

DSC291: Machine Learning with Few Labels

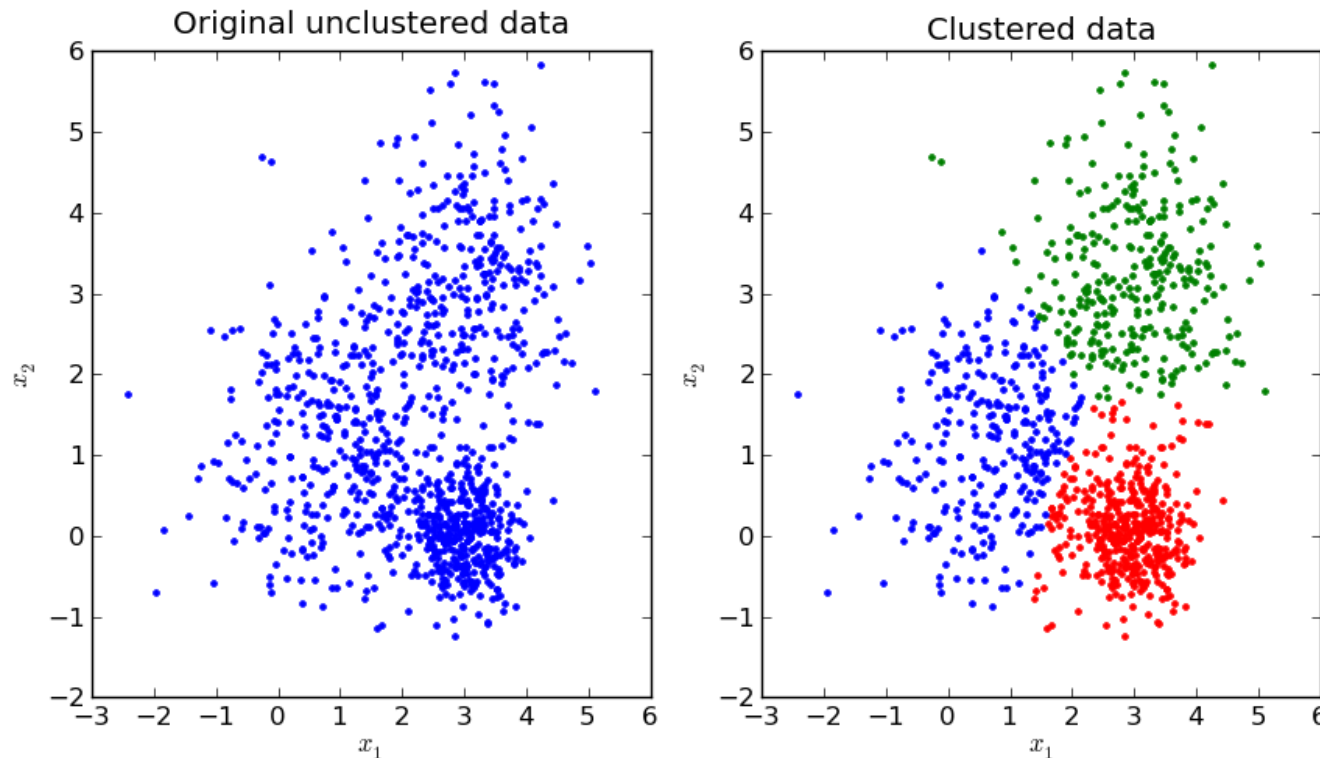
Unsupervised Learning

Zhiting Hu

Lecture 4, April 10, 2025

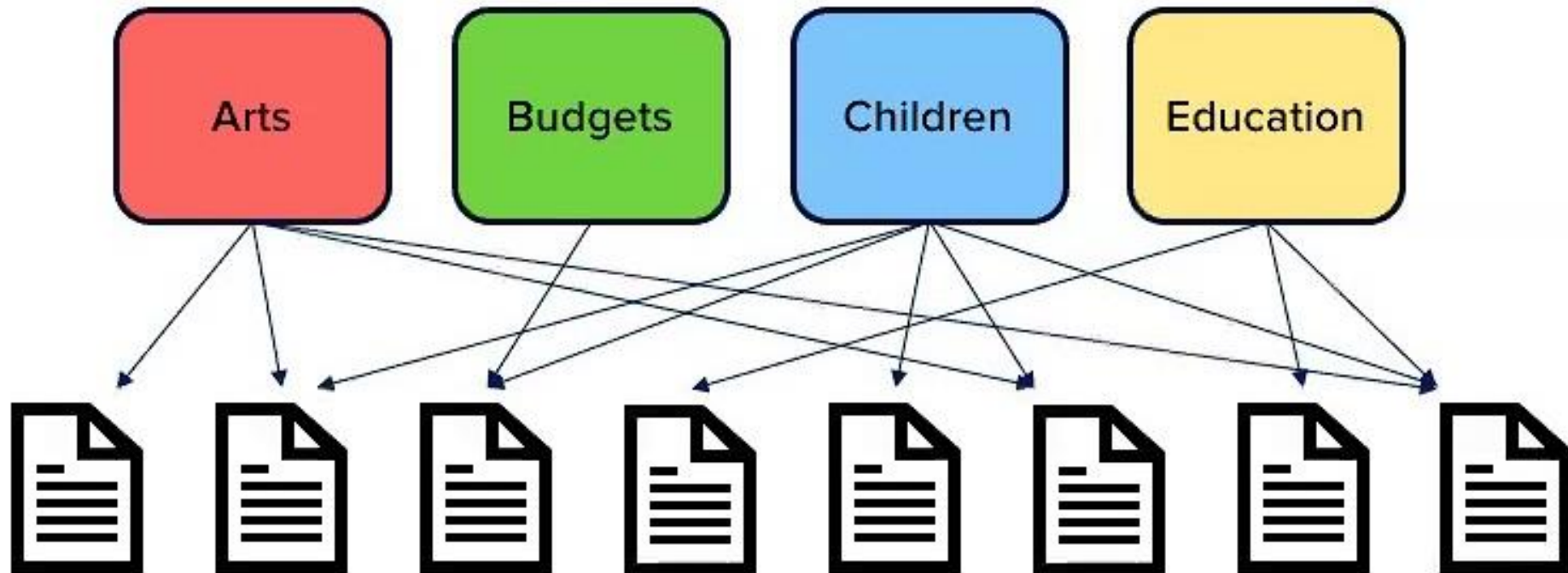
Unsupervised Learning

- Each data instance is partitioned into two parts:
 - observed variables \mathbf{x}
 - latent (unobserved) variables \mathbf{z}
- Want to learn a model $p_{\theta}(\mathbf{x}, \mathbf{z})$



Unsupervised Learning

- Each data instance is partitioned into two parts:
 - observed variables \mathbf{x}
 - latent (unobserved) variables \mathbf{z}
- Want to learn a model $p_{\theta}(\mathbf{x}, \mathbf{z})$



