

## \* Representative-Based Clustering

- For each cluster, there is a point (or parameter vector) that summarizes it (common choice: mean)
- Given a dataset of  $N$  instances and desired # of clusters  $k$ , this class of algorithms generates a partition  $C = \{C_1, \dots, C_k\}$

$$M_i = \frac{1}{n_i} \cdot \sum_{x_j \in C_i} x_j$$

### • Brute-Force Approach

- Generate all possible clustering  $C = \{C_1, \dots, C_k\} \rightarrow$  Then select the best one (But there are  $O(2^N/k!)$  partitions)

### • K-means Algorithm

- Greedy iterative approach to find a clustering that minimizes the SSE objective. (Then it can converge to local optimal instead of global one)

$$SSE(C) = \sum_{i=1}^k \sum_{x_j \in C_i} \|x_j - \mu_i\|^2 \quad \text{Goal is to find } C^* = \arg \min_C SSE(C)$$

```
Algorithm 13.1: K-means Algorithm
K-MEANS(D, k, ε):
1 t ← 0
2 Randomly initialize k centroids:  $\mu_1^t, \mu_2^t, \dots, \mu_k^t \in \mathbb{R}^d$ 
3 repeat
4   t ← t + 1
5    $C_t \leftarrow \emptyset$  for all  $i = 1, \dots, k$ 
// Cluster Assignment Step
6   foreach  $x_i \in D$  do
7      $i^* \leftarrow \operatorname{argmin}_i \|\|x_i - \mu_i^{t-1}\|^2\|$ 
8      $C_{i^*} \leftarrow C_{i^*} \cup \{x_i\}$  // Assign  $x_i$  to closest centroid
// Centroid Update Step
9   foreach  $i = 1, \dots, k$  do
10     $\mu_i^t \leftarrow \frac{1}{|C_i|} \sum_{x_j \in C_i} x_j$ 
11 until  $\sum_{i=1}^k \|\mu_i^t - \mu_i^{t-1}\|^2 \leq \epsilon$ 
```

→ Algorithm Complexity;

$O(t n k d)$  time  
 $n$  points,  $k$  clusters,  
 $d$  operations  
 $t$  iterations

### » Centroid Initialization:

Solution 1) Pick points that are as far away from one another as possible

Solution 2) Pick first point random → while there are fewer than  $k$  points; add the point whose min. distance from the selected points is as large as possible

Solution 3) Cluster a sample of data (hierarchically?) →  $k$  clusters → pick a point from each cluster, perhaps that point closest to the centroid of the cluster

! Selecting best initial centroid is difficult:  $\frac{\text{ways to select one centroid for each cluster}}{\text{ways to select } k \text{ centroids}} = \frac{k! \cdot n^k}{(kn)!} = \frac{k!}{(kn-k)!}$

for 10 clusters:  $\frac{10!}{100!} \sim 0.00086$

```
Algorithm 3 Bisecting K-means Algorithm.
1: Initialize the list of clusters to contain the cluster containing all points.
2: repeat
3:   Select a cluster from the list of clusters
4:   for i = 1 to number_of_iterations do
5:     Bisect the selected cluster using basic K-means
      → "selected cluster"
      → "orthogonal to hyperplane"
6:   end for
7:   Add the two clusters from the bisection with the lowest SSE to the list of clusters.
8: until Until the list of clusters contains K clusters
```

### » Dealing with Initial Centroids Issue:

- Multiple runs helps → but probability is not on your side.
- Sample and use another clustering method (hierarchical?) to determine initial centroids.
- Select more than  $k$  initial centroids → then select among these
- Postprocessing
- Bisecting K-means → not sensitive to initialization issues.

