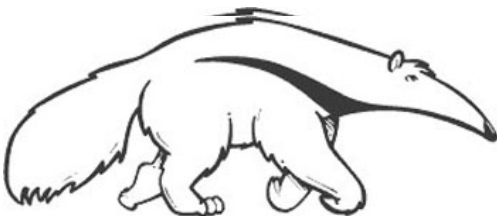


Latent Variable Models

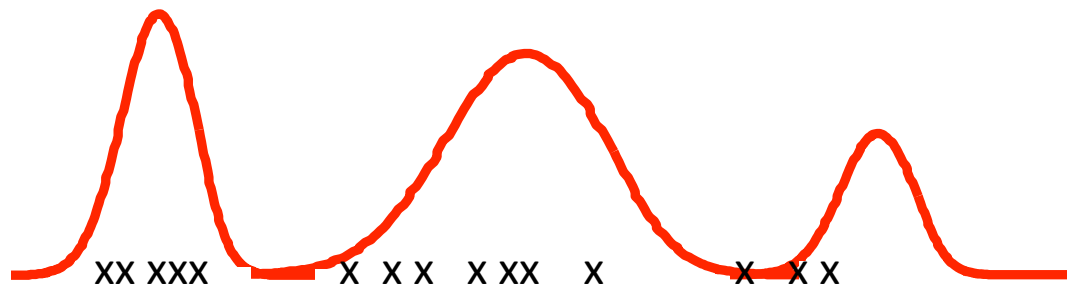
Learning in Graphical Models

Prof. Alexander Ihler



Latent variable models

- Never observe some variables?
- Ex: Gaussian mixture models
- Probability distribution: $p(x) = \sum_c \pi_c \mathcal{N}(x ; \mu_c, \sigma_c)$



Latent variable models

- Never observe some variables?
- Ex: Gaussian mixture models
- Probability distribution: $p(x) = \sum_c \pi_c \mathcal{N}(x ; \mu_c, \sigma_c)$
- Equivalent “latent variable” form:

$$p(z = c) = \pi_c$$

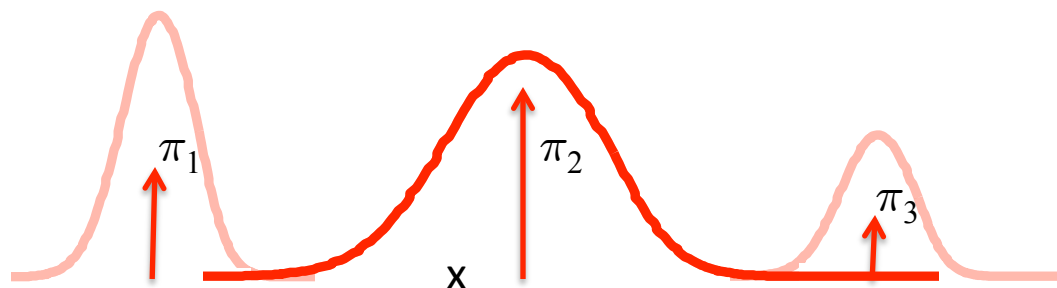
Select a mixture component with probability π

$$p(x|z = c) = \mathcal{N}(x ; \mu_c, \sigma_c)$$

Sample from that component's Gaussian

“Latent assignment” z :
we observe x , but z is hidden

$p(x)$ = marginal over x



Learning mixture models

- Maximum likelihood? $p(x) = \sum_c \pi_c \mathcal{N}(x ; \mu_c, \sigma_c)$

- Observe iid samples $D = \{x^{(1)} \dots x^{(m)}\}$

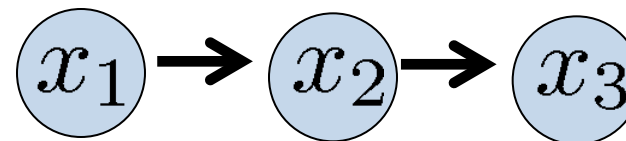
$$\mathcal{L}_X(\theta = \{\pi, \mu, \sigma^2\}) = \sum_j \log \left[\sum_c \pi_c h(\sigma_c) \exp(-(x^{(j)} - \mu_c)^2 / \sigma_c^2) \right]$$

- Gradient descent?

$$\frac{\partial \mathcal{L}_X(\theta)}{\partial \pi_c} = \dots \quad \frac{\partial \mathcal{L}_X(\theta)}{\partial \mu_c} = \dots \quad \frac{\partial \mathcal{L}_X(\theta)}{\partial \sigma_c} = \dots$$

- This can be very slow...

Recall ML for Bayes Nets



- Fully observed:

Likelihood decomposes; closed form MLE

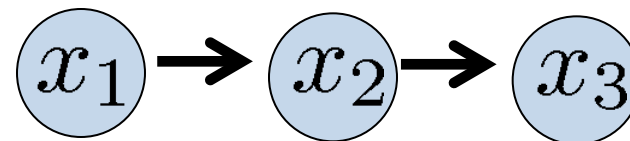
$$\mathcal{L} = \sum_i \log [p(x_1^i) p(x_2^i | x_1^i) p(x_3^i | x_2^i)]$$

$$D = \begin{matrix} & x_1 & x_2 & x_3 \\ \begin{bmatrix} x^1 \\ x^2 \\ x^3 \\ x^4 \\ \vdots \end{bmatrix} & = & \begin{bmatrix} x_1^1 & x_2^1 & x_3^1 \\ x_1^2 & x_2^2 & x_3^2 \\ x_1^3 & x_2^3 & x_3^3 \\ x_1^4 & x_2^4 & x_3^4 \\ \vdots & \vdots & \vdots \end{bmatrix} \end{matrix}$$

What if...