Unsupervised Learning

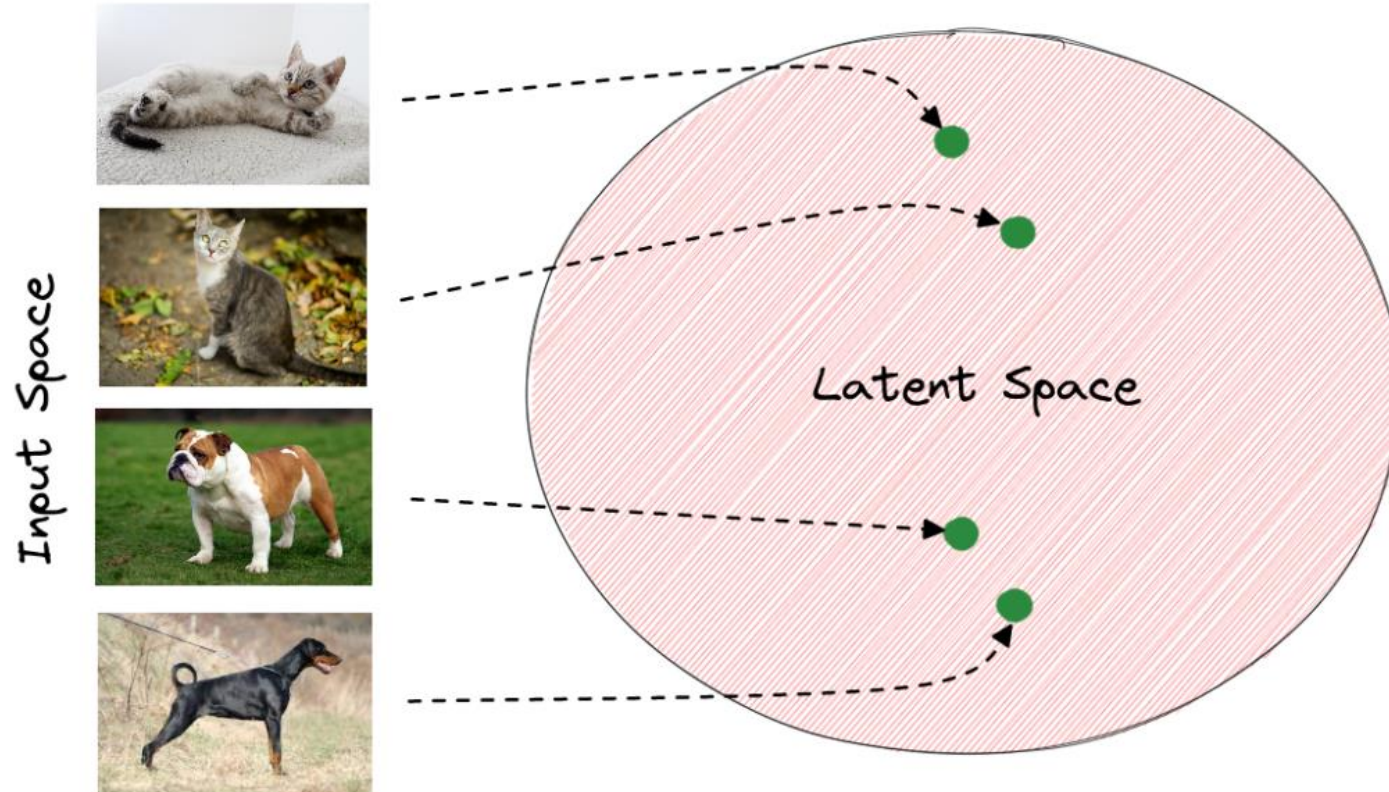
- Each data instance is partitioned into two parts:
 - observed variables \mathbf{x}
 - latent (unobserved) variables \mathbf{z}
- Want to learn a model $p_{\theta}(\mathbf{x}, \mathbf{z})$

Input Space



Unsupervised Learning

- Each data instance is partitioned into two parts:
 - observed variables \mathbf{x}
 - latent (unobserved) variables \mathbf{z}
- Want to learn a model $p_{\theta}(\mathbf{x}, \mathbf{z})$



Why is Unsupervised Learning Harder?

- **Complete log likelihood:** if both x and z can be observed, then

$$\ell_c(\theta; x, z) = \log p(x, z|\theta) = \log p(z|\theta_z) + \log p(x|z, \theta_x)$$

- Decomposes into a sum of factors, the parameter for each factor can be estimated separately

Now z is not observed:

- **Incomplete (or marginal) log likelihood:** with z unobserved, our objective becomes the log of a marginal probability:

$$\ell(\theta; x) = \log p(x|\theta)$$

Why is Unsupervised Learning Harder?

- **Complete log likelihood:** if both \mathbf{x} and \mathbf{z} can be observed, then

$$\ell_c(\theta; \mathbf{x}, \mathbf{z}) = \log p(\mathbf{x}, \mathbf{z} | \theta) = \log p(\mathbf{z} | \theta_z) + \log p(\mathbf{x} | \mathbf{z}, \theta_x)$$

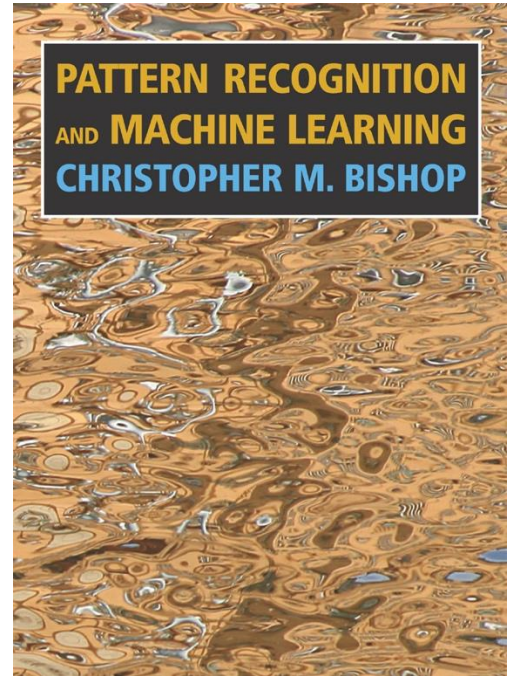
- Decomposes into a sum of factors, the parameter for each factor can be estimated separately

Now \mathbf{z} is not observed:

- **Incomplete (or marginal) log likelihood:** with \mathbf{z} unobserved, our objective becomes the log of a marginal probability:

$$\ell(\theta; \mathbf{x}) = \log p(\mathbf{x} | \theta) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z} | \theta)$$

- All parameters become coupled together
- In other models when \mathbf{z} is complex (continuous) variables (as we'll see later), marginalization over \mathbf{z} is intractable.