\* Centroids are updated after all points are assigned to a centroid (alternative) → update centroids after each assignment  
 (more expensive, introduces order dependency, never get empty cluster)

### » Preprocessing

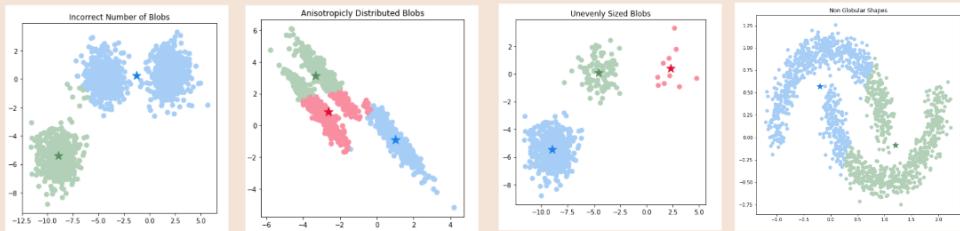
- Normalize the data
- Eliminate outliers

### » Post processing

- Eliminate small clusters that may represent outliers
- Split 'loose' clusters (clusters with high SSE)
- Merge clusters that are 'close' and have relatively low SSE

## » Limitations of K-means:

- K-means has problems when clusters are differing
  - non-globular shape
  - size
  - densities
- It also has problems with outliers



## » K-means Clustering Summary:

### Strength

- Efficient
- Often terminates at local optimal
- Global optimal can be found using; deterministic annealing, genetic algorithms

### Weakness

- Applicable only when mean defined (categorical data?)
- Need to specify  $k$  in advance
- Unable to handle noise, outliers
- Not suitable for non-convex shapes

### Advantages

- Simple, understandable
- Items automatically assigned to clusters

### Disadvantages

- Must pick  $k$  before
- All items forced into a cluster
- Too sensitive to outliers

## » Variants of K-means:

- K-means differ in
  - selection of initial  $K$  means
  - dissimilarity calculations
  - strategies to calculate cluster mean

### Handling categorical data: K-modes

- Replace means with modes
- Use new dissimilarity measure to deal with categorical objects
- Use a frequency-based method to update modes.
- Mixture of numerical and categorical data: K-prototype method

### How do we choose $k$ ?

- Knee and Elbow analysis → Not enough!

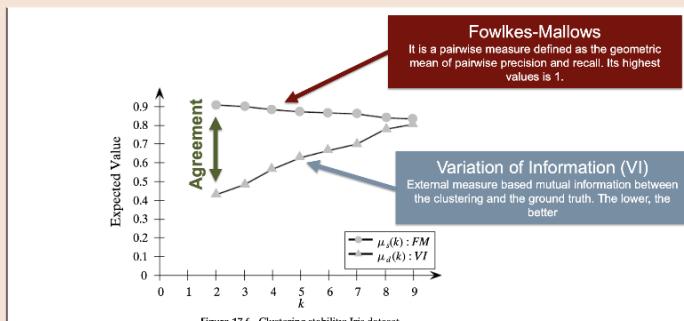
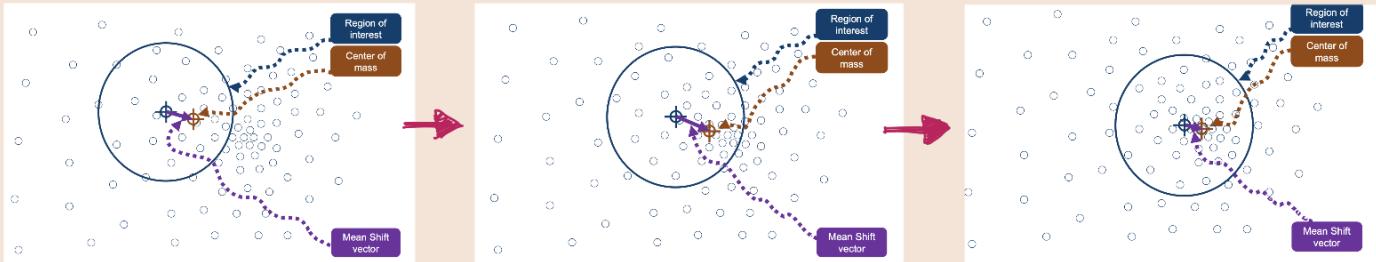


Figure 17.6. Clustering stability: Iris dataset.  
Clustering stability for the Iris dataset, with  $n = 150$ , using the K-means algorithm. It uses  $t = 500$  bootstrap samples. For each dataset  $D_i$ , and each value of  $k$ , it run K-means with 100 initial starting configurations and select the best clustering.