- Not fully observed?

$$D = \begin{bmatrix} x^1 \\ x^2 \\ x^3 \\ x^4 \\ \vdots \end{bmatrix} = \begin{array}{ccc} & x_1 & x_2 & x_3 \\ \begin{bmatrix} x^1 \\ x^2 \\ x^3 \\ x^4 \\ \vdots \end{bmatrix} & \begin{bmatrix} x_1^1 \\ x_1^2 \\ ? \\ x_1^4 \\ \vdots \end{bmatrix} & \begin{bmatrix} ? \\ x_2^2 \\ x_2^3 \\ ? \\ \vdots \end{bmatrix} & \begin{bmatrix} x_3^1 \\ x_3^2 \\ x_3^3 \\ ? \\ \vdots \end{bmatrix} \end{array}$$

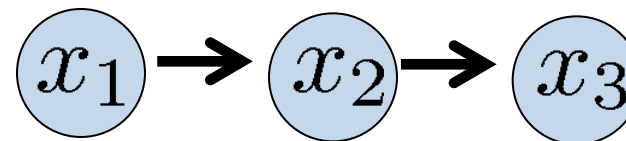


$$\mathcal{L} = \sum_i \log [p(x_1^i) p(x_2^i | x_1^i) p(x_3^i | x_2^i)]$$

$$\mathcal{L} = \log [p(x_1^1) p(? | x_1^1) p(x_3^1 | x_2^1)] + \dots$$

What if...

- Not fully observed?



$$\mathcal{L} = \sum_i \log [p(x_1^i) p(x_2^i|x_1^i) p(x_3^i|x_2^i)]$$

$$D = \begin{bmatrix} x^1 \\ x^2 \\ x^3 \\ x^4 \\ \vdots \end{bmatrix} = \begin{bmatrix} x_1^1 & ? & x_3^1 \\ x_1^2 & x_2^2 & x_3^2 \\ ? & x_2^3 & x_3^3 \\ x_1^4 & ? & ? \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$$\begin{aligned} \mathcal{L} &= \log [p(x_1^1) p(?|x_1^1) p(x_3^1|x_2^1)] + \dots \\ &= \log \left[\sum_{x_2} p(x_1^1) p(x_2|x_1^1) p(x_3^1|x_2) \right] + \dots \\ &= \sum_i \log \sum_{x_j \in S^i} p(x^i) \end{aligned}$$

No longer decomposes nicely into individual parts...
Hard to estimate in closed form.

Solutions?

1. Optimize directly anyway (gradient ascent, etc.)
2. Discard data with missing entries
3. Fill in (“impute”) missing entries somehow

Complete log-likelihood

- Gaussian mixture model, log-likelihood

$$\mathcal{L}_X(\theta = \{\pi, \mu, \sigma^2\}) = \sum_j \log \left[\sum_k \pi_k h(\sigma) \exp((x^{(j)} - \mu)^2 / \sigma^2) \right]$$

- If we observed z : the “complete data likelihood”

$$\begin{aligned} \mathcal{L}_{XZ}(\theta) &= \sum_j \log p(z^{(j)}) p(x^{(j)} | z^{(j)}) \\ &= \sum_j \log \pi_{z^{(j)}} + (x^{(j)} - \mu_{z^{(j)}})^2 / \sigma_{z^{(j)}}^2 + \dots \end{aligned}$$

- Now, given z , this is exponential family – easy:

$$\hat{\pi}_k = \frac{1}{m} \sum \delta(z^{(j)} = k) = \frac{m_k}{m} \quad \hat{\mu}_k = \frac{1}{m_k} \sum_{z^j=k} x^{(j)} \quad \dots$$

Learning with hidden data

- “Estimate” z , then optimize complete LL
- Example: *k-means-like* algorithm (“hard EM”)
 - Find “best” z -values, then best means; repeat
- Expectation-Maximization
 - Instead of fixing z 's, take “soft assignment”
 - Use $q(z=1) \dots q(z=k)$ (as if “partial” observations)

Expectation-Maximization

- Recall the K-L divergence between q & p :

$$D(\hat{q}(x) \| p(x)) = \mathbb{E}_q \left[\log \frac{q(x)}{p(x)} \right]$$

- $D(q \| p) \geq 0$, $= 0$ iff $p=q$ a.e.

