

# Lecture 18.

## Gaussian Mixture Model.

## Expectation Maximization.

COMP90051 Statistical Machine Learning

Semester 2, 2019  
Lecturer: Ben Rubinstein



THE UNIVERSITY OF  
MELBOURNE

# This lecture

- Unsupervised learning
  - \* Diversity of problems
- Gaussian mixture model (GMM)
  - \* A probabilistic approach to clustering
  - \* The GMM model
  - \* GMM clustering as an optimisation problem
- The Expectation Maximization (EM) algorithm

# Unsupervised Learning

A large branch of ML that concerns  
with learning the structure of the  
data in the absence of labels

# Previously: Supervised learning

- Supervised learning: Overarching aim is making predictions from data
- We studied methods such as random forest, ANN and SVM in the context of this aim
- We had instances  $\mathbf{x}_i \in \mathbf{R}^m$ ,  $i = 1, \dots, n$  and corresponding labels  $y_i$  as inputs, and the aim was to predict labels for new instances
- Can be viewed as a function approximation problem, but with a big caveat: ability to generalise is critical
- Bandits: a setting of partial supervision

# Now: Unsupervised learning

- Next few lectures: unsupervised learning methods
- In unsupervised learning, there is no dedicated variable called a “label”
- Instead, we just have a set of points  $\mathbf{x}_i \in \mathbf{R}^m$ ,  $i = 1, \dots, n$
- The aim of unsupervised learning is to **explore the structure** (patterns, regularities) of the data
- The aim of “exploring the structure” is vague

# Unsupervised learning tasks

- Diversity of tasks fall into unsupervised learning category
  - \* Clustering (now)
  - \* Dimensionality reduction (soon)
  - \* Learning parameters of probabilistic models (before/now)
- Applications and related tasks are numerous :
  - \* Market basket analysis. E.g., use supermarket transaction logs to find items that are frequently purchased together
  - \* Outlier detection. E.g., find potentially fraudulent credit card transactions
  - \* Often unsupervised tasks in (supervised) ML pipelines

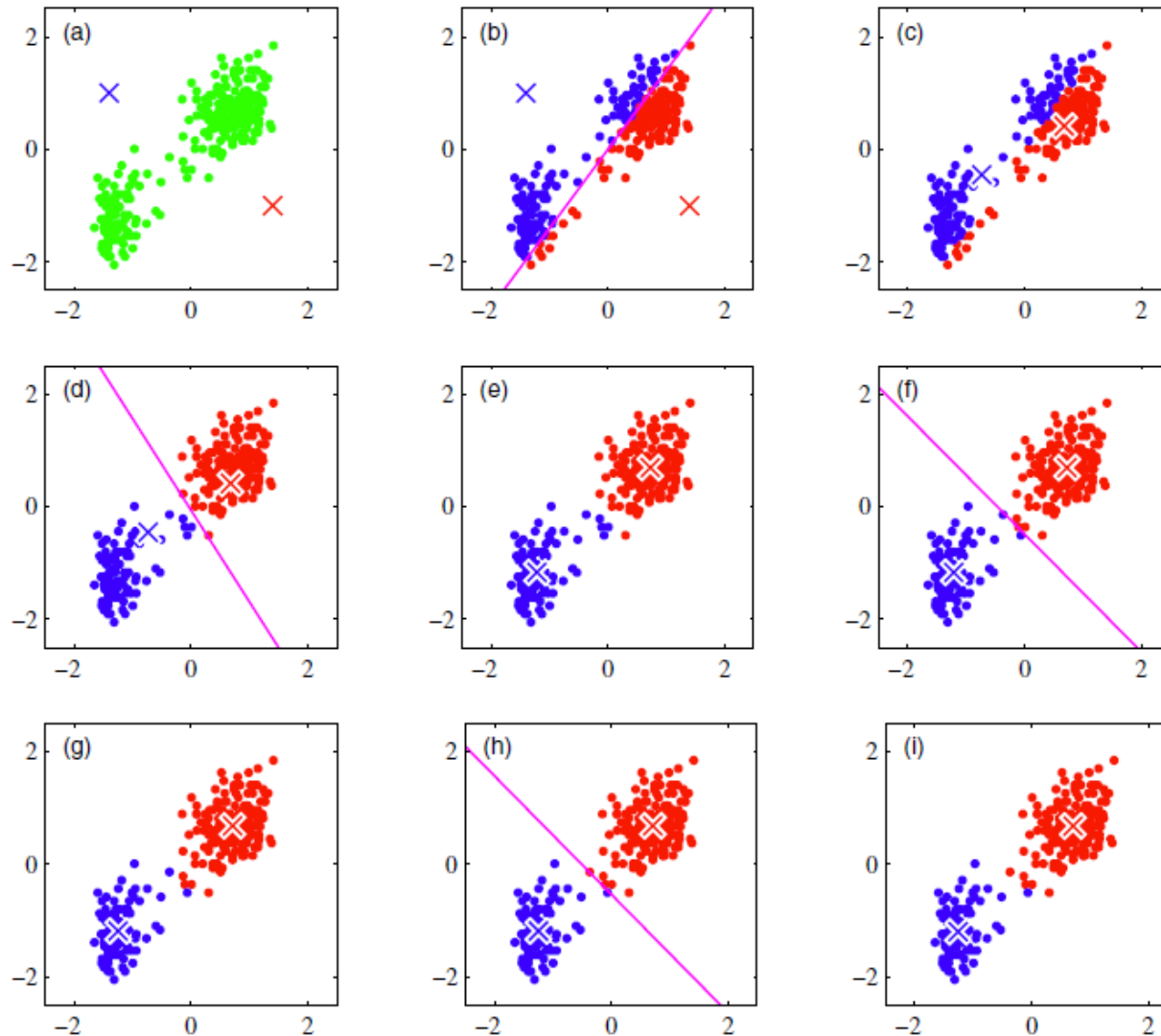
# Refresher: k-means clustering

1. Initialisation: choose  $k$  cluster **centroids** randomly
2. Update:
  - a) **Assign points** to the nearest\* centroid
  - b) **Compute centroids** under the current assignment
3. Termination: if no change then **stop**
4. Go to **Step 2**

\*Distance represented by choice of metric typically  $L_2$

Still one of the most popular data mining algorithms.