Expectation Maximization (EM)

9	Mixture Models and EM	423
9.1	<i>K</i> -means Clustering	424
9.1.1	Image segmentation and compression	428
9.2	Mixtures of Gaussians	430
9.2.1	Maximum likelihood	432
9.2.2	EM for Gaussian mixtures	435
9.3	An Alternative View of EM	439
9.3.1	Gaussian mixtures revisited	441
9.3.2	Relation to <i>K</i> -means	443
9.3.3	Mixtures of Bernoulli distributions	444
9.3.4	EM for Bayesian linear regression	448
This class → 9.4	The EM Algorithm in General	450
	Exercises	455

Expectation Maximization (EM): Intuition

Expectation Maximization (EM): Intuition

- Supervised MLE is easy:

$$\max_{\theta} \ell_c(\theta; \mathbf{x}, \mathbf{z}) = \log p(\mathbf{x}, \mathbf{z} | \theta)$$

- Observe both \mathbf{x} and \mathbf{z}

- Unsupervised MLE is hard:

$$\max_{\theta} \ell(\theta; \mathbf{x}) = \log p(\mathbf{x} | \theta) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z} | \theta)$$

- Observe only \mathbf{x}

- EM, intuitively:

E-step: $q(\mathbf{z} | \mathbf{x}) = p(\mathbf{z} | \mathbf{x}, \theta)$

We don't actually observe q , let's estimate it

M-step: $\max_{\theta} \mathbb{E}_{q(\mathbf{z} | \mathbf{x})} [\log p(\mathbf{x}, \mathbf{z} | \theta)]$

Let's "pretend" we also observe \mathbf{z} (its distribution)

Expectation Maximization (EM): Intuition

- Supervised MLE is easy:

$$\max_{\theta} \ell_c(\theta; \mathbf{x}, \mathbf{z}) = \log p(\mathbf{x}, \mathbf{z} | \theta)$$


- Observe both \mathbf{x} and \mathbf{z}

- Unsupervised MLE is hard:

$$\max_{\theta} \ell(\theta; \mathbf{x}) = \log p(\mathbf{x} | \theta) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z} | \theta)$$

- Observe only \mathbf{x}

- EM, intuitively:



E-step: $q^{t+1}(\mathbf{z} | \mathbf{x}) = p(\mathbf{z} | \mathbf{x}, \theta^t)$

We don't actually observe \mathbf{z} , let's estimate it

M-step: $\max_{\theta} \mathbb{E}_{q^{t+1}(\mathbf{z} | \mathbf{x})} [\log p(\mathbf{x}, \mathbf{z} | \theta)]$

Let's "pretend" we also observe \mathbf{z} (its distribution)

This is an iterative process

Expectation Maximization (EM): Intuition

- Supervised MLE is easy:

$$\max_{\theta} \ell_c(\theta; \mathbf{x}, \mathbf{z}) = \log p(\mathbf{x}, \mathbf{z} | \theta)$$

- Observe both \mathbf{x} and \mathbf{z}

- Unsupervised MLE is hard:

$$\max_{\theta} \ell(\theta; \mathbf{x}) = \log p(\mathbf{x} | \theta) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z} | \theta)$$

- Observe only \mathbf{x}

- EM, intuitively:



E-step: $q^{t+1}(\mathbf{z} | \mathbf{x}) = p(\mathbf{z} | \mathbf{x}, \theta^t)$

We don't actually observe q , let's estimate it

M-step: $\max_{\theta} \mathbb{E}_{q(\mathbf{z} | \mathbf{x})} [\log p(\mathbf{x}, \mathbf{z} | \theta)]$

Let's "pretend" we also observe \mathbf{z} (its distribution)

This is an iterative process

Expectation Maximization (EM): Intuition

- Supervised MLE is easy:

$$\max_{\theta} \ell_c(\theta; \mathbf{x}, \mathbf{z}) = \log p(\mathbf{x}, \mathbf{z} | \theta)$$

