

# DSC291: Machine Learning with Few Labels

## Unsupervised Learning

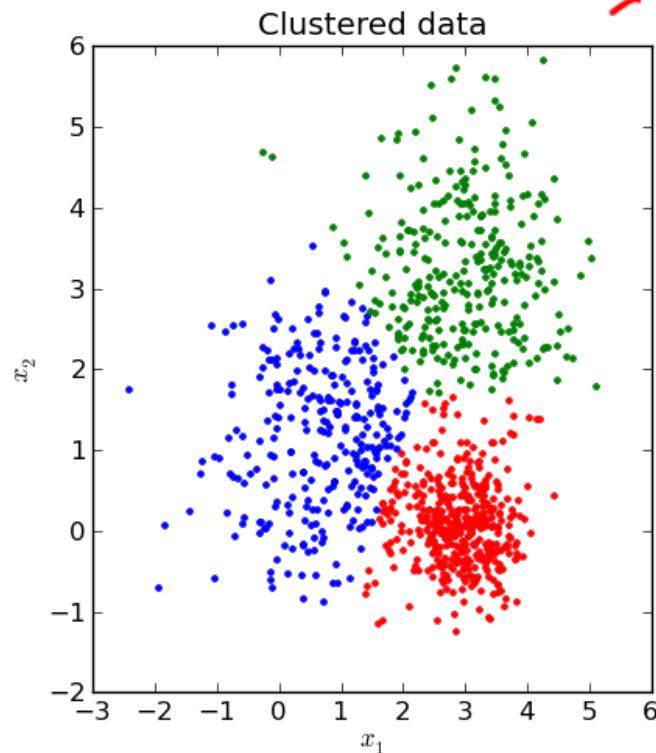
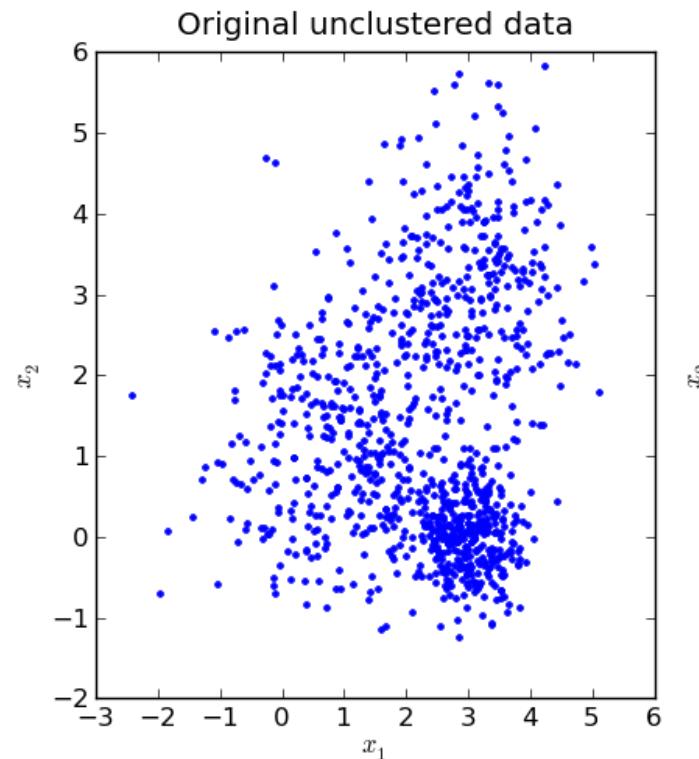
Zhitong Hu

