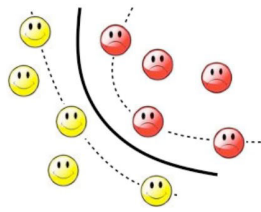


# Machine Learning and Data Mining (COMP 5318)

## Clustering and Expectation-Maximisation

Nguyen Hoang Tran



# Clustering

C. Bishop, *Pattern Recognition and Machine Learning*,  
Chapter 9: Mixture Models and EM  
Springer New York, 2006

K.P. Murphy, *Machine Learning: a Probabilistic Perspective*,  
Chapters 11 and 25, Massachusetts Institute of Technology, 2006

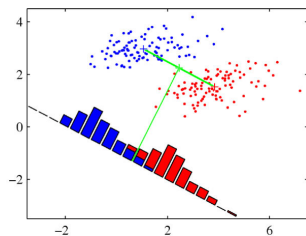
1

2

## Types of Learning

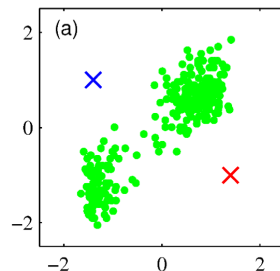


### Supervised Learning



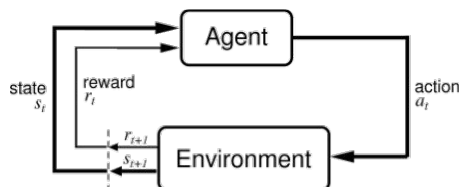
Learning  
input-output  
from examples  
  
Regression  
Classification

### Unsupervised Learning



Learning  
underlying  
structure  
  
Clustering  
Density  
Estimation

### Reinforcement Learning



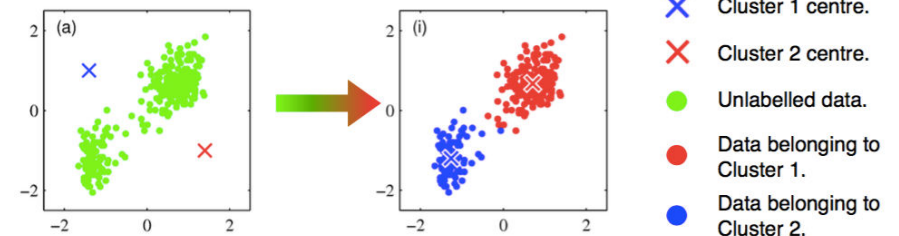
Learning policy  
from state-action-reward  
sequences.  
  
Learning behaviour

3

## Clustering



### Process of grouping similar objects together



Learn a set of clusters and assign data to a specific cluster.

**Deterministic:** Hard assignment to each cluster (K-means).  
**Probabilistic:** Model assignment as a discrete latent variable.  
(Mixtures of Gaussians, Dirichlet Process)

4

# Clustering

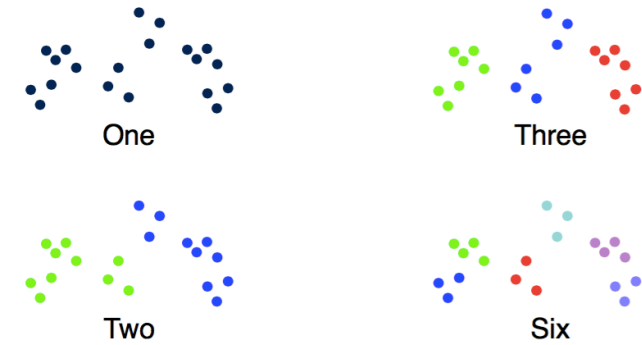
How many clusters?



5

# Clustering

How many clusters?

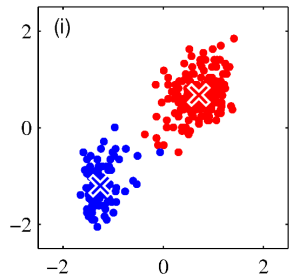


Presence of ambiguous solutions.

6

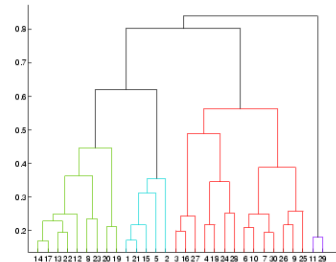
## Types of Clustering

### Partition Clustering



Partition the objects into disjoint sets. Faster to create.  
Sensible to initial conditions.  
Model selection for K.

### Hierarchical Clustering



Nested tree of partitions.  
Slower to create.  
Often more useful.  
Do not require knowing the number of clusters.

