

# Parameter learning with EM for discrete BNs

Pablo M. Olmos, [olmos@tsc.uc3m.es](mailto:olmos@tsc.uc3m.es)

Course on Bayesian Networks, November 2016

# Index

Learning parameters with full observations

Learning with partial observations

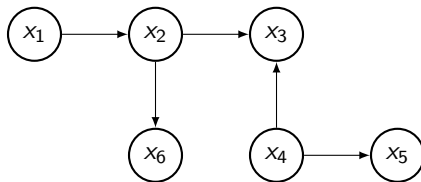
## Section 1

Learning parameters with full observations

# Motivation

- ▶ Let  $X_j \in \mathcal{X}$ ,  $j = 1, \dots, 5$ , be discrete R.V., where  $K \doteq |\mathcal{X}|$
- ▶ Consider the following BN:

$$p(\mathbf{x}) = p(x_1)p(x_2|x_1)p(x_4)p(x_3|x_2, x_4)p(x_5|x_4)p(x_6|x_2)$$



- ▶ **BN structure is known, CPD tables are unknown.**
- ▶ Our goal is to **estimate the CPD tables from  $N$  independent samples  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$** , drawn from  $p(\mathbf{x})$ .

# Log-Likelihood of data (I)

- ▶ Let  $\mathbf{x} \in \mathcal{X}^V$  be a discrete R.V. such that  $p(\mathbf{x}) = \prod_{t=1}^V p(x_t | \mathbf{x}_{\text{pa}(t)})$ .
- ▶ Consider  $N$  independent samples  $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$ . For each sample, we can write each  $p(x_t | \mathbf{x}_{\text{pa}(t)})$  term,  $t = 1, \dots, V$  as follows:

$$p(x_{it} | \mathbf{x}_{i,\text{pa}(t)}) = \prod_{c=1}^{K_{\text{pa}(t)}} \prod_{k=1}^{K_t} \theta_{tck}^{\mathbb{1}[x_{it}=k, \mathbf{x}_{i,\text{pa}(t)}=c]}$$

- ▶  $\prod_{c=1}^{K_{\text{pa}(t)}}$  → product over all possible values of  $\mathbf{x}_{i,\text{pa}(t)}$ .
- ▶  $\prod_{k=1}^K$  → product over all possible values of  $x_{it}$ .
- ▶  $\theta_{tck}^{\mathbb{1}[x_{it}=k, \mathbf{x}_{i,\text{pa}(t)}=c]}$  → equal to  $\theta_{tck}$  if  $x_{it} = k$ , and  $\mathbf{x}_{i,\text{pa}(t)} = c$ . Thus,

$$p(x_t = k | \mathbf{x}_{\text{pa}(t)} = c) \doteq \theta_{tck}$$

- ▶ Note that  $\sum_{k=1}^{K_t} \theta_{tck} = 1$ .

# Log-Likelihood of data (II)

- ▶ The log-likelihood of the complete data is given by

$$\log p(\mathcal{D}|\boldsymbol{\theta}) = \sum_{t=1}^V \sum_{c=1}^{K_{\text{pa}(t)}} \sum_{k=1}^{K_t} N_{tck} \log \theta_{tck}$$

where

$$N_{tck} = \sum_{i=1}^N \mathbb{1}[x_{it} = k, \mathbf{x}_{i,\text{pa}(t)} = c]$$

are the empirical counts.

- ▶ Log-likelihood is maximized if we take

$$\hat{\theta}_{tck} = \frac{N_{tck}}{\sum_{k'=1}^{K_t} N_{tck'}}$$

- ▶ **ML solution is very simple! We calculate frequencies!**

# MAP solution

- ▶ Dirichlet prior over  $\theta_{tc} = [\theta_{tc1}, \theta_{tc2}, \dots, \theta_{tcK_t}]$ :

$$\theta_{tc} \sim \text{Dir}[\alpha_{tc1}, \alpha_{tc2}, \dots, \alpha_{tcK_t}]$$

- ▶ It is easy to show that  $p(\theta_{tc}|\mathcal{D})$  is another Dirichlet distribution with parameters

$$\theta_{tc} \sim \text{Dir}[\alpha_{tc1} + N_{tc1}, \alpha_{tc2} + N_{tc2}, \dots, \alpha_{tcK_t} + N_{tcK_t}]$$

- ▶ The mean of  $\theta$  w.r.t. the posterior distribution is

$$\mathbb{E}_{p(\theta_{tc}|\mathcal{D})}[\theta_{tck}] = \frac{\alpha_{tck} + N_{tck}}{\sum_{k'=1}^{K_t} \alpha_{tck'} + N_{tck'}}$$

- ▶  $\theta_{tck}$  act as *pseudocounts*, avoiding to assign zero probability to unobserved outcomes in our data

