

Dr.Ramez Alkhatib

INTRODUCTION-What is clustering?



 Clustering is the <u>classification</u> of objects into different groups, or more precisely, the <u>partitioning</u> of a <u>data set</u> into <u>subsets</u> (clusters), so that the data in each subset (ideally) share some common trait - often according to some defined <u>distance measure</u>.

Types of clustering:

- ers
- 1. <u>Hierarchical algorithms</u>: these find successive clusters using previously established clusters.
 - 1. <u>Agglomerative ("bottom-up")</u>: Agglomerative algorithms begin with each element as a separate cluster and merge them into successively larger clusters.
 - 2. <u>Divisive ("top-down")</u>: Divisive algorithms begin with the whole set and proceed to divide it into successively smaller clusters.
- **2.** <u>Partitional clustering</u>: Partitional algorithms determine all clusters at once. They include:
 - *K*-means and derivatives
 - Fuzzy *c*-means clustering
 - QT clustering algorithm

Common Distance measures:



• *Distance measure* will determine how the *similarity* of two elements is calculated and it will influence the shape of the clusters.

They include:

1. The Euclidean distance (also called 2-norm distance) is given by:

$$d(x, y) = \sum_{i=1}^{n} |x_i - y_i|$$

2. The <u>Manhattan distance</u> (also called taxicab norm or 1norm) is given by:

$$d(x, y) = \sqrt[2]{\sum_{i=1}^{p} |x_i - y_i|^2}$$

3.The <u>maximum norm</u> is given by:



 $d(x, y) = \max_{1 \le i \le p} |x_i - y_i|$

- 4. The <u>Mahalanobis distance</u> corrects data for different scales and correlations in the variables.
- 5. Inner product space: The angle between two vectors can be used as a distance measure when clustering high dimensional data
- 6. <u>Hamming distance</u> (sometimes edit distance) measures the minimum number of substitutions required to change one member into another.

K-MEANS CLUSTERING



- The k-means algorithm is an algorithm to <u>cluster</u> *n* objects based on attributes into *k* <u>partitions</u>, where *k < n*.
- It is similar to the <u>expectation-maximization</u> <u>algorithm</u> for mixtures of <u>Gaussians</u> in that they both attempt to find the centers of natural clusters in the data.
- It assumes that the object attributes form a <u>vector</u>
 <u>space</u>.



 An algorithm for partitioning (or clustering) N data points into K disjoint subsets S_j containing data points so as to minimize the sum-of-squares criterion

$$J = \sum_{j=1}^{K} \sum_{n \in \mathcal{S}_j} |x_n - \mu_j|^2,$$

where x_n is a vector representing the the nth data point and u_j is the <u>geometric centroid</u> of the data points in S_j.



- Simply speaking k-means clustering is an algorithm to classify or to group the objects based on attributes/features into K number of group.
- K is positive integer number.
- The grouping is done by minimizing the sum of squares of distances between data and the corresponding cluster centroid.





- Step 1: Begin with a decision on the value of k = number of clusters.
- Step 2: Put any initial partition that classifies the data into k clusters. You may assign the training samples randomly,or systematically as the following:
 - 1. Take the first k training sample as singleelement clusters
 - 2. Assign each of the remaining (N-k) training sample to the cluster with the nearest centroid. After each assignment, recompute the centroid of the gaining cluster.



- <u>Step 3:</u> Take each sample in sequence and compute its <u>distance</u> from the centroid of each of the clusters. If a sample is not currently in the cluster with the closest centroid, switch this sample to that cluster and update the centroid of the cluster gaining the new sample and the cluster losing the sample.
- <u>Step 4</u>. Repeat step 3 until convergence is achieved, that is until a pass through the training sample causes no new assignments.

Real-Life Numerical Example of K-Means Clustering

We have 4 medicines as our training data points object and each medicine has 2 attributes. Each attribute represents coordinate of the object. We have to determine which medicines belong to cluster 1 and which medicines belong to the other cluster.