7

## Clustering

Dataset  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  with N Observations

Each data point is D dimension

**Goal:** Partition dataset into K clusters. (For now, assume K is given)

$\mu_k = (\mu_1, \dots, \mu_D)$  : Centroid for each cluster  $k \in 1, \dots, K$

Binary indicator variables

$$r_{nk} = \begin{cases} 1, & \text{if datapoint } n \text{ belongs to cluster } k \\ 0, & \sim \end{cases}$$

If  $\mathbf{x}_n$  is assigned to cluster k, then  $r_{nk} = 1 \wedge r_{nj} = 0 \forall j \neq k$

8

# K-Means



Objective function:

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

Represents the sum of the squares of the distances of each datapoint to its assigned centroid vector.

Goal: Find  $\{\boldsymbol{\mu}_k\}$  and  $\{r_{nk}\}$  that minimise J.

$$\{r_{nk}, \boldsymbol{\mu}_k\}^* = \underset{\{r_{nk}, \boldsymbol{\mu}_k\}}{\operatorname{argmin}} \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

9

# K-Means



Iterative solution to minimise J:

1. Data Preprocessing
2. Initialise  $\{\boldsymbol{\mu}_k\}$
3. Repeat 4 and 5 until convergence or Max Iterations
4. Minimise J w.r.t.  $\{r_{nk}\}$  keeping  $\{\boldsymbol{\mu}_k\}$
5. Minimise J w.r.t.  $\{\boldsymbol{\mu}_k\}$  keeping  $\{r_{nk}\}$

10

# K-Means



$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

Line 4: Optimise w.r.t  $r_{nk}$

Each data point is independent, so we can optimise for each n separately:

$$r_{nk} = \begin{cases} 1 & \text{if } k = \underset{j}{\operatorname{argmin}} \|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2 \\ 0 & \text{Otherwise} \end{cases}$$

Assign each data point to its closest centroid.

# K-Means



$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

Line 5: Optimise w.r.t  $\boldsymbol{\mu}_k$

$$\frac{\partial J}{\partial \boldsymbol{\mu}_k} = 0$$

$$2 \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0$$

$$\sum_{n=1}^N r_{nk} \mathbf{x}_n = \sum_{n=1}^N r_{nk} \boldsymbol{\mu}_k$$

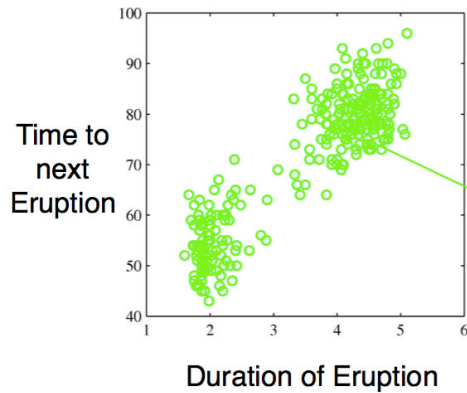
$$\frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}} = \boldsymbol{\mu}_k$$

Set  $\boldsymbol{\mu}_k$  equal to the mean of all data points  $\mathbf{x}_n$  assigned to cluster k.

# K-means

# K-Means Example

Hydrothermal Geyser: Old Faithful



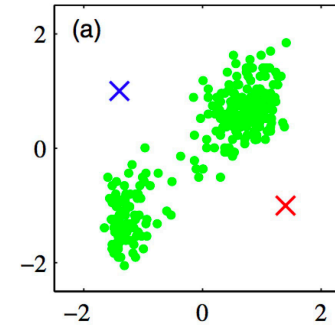
Single Eruption  
 $N = 272$  Observations

13

# K-Means Example

Number of clusters:  $K = 2$

- 1 Data Preprocessing
- 2 Initialise  $\{\mu_k\}$
- 3 Repeat until convergence or Max Iterations
- 4 Minimise  $J$  w.r.t.  $\{r_{nk}\}$  keeping  $\{\mu_k\}$  fixed.
- 5 Minimise  $J$  w.r.t.  $\{\mu_k\}$  keeping  $\{r_{nk}\}$  fixed.



Each dimension has zero mean and unit standard deviation.

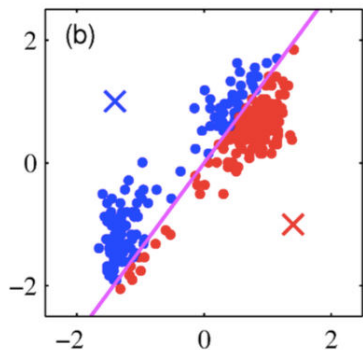
Better initialisation: Choose  $\{\mu_k\}$  as average of a random subset.

14

# K-Means Example

Number of clusters:  $K = 2$

- 1 Data Preprocessing
- 2 Initialise  $\{\mu_k\}$
- 3 Repeat until convergence or Max Iterations
- 4 Minimise  $J$  w.r.t.  $\{r_{nk}\}$  keeping  $\{\mu_k\}$  fixed.
- 5 Minimise  $J$  w.r.t.  $\{\mu_k\}$  keeping  $\{r_{nk}\}$  fixed.