- Also the empirical distribution $\hat{p}(x) = \frac{1}{m} \sum_j \delta(x = x^{(j)})$

$$\begin{aligned} D(\hat{p}(x) \| p(x; \theta)) &= \frac{1}{m} \sum_j \delta(x = x^{(j)}) (\log \hat{p}(x) - \log p(x; \theta)) \\ &= H(\hat{p}) - \frac{1}{m} \mathcal{L}(\theta) \end{aligned}$$

Expectation-Maximization

$$\begin{aligned} D(\hat{p} \| p(x; \theta)) &= -\mathbb{E}_{\hat{p}} \log p(x; \theta) - H(\hat{p}) \\ &= -\mathbb{E}_{\hat{p}} \left[\log \sum_z p(x, z; \theta) \right] - H(\hat{p}) \\ &= -\mathbb{E}_{\hat{p}} \left[\log \sum_z q(z|x) \frac{p(x, z; \theta)}{q(z|x)} \right] - H(\hat{p}) \\ &= -\mathbb{E}_{\hat{p}} \left[\log \mathbb{E}_q \left[\frac{p(x, z; \theta)}{q(z|x)} \right] \right] - H(\hat{p}) \\ &\leq -\mathbb{E}_{\hat{p}} \mathbb{E}_q \left[\log \frac{p(x, z; \theta)}{q(z|x)} \right] - H(\hat{p}) \quad (\text{Jensen's ineq.}) \\ &= D(\hat{p}(x) \cdot q(z|x) \| p(x, z; \theta)) \end{aligned}$$

$-\frac{1}{m} \mathcal{L}_X + \hat{H}_X$

Expectation-Maximization

$$\max_{\theta} \max_q -D(\hat{p}(x)q(z|x) \| p(x, z; \theta))$$

- E-step: max with respect to q
 - Optimal $q(z|x) = p(z|x, \theta)$

$$\mathbb{E}_{\hat{p}} \mathbb{E}_q \left[\log \frac{p(x|\theta)p(z|x, \theta)}{\hat{p}(x)q(z|x)} \right] \xrightarrow{q=p(z|x, \theta)} \mathbb{E}_{\hat{p}} \left[\log \frac{p(x|\theta)}{\hat{p}(x)} \right] = \frac{1}{m} \mathcal{L}_x$$

- M-step: max with respect to θ
 - Given weights of $q(z|x)$

$$\mathbb{E}_q \mathbb{E}_{\hat{p}} \left[\log \frac{p(x; \theta)p(z|x; \theta)}{\hat{p}(x)q(z|x)} \right] = \mathbb{E}_q \mathbb{E}_{\hat{p}} \left[\log p(z; \theta)p(x|z; \theta) \right] + \text{const.}$$

(Expected Complete Log-Likelihood)