- Observe both \mathbf{x} and \mathbf{z}

- Unsupervised MLE is hard:

$$\max_{\theta} \ell(\theta; \mathbf{x}) = \log p(\mathbf{x} | \theta) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z} | \theta)$$

- Observe only \mathbf{x}

- EM, intuitively

$$\ell(\theta; \mathbf{x}) \geq \mathbb{E}_{q(\mathbf{z} | \mathbf{x})} [\log p(\mathbf{x}, \mathbf{z} | \theta)] + H(q)$$



M-step: $\max_{\theta} \mathbb{E}_{q(\mathbf{z} | \mathbf{x})} [\log p(\mathbf{x}, \mathbf{z} | \theta)]$

only observe q , let's estimate it

Let's "pretend" we also observe \mathbf{z} (its distribution)

This is an iterative process

Expectation Maximization (EM): Intuition

- **Question:** show that $\ell(\theta; \mathbf{x}) \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}, \mathbf{z}|\theta)] + H(q)$

Expectation Maximization (EM): Intuition

- **Question:** show that $\ell(\theta; \mathbf{x}) \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}, \mathbf{z}|\theta)] + H(q)$
$$= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right]$$

Expectation Maximization (EM): Intuition

- **Question:** show that $\ell(\theta; \mathbf{x}) \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}, \mathbf{z}|\theta)] + H(q)$
$$= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right]$$

- **Hint:** first show that

$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta))$$

- Since KL divergence is non-negative, we arrive at the conclusion

Expectation Maximization (EM): Intuition

- **Question:** show that $\ell(\theta; \mathbf{x}) \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}, \mathbf{z}|\theta)] + H(q)$
$$= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right]$$

- **Hint:** first show that

$$\begin{aligned} \ell(\theta; \mathbf{x}) &= \boxed{\mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right]} + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta)) && \text{Evidence Lower Bound (ELBO)} \\ &= -\boxed{F(q, \theta)} + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta)) \end{aligned}$$

Variational free energy

$$F(q, \theta) = - \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}, \mathbf{z}|\theta)] - H(q)$$

Lower Bound and Free Energy

- Variational free energy:

$$F(q, \theta) = - \mathbb{E}_{q(z|x)} [\log p(\mathbf{x}, \mathbf{z}|\theta)] - H(q)$$

- The EM algorithm is coordinate-decent on F

- At each step t :

- E-step: $q^{t+1} = \arg \min_q F(q, \theta^t) \Rightarrow q^{t+1}(z|x) = p(z|x, \theta^t)$

- M-step: $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta) \Rightarrow \max_{\theta} \mathbb{E}_{q^{t+1}(z|x)} [\log p(\mathbf{x}, \mathbf{z}|\theta)]$


Lower Bound and Free Energy

- Variational free energy:

$$F(q, \theta) = - \mathbb{E}_{q(z|x)} [\log p(x, z|\theta)] - H(q)$$

- The EM algorithm is coordinate-decent on F

- At each step t :

- E-step: $q^{t+1} = \arg \min_q F(q, \theta^t)$  $\Rightarrow q^{t+1}(z|x) = p(z|x, \theta^t)$

- M-step: $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta)$ $\Rightarrow \max_{\theta} \mathbb{E}_{q^{t+1}(z|x)} [\log p(x, z|\theta)]$

E-step: minimization of $F(q, \theta)$ w.r.t q

- **Question:** show that that optimal solution of E-step is

$$q^{t+1} = \operatorname{argmin}_q F(q, \theta^t) = p(\mathbf{z}|\mathbf{x}, \theta^t)$$

- I.e., the posterior distribution over the latent variables given the data and the current parameters.