Lecture 4, April 10, 2025

# Unsupervised Learning

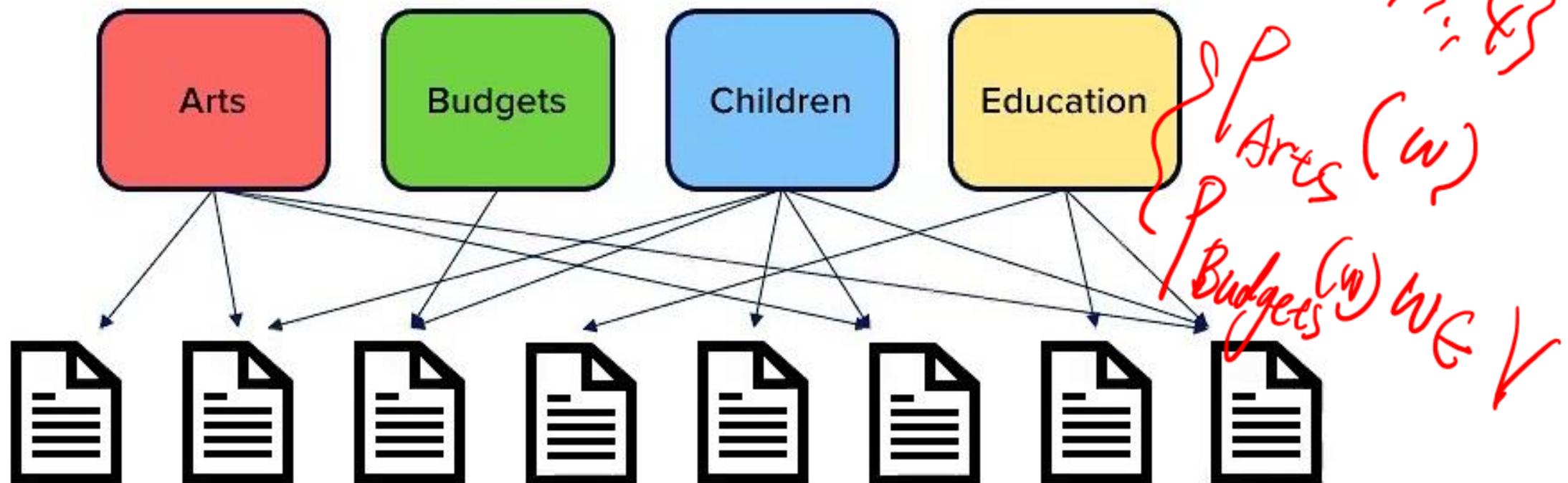
- Each data instance is partitioned into two parts:
  - observed variables  $x$
  - latent (unobserved) variables  $z$
- Want to learn a model  $p_\theta(x, z)$