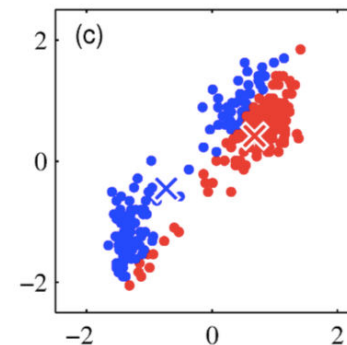
Each data point is assigned to the closest cluster centre.

15

# K-Means Example

Number of clusters:  $K = 2$

- 1 Data Preprocessing
- 2 Initialise  $\{\mu_k\}$
- 3 Repeat until convergence or Max Iterations
- 4 Minimise  $J$  w.r.t.  $\{r_{nk}\}$  keeping  $\{\mu_k\}$  fixed.
- 5 Minimise  $J$  w.r.t.  $\{\mu_k\}$  keeping  $\{r_{nk}\}$  fixed.



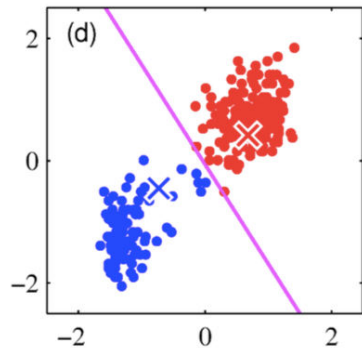
Re-compute each cluster centre to be the mean of the points previously assigned.

16

# K-Means Example

Number of clusters:  $K = 2$

- 1 Data Preprocessing
- 2 Initialise  $\{\mu_k\}$
- 3 Repeat until convergence or Max Iterations
- 4 Minimise  $J$  w.r.t.  $\{r_{nk}\}$  keeping  $\{\mu_k\}$  fixed.
- 5 Minimise  $J$  w.r.t.  $\{\mu_k\}$  keeping  $\{r_{nk}\}$  fixed.



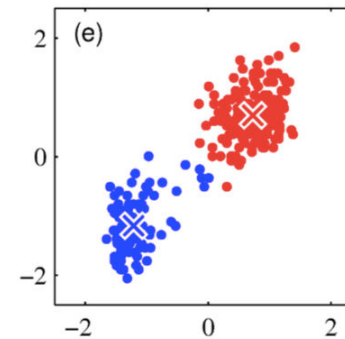
Each data point is assigned to the closest cluster centre.

17

# K-Means Example

Number of clusters:  $K = 2$

- 1 Data Preprocessing
- 2 Initialise  $\{\mu_k\}$
- 3 Repeat until convergence or Max Iterations
- 4 Minimise  $J$  w.r.t.  $\{r_{nk}\}$  keeping  $\{\mu_k\}$  fixed.
- 5 Minimise  $J$  w.r.t.  $\{\mu_k\}$  keeping  $\{r_{nk}\}$  fixed.

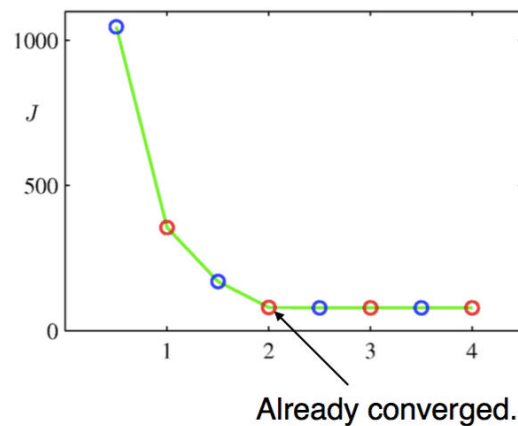


Re-compute each cluster centre to be the mean of the points previously assigned.

18

# K-Means Example

Plot of the cost function for each iteration.



19

# K-Means Example 2

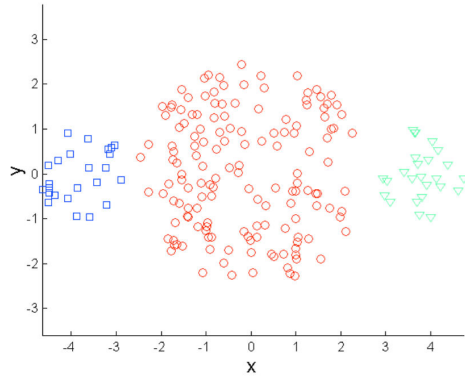
Image segmentation and compression.



