# LEARNING BAYESIAN NETWORKS

$$p(x) = \prod_i p(x_i \mid pa(x_i))$$

$p(x_i \mid pa(x_i), \vartheta)$, estimate by ML $\vartheta$ from data

$x_1 \bigcirc \qquad \bigcirc x_2$

$\bigcirc x_3$

$$P(x_1, x_2, x_3) = \underline{p(x_3 \mid x_1, x_2)} \, p(x_1) \, p(x_2)$$

$$x_i \in \{0, 1\}$$

$$P(x_3 = 1 \mid x_1 = 0, x_2 = 0) = \vartheta_{00} \longrightarrow \quad \frac{\#(x_1 = 0, x_2 = 0, x_3 = 1)}{\#(x_1 = 0, x_2 = 0)} \quad ML.$$

$$P(x_3 = 1 \mid x_1 = 0, x_2 = 1) = \vartheta_{01}$$

$$P(x_3 = 1 \mid x_1 = 1, x_2 = 0) = \vartheta_{10}$$

$$P(x_3 = 1 \mid x_1 = 1, x_2 = 1) = \vartheta_{11}$$

What if some variables are LATENT (not-observed)?

$X$ = observed

$Z$ = latent

$p(x, z \mid \vartheta)$ and we want to estimate $\vartheta$ by ML.

$$p(x \mid \vartheta) = \sum_z p(x, z \mid \vartheta)$$

$x_1, \ldots, x_N$ observation

$z_1, \ldots, z_N$ latent states of observations

$$p(\underline{x}, \underline{z} \mid \vartheta) = \prod_n p(x_n, z_n \mid \vartheta) \qquad \log p(\underline{x}, \underline{z}) = \sum_n p(x_n, z_n)$$

$$\log p(\underline{x} \mid \vartheta) \neq \sum \log p(x_n \mid \vartheta)$$

# ELBO

$P(x,z)$

$x$ observed $\quad x_1 \ldots x_N$

$z$ latent $\quad z_1 \ldots z_N$

$P(x \mid \vartheta) = \sum\limits_{z} P(x,z \mid \vartheta)$

$\swarrow$ function to optimize.

$\operatorname*{argmax}\limits_{\vartheta} \; P(\underline{x} \mid \vartheta)$

$P(x,z \mid \vartheta) = P(x \mid \vartheta) \, P(z \mid x, \vartheta)$

VARIATIONAL APPROXMATION $q(z)$ of $P(z \mid x, \vartheta)$

$$KL\big[q \| P\big] = KL\big[q(z) \| P(z \mid x, \vartheta)\big] = \mathbb{E}_{q(z)}\left[- \log \frac{P(z \mid x, \vartheta)}{q(z)}\right]$$

$$= -\sum_{z} q(z) \left[\log \frac{P(z \mid x, \vartheta)}{q(z)} + \log P(x \mid \vartheta) - \log P(x \mid \vartheta)\right]$$

$$= -\sum_{z} q(z) \log \frac{P(x,z \mid \vartheta)}{q(z)} + \log P(x \mid \vartheta)$$

$$\mathcal{L}(q, \vartheta) = \sum_{z} q(z) \cdot \log \frac{P(x,z \mid \vartheta)}{q(z)}$$

$$\log p(x|\vartheta) = \mathcal{L}(q, \vartheta) + KL[q(z) \| p(z|x, \vartheta)]$$

$$\mathcal{L}(q, \vartheta) = \sum_z q(z) \cdot \log \frac{p(x, z|\vartheta)}{q(z)}$$

$\uparrow$

EVIDENCE LOWER BOUND (ELBO)

$$KL[q\|p] \geq 0 \implies \mathcal{L}(q, \vartheta) \leq \log(x|\vartheta)$$

variational distribution

# EXPECTATION MAXIMIZATION

$$\log p(\underline{x}|\vartheta) = \mathcal{L}(q,\vartheta) + KL[q(z) \| p(z|\underline{x},\vartheta)]$$

$$\mathcal{L}(q,\vartheta) = \underbrace{\mathbb{E}_q[\log p(\underline{x},z|\vartheta)]}_{\text{energy}} + \underbrace{\mathbb{E}_q[-\log q(z)]}_{H(q) \text{ entropy}} = \mathbb{E}_q\left[\log \frac{p(x,z|\vartheta)}{q(z)}\right]$$