- **Hint:** use the fact

$$\ell(\theta^t; \mathbf{x}) = -F(q, \theta^t) + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta^t))$$



Independent of q



≥ 0

- $F(q, \theta^t)$ is minimized when $\text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta^t)) = 0$, which is achieved only when $q(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x}, \theta^t)$


Lower Bound and Free Energy


- Variational free energy:

$$F(q, \theta) = - \mathbb{E}_{q(z|x)} [\log p(x, z|\theta)] - H(q)$$

- The EM algorithm is coordinate-decent on F

- At each step t :

- E-step: $q^{t+1} = \arg \min_q F(q, \theta^t)$  $\Rightarrow q^{t+1}(z|x) = p(z|x, \theta^t)$

- M-step: $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta)$  $\Rightarrow \max_{\theta} \mathbb{E}_{q^{t+1}(z|x)} [\log p(x, z|\theta)]$

EM Algorithm: Quick Summary

- Observed variables \mathbf{x} , latent variables \mathbf{z}
- To learn a model $p(\mathbf{x}, \mathbf{z}|\theta)$, we want to maximize the marginal log-likelihood

○ But it's too difficult $\ell(\theta; \mathbf{x}) = \log p(\mathbf{x}|\theta) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}|\theta)$

- EM algorithm:
 - maximize a lower bound of $\ell(\theta; \mathbf{x})$
 - Or equivalently, minimize an upper bound of $-\ell(\theta; \mathbf{x})$
- Key equation:

$$\begin{aligned} \ell(\theta; \mathbf{x}) &= \underbrace{\mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right]}_{\text{Evidence Lower Bound (ELBO)}} + \text{KL}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}, \theta)) \\ &= -\underbrace{F(q, \theta)}_{\text{Variational free energy}} + \text{KL}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}, \theta)) \end{aligned}$$

Variational free energy

EM Algorithm: Quick Summary

- The EM algorithm is coordinate-decent on $F(q, \theta)$
 - E-step: $q^{t+1} = \arg \min_q F(q, \theta^t) = p(\mathbf{z}|\mathbf{x}, \theta^t)$
 - the posterior distribution over the latent variables given the data and the current parameters
 - M-step: $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta) = \operatorname{argmax}_{\theta} \sum_{\mathbf{z}} q^{t+1}(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{z}|\theta)$

$$\begin{aligned}\ell(\theta; \mathbf{x}) &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta)) \\ &= -F(q, \theta) + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta))\end{aligned}$$

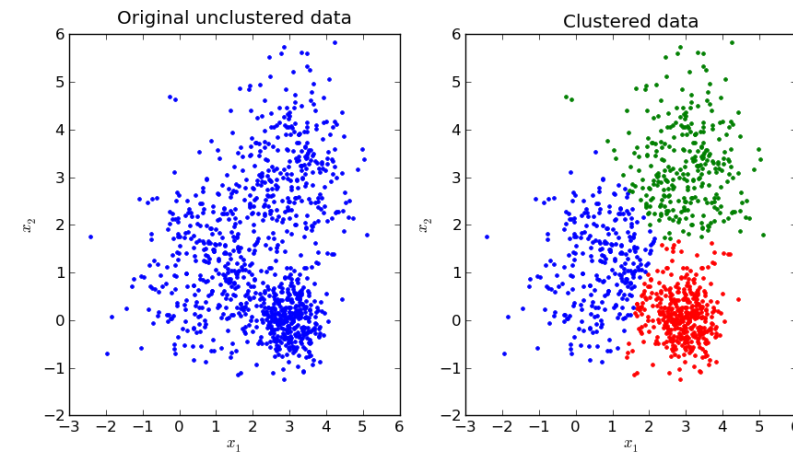
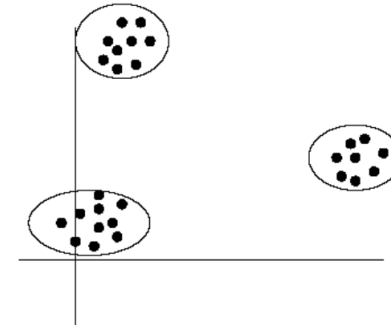
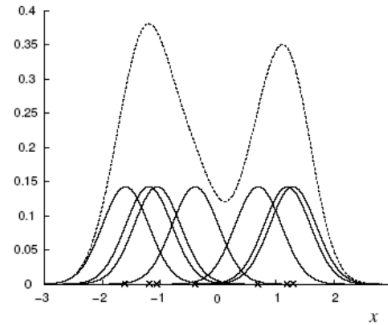
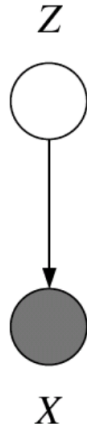
Example: Gaussian Mixture Models

- Consider a mixture of K Gaussian components:

$$p(x_n|\mu, \Sigma) = \sum_k \pi_k N(x, | \mu_k, \Sigma_k)$$

mixture proportion

mixture component



- This model can be used for unsupervised clustering.
 - This model (fit by AutoClass) has been used to discover new kinds of stars in astronomical data, etc.