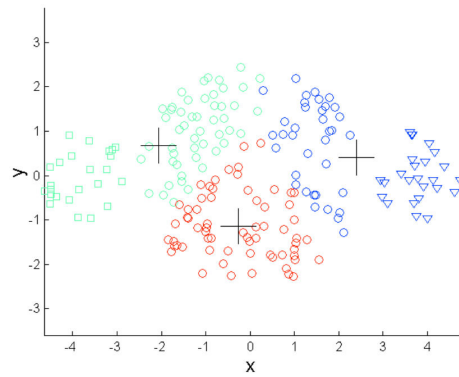
20

# K-Means Limitations

Differing Sizes:



Original Points

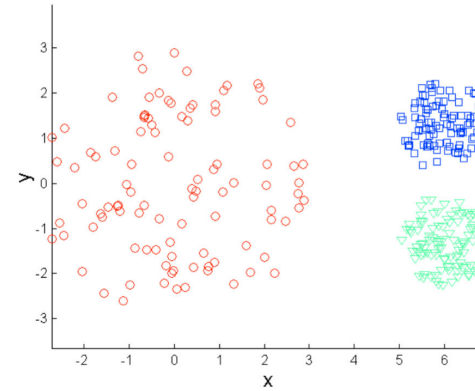


K-Means 3 Clusters

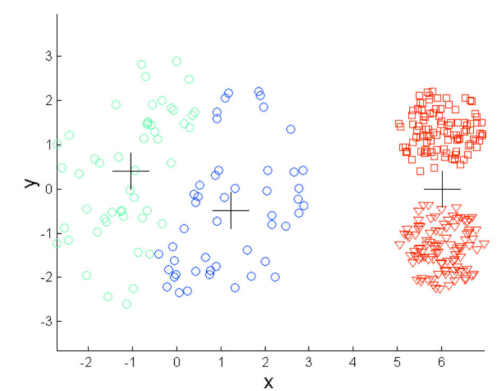
21

# K-Means Limitations

Differing Density:



Original Points

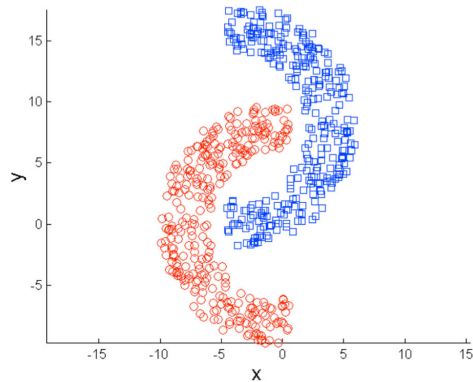


K-Means 3 Clusters

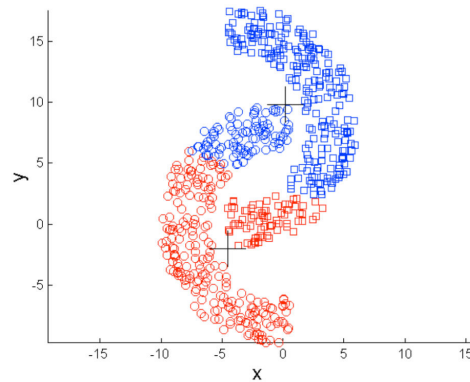
22

# K-Means Limitations

Non-Globular Shapes:



Original Points

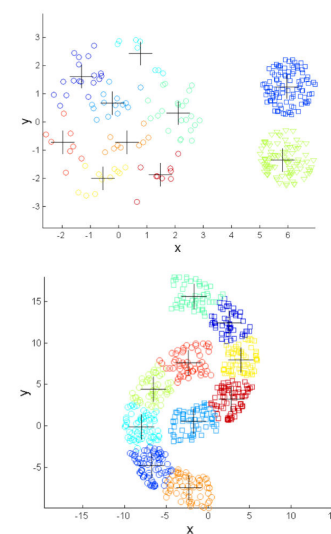
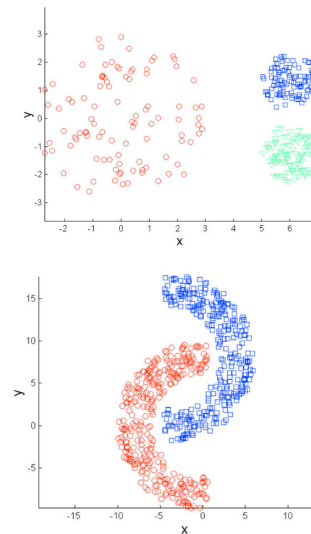


K-Means 2 Clusters

23

# Overcome K-Means Limitations

Use large number of clusters.