GOAL $\vartheta_{ML} = \arg\max\limits_{\vartheta} \log p(\underline{x}|\vartheta)$ , intractable. THEN max $\mathcal{L}(q,\vartheta)$

E-step: maximise $\mathcal{L}(q,\vartheta)$ wrt $q(z)$, with $\vartheta$ FIXED to $\vartheta_{old}$

$q$ is maximized iff $KL[q \| p(z|x,\vartheta_{old})] = 0$

iff $q_{new}(z) = q(z) = p(\underline{z}|\underline{x},\vartheta_{old})$ ⬅

then compute $\mathbb{E}_{q_{new}}[\log p(\underline{x},\bar{z}|\vartheta)]$

$$P(\underline{z}|\underline{x},\vartheta) = \prod_{i=1}^{N} p(z_i|x_i,\vartheta)$$

$\text{\textit{s}} \ x_1 \ldots x_N$ are iid

$$P(\underline{z}|\underline{x},\vartheta) = \frac{\prod p(z_i,x_i|\vartheta)}{\prod \sum\limits_{z_i} p(z_i,x_i|\vartheta)}$$

$$P(\underline{z}|\underline{x},\vartheta) = \frac{p(\underline{x},\underline{z}|\vartheta)}{\sum\limits_{z} p(\underline{x},\underline{z}|\vartheta)} = \frac{\prod p(z_i,x_i|\vartheta)}{\sum\limits_{z} \prod p(z_i,x_i|\vartheta)}$$

M-step    maximize $\mathcal{L}(q, \theta)$ keeping $q$ fixed to $q_{new} = p(z \mid x, \theta_{old})$
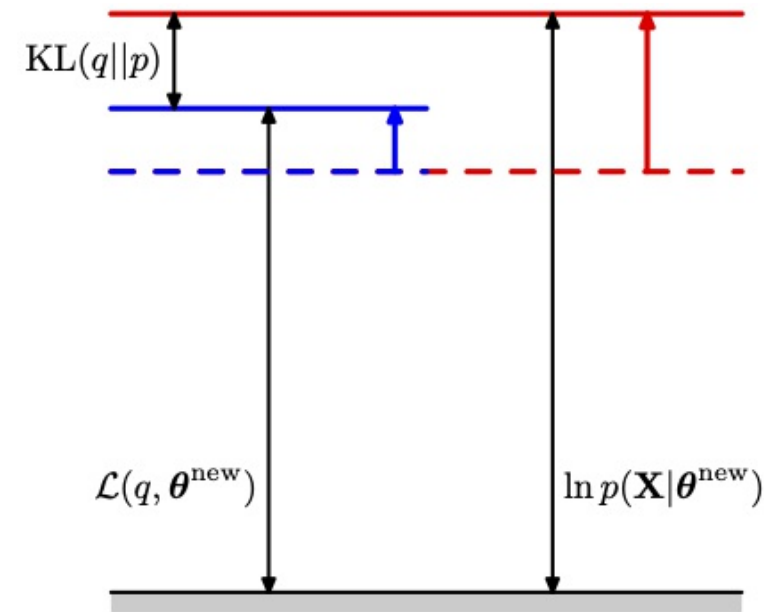
$\uparrow\downarrow$

maximize $\mathbb{E}_{q_{new}}[\log p(z, x \mid \theta)]$

Then we have $\theta_{new} = \arg\max_{\theta} \mathbb{E}_{q_{new}}[\log p(x, z \mid \theta)]$