Example: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components:

- Z is a latent class indicator vector:

$$p(z_n) = \text{multi}(z_n : \pi) = \prod_k (\pi_k)^{z_n^k}$$

- X is a conditional Gaussian variable with a class-specific mean/covariance

$$p(x_n | z_n^k = 1, \mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)\right\}$$

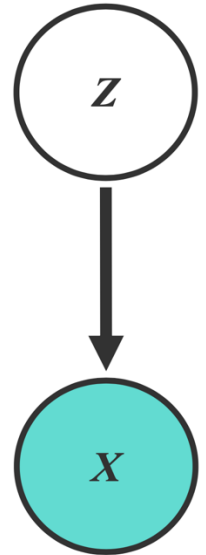
- The likelihood of a sample:

Parameters to be learned:

$$\begin{aligned} p(x_n | \mu, \Sigma) &= \sum_k p(z^k = 1 | \pi) p(x, | z^k = 1, \mu, \Sigma) \\ &= \sum_{z_n} \prod_k \left((\pi_k)^{z_n^k} N(x_n : \mu_k, \Sigma_k)^{z_n^k} \right) = \sum_k \pi_k N(x, | \mu_k, \Sigma_k) \end{aligned}$$

mixture proportion

mixture component



Example: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components
- The expected complete log likelihood

$$\begin{aligned}\mathbb{E}_q [\ell_c(\boldsymbol{\theta}; x, z)] &= \sum_n \mathbb{E}_q [\log p(z_n | \pi)] + \sum_n \mathbb{E}_q [\log p(x_n | z_n, \mu, \Sigma)] \\ &= \sum_n \sum_k \mathbb{E}_q [z_n^k] \log \pi_k - \frac{1}{2} \sum_n \sum_k \mathbb{E}_q [z_n^k] \left((x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) + \log |\Sigma_k| + C \right)\end{aligned}$$

- E-step: computing the posterior of z_n given the current estimate of the parameters (i.e., π, μ, Σ)

$$p(z_n^k = 1 | x, \mu^{(t)}, \Sigma^{(t)}) = \frac{\pi_k^{(t)} N(x_n, | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_i \pi_i^{(t)} N(x_n, | \mu_i^{(t)}, \Sigma_i^{(t)})}$$

$\nearrow p(z_n^k = 1, x, \mu^{(t)}, \Sigma^{(t)})$
 $\searrow p(x, \mu^{(t)}, \Sigma^{(t)})$

Example: Gaussian Mixture Models (GMMs)

- M-step: computing the parameters given the current estimate of Z_n

$$\begin{aligned}\pi_k^* &= \arg \max \langle l_c(\boldsymbol{\theta}) \rangle, & \Rightarrow \quad \frac{\partial}{\partial \pi_k} \langle l_c(\boldsymbol{\theta}) \rangle = 0, \forall k, \quad \text{s.t. } \sum_k \pi_k = 1 \\ & & \Rightarrow \quad \pi_k^* = \frac{\sum_n \langle Z_n^k \rangle_{q^{(t)}}}{N} = \frac{\sum_n \tau_n^{k(t)}}{N} = \frac{\langle n_k \rangle}{N}\end{aligned}$$

$$\mu_k^* = \arg \max \langle l(\boldsymbol{\theta}) \rangle, \quad \Rightarrow \quad \mu_k^{(t+1)} = \frac{\sum_n \tau_n^{k(t)} x_n}{\sum_n \tau_n^{k(t)}}$$