# Refresher: k-means clustering



Requires specifying the number of clusters in advance

Measures “dissimilarity” using Euclidean distance

Finds “spherical” clusters

An iterative optimization procedure

Data: Old Faithful  
Geyser Data: waiting time between eruptions and the duration of eruptions



# Gaussian Mixture Model

A probabilistic view of clustering

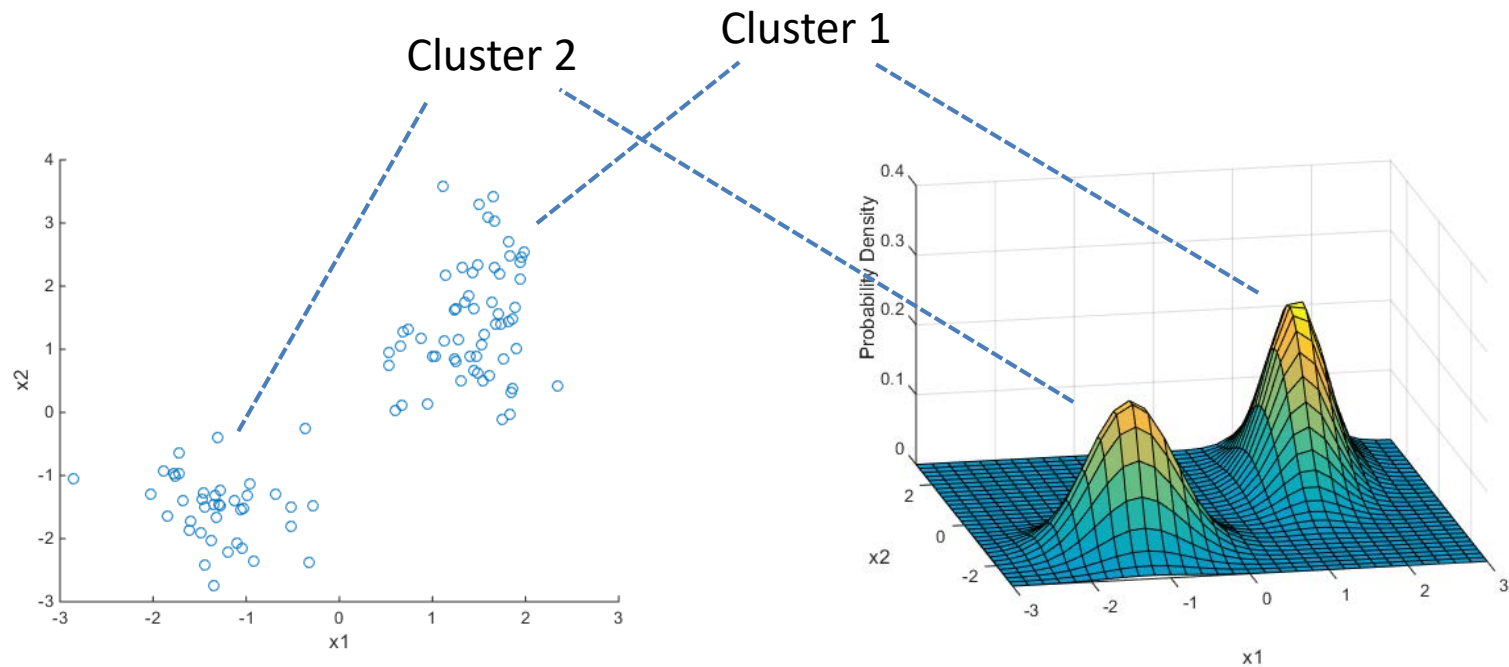
# Modelling uncertainty in data clustering

- k-means clustering assigns each point to exactly one cluster
- Similar to k-means, a probabilistic mixture model requires the user to choose the number of clusters in advance
- Unlike k-means, the probabilistic model gives us a power to express **uncertainly about the origin** of each point
  - \* Each point originates from cluster  $c$  with probability  $w_c$ ,  $c = 1, \dots, k$
- That is, each point still originates from one particular cluster (aka component), but we are not sure from which one
- Next
  - \* Individual components modelled as Gaussians
  - \* Fitting illustrates general Expectation Maximization (EM) algorithm

# Clustering: probabilistic model

Data points  $x_i$  are independent and identically distributed (i.i.d.) samples from a **mixture** of  $K$  distributions (components)

Each component in the mixture is what we call a cluster



In principle, we can adopt any probability distribution for the **components**, however, the normal distribution is a common modelling choice → Gaussian Mixture Model

# Normal (aka Gaussian) distribution

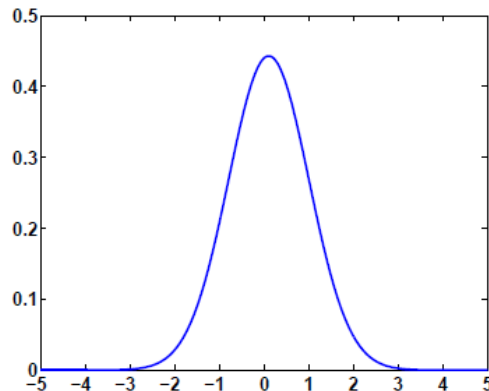
- Recall that a 1D Gaussian is

$$\mathcal{N}(x|\mu, \sigma) \equiv \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

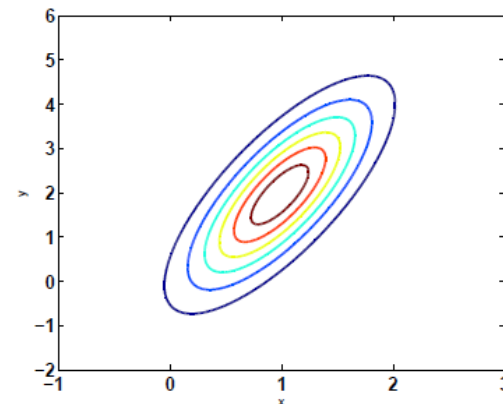
- And a  $d$ -dimensional Gaussian is

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \equiv (2\pi)^{-\frac{d}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$$

- \*  $\boldsymbol{\Sigma}$  is a PSD symmetric  $d \times d$  matrix, the **covariance matrix**
- \*  $|\boldsymbol{\Sigma}|$  denotes determinant
- \* No need to memorize the full formula.



(a) 1-Dim

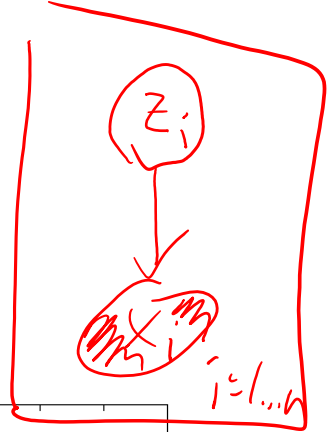


(b) 2-Dim

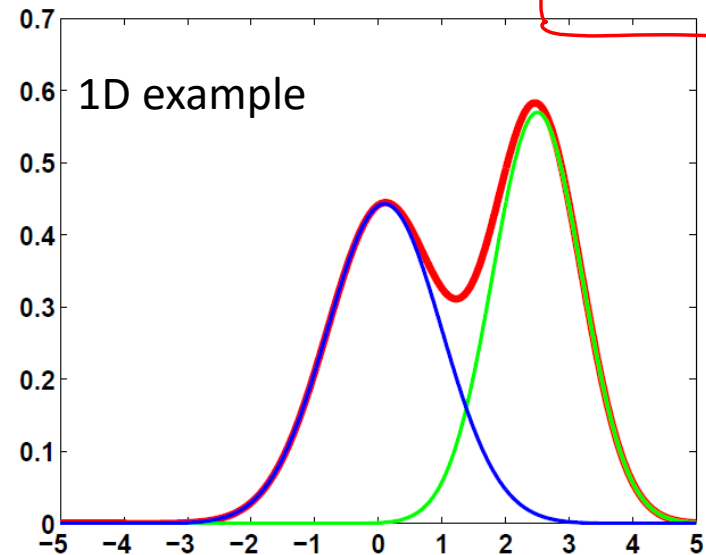
# Gaussian mixture model (GMM)