Iterate E and M steps until convergence (when log-like or $\|\theta_{old} - \theta_{new}\| \leq \varepsilon$)
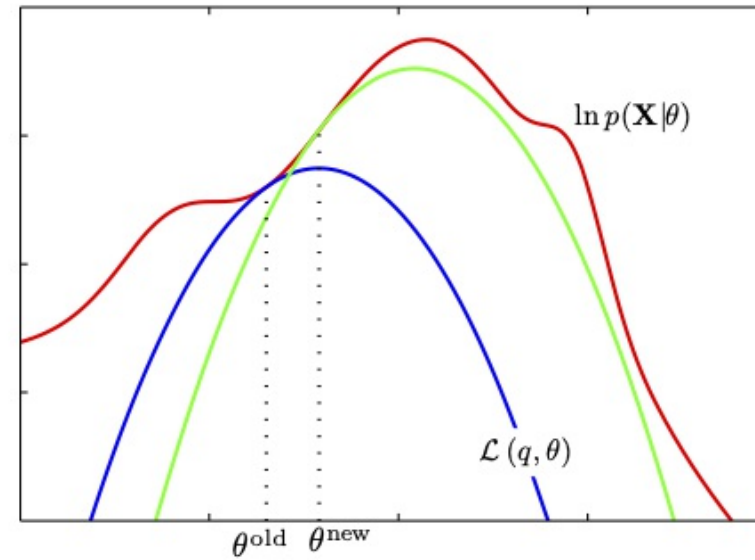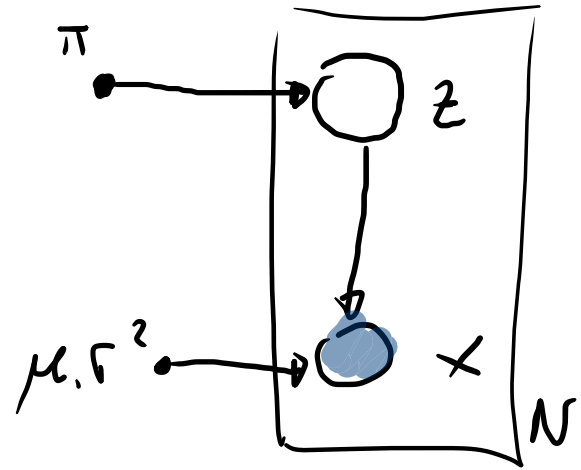


E-STEP

M-STEP

**Figure 9.14** The EM algorithm involves alternately computing a lower bound on the log likelihood for the current parameter values and then maximizing this bound to obtain the new parameter values. See the text for a full discussion.



$\ln p(\mathbf{X}|\theta)$

$\mathcal{L}(q, \theta)$

$\theta^{\mathrm{old}}$  $\theta^{\mathrm{new}}$

EM converges to
a local optimum of
$\log p(\mathbf{x}|\theta)$

# MIXTURE OF GAUSSIANS



$z$ disrete : $z_1 - z_K$

$z_j \in \{0,1\}$, $\sum_\delta z_\delta = 1$

$\theta = (\pi, \mu, \sigma^2)$

$\underline{z} = (z_{i j})$, $i = 1 \to N$, $j = 1 \dots K$   not observed

$\underline{x} = x_i$ , $i = 1 \dots N$   observed

$$P(x, z \mid \theta) = \prod_{j=1}^{K} \pi_\delta^{z_\delta} \cdot \mathcal{N}(x \mid \mu_j, \sigma_\delta^2)^{z_j}$$

$$P(x \mid \theta) = \sum_{j=1}^{K} \pi_\delta \, \mathcal{N}(x \mid \mu_\delta, \sigma_\delta^2)$$

$$P(z = j \mid x, \theta) = \frac{\pi_\delta \, \mathcal{N}(x \mid \mu_j, \sigma_\delta^2)}{\sum_{i=1}^{K} \pi_i \, \mathcal{N}(x \mid \mu_i, \sigma_i^2)}$$

$$P(z \mid \theta) = \prod_\delta \pi_1^{z_\delta}$$

$$P(\underline{z} \mid \underline{x}, \theta) \propto \prod_{n=1}^{N} \prod_{j=1}^{K} \pi_\delta^{z_{n j}} \, \mathcal{N}(x_n \mid \mu_\delta, \sigma_\delta^2)^{z_{n j}}$$