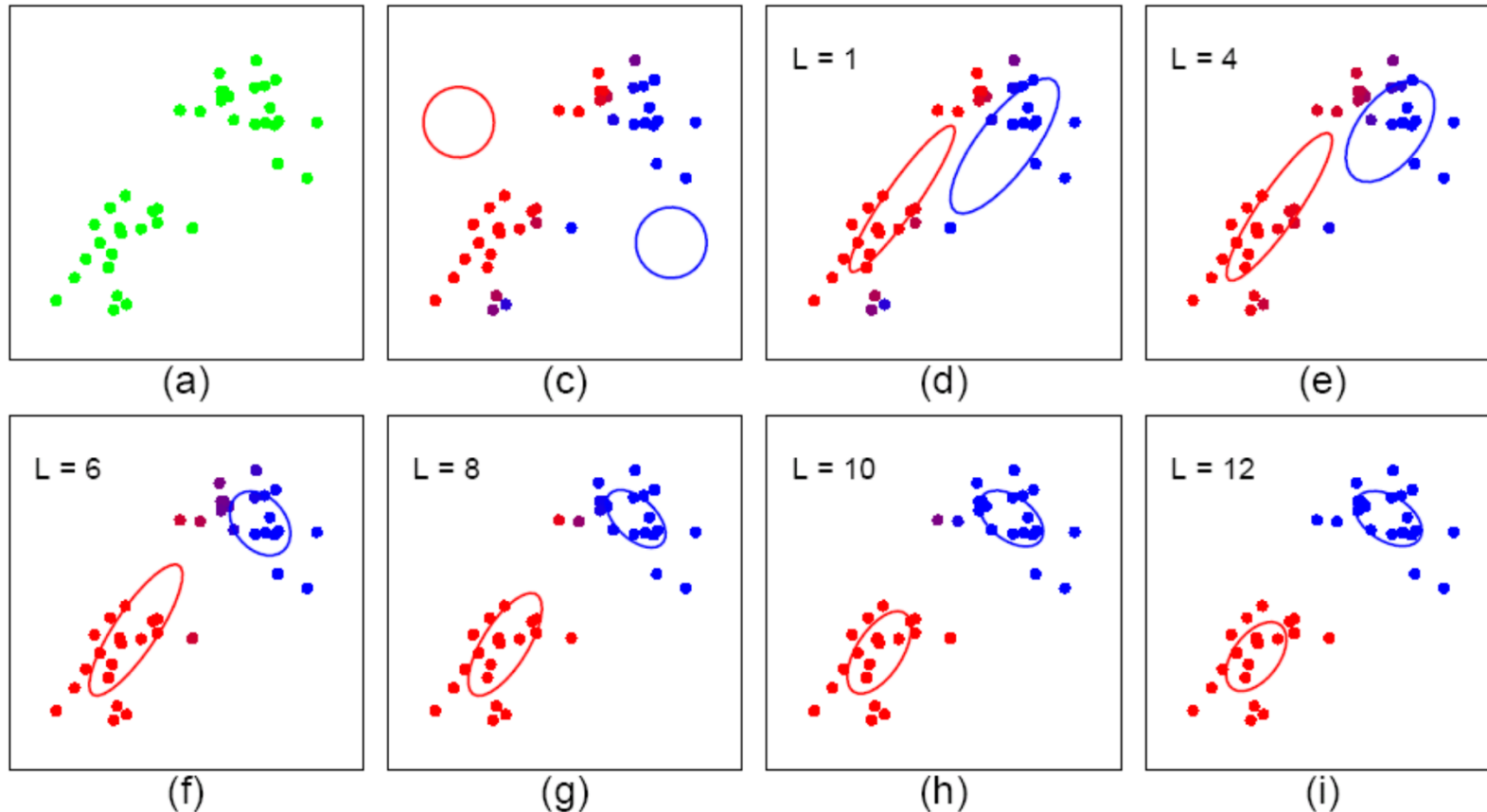
$$\Sigma_k^* = \arg \max \langle l(\boldsymbol{\theta}) \rangle, \quad \Rightarrow \quad \Sigma_k^{(t+1)} = \frac{\sum_n \tau_n^{k(t)} (x_n - \mu_k^{(t+1)})(x_n - \mu_k^{(t+1)})^T}{\sum_n \tau_n^{k(t)}}$$

Fact :

$$\begin{aligned}\frac{\partial \log |\mathbf{A}^{-1}|}{\partial \mathbf{A}^{-1}} &= \mathbf{A}^T \\ \frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{A}} &= \mathbf{x} \mathbf{x}^T\end{aligned}$$

Example: Gaussian Mixture Models (GMMs)

- Start: “guess” the centroid μ_k and covariance Σ_k of each of the K clusters
- Loop:



Summary: EM Algorithm

- A way of maximizing likelihood function for latent variable models. Finds MLE of parameters when the original (hard) problem can be broken up into two (easy) pieces
 - Estimate some “missing” or “unobserved” data from observed data and current parameters.
 - Using this “complete” data, find the maximum likelihood parameter estimates.
- Alternate between filling in the latent variables using the best guess (posterior) and updating the parameters based on this guess:
 - E-step:
 - M-step: $q^{t+1} = \arg \min_q F(q, \theta^t)$
 $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta)$

Questions?