- Gaussian mixture distribution (for one data point):

$$P(\mathbf{x}) \equiv \sum_{j=1}^k w_j \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \equiv \sum_{j=1}^k P(C_j) P(\mathbf{x} | C_j)$$



- $P(\mathbf{x} | C_j) \equiv \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$  is class/component conditional density (modeled as a Gaussian) for class  $j$
- Here  $P(C_j) \geq 0$  and  $\sum_{j=1}^k P(C_j) = 1$
- Parameters of the model are  $P(C_j)$ ,  $\boldsymbol{\mu}_j$ ,  $\boldsymbol{\Sigma}_j$ ,  $j = 1, \dots, k$



Mixture and individual component densities are re-scaled for visualisation purposes

Figure: Bishop

**Consider a GMM with five components for 3D data. How many independent scalar parameters does this model have?**

$$49 = 6 \times 5 + 3 \times 5 + 4$$

$$50 = 6 \times 5 + 3 \times 5 + 5$$

$\sum w_c$

$$65 = 9 \times 5 + 3 \times 5 + 5$$

$\mu_c$

# Clustering as model estimation

- Given a set of data points, we assume that data points are generated by a GMM
  - \* Each point in our dataset originates from the  $j$ -th normal distribution component with probability  $w_j = P(C_j)$
- Clustering now amounts to finding parameters of the GMM that “best explain” the observed data
- Call upon old friend **MLE** principle to find parameter values that maximise  $p(\mathbf{x}_1, \dots, \mathbf{x}_n)$

# Fitting the GMM

- Modelling the data points as independent, aim is to find  $\mathbf{P}(\mathbf{C}_j)$ ,  $\boldsymbol{\mu}_j$ ,  $\boldsymbol{\Sigma}_j$ ,  $j = 1, \dots, k$  that maximise
$$P(\mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{i=1}^n \sum_{j=1}^k P(C_j) P(\mathbf{x}_i | C_j)$$
where  $P(\mathbf{x} | \mathbf{C}_j) \equiv \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$   
Can be solved analytically?
- Taking the derivative of this expression is pretty awkward, **try the usual log trick**

$$\log P(\mathbf{x}_1, \dots, \mathbf{x}_n) = \sum_{i=1}^n \log \left( \sum_{j=1}^k P(C_j) P(\mathbf{x}_i | \mathbf{C}_j) \right)$$


→ Expectation-Maximisation (EM)



# Expectation Maximisation Algorithm

For a moment, let's put GMM problem  
aside – to come back to later.

# Motivation of EM

- Consider a parametric probabilistic model  $p(\mathbf{X}|\boldsymbol{\theta})$ , where  $\mathbf{X}$  denotes data and  $\boldsymbol{\theta}$  denotes a vector of parameters
- According to MLE, we need to maximise  $p(\mathbf{X}|\boldsymbol{\theta})$  as a function of  $\boldsymbol{\theta}$ 
  - \* equivalently maximise  $\log p(\mathbf{X}|\boldsymbol{\theta})$
- There can be a couple of issues with this task 
  1. Sometimes we **don't observe** some of the variables needed to compute the log likelihood
    - \* Example: GMM cluster membership is not known in advance
  2. Sometimes the form of the log likelihood is **inconvenient** to work with
    - \* Example: taking a derivative of GMM log likelihood results in a cumbersome equation

# MLE vs EM

- MLE is a frequentist *principle* that suggests that given a dataset, the “best” parameters to use are the ones that maximise the probability of the data
  - \* MLE is a way *to formally pose* the problem
- EM is an *algorithm*
  - \* EM is a way *to solve* the problem posed by MLE
  - \* Especially convenient under unobserved latent variables
- MLE can be found by other methods such as gradient descent (but gradient descent is not always the most convenient method)

# Expectation-Maximisation (EM) Algorithm

- Initialisation Step:

- \* Initialize  $K$  clusters:  $C_1, \dots, C_K$   
 $(\mu_j, \Sigma_j)$  and  $P(C_j)$  for each cluster  $j$ .

- Iteration Step:

- \* Estimate the cluster of each datum

$$p(C_j | x_i)$$

 Expectation

- \* Re-estimate the cluster parameters

 Maximisation

$$(\mu_j, \Sigma_j), p(C_j) \text{ for each cluster } j$$

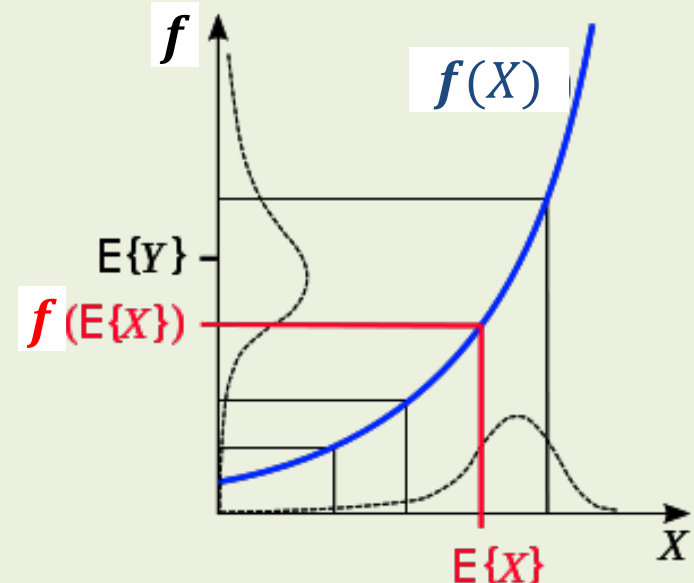
# EM for GMM and generally

- EM is a general approach, goes beyond GMMs
  - \* Purpose: Implement **MLE under latent variables  $\mathbf{Z}$**  ('latent' is fancy for 'missing')
- What are variables, parameters in GMMs?
  - \* Variables: Point locations  $\mathbf{X}$  and cluster assignments  $\mathbf{Z}$ 
    - let  $z_i$  denote true cluster membership for each point  $x_i$ , computing the likelihood with known values  $\mathbf{z}$  is simplified (see next section)
  - \* Parameters:  $\boldsymbol{\theta}$  are cluster locations and scales
- What is EM really doing?
  - \* **Coordinate ascent** on a lower bound on the log-likelihood
    - M-step: ascent in modeled parameters  $\boldsymbol{\theta}$
    - E-step: ascent in the marginal likelihood  $P(\mathbf{Z})$
  - \* Each step moves towards a **local** optimum
  - \* Can get stuck, can need **random restarts**