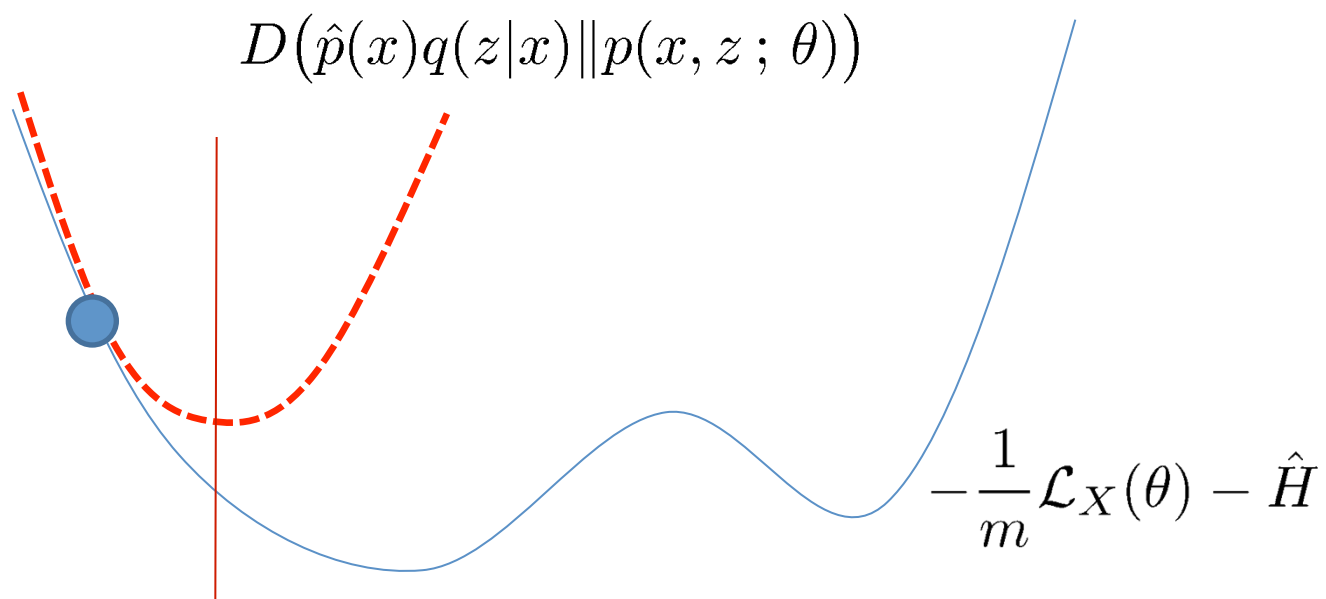
Expectation-Maximization

$$\max_{\theta} \max_q -D(\hat{p}(x)q(z|x) || p(x, z; \theta))$$

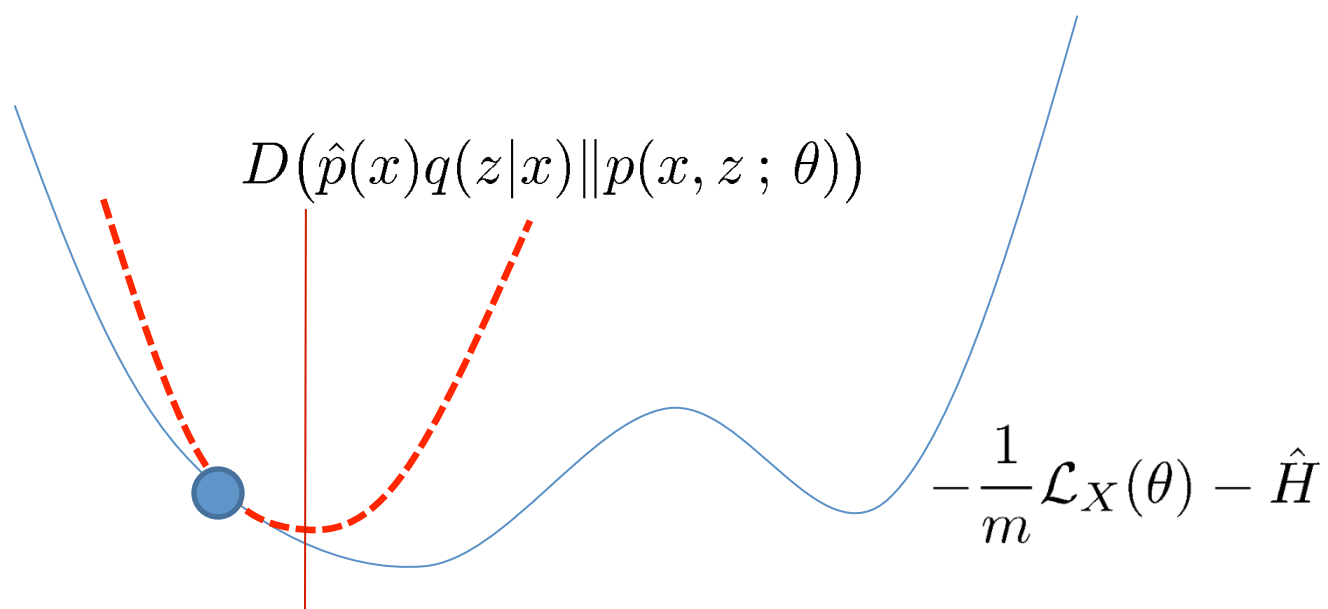
- E-step: max with respect to q
 - Optimal $q(z) = p(z | x, \theta)$
- M-step: max with respect to θ
 - Expected complete log-likelihood
- Coordinate ascent
 - Always converges; maximizes lower bound on LL
 - Touches (equality) after E-step

Note: if we restrict $q(z)$ to 0/1, this is *hard EM* (k-means-like)

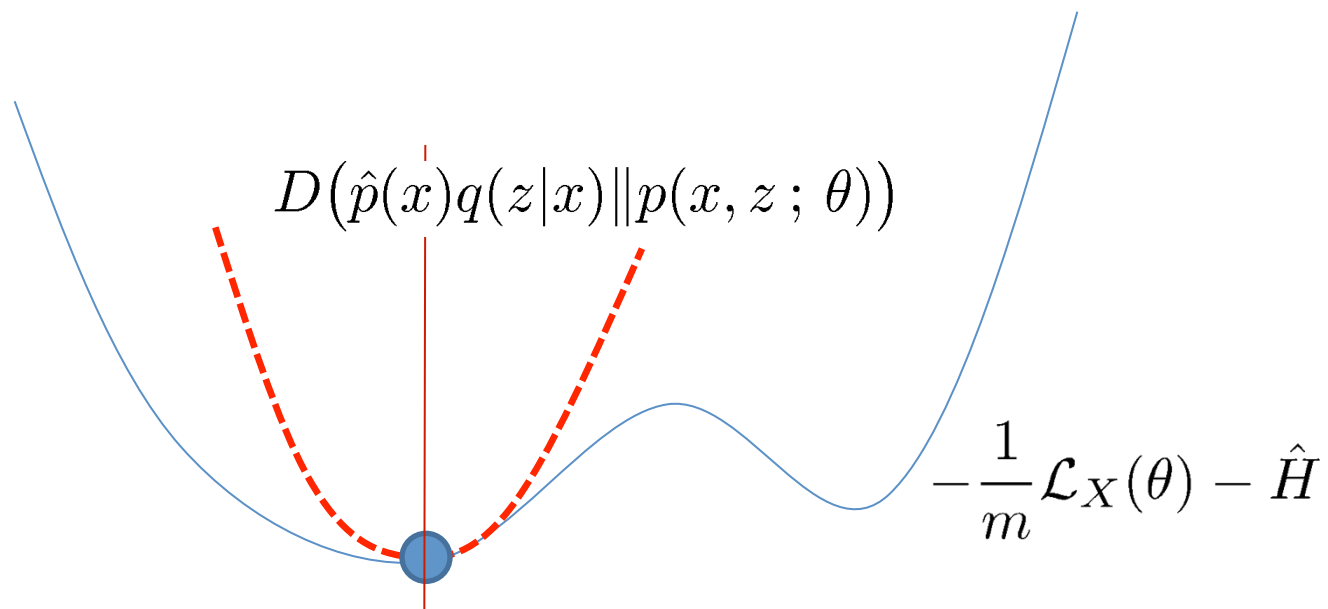
Expectation-Maximization



Expectation-Maximization



Expectation-Maximization



Expectation-Maximization

- Intuition:
 - Instead of fixing z 's, take “soft assignment”
 - Use $q(z=1) \dots q(z=k)$ (as if “partial” observations)
- Derivation: optimize the expected complete LL