## • Mean-Shift Algorithm

- It searches for the mode of a data distribution (nonparametric, iterative, versatile algorithm)



### ⇒ Mean-Shift Pseudocode

- 1) Choose a search window size
- 2) For each point;
  - a) Center a window on that point
  - b) Compute mean in the search window
  - c) Center the search window at new mean location
  - d) Repeat (b,c) until convergence
- 3) Assign points that lead to nearby modes to the same cluster.

### ⇒ K-Means Clustering vs. Mean-Shift Clustering

- Assumes # of clusters already known and clusters shaped spherically

- Sensitive to initialization
  - Sensitive to outliers
  - Fast and has complexity  $O(knT)$
- # of cluster  
# of points

### Mean-Shift Clustering

- Does not assume anything about # of clusters, instead, # of modes give the # of clusters.  
Since it based on density estimation, it can handle arbitrarily shaped clusters

- Robust to initializations
  - Less-sensitive to outliers
  - Computationally expensive, has time complexity  $O(T n^2)$
- # of iterations

### ⇒ Mean-Shift Clustering Strength & Weakness

#### Strength:

- Does not assume prior shape on clusters
- Can handle arbitrary feature spaces
- Only window size to choose
- Window size has a physical meaning

#### Weakness

- Window size selection is not trivial (Inappropriate window size can cause modes to be merged, or generate additional shallow modes → Adaptive window size can help)
- Not suitable for high-dimensional features

## • Expectation - Maximization (EM) clustering

- We assume that each cluster  $C_i$  is characterized by a multivariate normal distribution and thus, identified by the mean vector  $\mu_i$  and covariance matrix  $\Sigma_i$
- A clustering identified by a vector of parameter  $\Theta = \{\mu_i, \Sigma_i, \underbrace{P(C_i)}_{\text{prior prob. of cluster } C_i}\}$   $\rightarrow \sum P(C_i) = 1$
- Goal of Maximum Likelihood Estimation (MLE) is choose the parameters  $\Theta$  that maximize the likelihood  $\Theta^* = \arg \max_{\Theta} P(D | \Theta)$
- General Idea:
  - Starts with an initial estimate of parameter vector  $\Theta$
  - Iteratively rescores the patterns against mixture density produced by the parameter vector  $\Theta$
  - The rescored patterns are used to update the parameters
  - Patterns belonging to some cluster if they're placed by their scores in particular component.

### » EM in One Direction

- Consider a dataset consisting of single attribute  $X = \{x_1, \dots, x_n\}$
- Clusters represented using Univariate Normal ;  $f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
- Each cluster represented by the parameters  $\{\mu_i, \sigma_i^2, P(C_i)\}$
- Initialize parameters;  $\mu_i$  selected uniformly at random,  $\sigma_i^2 = 1$ ,  $P(C_i) = \gamma_k$   $\rightarrow$  # of clusters
- For each cluster, using current estimate of parameters, we can compute posterior probabilities

$$P(C_i | x_j) = \frac{P(x_j | \mu_i, \sigma_i^2) \cdot P(C_i)}{\sum_{k=1}^K f(x_j | \mu_k, \sigma_k^2) \cdot P(C_k)}$$

→ denote;  $w_{ij} = P(C_i | x_j)$   
 $w_i : \text{weight vector for } C_i \text{ over all } n \text{ points}$   
 $(w_{i1}, \dots, w_{in})^T$

- Given all posterior values  $w_{ij}$ , maximization step computes MLE of cluster parameters  $(\mu_i, \sigma_i^2, P(C_i))$ 
  - Mean updated as:  $\mu_i = \frac{\sum_{j=1}^n w_{ij} x_j}{\sum_{j=1}^n w_{ij}}$
  - $\sigma_i^2 = \frac{\sum_{j=1}^n w_{ij} (x_j - \mu_i)^2}{\sum_{j=1}^n w_{ij}}$
  - $P(C_i) = \frac{\sum_{j=1}^n w_{ij}}{\sum_{j=1}^n \sum_{k=1}^n w_{kj}} = \frac{\sum_{j=1}^n w_{ij}}{n} = 1$
- Expectation and Maximization steps are repeated until convergence