# Needed tool: Jensen's inequality

- Compares effect of averaging before and after applying a **convex function**:  

$$f(\text{Average}(\mathbf{x})) \leq \text{Average}(f(\mathbf{x}))$$
- Example:
  - \* Let  $f$  be some convex function, such as  $f(x) = x^2$
  - \* Consider  $\mathbf{x} = [1, 2, 3, 4, 5]'$ , then  $f(\mathbf{x}) = [1, 4, 9, 16, 25]'$
  - \* Average of input  $\text{Average}(\mathbf{x}) = 3$
  - \*  $f(\text{Average}(\mathbf{x})) = 9$
  - \* Average of output  $\text{Average}(f(\mathbf{x})) = 12.4$
- Proof follows from the definition of convexity
  - \* Proof by induction
- General statement:
  - \* If  $\mathbf{X}$  random variable,  $f$  is a convex function
  - \*  $f(\mathbb{E}[\mathbf{X}]) \leq \mathbb{E}[f(\mathbf{X})]$



# Putting the latent variables to use

- We want to maximise  $\log p(\mathbf{X}|\boldsymbol{\theta})$ . We don't observe  $\mathbf{Z}$  (here discrete), but can introduce it nonetheless.
- $\log p(\mathbf{X}|\boldsymbol{\theta}) = \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ 
  - ← Marginalisation (here  $\sum_{\mathbf{Z}} \dots$  iterates over all possible values of  $\mathbf{Z}$ )
- $= \log \sum_{\mathbf{Z}} \left( p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \frac{p(\mathbf{Z})}{p(\mathbf{Z})} \right)$ 
  - ← Need  $\mathbf{Z}$  to have non-zero marginal
- $= \log \sum_{\mathbf{Z}} \left( p(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{p(\mathbf{Z})} \right)$
- $= \log \mathbb{E}_{\mathbf{Z}} \left[ \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{p(\mathbf{Z})} \right]$
- $\geq \mathbb{E}_{\mathbf{Z}} \left[ \log \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{p(\mathbf{Z})} \right]$ 
  - ← Jensen's inequality holds since  $\log(\dots)$  is a concave function
- $= \mathbb{E}_{\mathbf{Z}} [\log p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})] - \mathbb{E}_{\mathbf{Z}} [\log p(\mathbf{Z})]$

# Maximising the lower bound (1/2)

- $\log p(\mathbf{X}|\boldsymbol{\theta}) \geq \mathbb{E}_{\mathbf{Z}}[\log p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})] - \mathbb{E}_{\mathbf{Z}}[\log p(\mathbf{Z})]$
- The right hand side (RHS) is a **lower bound** on the original log likelihood
  - \* This holds for any  $\boldsymbol{\theta}$  and any non zero  $p(\mathbf{Z})$
- Intuitively, we want to push the lower bound up
- This lower bound is a function of **two “variables”  $\boldsymbol{\theta}$  and  $p(\mathbf{Z})$** . We want to maximise the RHS as a function of these two “variables”
- It is hard to optimise with respect to both at the same time, so EM resorts to an iterative procedure



# Maximising the lower bound (2/2)

- $\log p(\mathbf{X}|\boldsymbol{\theta}) \geq \mathbb{E}_{\mathbf{Z}}[\log p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})] - \mathbb{E}_{\mathbf{Z}}[\log p(\mathbf{Z})]$

- EM is essentially **coordinate ascent**:

- \* Fix  $\boldsymbol{\theta}$  and optimise the lower bound for  $p(\mathbf{Z})$
- \* Fix  $p(\mathbf{Z})$  and optimise for  $\boldsymbol{\theta}$

we will  
prove this  
shortly

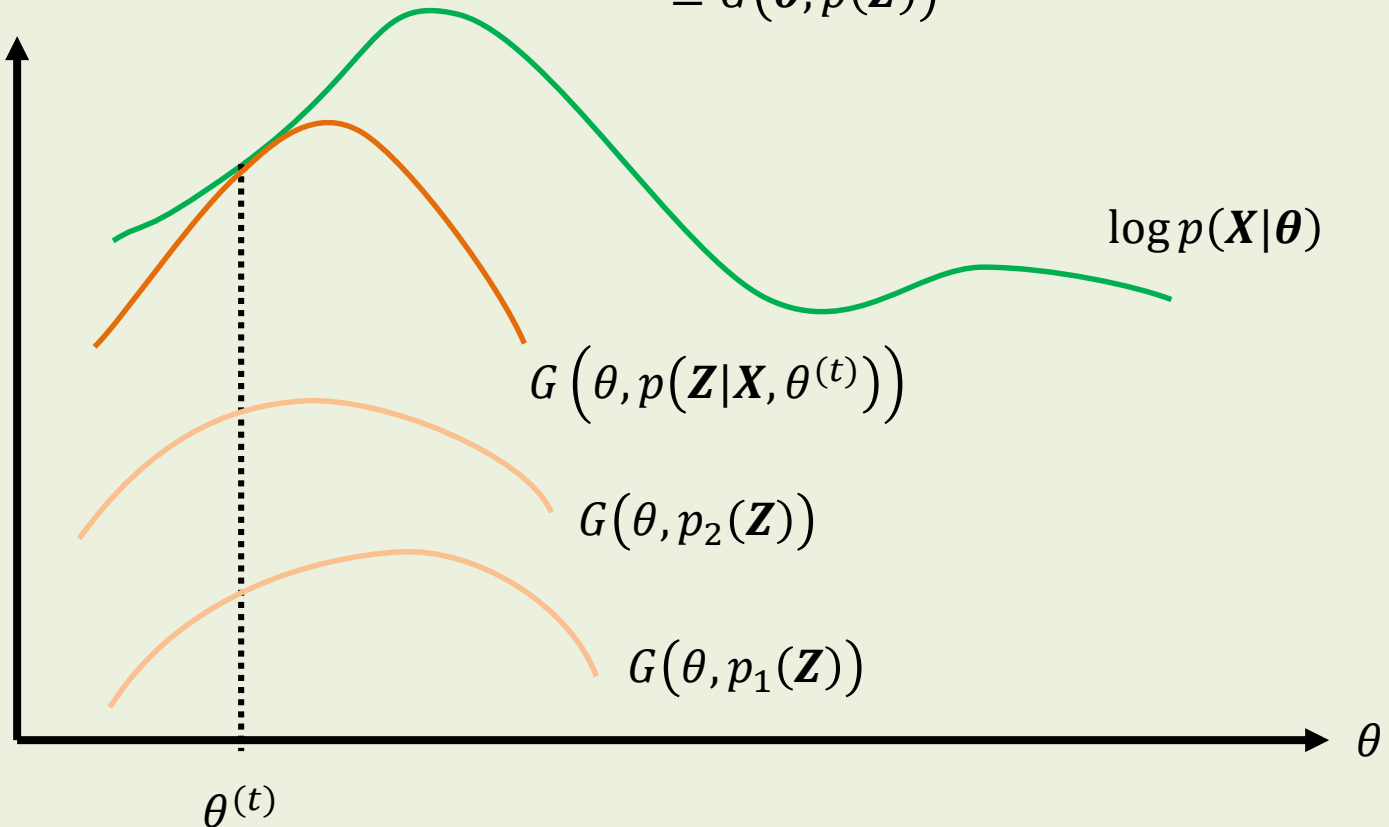
- The convenience of EM comes from the following