$$\max_{\theta} \mathbb{E}_{q(Z)} [\mathcal{L}(X, Z)] = \sum_Z q(Z) \underline{\mathcal{L}(X, Z)}$$

= log $p(x, z)$

$p(x, z)$ is structured (factors into a product)

=> log p decomposes into a sum of log terms

Optimizing the ECLL

- Usually not harder than standard ML estimation for the CLL

$$\max_{\theta} \mathbb{E}_{q(Z)} [\mathcal{L}_{XZ}(\theta)] = \sum_Z q(Z) \mathcal{L}_{XZ}(\theta)$$

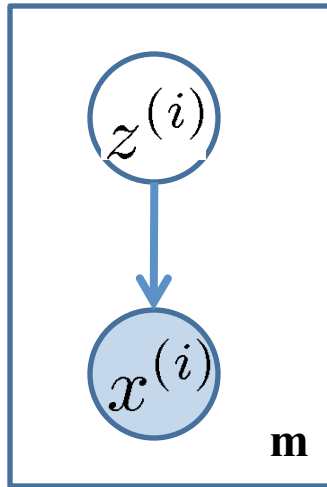
$$\frac{\partial}{\partial \theta_i} \mathcal{L}_{XZ}(\theta) = \sum_j \sum_k q(z^{(j)} = k) \frac{\partial}{\partial \theta_i} \left[\log \pi_k - \frac{1}{2} (x^{(j)} - \mu_k)^2 / \sigma_k^2 \right]$$

π, μ, σ **Weight, W_{jk}** $u_i(x^{(j)}, z = k) - \mathbb{E}[u_i(x^{(j)}, z = k)]$

Expectation-Maximization for GMMs

- E-step: compute the expected values of z 's
 - Weights $w_{jk} = p(z^{(j)} = k | x^{(j)}, \theta)$
$$= p(z^{(j)} = k, x^{(j)} | \theta) / \sum_{k'} p(z^{(j)} = k', x^{(j)} | \theta)$$
- M-step: maximize θ given weights (closed form)
 - $\hat{\pi}_k = \frac{1}{m} \sum w_{jk} = \frac{W_k}{m}$
 - $\hat{\mu}_k = \frac{1}{W_k} \sum_j w_{jk} x^{(j)}$
 - $\hat{\sigma}_k^2 = \frac{1}{W_k} \sum_j w_{jk} (x^{(j)} - \hat{\mu}_k)^2$

Three Mixture Models



Probabilistic PCA

$$\prod_i p(z^{(i)}) = \prod \mathcal{N}(0, I)$$

$$\prod_i p(x^{(i)} | z^{(i)}) = \mathcal{N}(Wz + \mu, \sigma^2 I)$$

Gaussian Mixture Models

$$\prod_i p(z^{(i)}) = [\pi_1 \dots \pi_K]$$

$$\prod_i p(x^{(i)} | z^{(i)}) = \mathcal{N}(\mu_z, \Sigma_z)$$

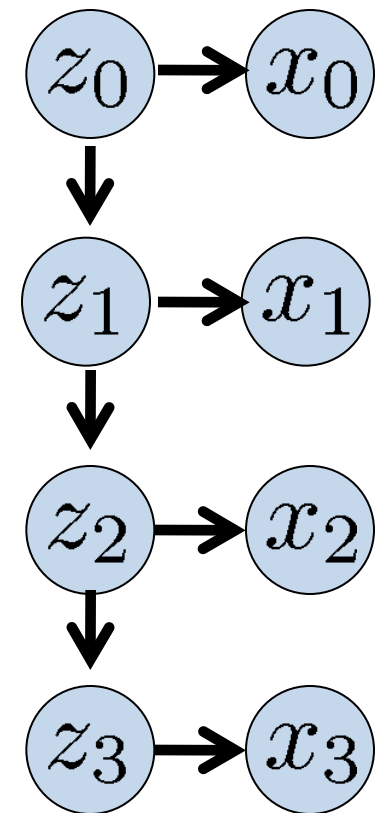
Factor Analysis

$$\prod_i p(z^{(i)}) = \prod \mathcal{N}(0, \Lambda)$$

$$\prod_i p(x^{(i)} | z^{(i)}) = \mathcal{N}(Wz + \mu, \sigma^2 I)$$

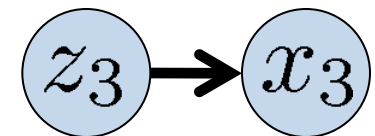
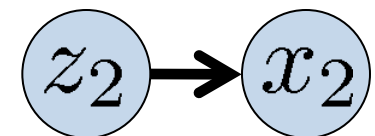
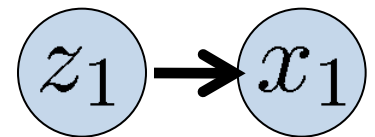
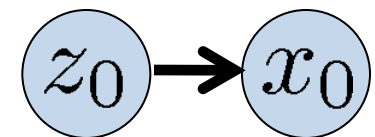
Hidden Markov models

