\sum_i I_{\beta}(X_i, P_{X_i}) - \sum_i H_{\beta}(X_i) \right]$$

Indep. of G
Constant!

→ Good news: score seems reasonable; high mutual info between X_i & P_{X_i} is good!

→ Unfortunately $I_{\beta}(X, Y) \geq 0$ [see HWO in ECE 4984/5984]

⇒ score is always non-negative

⇒ adding an edge never hurts!

⇒ Problem!