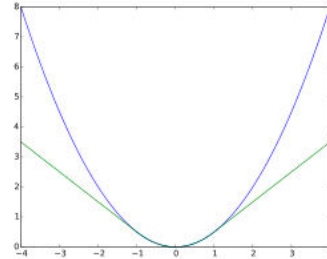


24



Generalise distance function:

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \mathcal{V}(\mathbf{x}_n, \mu_k)$$



Robustness to outliers.

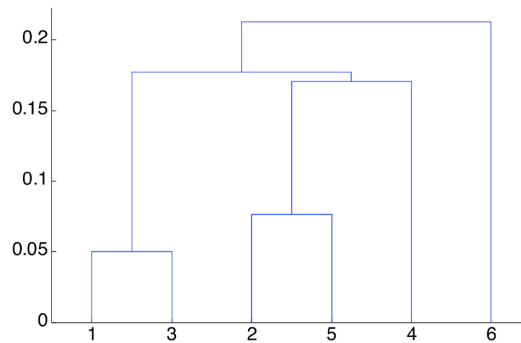
$$\mathcal{V}(\mathbf{x}_n, \mu_k) = \begin{cases} 1/2 \|\mathbf{x}_n - \mu_k\|_2^2, & \text{if } \|\mathbf{x}_n - \mu_k\| \leq \delta \\ \delta \|\mathbf{x}_n - \mu_k\|_1 - 1/2 \delta^2, & \text{otherwise} \end{cases}$$

25

## Hierarchical Clustering

Nested set of clusters organised as a hierarchical tree.

Any number of clusters can be obtained by 'cutting' the dendrogram.



Dendrogram

Uses a similarity matrix.

27

# Hierarchical Clustering

## Hierarchical Agglomerative Clustering

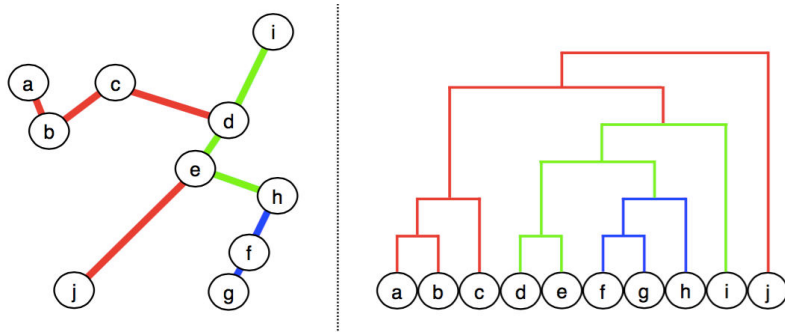
Simple clustering algorithm.  
Uses a inter cluster similarity measure.

1. Initialise: Every data point is a cluster.
2. Repeat until one cluster remains.
3. Compute distances between all clusters.
4. Merge closest clusters.
5. Update dendrogram.

28

# Hierarchical Agglomerative Clustering

Example:

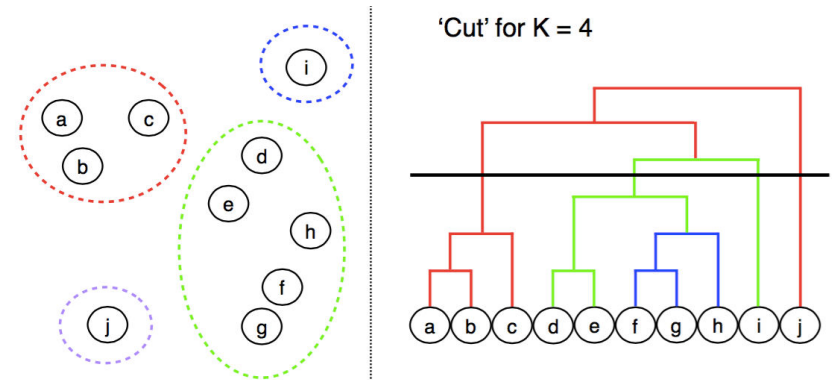


Dendrogram

29

# Hierarchical Agglomerative Clustering

Example:



Dendrogram

30

## Inter Cluster Similarity

### Nearest Neighbour

$$D_{\min}(C_i, C_j) = \min_{\mathbf{x} \in C_i, \mathbf{y} \in C_j} \|\mathbf{x} - \mathbf{y}\|^2$$

### Furthest Neighbour

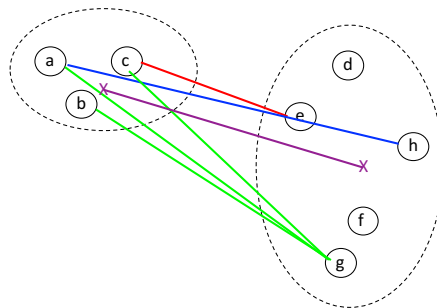
$$D_{\max}(C_i, C_j) = \max_{\mathbf{x} \in C_i, \mathbf{y} \in C_j} \|\mathbf{x} - \mathbf{y}\|^2$$