- For any point  $\boldsymbol{\theta}^*$ , it can be shown that setting  $p(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^*)$  makes the lower bound tight

- For any  $p(\mathbf{Z})$ , the second term does not depend on  $\boldsymbol{\theta}$
- When  $p(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^*)$ , the first term can usually be maximised as a function of  $\boldsymbol{\theta}$  in a closed-form
  - \* If not, then probably don't use EM

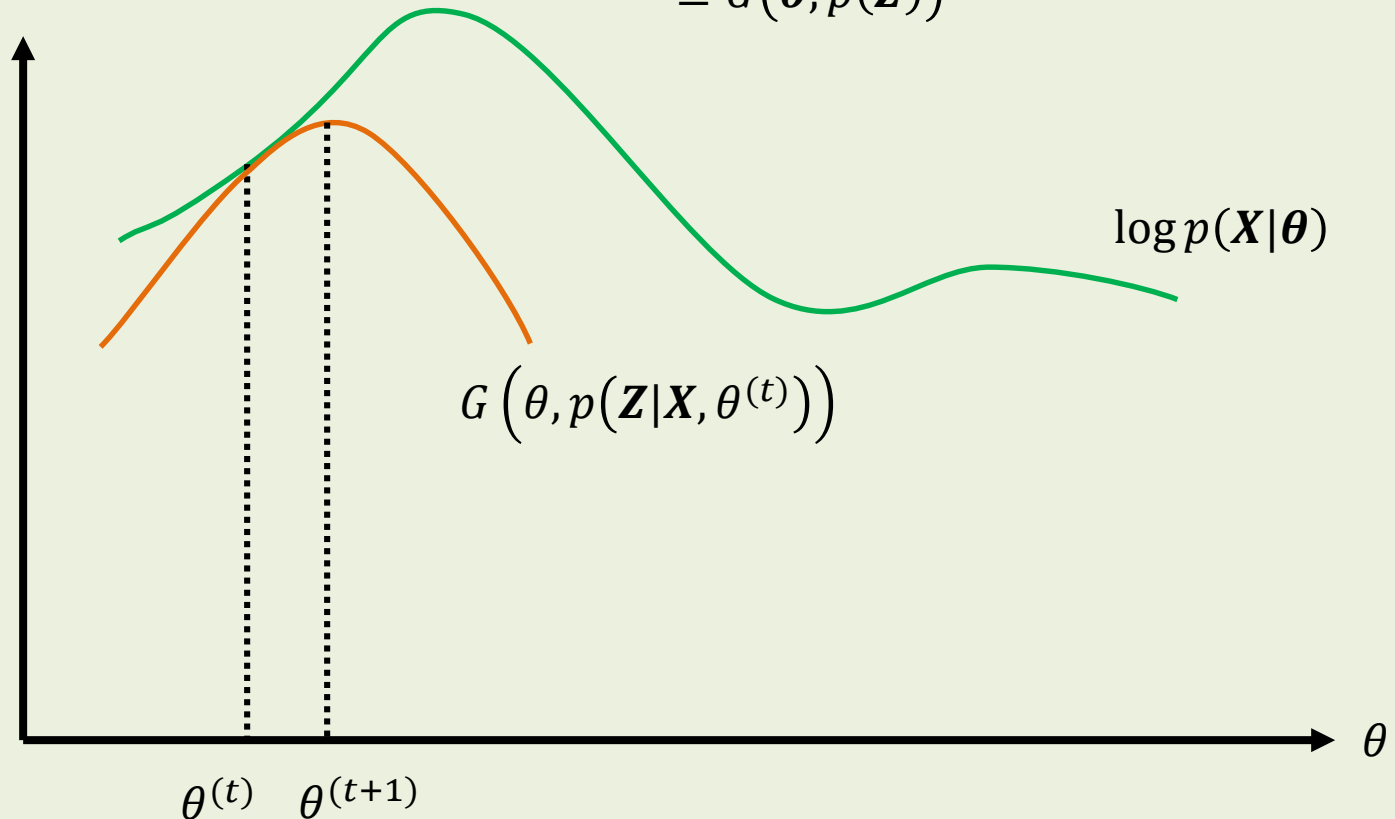
# Example (1/3)

$$\log p(\mathbf{X}|\boldsymbol{\theta}) \geq \underbrace{\mathbb{E}_{\mathbf{Z}}[\log p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})] - \mathbb{E}_{\mathbf{Z}}[\log p(\mathbf{Z})]}_{\equiv G(\boldsymbol{\theta}, p(\mathbf{Z}))}$$