Object	Attribute1 (X): weight index	Attribute 2 (Y): pH
Medicine A	1	1
Medicine B	2	1
Medicine C	4	3
Medicine D	5	4



<u>Step 1:</u>

- Initial value of centroids : Suppose we use medicine A and medicine B as the first centroids.
- Let and c_1 and c_2 denote the coordinate of the centroids, then $c_1=(1,1)$ and $c_2=(2,1)$



Objects-Centroids distance : we calculate the distance between cluster centroid to each object. Let us use Euclidean distance, then we have distance matrix at iteration 0 is

$$\mathbf{D}^{0} = \begin{bmatrix} 0 & 1 & 3.61 & 5 \\ 1 & 0 & 2.83 & 4.24 \end{bmatrix} \quad \begin{array}{c} \mathbf{c}_{1} = (1,1) & group - 1 \\ \mathbf{c}_{2} = (2,1) & group - 2 \\ \end{array}$$

$$\begin{bmatrix} 1 & 2 & 4 & 5 \\ 1 & 1 & 3 & 4 \end{bmatrix} \quad \begin{array}{c} X \\ Y \\ \end{array}$$



- Each column in the distance matrix symbolizes the object.
- The first row of the distance matrix corresponds to the distance of each object to the first centroid and the second row is the distance of each object to the second centroid.
- For example, distance from medicine C = (4, 3) to the first centroid $r_1 = (1,1)$ is $\sqrt{(4-1)^2 + (3-1)^2} = 3.61$ and its distance to the second centroid is , $r_2 = (2,1)$ is $\sqrt{(4-2)^2 + (3-1)^2} = 2.83$ etc.

<u>Step 2:</u>

- **Objects clustering** : We assign each object based on the minimum distance.
- Medicine A is assigned to group 1, medicine B to group 2, medicine C to group 2 and medicine D to group 2.
- The elements of Group matrix below is 1 if and only if the object is assigned to that group.

$$\mathbf{G}^{0} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 \end{bmatrix} \quad \begin{array}{c} group - 1 \\ group - 2 \\ A & B & C & D \end{array}$$



- Iteration-1, Objects-Centroids distances : The next step is to compute the distance of all objects to the new centroids.
- Similar to step 2, we have distance matrix at iteration 1 is

$$\mathbf{D}^{1} = \begin{bmatrix} 0 & 1 & 3.61 & 5 \\ 3.14 & 2.36 & 0.47 & 1.89 \end{bmatrix} \begin{array}{c} \mathbf{c}_{1} = (1,1) & group - 1 \\ \mathbf{c}_{2} = (\frac{11}{3}, \frac{8}{3}) & group - 2 \\ A & B & C & D \\ \begin{bmatrix} 1 & 2 & 4 & 5 \\ 1 & 1 & 3 & 4 \end{bmatrix} \begin{array}{c} X \\ Y \end{array}$$



Iteration-1, Objects
 <u>clustering</u>: Based on the new distance matrix, we move the medicine B to Group 1 while all the other objects remain. The Group matrix is shown below

$$\mathbf{G}^{1} = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix} \quad group - 1$$

$$A \quad B \quad C \quad D$$

• <u>Iteration 2, determine</u> <u>centroids</u>: Now we repeat step 4 to calculate the new centroids coordinate based on the clustering of previous iteration. Group1 and group 2 both has two members, thus the new centroids are $c_1 = (\frac{1+2}{2}, \frac{1+1}{2}) = (1\frac{1}{2}, 1)$ and $c_2 = (\frac{4+5}{2}, \frac{3+4}{2}) = (4\frac{1}{2}, 3\frac{1}{2})$





Iteration-2, Objects-Centroids distances : Repeat step 2 again, we have new distance matrix at iteration 2 as

$$\mathbf{D}^{2} = \begin{bmatrix} 0.5 & 0.5 & 3.20 & 4.61 \\ 4.30 & 3.54 & 0.71 & 0.71 \end{bmatrix} \quad \mathbf{c}_{1} = (1\frac{1}{2}, 1) \quad group - 1 \\ \mathbf{c}_{2} = (4\frac{1}{2}, 3\frac{1}{2}) \quad group - 2 \\ A \quad B \quad C \quad D \\ \begin{bmatrix} 1 & 2 & 4 & 5 \\ 1 & 1 & 3 & 4 \end{bmatrix} \quad X \\ \begin{bmatrix} 1 & 1 & 3 & 4 \end{bmatrix} \quad Y \end{bmatrix}$$

Iteration-2, Objects clustering: Again, we assign each object based on the minimum distance.

$$\mathbf{G}^{2} = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix} \quad group - 1$$

$$A \quad B \quad C \quad D$$

- We obtain result that G² = G¹. Comparing the grouping of last iteration and this iteration reveals that the objects does not move group anymore.
- Thus, the computation of the k-mean clustering has reached its stability and no more iteration is needed..



