

# The EM algorithm for Bayesian networks

**Aarhus University, Fall 2003, Lecture 8**

Steffen L. Lauritzen, Aalborg University

# Entropy

The *entropy* of a discrete probability distribution  $P$  is

$$\text{Ent}(P) = - \sum_x p(x) \log p(x).$$

Entropy is a measure of *spread* of the distribution and it is always positive.

The entropy is never larger than the entropy of the uniform distribution:

Let  $P_u(x) = 1/|\mathcal{X}|$ , then it holds that

$$0 \leq \text{Ent}(P) \leq \text{Ent}(P_u) = \log |\mathcal{X}|.$$

Proof on next overhead.

## Uniform distribution has maximal entropy

The information inequality

$$\sum_x p(x) \log p(x) \geq \sum_x p(x) \log q(x)$$

yields

$$\begin{aligned} \text{Ent}(P) &= - \sum_x p(x) \log p(x) \\ &\leq - \sum_x p(x) \log \frac{1}{|\mathcal{X}|} \\ &= - \sum_x \frac{1}{|\mathcal{X}|} \log \frac{1}{|\mathcal{X}|} = \text{Ent}(P_u). \end{aligned}$$

## Kullback-Leibler divergence

The *KL divergence* between  $P$  and  $Q$  is

$$KL(P : Q) = \sum_x p(x) \log \frac{p(x)}{q(x)}.$$

Also known as *relative entropy* of  $Q$  with respect to  $P$ .

Information inequality says that

$KL(P : Q) \geq 0$  and  $KL(P : Q) = 0$  if and only if  $P = Q$ ,

so KL divergence defines an (asymmetric) distance measure between probability distributions.

## Incomplete observations

Bayesian network with conditional probability distributions only partially known:

$$p(x) = \prod_{v \in V} p(x_v \mid x_{\text{pa}(v)}, \theta)$$

where  $\theta \in \Theta \subseteq \mathcal{R}^k$  is unknown parameter.

Instead of *complete data*  $(x^1, \dots, x^n)$ , only *incomplete data*  $(x_{A_1}^1, \dots, x_{A_n}^n)$  available, where  $A_i \subseteq V$ .

Example: paternity cases. Unknown parameters: gene frequencies, probability of paternity, possibly mutation rates.

## EM algorithm

Complete data  $x$ , incomplete data (observed)  $y = g(x)$ .

*Complete data log-likelihood:*

$$l(\theta) = \log L(x | \theta) = \log p(x | \theta).$$

The *marginal log-likelihood* is

$$l_y(\theta) = \log L(\theta | y) = \log p(y | \theta).$$

Wish to maximize  $l_y$  in  $\theta$  but  $l_y$  is unpleasant:

$$l_y(\theta) = \log \sum_{x:g(x)=y} p(x | \theta).$$

However, we assume that we know how to maximize  $l$ .  
How can this be exploited?

## E-step and M-step

We let  $\theta^*$  be arbitrary but fixed.

The **E-step** calculates *expected complete data log-likelihood*  $q(\theta | \theta^*)$ :

$$q(\theta | \theta^*) = \mathbf{E}_{\theta^*} \{l(\theta) | y\} = \sum_{x:g(x)=y} p(x | y, \theta^*) \log p(x | \theta).$$

The **M-step** maximizes  $q(\cdot | \theta^*)$  for fixed  $\theta^*$ :

The algorithm alternates between an E-step and an M-step.

*After an E-step and subsequent M-step, the likelihood function has never decreased, as we shall now show.*

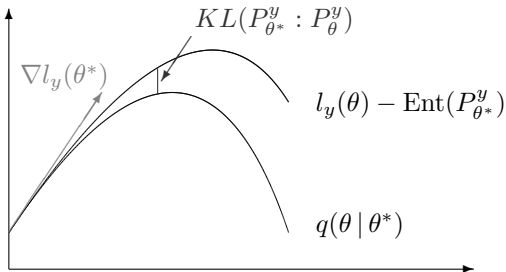
## EM algorithm

Since  $p(x | y, \theta) = \chi_{g(x)}(y)p(x | \theta)/p(y | \theta)$  we have

$$\begin{aligned}q(\theta | \theta^*) &= \sum_x p(x | y, \theta^*) \log\{p(y | \theta)p(x | y, \theta)\} \\&= \log p(y | \theta) + \sum_x p(x | y, \theta^*) \log p(x | y, \theta) \\&= l_y(\theta) - \sum_x p(x | y, \theta^*) \log p(x | y, \theta^*) \\&\quad - \sum_x p(x | y, \theta^*) \log \frac{p(x | y, \theta^*)}{p(x | y, \theta)} \\&= l_y(\theta) - \text{Ent } P_{\theta^*}^y - KL(P_{\theta^*}^y : P_{\theta}^y).\end{aligned}$$