$$\begin{aligned} \vec{x} &= (x_1, x_2) \\ Z \in \{1, 2, 3\} \end{aligned}$$



# Unsupervised Learning

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



# Unsupervised Learning

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

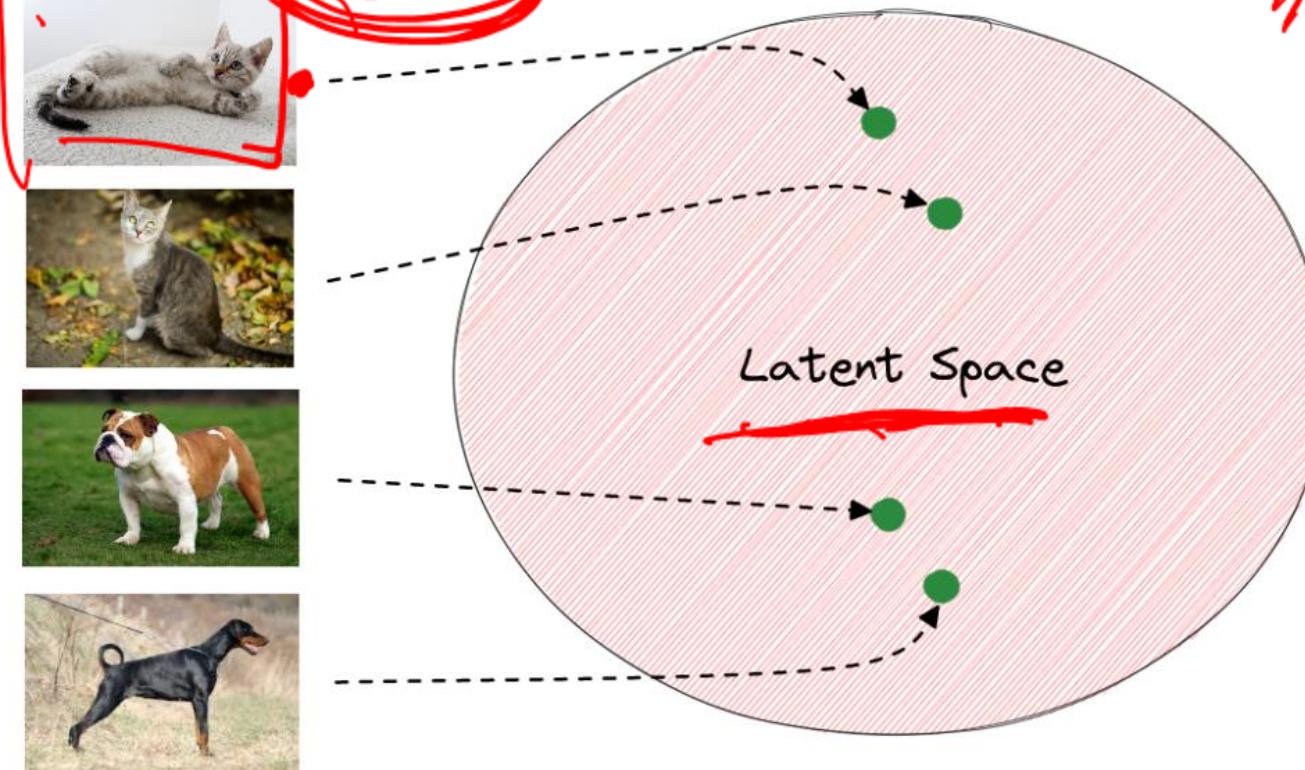


# Unsupervised Learning

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

Latent-space  
model

deep generative  
models



~~VAE GAN~~

~~Diffusion~~

$p_\theta(x|z)p(z)$

$x = \text{image}$

$z \in \mathbb{R}^{100}$

$$p_\theta(x^+ | z)$$

$$P(x, z|\theta) = P(x|z, \theta) \cdot P(z|\theta)$$

*Supervised, complete data  
(x, 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 parameters 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:

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

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

$$P(x, z|\theta) \\ P(x|\theta) = \sum_z P(x, z|\theta)$$

$z \in \mathbb{R}^{100}$

$$\nabla_{\theta_z} \int_{\mathcal{Z}} P(x, z | \theta) dz$$

# Why is Unsupervised Learning Harder?

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

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

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

Now  $z$  is not observed:

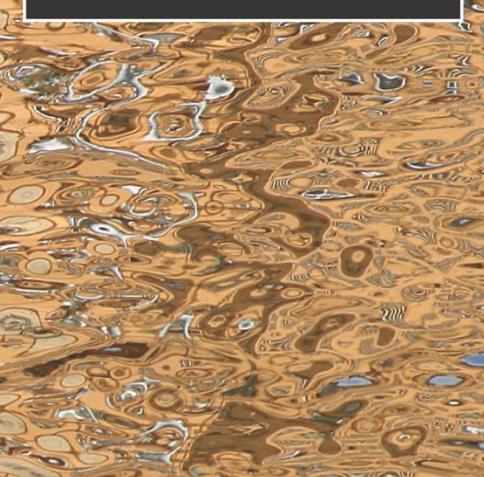
- **Incomplete (or marginal) log likelihood:** with  $z$  unobserved, our objective becomes the log of a marginal probability:

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

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

$$\nabla_{\theta_z} \ell(\theta; x) = \nabla_{\theta_z} \log \sum_z P(x, z | \theta_z, \theta_x)$$

$$= \frac{\nabla_{\theta_z} \sum_z P(x, z | \theta_z, \theta_x)}{\sum_z P(x, z | \theta_z, \theta_x)}$$