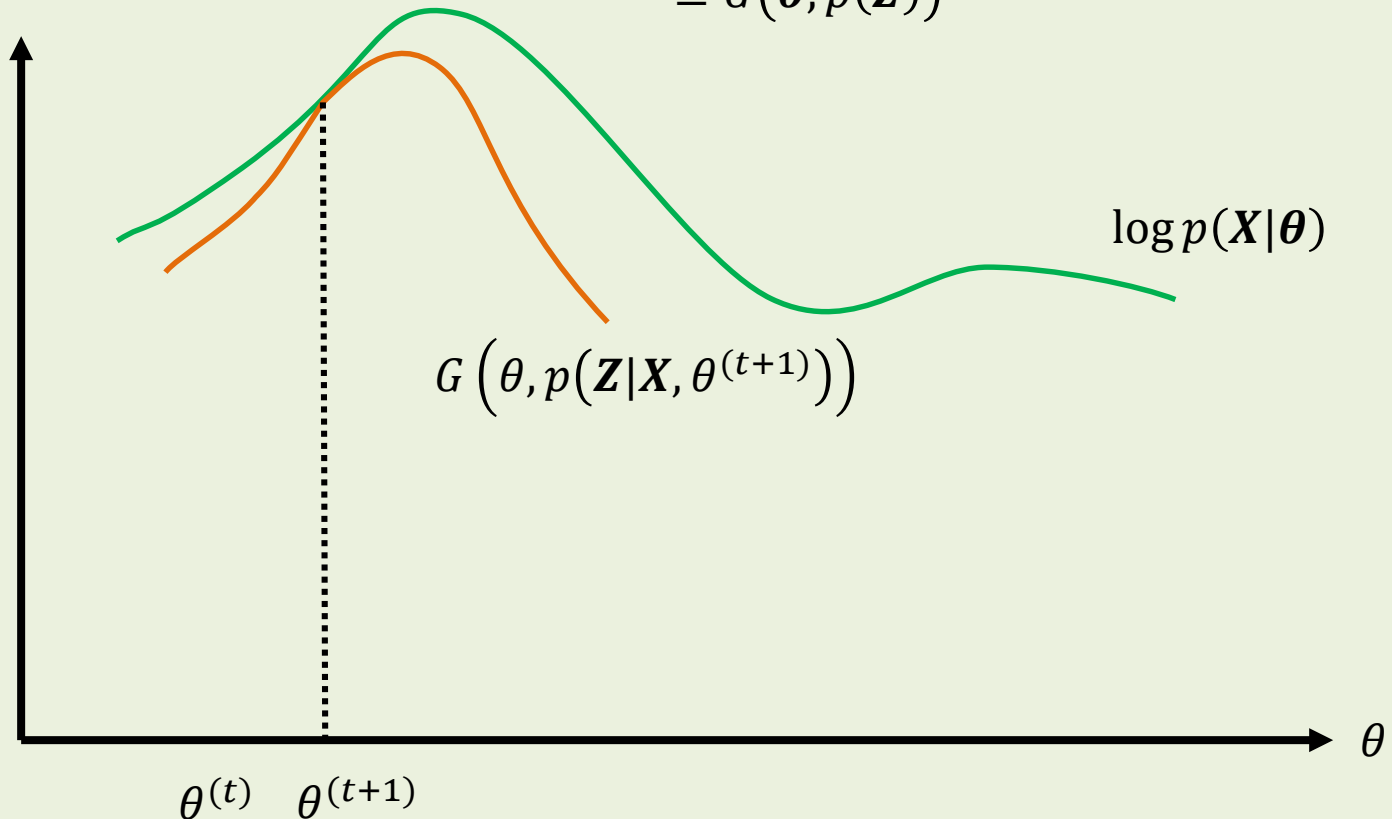
## Example (2/3)

$$\log p(\mathbf{X}|\boldsymbol{\theta}) \geq \underbrace{\mathbb{E}_{\mathbf{Z}}[\log p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})] - \mathbb{E}_{\mathbf{Z}}[\log p(\mathbf{Z})]}_{\equiv G(\boldsymbol{\theta}, p(\mathbf{Z}))}$$



## Example (3/3)

$$\log p(\mathbf{X}|\boldsymbol{\theta}) \geq \underbrace{\mathbb{E}_{\mathbf{Z}}[\log p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})] - \mathbb{E}_{\mathbf{Z}}[\log p(\mathbf{Z})]}_{\equiv G(\boldsymbol{\theta}, p(\mathbf{Z}))}$$



# EM as iterative optimisation

1. Initialisation: choose (random) initial values of  $\theta^{(1)}$
2. Update:
  - \* **E-step**: compute  $Q(\theta, \theta^{(t)}) \equiv \mathbb{E}_{\mathbf{Z}|\mathbf{X}, \theta^{(t)}} [\log p(\mathbf{X}, \mathbf{Z}|\theta)]$
  - \* **M-step**:  $\theta^{(t+1)} = \operatorname{argmax}_{\theta} Q(\theta, \theta^{(t)})$
3. Termination: if no change then stop
4. Go to Step 2

This algorithm will eventually stop (converge), but the resulting estimate can be only a local maximum

# Maximising the lower bound (2/2)

- $\log p(\mathbf{X}|\boldsymbol{\theta}) \geq \mathbb{E}_{\mathbf{Z}}[\log p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})] - \mathbb{E}_{\mathbf{Z}}[\log p(\mathbf{Z})]$
- EM is essentially coordinate descent:
  - \* Fix  $\boldsymbol{\theta}$  and optimise the lower bound for  $p(\mathbf{Z})$
  - \* Fix  $p(\mathbf{Z})$  and optimise for  $\boldsymbol{\theta}$
- The convenience of EM follows from the following
- For any point  $\boldsymbol{\theta}^*$ , it can be shown that setting  $p(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^*)$  makes the lower bound tight
- For any  $p(\mathbf{Z})$ , the second term does not depend on  $\boldsymbol{\theta}$
- When  $p(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^*)$ , the first term can usually be maximised as a function of  $\boldsymbol{\theta}$  in a closed-form
  - \* If not, then probably don't use EM

we will  
prove this  
now



# Putting the latent variables in use

- We want to maximise  $\log p(\mathbf{X}|\boldsymbol{\theta})$ . We don't know  $\mathbf{Z}$ , but consider an arbitrary non-zero distribution  $p(\mathbf{Z})$

- $\log p(\mathbf{X}|\boldsymbol{\theta}) = \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$

← Rule of marginal distribution  
(here  $\sum_{\mathbf{Z}} \dots$  iterates over all possible values of  $\mathbf{Z}$ )

- $= \log \sum_{\mathbf{Z}} \left( p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \frac{p(\mathbf{Z})}{p(\mathbf{Z})} \right)$

- $= \log \sum_{\mathbf{Z}} \left( p(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{p(\mathbf{Z})} \right)$

- $= \log \mathbb{E}_{\mathbf{Z}} \left[ \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{p(\mathbf{Z})} \right]$

← Jensen's inequality holds since  $\log(\dots)$  is a concave function

- $\geq \mathbb{E}_{\mathbf{Z}} \left[ \log \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{p(\mathbf{Z})} \right]$

- $= \mathbb{E}_{\mathbf{Z}} [\log p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})] - \mathbb{E}_{\mathbf{Z}} [\log p(\mathbf{Z})]$

# Setting a tight lower bound (1/2)

- $\log p(\mathbf{X}|\boldsymbol{\theta}) \geq \mathbb{E}_{\mathbf{Z}} \left[ \log \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{p(\mathbf{Z})} \right]$ 
  - $= \mathbb{E}_{\mathbf{Z}} \left[ \log \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})p(\mathbf{X}|\boldsymbol{\theta})}{p(\mathbf{Z})} \right]$  ← Chain rule of probability
  - $= \mathbb{E}_{\mathbf{Z}} \left[ \log \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{p(\mathbf{Z})} + \log p(\mathbf{X}|\boldsymbol{\theta}) \right]$
  - $= \mathbb{E}_{\mathbf{Z}} \left[ \log \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{p(\mathbf{Z})} \right] + \mathbb{E}_{\mathbf{Z}}[\log p(\mathbf{X}|\boldsymbol{\theta})]$  ← Linearity of  $\mathbb{E}[\cdot]$
  - $= \mathbb{E}_{\mathbf{Z}} \left[ \log \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{p(\mathbf{Z})} \right] + \log p(\mathbf{X}|\boldsymbol{\theta})$  ←  $\mathbb{E}[\cdot]$  of a constant
- $\log p(\mathbf{X}|\boldsymbol{\theta}) \geq \mathbb{E}_{\mathbf{Z}} \left[ \log \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{p(\mathbf{Z})} \right] + \log p(\mathbf{X}|\boldsymbol{\theta})$