### Group Average

$$D_{\text{avg}}(C_i, C_j) = \frac{1}{|C_i||C_j|} \sum_{\mathbf{x} \in C_i} \sum_{\mathbf{y} \in C_j} \|\mathbf{x} - \mathbf{y}\|^2$$

### Centroid Distance

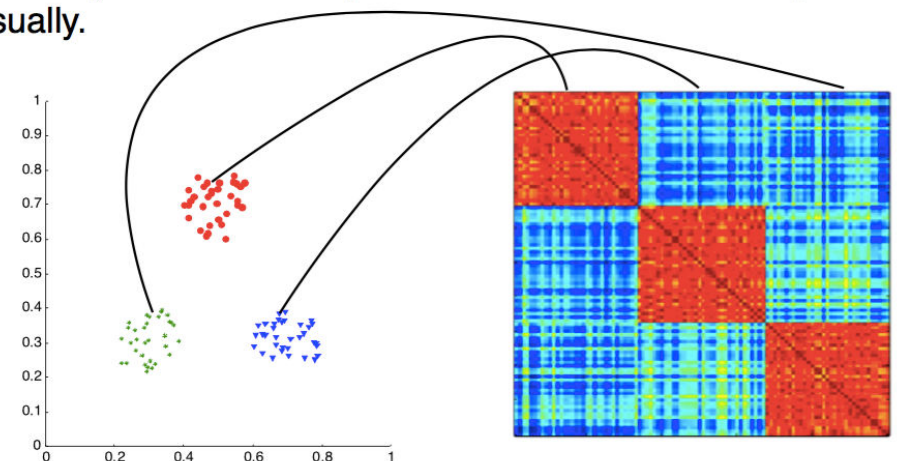
$$D_{\text{means}}(C_i, C_j) = \|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|$$



31

## Cluster Validation

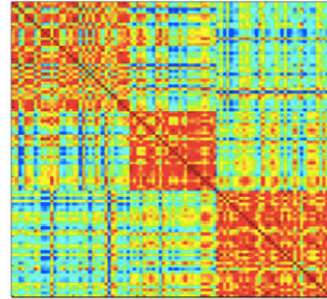
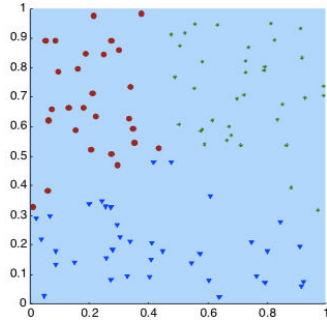
Similarity matrix with respect to cluster labels and inspect visually.



32



Random data clusters are not well defined.



33

# Probabilistic Approach to Clustering

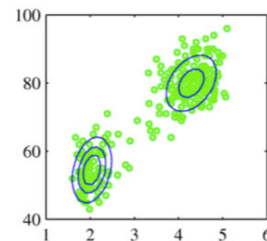
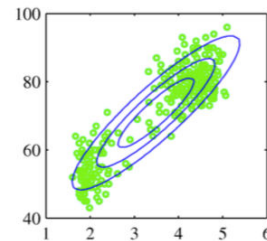
34

## Mixture of Gaussians

Gaussian mixture distribution with K components.

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$0 \leq \pi_k \leq 1 \quad \sum_{k=1}^K \pi_k = 1$$



Mixture models provide a probabilistic framework for clustering.

35

## Mixture of Gaussians

Let us introduce a latent random variable

$$\mathbf{z} = \{z_k\}_{k=1, \dots, K} \quad z_k \in \{0, 1\} \quad \sum_{k=1}^K z_k = 1$$

$\mathbf{z}$  has K possible states.

$$p(z_k = 1) = \pi_k \quad p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$$

$$p(\mathbf{x} | z_k = 1) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad p(\mathbf{x} | \mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}$$

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

36

Let us apply Bayes theorem and infer the value of the latent variable.

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$$

$$\gamma(z_k) \equiv p(z_k = 1|\mathbf{x}) = \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{p(\mathbf{x})}$$

$$= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

$\gamma(z_k)$  is called *responsibility* that component  $k$  takes for explaining  $\mathbf{x}$ .

$\pi_k$  is the prior probability of component  $k$ .

$\gamma(z_k)$  is the posterior probability after  $\mathbf{x}$  is observed.

37

## Expectation Maximisation (EM)

Elegant and powerful method for finding MLE or MAP solutions for models with latent variables.

Intuition: If we knew what cluster each point belonged to (i.e. the  $\mathbf{z}$  variables), we could partition the data and find the MLE for each cluster separately.