# Expectation Maximization (EM)

|           |   |            |
|-----------|---|------------|
| <b>9</b>  | <b>Mixture Models and EM</b>                  | <b>423</b> |
| 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        |
| 9.4       | The EM Algorithm in General . . . . .         | 450        |
| Exercises | . . . . .                                     | 455        |

This class

EM → VAEs → diffusion

# Expectation Maximization (EM): Intuition

# Expectation Maximization (EM): Intuition

- Supervised MLE is easy:

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

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

(c)

- Unsupervised MLE is hard:

- Observe only  $\mathbf{x}$

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

- EM, intuitively:

**E-step:**

$$q(\mathbf{z} | \mathbf{x}) = \underline{p(\mathbf{z} | \mathbf{x}, \theta)}$$

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)

$$\rightarrow \mathcal{E}(\mathcal{C}(\mathbf{x}))$$

# Expectation Maximization (EM): Intuition

- Supervised MLE is easy:
  - Observe both  $x$  and  $z$
- Unsupervised MLE is hard:
  - Observe only  $x$
- EM, intuitively:

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

*Posterior,*

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

$$q^*(z|x) = p(z|x, \theta^*)$$

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

$$q(z|x)$$



E-step:

$$q^{t=1}(z|x)$$

$$q^{t=2}(z|x)$$

$$q^{t+1}(z|x) = p(z|x, \theta^t)$$

M-step:

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

$$\theta^{t=1}$$

$$\theta^{t=2}$$

...

This is an iterative process

Let's "pretend" we also observe  $z$  (its distribution)

$$\rightarrow \theta^* \quad q^*(z|x)$$

# Expectation Maximization (EM): Intuition

- Supervised MLE is easy:
  - Observe both  $\mathbf{x}$  and  $\mathbf{z}$
- Unsupervised MLE is hard:
  - Observe only  $\mathbf{x}$
- EM, intuitively:

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

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



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

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

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

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

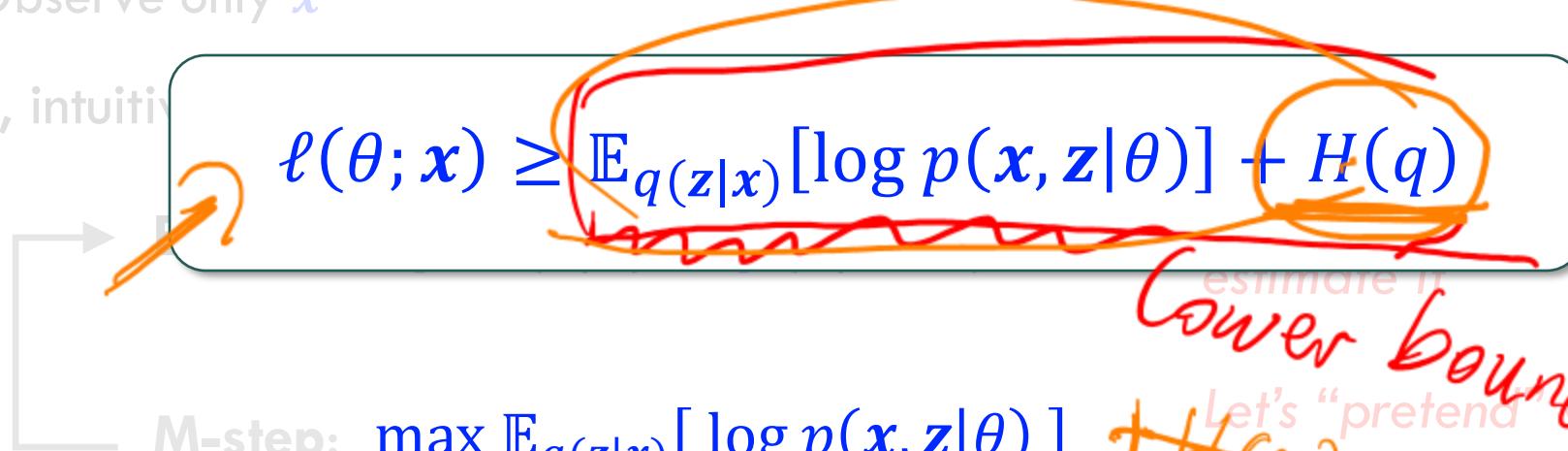
This is an iterative process

# Expectation Maximization (EM): Intuition

- ~~Supervised MLE is easy:~~  
o Observe both  $x$  and  $z$
- Unsupervised MLE is hard:  
o Observe only  $x$
- EM, intuitively