# Setting a tight lower bound (2/2)

Ultimate aim:  
maximise this

Lower bound of what  
we want to maximise

$$\log p(\mathbf{X}|\boldsymbol{\theta}) \geq \underbrace{\mathbb{E}_{\mathbf{Z}} \left[ \log \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{p(\mathbf{Z})} \right]}_{\text{Lower bound}} + \log p(\mathbf{X}|\boldsymbol{\theta})$$

First, note that this term\*  $\leq 0$

Second, note that if  $p(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ , then

$$\mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})} \left[ \log \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})} \right] = \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})} [\log 1] = 0$$

For any  $\boldsymbol{\theta}^*$ , setting  $p(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^*)$  maximises the lower bound on  $\log p(\mathbf{X}|\boldsymbol{\theta}^*)$  and makes it tight

\*Negative Kullback-Leibler divergence between  $p(\mathbf{Z})$  and  $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$

# Estimating Parameters of Gaussian Mixture Model

A classical application of the  
Expectation-Maximisation algorithm

# Latent variables of GMM

- Let  $z_1, \dots, z_n$  denote **true origins** of the corresponding points  $\mathbf{x}_1, \dots, \mathbf{x}_n$ . Each  $z_i$  is a discrete variable that takes values in  $1, \dots, k$ , where  $k$  is a number of clusters

- Now compare the original log likelihood

$$\log p(\mathbf{x}_1, \dots, \mathbf{x}_n) = \sum_{i=1}^n \log \left( \sum_{c=1}^k w_c \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \right)$$

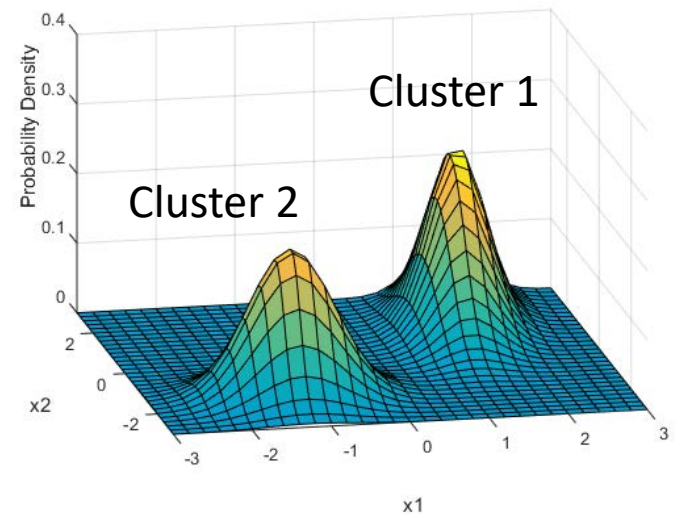
- With **complete log likelihood** (if we knew  $\mathbf{z}$ )

$$\log p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}) = \sum_{i=1}^n \log \left( w_{z_i} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i}) \right)$$

- Recall that taking a log of a normal density function results in a **tractable** expression

# Handling uncertainty about $\mathbf{z}$

- We cannot compute complete log likelihood because we don't know  $\mathbf{z}$
- EM algorithm handles this uncertainty replacing  $\log p(\mathbf{X}, \mathbf{z} | \boldsymbol{\theta})$  with expectation  $\mathbb{E}_{\mathbf{z} | \mathbf{X}, \boldsymbol{\theta}^{(t)}} [\log p(\mathbf{X}, \mathbf{z} | \boldsymbol{\theta})]$
- This in turn requires the distribution of  $p(\mathbf{z} | \mathbf{X}, \boldsymbol{\theta}^{(t)})$  given current parameter estimates
- Assuming that  $z_i$  are pairwise independent, we need  $P(z_i = c | \mathbf{x}_i, \boldsymbol{\theta}^{(t)})$
- E.g., suppose  $\mathbf{x}_i = (-2, -2)$ . What is the probability that this point originated from Cluster 1



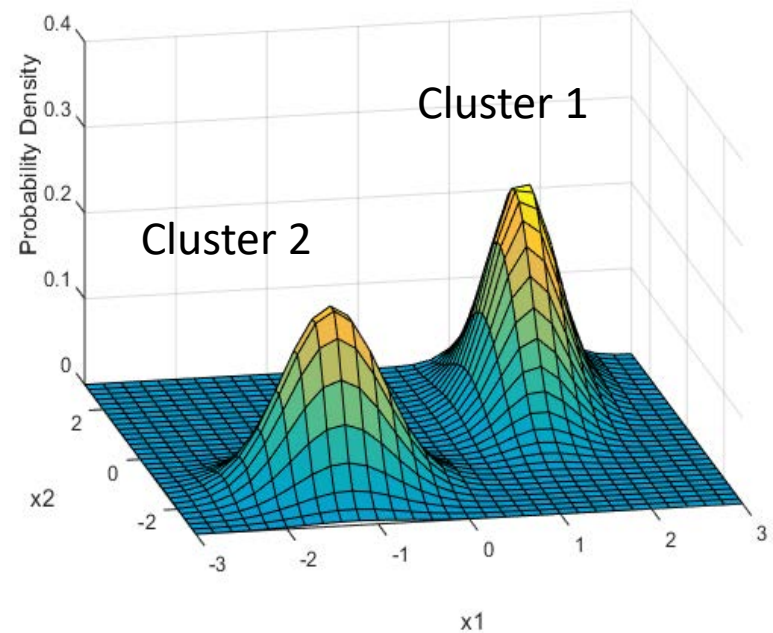
# E-step: Cluster responsibilities

- Setting latent  $Z$  as originating cluster, yields (via Bayes rule)

$$P(z_i = c | \mathbf{x}_i, \boldsymbol{\theta}^{(t)}) = \frac{w_c \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)}{\sum_{l=1}^k w_l \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)}$$

- This probability is called *responsibility* that cluster  $c$  takes for data point  $i$

$$r_{ic} \equiv P(z_i = c | \mathbf{x}_i, \boldsymbol{\theta}^{(t)})$$



# Expectation step for GMM

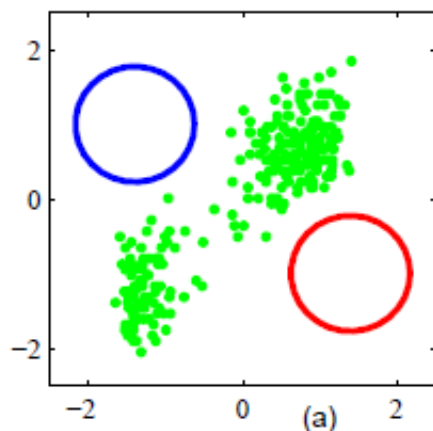
To simplify notation, we denote  $\mathbf{x}_1, \dots, \mathbf{x}_n$  as  $\mathbf{X}$