$$\mathbb{E}_{p(z|x)}\left[\bar{z}_{nj}\right] = P(z_n = j \mid x_n, \vartheta) = \frac{\pi_j \mathcal{N}(x_n \mid \mu_0, \sigma_j^2)}{\sum_i \pi_i \mathcal{N}(x_n \mid \mu_i, \sigma_i^2)} = \underbrace{\gamma(z_{nj})}_{\text{RESPONSIBILITY}}$$

$$\log p(\underline{x}, \underline{z} \mid \vartheta) = \sum_{n=1}^{N} \sum_{j=1}^{K} z_{nj} \left[ \log \pi_j + \log \mathcal{N}(x_n \mid \mu_0, \sigma_j^2) \right]$$

$$\mathbb{E}_{p(\underline{z}|x,\vartheta)}\left[\log p(\underline{x}, \underline{z} \mid \vartheta)\right] = \sum_{n=1}^{N} \sum_{j=1}^{K} \underbrace{\mathbb{E}[z_{nj}]}_{\gamma(z_{nj})} \left[ \log \pi_\vartheta + \log \mathcal{N}(x_n \mid \mu_0, \sigma_0^2) \right] \quad \leftarrow \text{E step.}$$

$$\begin{cases} \mu_j^{\text{new}} = \frac{1}{N_0} \sum_n \gamma(z_{nj}) x_n \\[2mm] \overline{\Sigma}_j^{\text{new}} = \frac{1}{N_j} \sum_n \gamma(z_{nj}) (x_n - \mu_j^{\text{new}})^T (x_n - \mu_0^{\text{new}}) \quad \text{M step} \\[2mm] \pi_j^{\text{new}} = N_j / N \quad, \quad N_j = \sum_{n=1}^{N} \gamma(z_{nj}) \quad, \quad \sum_j N_j = N \end{cases}$$

# EM FOR BAYESIAN NETWORKS

$$p(x) = \prod_i p(x_i \mid pa(x_i), \vartheta_i) \qquad x = (v, z), \quad \vartheta = (\vartheta_i)_{i=1 \dots m}$$

$$p(x) = p(v, z \mid \vartheta)$$

$$\rightsquigarrow p(z \mid v = \hat{v}, \vartheta) \quad \text{for fixed } \vartheta$$

$$\underline{v} = (v_1, \dots v_N) \text{ observations of } v \qquad \overset{E\text{-step}}{\swarrow}$$

$$q^n(z) = p(z \mid v_n, \vartheta) \quad \rightsquigarrow \quad \underline{q^n(x)} = p(z \mid v_n, \vartheta) \, \delta(v, v_n)$$

ENERGY for M-step

$$\sum_n \mathbb{E}_{q^n}\left[ \log p(v_n, z_n \mid \vartheta) \right] = \sum_n \sum_i \mathbb{E}_{q^n}\left[ \log p(x_i^n \mid pa(x_i^n), \vartheta_i) \right]$$

then optimize $\displaystyle\sum_n \mathbb{E}_{q^n}\left[ \log p(x_i \mid pa(x_i), \vartheta_i) \right]$ over $\vartheta_i$ for each $i$

$t \bigcirc \longrightarrow \bigcirc v$     $, t, v, w \in \{0, 1\}$     $P(t=1) = \vartheta_t$     $x = (v, w, z)$