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

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



This is an iterative process

Only observe  $q$ , let's estimate it

Let's "pretend" we also observe  $z$  (its distribution)

# Expectation Maximization (EM): Intuition

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

## Expectation Maximization (EM): Intuition

$$-\mathbb{E}_{q(z|x)} [\log q(z|x)]$$

- **Question:** show that

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

# Expectation Maximization (EM): Intuition

- **Question:** show that

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

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

- **Hint:** first show that

$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(z|x)} \left[ \log \frac{p(\mathbf{x}, z|\theta)}{q(z|x)} \right] + \text{KL}(q(z|x) || p(z|x, \theta)) \geq 0$$

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

# Expectation Maximization (EM): Intuition

- **Question:** show that

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

- **Hint:** first show that

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

Variational free energy

Evidence Lower Bound (ELBO)

$$= -F(q, \theta) + \text{KL}(q(z|x) \parallel p(z|x, \theta))$$

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

# Lower Bound and Free Energy

- Variational free energy:

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

- The EM algorithm is coordinate-decent on  $F$

- At each step  $t$ :

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

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

# 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:  $\underset{\longleftarrow}{q^{t+1}} = \arg \min_q F(q, \theta^t)$

$$\stackrel{?}{\rightarrow} q^{t+1}(z|x) = p(z|x, \theta^t)$$

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

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

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

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

- I.e., the posterior distribution over the latent variables given the data and the current parameters.
- **Hint:** use the fact

$$\begin{aligned} F &= -\ell + \text{KL} \\ \ell(\theta^t; x) &= -F(q, \theta^t) + \text{KL}(q(z|x) \parallel p(z|x, \theta)) \\ &\quad \downarrow \\ &\quad \text{Independent of } q \\ &\quad \downarrow \\ &\geq 0 \end{aligned}$$

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

# Lower Bound and Free Energy

- Variational free energy:

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

- 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:  $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta^t)$

?  $\rightarrow \max_{\theta} \mathbb{E}_{q^{t+1}(z|x)} [\log p(x, z|\theta)]$

# EM Algorithm: Quick Summary

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

- But it's too difficult  $\ell(\theta; x) = \log p(x|\theta) = \log \sum_z p(x, z|\theta)$

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

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

Variational free energy

## EM Algorithm: Quick Summary

- The EM algorithm is coordinate-decent on  $F(q, \theta)$

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

- the posterior distribution over the latent variables given the data and the current parameters

- M-step:  $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta^t) = \operatorname{argmax}_{\theta} \sum_z q^{t+1}(z|x) \log p(x, z|\theta)$