- A Markov model, but in which we can not observe the state (just something indirect)
- Hidden variables Z_t
- Observed variables X_t
- “Emission probability distribution”
 - $p(X_t = x \mid Z_t = k)$



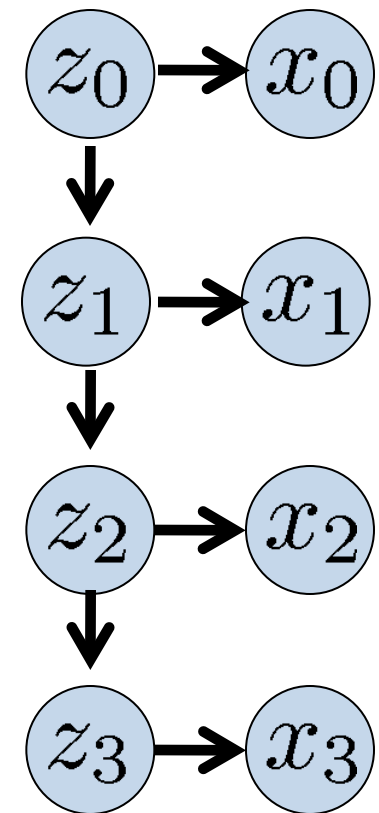
Hidden Markov models

- A Markov model, but in which we can not observe the state (just something indirect)
- Hidden variables Z_t
- Observed variables X_t
- “Emission probability distribution”
 - $p(X_t = x \mid Z_t = k)$
- Special case: Z_t independent of Z_{t-1}
 - Discrete distribution over z 's
 - Given z , a distribution over obs. x
 - A mixture model, e.g. in clustering



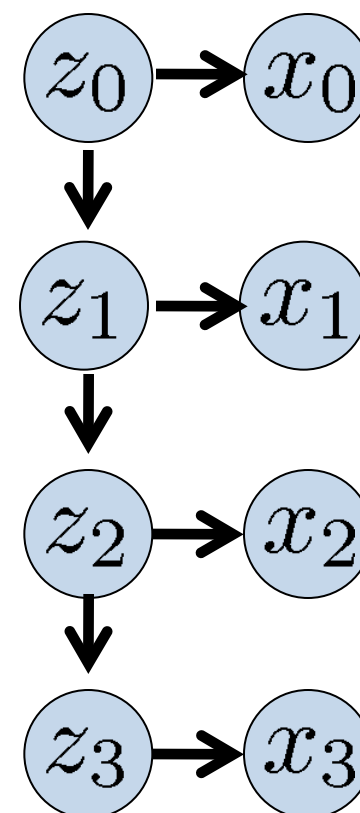
Hidden Markov models

- A Markov model, but in which we can not observe the state (just something indirect)
- Hidden variables Z_t
- Observed variables X_t
- “Emission probability distribution”
 - $p(X_t = x \mid Z_t = k)$



EM algorithm

- E-step: compute $q(Z) = p(Z|X; \theta)$
- M-step: max over θ of $E[\log p(Z, X; \theta)]$
- Sufficient statistics:
 - # times see $j \rightarrow i$ for each conditional
- Depends on $q(z_{t+1} = i, z_t = j)$
 $q(z_t = i, x_t = j)$



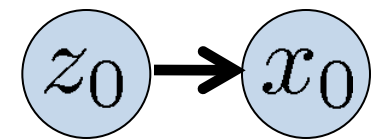
Inference in HMMs

- Need to compute $q(z_t) = p(z_t \mid x_1 \dots x_n)$
- Two-pass algorithm:
 - Forward pass: $p(z_t = i \mid x_1 \dots x_t)$
 - Backward pass: $p(x_{t+1} \dots x_n \mid z_t = i)$
- $p(z_0=i) \propto p(z_0=i)$

z_0

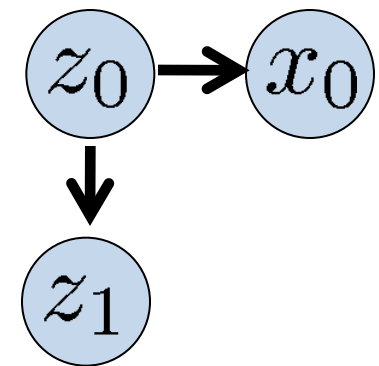
Inference in HMMs

- Need to compute $q(z_t) = p(z_t \mid x_1 \dots x_n)$
- Two-pass algorithm:
 - Forward pass: $p(z_t = i \mid x_1 \dots x_t)$
 - Backward pass: $p(x_{t+1} \dots x_n \mid z_t = i)$
- $p(z_0 = i \mid x_0) \propto p(z_0 = i) p(x_0 \mid z_0 = i)$