39

# Expectation Maximisation

## EM for Gaussian Mixtures

38

## EM for Gaussian Mixtures

Likelihood function:

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

Conditions to be satisfied at maximum likelihood:

$$0 = \frac{\partial \ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\mu}_k}$$

$$0 = - \sum_{n=1}^N \frac{\pi_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\underbrace{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}_{\gamma(z_{nk})}} \boldsymbol{\Sigma}_k (\mathbf{x}_n - \boldsymbol{\mu}_k)$$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad N_k = \sum_{n=1}^N \gamma(z_{nk})$$

$N_k$  is the effective number of points assigned to cluster  $k$ .

40

$$0 = \frac{\partial \ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\Sigma}_k}$$

$$\Rightarrow \boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

$$0 = \frac{\partial \ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \pi_k}$$

$$\Rightarrow \pi_k = \frac{N_k}{N}$$

41

1 Initialise means  $\boldsymbol{\mu}_k$ , covariances  $\boldsymbol{\Sigma}_k$  and mixing coefficients  $\pi_k$ .

2 E-step

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

3 M-step

$$\boldsymbol{\mu}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})(\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})^T$$

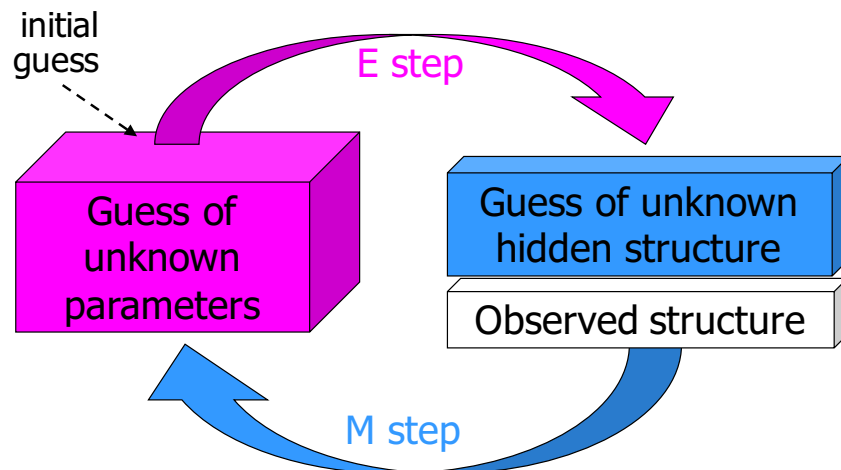
$$\pi_k^{\text{new}} = \frac{N_k}{N}$$

4 Eval Likelihood

$$\ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

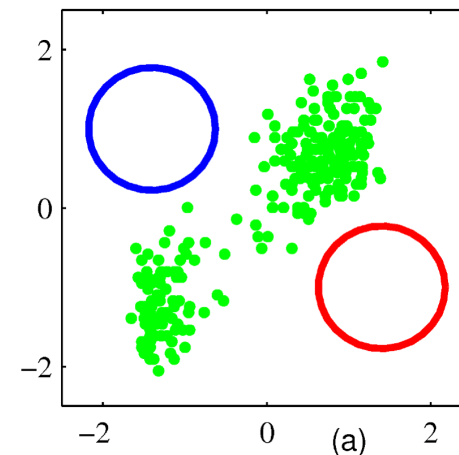
42

## EM Algorithm



43

## EM Example

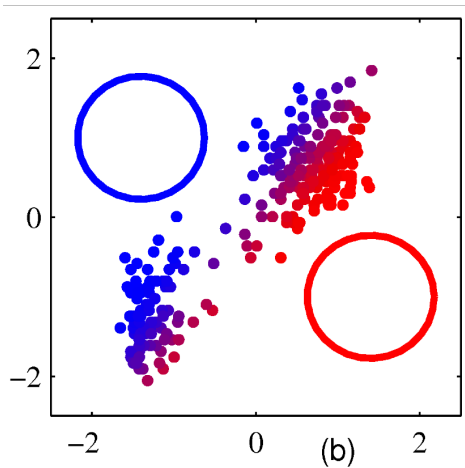


Initial values for mean vectors  
(same as K-means example).

Diagonal covariance matrices  
(showing one std contour).

44

## EM Example



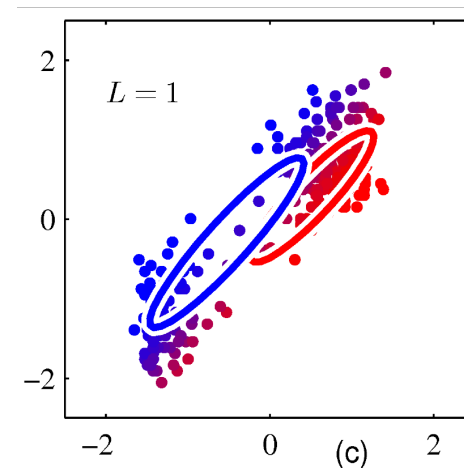
Initial E step.