$\qquad \searrow \qquad \swarrow$

$\qquad \bigcirc w$     $P(v=1) = \vartheta_v$

$\qquad P(w = 1 \mid t = a, v = b) = \vartheta_{wab}$   $, a, b \in \{0, 1\}$

$(v_1, w_1), \longrightarrow (v_N, w_N)$ observations

E-step

$q^n(z) = P(z \mid v = v_n, w = w_n, \vartheta)$     $q^n(x) = P(t \mid v = v_n, w = w_n, \vartheta) \, \delta(v, v_n) \delta(w, w_n)$

$$\sum_n \mathbb{E}_{q^n} \left[ \log \underbrace{P(z^n \mid \vartheta_z)}_{\vartheta_z \;\&\; z^n = 1} \right] = \boxed{\sum_n} \log \vartheta_z \, q^n(z=1) + \log(1 - \vartheta_z) \cdot q^n(z=0)$$

$$\vartheta_z = \frac{\sum_n q^n(z=1)}{\sum_n q^n(z=1) + \sum_n q^n(z=0)} = \frac{1}{N} \sum_n q^n(z=1)$$
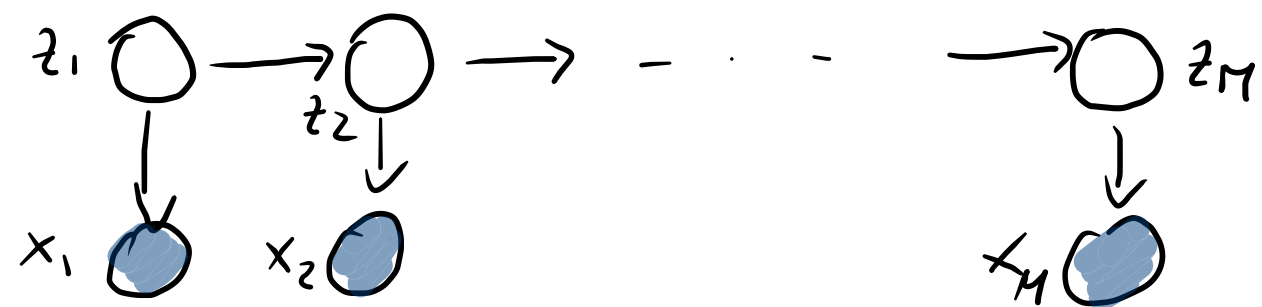
$$\sum_n \mathbb{E}_{q^n}\left[\log p(w_n \mid z, v_n, \vartheta_w)\right]$$

for $z=0$, $v=1$ $(\vartheta_{w01})$:

$$\sum_{n:\, w_n=1,\, v_n=1} q^n(z=0)\cdot \log \vartheta_{w01} + \sum_{n:\, w_n=0,\, v_n=1} q^n(z=0) \log\left(1-\vartheta_{w01}\right)$$

$$\Rightarrow \vartheta_{w01} = \frac{\sum_n' \mathbb{I}(w_n=1)\,\mathbb{I}(v_n=1)\,q^n(z=0)}{\sum_n \mathbb{I}(w_n=1)\,\mathbb{I}(v_n=1)\,q^n(z=0) + \sum_n \mathbb{I}(w_n=0)\,\mathbb{I}(v_n=1)\,q^n(z=0)}$$

# EM FOR HMM (Baum-Welch)

$z \in \{1 \to K\}$

$P(z_1 = i) = \pi_i$

$P(z_i = j \mid z_{i-1} = k) = A_{kj}$

$P(x_i \mid z_1 = k) = p(x_i \mid \phi_k)$

$\vartheta = (\pi, A, \phi)$

$$\boxed{\pi_k = \frac{\sum_n q^n(z_{1k})}{\sum_j \sum_n q^n(z_{1j})}} \quad A$$

$$\underline{x} = x^1 \, , \ldots \, x^N$$

$$(x_1^1 \, , \ldots \, x_M^1)$$

E step $\quad q^n(z) = \underbrace{P(z \mid x^n, \vartheta)}_{\text{message passing.}} \; \forall n$

M step $\quad E(\vartheta) = \sum_{n=1}^{N} \left[ \sum_{k=1}^{K} q^n(z_{1k}) \ln \pi_k + \sum_{i=2}^{M} \sum_{j,k=1}^{K} q^n(z_{i-1j}, z_{ik}) \ln A_{jk} \right.$

$$\left. + \sum_{i=1}^{M} \sum_{k=1}^{K} q^n(z_{ik}) \ln p(x_i^n \mid \phi_k) \right]$$