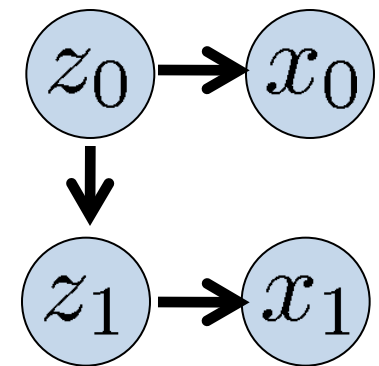
Inference in HMMs

- Need to compute $q(z_t) = p(z_t \mid x_1 \dots x_n)$
- Two-pass algorithm:
 - Forward pass: $p(z_t = i \mid x_1 \dots x_t)$
 - Backward pass: $p(x_{t+1} \dots x_n \mid z_t = i)$
- $p(z_0 = i \mid x_0) \propto p(z_0 = i) p(x_0 \mid z_0 = i)$
- $p(z_1 = i \mid x_0) \propto \sum_j p(z_1 = i \mid z_0 = j) p(z_0 = j \mid x_0)$



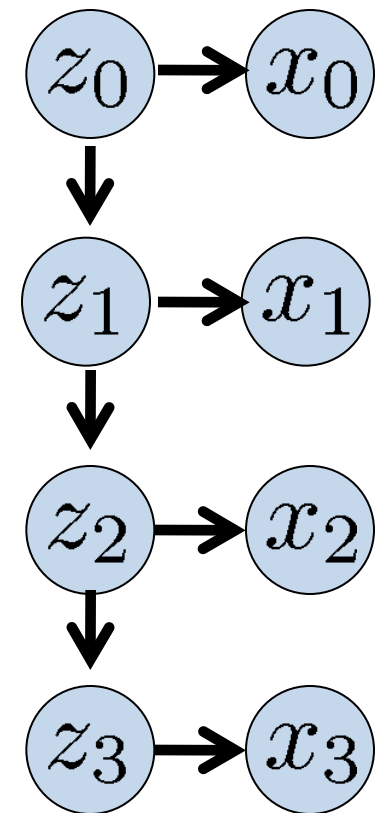
Inference in HMMs

- Need to compute $q(z_t) = p(z_t \mid x_1 \dots x_n)$
- Two-pass algorithm:
 - Forward pass: $p(z_t = i \mid x_1 \dots x_t)$
 - Backward pass: $p(x_{t+1} \dots x_n \mid z_t = i)$
- $p(z_0=i \mid x_0) \propto p(z_0=i) p(x_0 \mid z_0=i)$
- $p(z_1=i \mid x_0) \propto \sum_j p(z_1=i \mid z_0=j) p(z_0=j \mid x_0)$
- $p(z_1=i \mid x_0, x_1) \propto p(z_1=i \mid x_0) p(x_1 \mid z_1=i)$



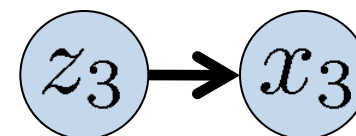
Inference in HMMs

- Need to compute $q(z_t) = p(z_t \mid x_1 \dots x_n)$
- Two-pass algorithm:
 - Forward pass: $p(z_t = i \mid x_1 \dots x_t)$
 - Backward pass: $p(x_{t+1} \dots x_n \mid z_t = i)$
- $p(z_0=i \mid x_0) \propto p(z_0=i) p(x_0 \mid z_0=i)$
- $p(z_1=i \mid x_0) \propto \sum_j p(z_1=i \mid z_0=j) p(z_0=j \mid x_0)$
- $p(z_1=i \mid x_0, x_1) \propto p(z_1=i \mid x_0) p(x_1 \mid z_1=i)$
- ...
- $p(z_n \mid x_0, \dots, x_n) \propto p(z_n \mid x_0 \dots x_{n-1}) p(x_n \mid z_n)$



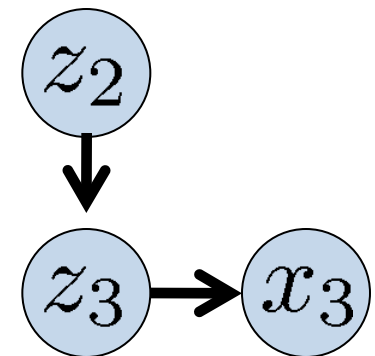
Inference in HMMs

- Need to compute $q(z_t) = p(z_t \mid x_1 \dots x_n)$
- Two-pass algorithm:
 - Forward pass: $p(z_t = i \mid x_1 \dots x_t)$
 - Backward pass: $p(x_{t+1} \dots x_n \mid z_t = i)$
- $p(x_n \mid z_n = i)$