We get the final grouping as the results as:

<u>Object</u>	<u>Feature1(X):</u> weight index	<u>Feature2</u> (Y): pH	<u>Group</u> (result)
Medicine A	1	1	1
Medicine B	2	1	1
Medicine C	4	3	2
Medicine D	5	4	2

K-Means Clustering Visual Basic Code

Sub kMeanCluster (Data() As Variant, numCluster As Integer) ' main function to cluster data into k number of Clusters ' input:

- ' + Data matrix (0 to 2, 1 to TotalData);
- 'Row 0 = cluster, 1 =X, 2= Y; data in columns
- ' + numCluster: number of cluster user want the data to be clustered
- ' + private variables: Centroid, TotalData
- 'ouput:
- ' o) update centroid
- ' o) assign cluster number to the Data (= row 0 of Data)

Dim i As Integer Dim j As Integer Dim X As Single Dim Y As Single Dim min As Single Dim cluster As Integer Dim d As Single Dim sumXY()

```
Dim isStillMoving As Boolean

isStillMoving = True

if totalData <= numCluster Then

'only the last data is put here because it designed to be interactive

Data(0, totalData) = totalData ' cluster No = total data

Centroid(1, totalData) = Data(1, totalData) ' X

Centroid(2, totalData) = Data(2, totalData) ' Y

Else

'calculate minimum distance to assign the new data

min = 10 ^ 10 'big number

X = Data(1, totalData)

Y = Data(2, totalData)

For i = 1 To numCluster
```



```
Do While isStillMoving
' this loop will surely convergent
'calculate new centroids
'1 =X, 2=Y, 3=count number of data
ReDim sumXY(1 To 3, 1 To numCluster)
For i = 1 To totalData
sumXY(1, Data(0, i)) = Data(1, i) + sumXY(1, Data(0, i))
sumXY(2, Data(0, i)) = Data(2, i) + sumXY(2, Data(0, i))
Data(0, i))
sumXY(3, Data(0, i)) = 1 + sumXY(3, Data(0, i))
Next i
For i = 1 To numCluster
Centroid(1, i) = sumXY(1, i) / sumXY(3, i)
Centroid(2, i) = sumXY(2, i) / sumXY(3, i)
Next i
'assign all data to the new centroids
isStillMoving = False
For i = 1 To totalData
min = 10 ^ 10 'big number
X = Data(1, i)
Y = Data(2, i)
For j = 1 To numCluster
d = dist(X, Y, Centroid(1, j), Centroid(2, j))
If d < min Then
min = d
cluster = j
End If
Next j
If Data(0, i) <> cluster Then
Data(0, i) = cluster
isStillMoving = True
End If
Next i
Loop
End If
End Sub
```



Weaknesses of K-Mean Clustering

- 1. When the numbers of data are not so many, initial grouping will determine the cluster significantly.
- 2. The number of cluster, K, must be determined before hand. Its disadvantage is that it does not yield the same result with each run, since the resulting clusters depend on the initial random assignments.
- 3. We never know the real cluster, using the same data, because if it is inputted in a different order it may produce different cluster if the number of data is few.
- It is sensitive to initial condition. Different initial condition may produce different result of cluster. The algorithm may be trapped in the *local optimum*.

Applications of K-Mean Clustering



- It is relatively *efficient and fast.* It computes result at O(tkn), where n is number of objects or points, k is number of clusters and t is number of iterations.
- k-means clustering can be applied to *machine learning or data mining*
- Used on acoustic data in speech understanding to convert waveforms into one of k categories (known as Vector Quantization or Image Segmentation).
- Also used for choosing color palettes on old fashioned graphical display devices and Image Quantization.

CONCLUSION



 K-means algorithm is useful for undirected knowledge discovery and is relatively simple.
 K-means has found wide spread usage in lot of fields, ranging from unsupervised learning of neural network, Pattern recognitions, Classification analysis, Artificial intelligence, image processing, machine vision, and many others.