- $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) \equiv \mathbb{E}_{\mathbf{z}|\mathbf{X}, \boldsymbol{\theta}^{(t)}} [\log p(\mathbf{X}, \mathbf{z}|\boldsymbol{\theta})]$
- $= \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{X}, \boldsymbol{\theta}^{(t)}) \log p(\mathbf{X}, \mathbf{z}|\boldsymbol{\theta})$
- $= \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{X}, \boldsymbol{\theta}^{(t)}) \sum_{i=1}^n \log w_{z_i} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i})$
- $= \sum_{i=1}^n \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{X}, \boldsymbol{\theta}^{(t)}) \log w_{z_i} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i})$
- $= \sum_{i=1}^n \sum_{c=1}^k r_{ic} \log w_{z_i} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i})$
- $= \sum_{i=1}^n \sum_{c=1}^k r_{ic} \log w_{z_i}$
- $+ \sum_{i=1}^n \sum_{c=1}^k r_{ic} \log \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i})$

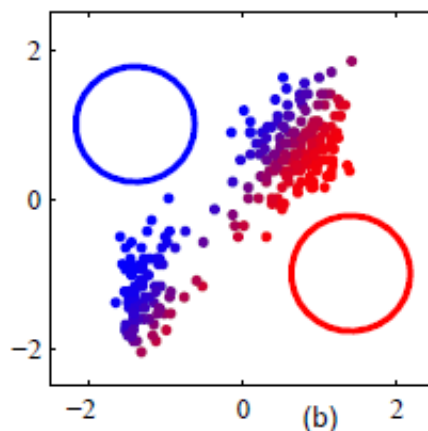
# Maximisation step for GMM

- In the maximisation step, take partial derivatives of  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})$  with respect to each of the parameters and set the derivatives to zero to obtain new parameter estimates
- $w_c^{(t+1)} = \frac{1}{n} \sum_{i=1}^n r_{ic}$
- $\boldsymbol{\mu}_c^{(t+1)} = \frac{\sum_{i=1}^n r_{ic} \mathbf{x}_i}{r_c}$   
 \* Here  $r_c \equiv \sum_{i=1}^n r_{ic}$
- $\boldsymbol{\Sigma}_c^{(t+1)} = \frac{\sum_{i=1}^n r_{ic} \mathbf{x}_i \mathbf{x}_i'}{r_c} - \boldsymbol{\mu}_c^{(t)} \left( \boldsymbol{\mu}_c^{(t)} \right)'$
- Note that these are the estimates for step  $(t + 1)$

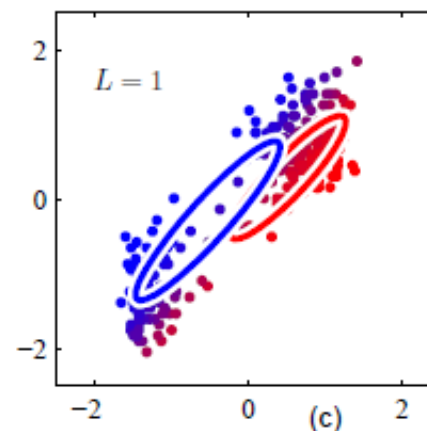
# Example of fitting Gaussian Mixture model



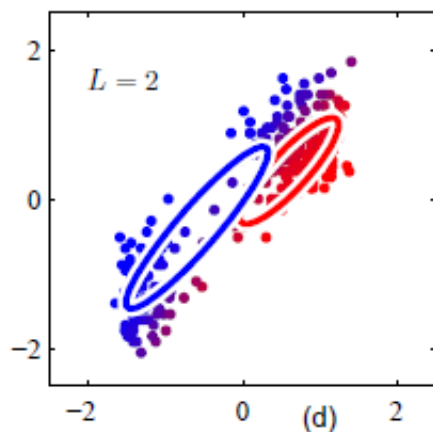
(a) Initial



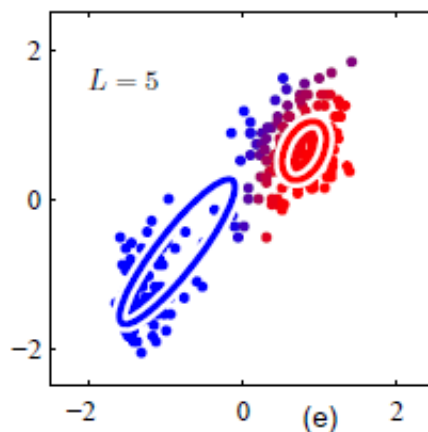
(b) E-step



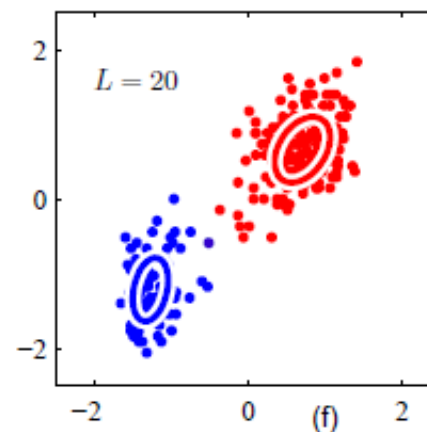
(c) M-step



(d) 2 cycles



(e) 5-cycles



(f) 20-cycles



# K-means as a EM for a restricted GMM

- Consider a GMM model in which all components have the same fixed probability  $w_c = 1/k$ , and each Gaussian has the same fixed covariance matrix  $\Sigma_c = \sigma^2 \mathbf{I}$ , where  $\mathbf{I}$  is the identity matrix
- In such a model, only component centroids  $\mu_c$  need to be estimated
- Next approximate a probabilistic cluster responsibility  $r_{ic} = P(z_i = c | \mathbf{x}_i, \mu_c^{(t)})$  with a deterministic assignment  $r_{ic} = 1$  if centroid  $\mu_c^{(t)}$  is closest to point  $\mathbf{x}_i$ , and  $r_{ic} = 0$  otherwise
- Such a formulation results in a E-step where  $\mu_c$  should be set as a centroid of points assigned to cluster  $c$
- In other words, **k-means algorithm is a EM algorithm for the restricted GMM model** described above!!!

# This lecture

- Unsupervised learning
  - \* Diversity of problems
- Gaussian mixture model (GMM)
  - \* A probabilistic approach to clustering
  - \* The GMM model
  - \* GMM clustering as an optimisation problem
- The Expectation Maximization (EM) algorithm
- Next lecture: More unsupervised with dim reduction