## Expected and complete data likelihood



$$l_y(\theta) - \text{Ent}(P_{\theta^*}^y) = q(\theta | \theta^*) + KL(P_{\theta^*}^y : P_{\theta}^y)$$

$$\nabla l_y(\theta^*) = \nabla q(\theta^* | \theta^*)$$

## Likelihood monotonicity of EM algorithm

Let  $\theta_0 = \theta^*$  and  $\theta_{n+1} = \arg \max_{\theta} q(\theta | \theta_n)$ .

Then

$$\begin{aligned} l_y(\theta_{n+1}) &= q(\theta_{n+1} | \theta_n) + \text{Ent}(P_{\theta_n}^y) + KL(P_{\theta_{n+1}}^y : P_{\theta_n}^y) \\ &\geq q(\theta_n | \theta_n) + \text{Ent}(P_{\theta_n}^y) = l_y(\theta_n). \end{aligned}$$

So likelihood never decreases. Note, this also holds if just  $q(\theta_{n+1} | \theta_n) \geq q(\theta_n | \theta_n)$ .

## E-step for Bayesian networks

The complete data likelihood is

$$\log p(x | \theta) = \sum_{i=1}^n \log p(x^i | \theta) = \sum_x n(x) \log p(x | \theta).$$

where  $n(x) = \#\{i : x^i = x\}$ . Using factorization we get

$$\begin{aligned} \log p(x | \theta) &= \sum_x \sum_v n(x) \log p(x_v | x_{\text{pa}(v)}, \theta) \\ &= \sum_v \sum_{x_{v \cup \text{pa}(v)}} n(x_{v \cup \text{pa}(v)}) \log p(x_v | x_{\text{pa}(v)}, \theta), \end{aligned}$$

with  $n(x_{v \cup \text{pa}(v)}) = \#\{i : x_{v \cup \text{pa}(v)}^i = x_{v \cup \text{pa}(v)}\}$ . So E-step equivalent to computing

$$n^*(x_{v \cup \text{pa}(v)}) = \mathbf{E}\{N(x_{v \cup \text{pa}(v)}) | y, \theta^*\}.$$

## Computing expected counts

We now get

$$\begin{aligned}n^*(x_{v \cup \text{pa}(v)}) &= \mathbf{E}\{N(x_{v \cup \text{pa}(v)}) \mid y, \theta^*\} \\&= \sum_i \mathbf{E}\{\chi_{x_{v \cup \text{pa}(v)}}(x_{v \cup \text{pa}(v)}^i) \mid y, \theta^*\} \\&= \sum_i \mathbf{E}\{\chi_{x_{v \cup \text{pa}(v)}}(x_{v \cup \text{pa}(v)}^i) \mid x_{A_i}^i, \theta^*\} \\&= \sum_i p(x_{v \cup \text{pa}(v)} \mid x_{A_i}^i, \theta^*).\end{aligned}$$

Each of the latter terms can be calculated by probability propagation as can the marginal likelihood function:

$$\log p(y \mid \theta) = \sum_i \log p(x_{A_i}^i \mid \theta).$$

## M-step for Bayesian networks

Note the similarity between the complete data likelihood and  $q$ :

$$\log p(x | \theta) = \sum_v \sum_{x_{v \cup \text{pa}(v)}} n(x_{v \cup \text{pa}(v)}) \log p(x_v | x_{\text{pa}(v)}, \theta)$$

whereas

$$q(\theta | \theta^*) = \sum_v \sum_{x_{v \cup \text{pa}(v)}} n^*(x_{v \cup \text{pa}(v)}) \log p(x_v | x_{\text{pa}(v)}, \theta).$$

So any algorithm which maximizes the complete data likelihood can be used to maximize  $q$  in the M-step.