## Section 2

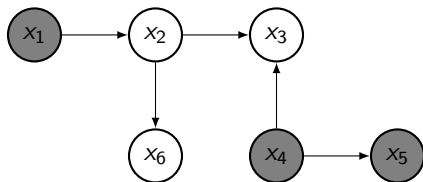
### Learning with partial observations



# Motivation

- ▶ Let  $X_j \in \mathcal{X}$ ,  $j = 1, \dots, 5$ , be discrete R.V., where  $K \doteq |\mathcal{X}|$
- ▶ Consider the following BN:

$$p(\mathbf{x}) = p(x_1)p(x_2|x_1)p(x_4)p(x_3|x_2, x_4)p(x_5|x_4)p(x_6|x_2)$$



- ▶ **BN structure is known, CPD tables are unknown.**
- ▶ Our goal is to **estimate the CPD tables from  $N$  independent samples  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$** , drawn from  $p(\mathbf{x})$ .
- ▶ Only a few elements of  $\mathbf{x}$  are observed!

# EM in a nutshell

- ▶  $p(\mathbf{y}_i, \mathbf{x}_i | \boldsymbol{\theta})$ , where  $\mathbf{y}_i$  is observed and  $\mathbf{x}_i$  is hidden,  $i = 1, \dots, N$ .
- ▶ Our goal is to estimate  $\boldsymbol{\theta}$  to maximize  $p(\mathcal{D} | \boldsymbol{\theta})$  (ML) or as the mode of the posterior distribution  $p(\boldsymbol{\theta} | \mathcal{D})$ . **Complex! We have to marginalize first over  $\mathbf{x}_i$ ,  $i = 1, \dots, N$ .**
- ▶ EM algorithm. Initialize  $\boldsymbol{\theta}$  to  $\boldsymbol{\theta}^0$ . For  $\ell = 1, 2, \dots$ 
  1. **E-step:**

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\ell-1}) &= \sum_{i=1}^N \int_{\mathbf{x}_i} p(\mathbf{x}_i | \mathbf{y}_i, \boldsymbol{\theta}^{\ell-1}) \log(p(\mathbf{x}_i, \mathbf{y}_i | \boldsymbol{\theta})) d\mathbf{x}_i \\ &= \sum_{i=1}^N \mathbb{E}_{p(\mathbf{x}_i | \mathbf{y}_i, \boldsymbol{\theta}^{\ell-1})} [\log(p(\mathbf{x}_i, \mathbf{y}_i | \boldsymbol{\theta}))] \end{aligned}$$

## 2. **M-step:**

$$\boldsymbol{\theta}^{\ell} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\ell-1}) \quad (\text{ML estimation})$$

$$\boldsymbol{\theta}^{\ell} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\ell-1}) + \log p(\boldsymbol{\theta}) \quad (\text{MAP estimation})$$

# EM for discrete BNs (I)

- ▶  $\mathbf{y}_i \rightarrow$  Set of observed variables for  $i$ -th data,  $i = 1, \dots, N$ .
- ▶ The log-likelihood of the complete data is given by

$$\log p(\{\mathbf{y}_i, \mathbf{x}_i\}_{i=1}^N \parallel \boldsymbol{\theta}) = \sum_{t=1}^V \sum_{c=1}^{K_{\text{pa}(t)}} \sum_{k=1}^{K_t} N_{tck} \log \theta_{tck}$$

where

$$N_{tck} = \sum_{i=1}^N \mathbb{1}[x_{it} = k, \mathbf{x}_{i,\text{pa}(t)} = \mathbf{c}]$$

## EM for discrete BNs (II)

► **E-step:**

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\ell-1}) = \sum_{i=1}^N \sum_{c=1}^{K_{\text{pa}(t)}} \sum_{k=1}^{K_t} \widetilde{N}_{tck} \log \theta_{tck}$$

where

$$\begin{aligned} \widetilde{N}_{tck} &= \sum_{i=1}^N \mathbb{E} [\mathbb{1}[x_{it} = k, \mathbf{x}_{i,\text{pa}(t)} = c]] \\ &= \sum_{i=1}^N p(x_{it} = k, \mathbf{x}_{i,\text{pa}(t)} = c | \mathbf{y}_i, \boldsymbol{\theta}^{\ell-1}) \end{aligned}$$

- Thus, given each pair  $(\mathbf{x}_i, \mathbf{y}_i)$ , we simply have to compute the marginal joint probabilities  $p(x_{it} = k, \mathbf{x}_{i,\text{pa}(t)} = c | \mathbf{y}_i)$ ,  $t = 1, \dots, V$ .
- We will use **Belief Propagation for this task!**

# EM for discrete BNs (III)

► **M-step:**

$$\hat{\theta}_{tck} = \frac{\widetilde{N}_{tck}}{\sum_{k'=1}^{K_t} \widetilde{N}_{tck'}} \quad (\text{ML estimation})$$

$$\hat{\theta}_{tck} = \frac{\alpha_{tck} + \widetilde{N}_{tck}}{\sum_{k'=1}^{K_t} \alpha_{tck'} + \widetilde{N}_{tck'}} \quad (\text{MAP estimation})$$