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

$$= -F(q, \theta) + \text{KL}(q(z|x) || p(z|x, \theta))$$

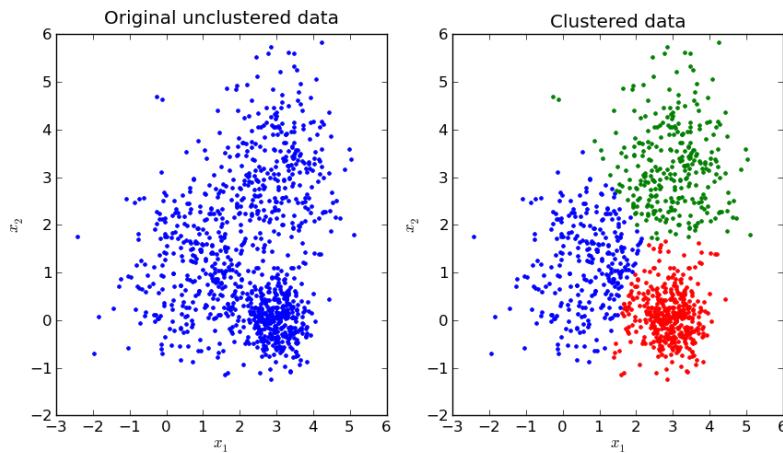
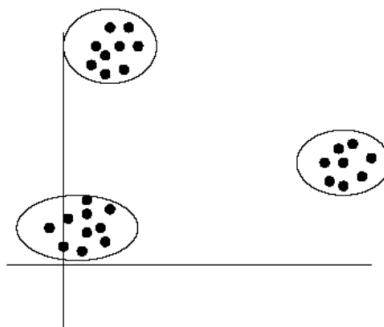
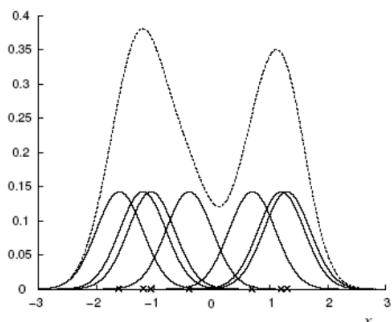
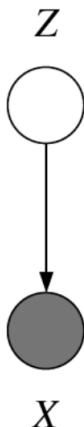
~~Surrogate objective~~

# Example: Gaussian Mixture Models

- Consider a mixture of K Gaussian components:

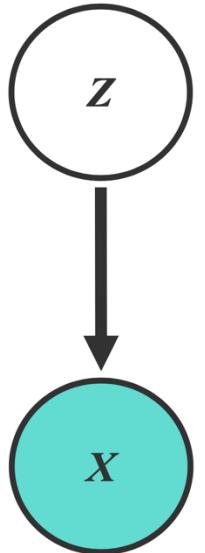
$$p(x_n | \mu, \Sigma) = \sum_k \pi_k N(x_n | \mu_k, \Sigma_k)$$

↑                                   ↑  
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:

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

mixture component   
mixture proportion 

# Example: Gaussian Mixture Models (GMMs)

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

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

- E-step: computing the posterior of  $z_n$  given the current estimate of the parameters (i.e.,  $\pi, \mu, \Sigma$ )

$$p(z_n^k = 1 | \mathbf{x}, \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)}) = \frac{\pi_k^{(t)} N(\mathbf{x}_n | \boldsymbol{\mu}_k^{(t)}, \boldsymbol{\Sigma}_k^{(t)})}{\sum_i \pi_i^{(t)} N(\mathbf{x}_n | \boldsymbol{\mu}_i^{(t)}, \boldsymbol{\Sigma}_i^{(t)})}$$

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

$p(\mathbf{x}, \boldsymbol{\mu}^{(t)}, \boldsymbol{\Sigma}^{(t)})$

# Example: Gaussian Mixture Models (GMMs)

- M-step: computing the parameters given the current estimate of  $z_n$

$$\pi_k^* = \arg \max \langle l_c(\theta) \rangle, \quad \Rightarrow \quad \frac{\partial}{\partial \pi_k} \langle l_c(\theta) \rangle = 0, \forall k, \quad \text{s.t.} \sum_k \pi_k = 1$$

$$\Rightarrow \pi_k^* = \left. \frac{\sum_n \langle z_n^k \rangle_{q^{(t)}}}{N} \right/ = \left. \frac{\sum_n \tau_n^{k(t)}}{N} \right/ = \left. \frac{\langle n_k \rangle}{N} \right/$$