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

Colour  
proportional to responsibilities.

45

## EM Example



M Step:

The means move towards the weighted average of dataset with respective ink colour (responsibilities).

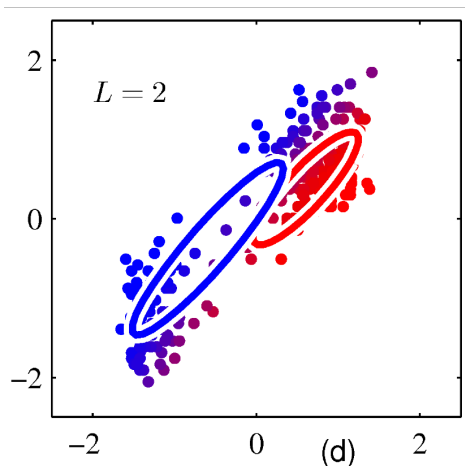
$$\boldsymbol{\mu}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

The covariance matrices adapt to the covariance of the respective ink.

$$\boldsymbol{\Sigma}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})^T$$

46

## EM Example



M Step:

The means move towards the weighted average of dataset with respective ink colour (responsibilities).

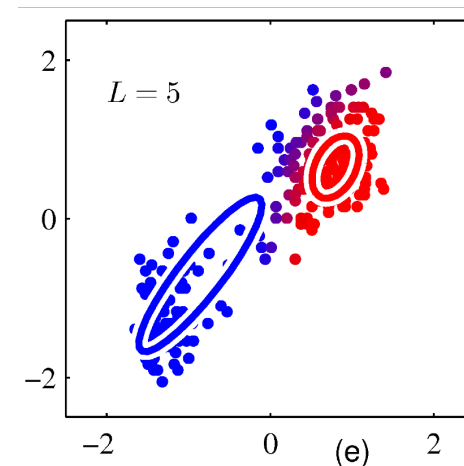
$$\boldsymbol{\mu}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

The covariance matrices adapt to the covariance of the respective ink.

$$\boldsymbol{\Sigma}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})^T$$

47

## EM Example



M Step:

The means move towards the weighted average of dataset with respective ink colour (responsibilities).

$$\boldsymbol{\mu}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

The covariance matrices adapt to the covariance of the respective ink.

$$\boldsymbol{\Sigma}_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{new}})^T$$

48

## EM Example

M Step:

The means move towards the weighted average of dataset with respective ink colour (responsibilities).

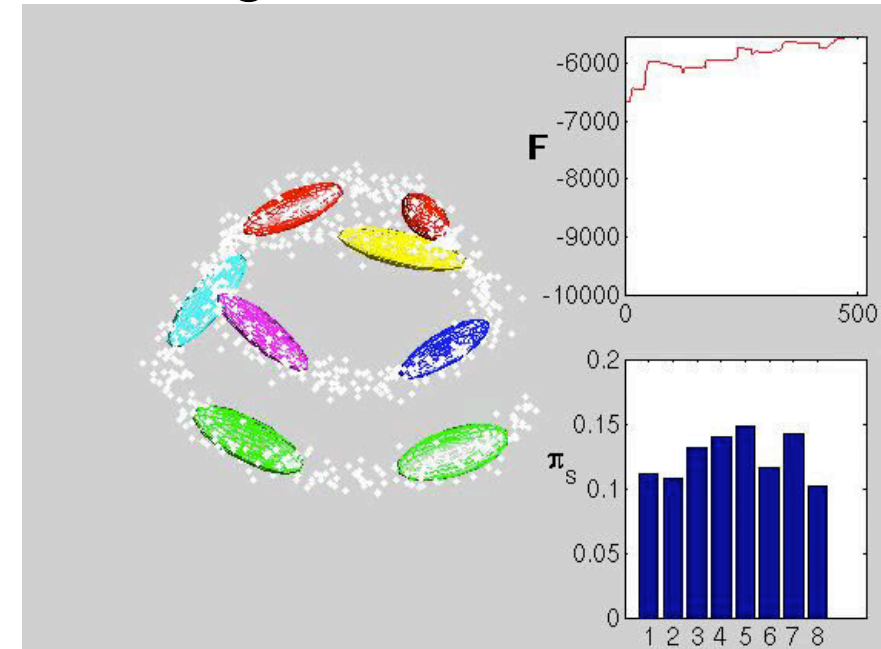
$$\mu_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

The covariance matrices adapt to the covariance of the respective ink.

$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k^{\text{new}}) (\mathbf{x}_n - \mu_k^{\text{new}})^T$$

49

## Clustering



COMP5318 - Nguyen Hoang Tran

50