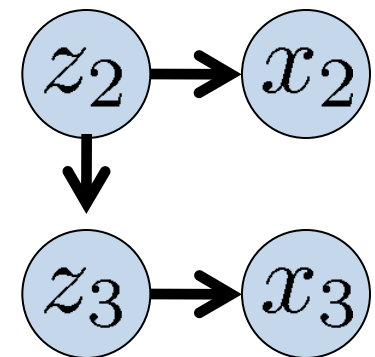
Inference in HMMs

- Need to compute $q(z_t) = p(z_t \mid x_1 \dots x_n)$
- Two-pass algorithm:
 - Forward pass: $p(z_t = i \mid x_1 \dots x_t)$
 - Backward pass: $p(x_{t+1} \dots x_n \mid z_t = i)$
- $p(x_n \mid z_n = i)$
- $p(x_n \mid z_{n-1} = i) \propto \sum_j p(z_n = j \mid z_{n-1} = i) p(x_n \mid z_n = j)$



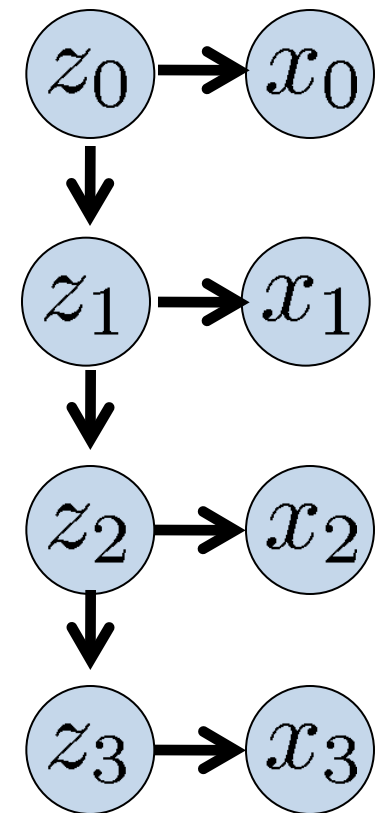
Inference in HMMs

- Need to compute $q(z_t) = p(z_t \mid x_1 \dots x_n)$
- Two-pass algorithm:
 - Forward pass: $p(z_t = i \mid x_1 \dots x_t)$
 - Backward pass: $p(x_{t+1} \dots x_n \mid z_t = i)$
- $p(x_n \mid z_n = i)$
- $p(x_n \mid z_{n-1} = i) \propto \sum_j p(z_n = j \mid z_{n-1} = i) p(x_n \mid z_n = j)$
- $p(x_n, x_{n-1} \mid z_{n-1} = i) \propto p(x_{n-1} \mid z_{n-1} = i) p(x_n \mid z_{n-1} = i)$



Inference in HMMs

- Need to compute $q(z_t) = p(z_t \mid x_1 \dots x_n)$
- Two-pass algorithm:
 - Forward pass: $p(z_t = i \mid x_1 \dots x_t)$
 - Backward pass: $p(x_{t+1} \dots x_n \mid z_t = i)$
- $p(x_n \mid z_n = i)$
- $p(x_n \mid z_{n-1} = i) \propto \sum_j p(z_n = j \mid z_{n-1} = i) p(x_n \mid z_n = j)$
- $p(x_n, x_{n-1} \mid z_{n-1} = i) \propto p(x_{n-1} \mid z_{n-1} = i) p(x_n \mid z_{n-1} = i)$
- ...



Forward-backward

- $O_t(i) = p(x_t = X_t \mid z_t = i)$
- $T(i,j) = p(z_t = i \mid z_{t-1}=j)$
- Recursion
 - $a_t = (T * a_{t-1}) .* O_t \propto p(z_t = i \mid x_1 \dots x_t)$
 - $b_t = T' * (b_{t+1} .* O_{t+1}) \propto p(x_{t+1} \dots x_N \mid z_t = i)$
 - $p(z_t = i) \propto a_t(i) b_t(i) \propto p(z_t = i \mid x_1 \dots x_N)$
 - $p(z_t = i, z_{t-1} = j \mid x_1 \dots x_N) \propto T(i,j) b_t(j) O_t(j) a_{t-1}(i)$
- Gaussian case: similar recursion, “Kalman filter / smoother”

EM for HMMs

- E-step: compute $q(z_t, z_{t-1}) = p(z_t, z_{t-1} | x_1 \dots x_N)$
 - Using forward-backward recursions
- M-step: maximize expected LL given $q(\cdot)$

$$\hat{T}_{ij} = \frac{\sum_t q(z_t = i, z_{t-1} = j)}{\sum_{t,k} q(z_t = k, z_{t-1} = j)} \quad (\mathbf{N}_{ij} / \mathbf{N}_j)$$

$$\hat{\pi}_i = q(z_0 = i)$$

$$\hat{O}_{iX} = \hat{p}(x_t = X | z_t = i) = \frac{\sum_{t:x_t=X} q(z_t = i)}{\sum_t q(z_t = i)} \quad (\mathbf{N}_{iX} / \mathbf{N}_i)$$

EM Variants

- To use EM, we need:
 - Complete log-likelihood easy to optimize
 - $q(z|x) = p(z | x ; \theta)$ efficiently computable
- Alternatives? What if $p(z|x)$ is difficult?
- Hard EM: assign “most likely” z
 - If maximizing $p(z | x ; \theta)$ is hard, we can approximate
- Stochastic EM: sample $z \sim p(z | x ; \theta)$
 - If sampling hard, approximate (e.g., using MCMC)
- Variational EM: approximate $q(z|x)$ directly
 - Replace $q(.)$ with a model that is easier to compute over