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

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

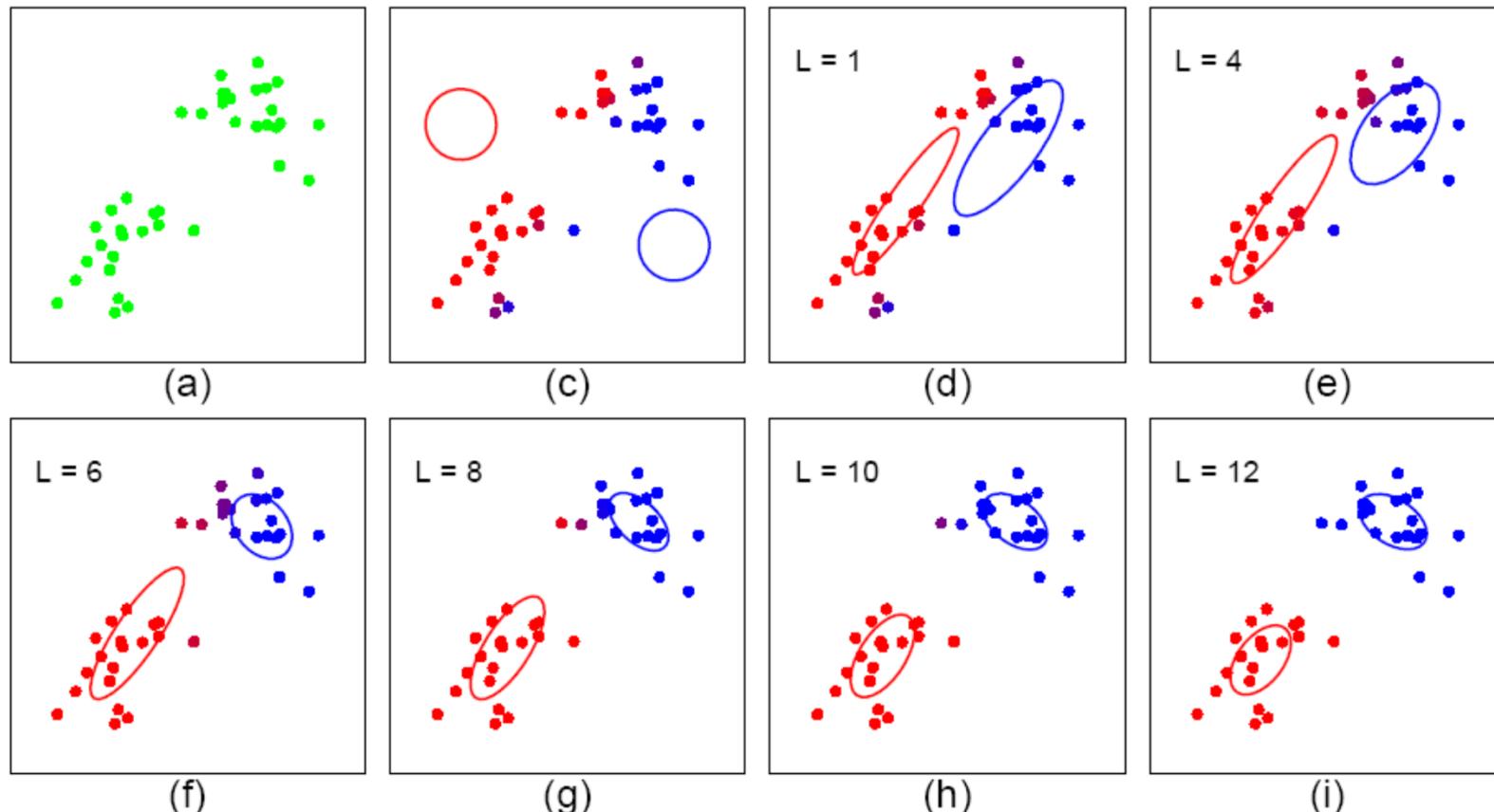
Fact:

$$\frac{\partial \log |A^{-1}|}{\partial A^{-1}} = A^T$$

$$\frac{\partial \mathbf{x}^T A \mathbf{x}}{\partial A} = \mathbf{x} \mathbf{x}^T$$

# Example: Gaussian Mixture Models (GMMs)

- Start: “guess” the centroid  $\mu_k$  and covariance  $\Sigma_k$  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?