

② CLUSTERING: Introduction

Clustering is the process of grouping the data into classes or clusters, so that the objects within a cluster have high similarity in comparison to one another but are very dissimilar to objects in other clusters.

Although classification is an effective means for distinguishing groups, it requires costly collection and labelling of large data set.

Clustering is an example of unsupervised learning, which do not depend on predefined classes and class labels.

Hence clustering is a form of learning by observation rather than learning by examples.

K-Means clustering :-

K Means is a clustering algorithm for finding groups in the data,

The groups of data (cluster) are represented by their centers.

For the sample $X = \{x^t\}_{t=1}^N$

Consider $m_i = x^t$ is an approximated value of x^t

The error in approximation should be as minimum as possible.

$$\text{The Error } E(\{m_i\}_{i=1}^K | X) = \sum_t \sum_i b_i^t \|x^t - m_i\|^2$$

$$\text{where } b_i^t = \begin{cases} 1 & \text{if } \|x^t - m_i\| = \min_j \|x^t - m_j\| \\ 0 & \text{otherwise} \end{cases}$$

* K-means algorithm iteratively calculate b_i^t for all x^t . $b_i^t = 1 \Rightarrow x^t$ belongs to group m_i

* with new x^t being added to m_i , b_i^t changes and needs to be recalculated.

* these steps are repeated until m_i stabilizes

K-means algorithm

1. Initialize $m_i, i = 1, \dots, K$ for example, to K random x^t .

2. Repeat

for all $x^t \in X$

$$b_i^t = \begin{cases} 1 & \text{if } \|x^t - m_i\| = \min_j \|x^t - m_j\| \\ 0 & \text{otherwise} \end{cases}$$

for all $m_i, i = 1, \dots, K$

$$m_i \leftarrow \frac{\sum_t b_i^t x^t}{\sum_t b_i^t}$$

Until m_i converge

Disadvantage

Prof. Rajitha V N

K-Means algorithm is a local search procedure and the final m_i highly depend on the initial m_i .

Methods to overcome this disadvantage

→ Considering randomly selected k instances as initial m_i

→ Calculating principal component, dividing its range, partitioning the data of take means of groups as initial centres etc.

Note: Best way is to initialize centres where there is data.

PROBLEM

* Given the dataset of medicines, group it to relevant medicine cluster. Apply k-means for $K=2$.

Solve

medicine	attrib1	attrib2
A	1	1
B	1	0
C	0	2
D	2	4
E	3	5

consider any two data as the first centroids.

Let us consider

medicine A = first cluster, $C_1 \Rightarrow$ group 1

medicine C = Second cluster, $C_2 \Rightarrow$ group 2

step 1

Considering cluster/group 1 of medicine A i.e. $C_1 = (1, 1)$
calculate Euclidean distance.

$$A \rightarrow \sqrt{(1-1)^2 + (1-1)^2} = 0$$

$$B \rightarrow \sqrt{(1-1)^2 + (0-1)^2} = 1$$

$$C \rightarrow \sqrt{(0-1)^2 + (2-1)^2} = 1.41$$

$$D \rightarrow \sqrt{(2-1)^2 + (4-1)^2} = 3.2$$

$$E \rightarrow \sqrt{(3-1)^2 + (5-1)^2} = 4.5$$

111th for $c_2 = (0, 2)$

$$A \rightarrow \sqrt{(1-0)^2 + (2-1)^2} = 1.4$$

$$B \rightarrow \sqrt{(1-0)^2 + (2-0)^2} = 2.2$$

$$C \rightarrow \sqrt{(0-0)^2 + (2-1)^2} = 1$$

$$D \rightarrow \sqrt{(2-0)^2 + (4-2)^2} = 2.8$$

$$E \rightarrow \sqrt{(3-0)^2 + (5-2)^2} = 4.2$$

The ^{new} table

	For c_1	For c_2	
A	0	1.4	→ belongs to c_1
B	1	2.2	→ " " "
C	1.41	0	→ " " c_2
D	3.2	2.8	→ " "
E	4.5	4.2	→ " "

Due to new grouping, the centroid also changes

$$\text{New } c_1 = \frac{1+1}{2}, \frac{1+0}{2}$$

$$\text{New } c_2 = \frac{0+2+3}{3}, \frac{2+4+5}{3}$$

$$c_1' = (1, 0.5)$$

$$c_2' = (1.7, 3.7)$$

step 2 Consider c_1' & c_2' and compute Euclidean distance for given data.

we get

	For c_1'	For c_2'	
A	0.5	2.7	→ belongs to c_1
B	0.5	3.7	→ " c_1
C	1.8	2.4	→ " c_1
D	3.6	0.5	→ " c_2
E	4.9	1.9	→ " c_2

New centroid

$$c_1'' = \frac{1+1+0}{3}, \frac{1+0+2}{3}$$

$$= (0.7, 1)$$

$$c_2'' = \frac{2+3}{2}, \frac{4+5}{2}$$

$$= (2.5, 4.5)$$

1) Apply ^{Means clustering} KNN on the given dataset for $k=2$.
 $k=2$, $D = \{2, 3, 4, 10, 11, 12, 20, 25, 30\}$

Solu Taking two means randomly (if not given)
Initialize $m_1 = 4$ $m_2 = 12$

$$D_1 = \{2, 3, 4\} \quad D_2 = \{10, 11, 12, 20, 25, 30\}$$

$$\text{Step 1: } m_1' = \frac{2+3+4}{3} = 3 \quad m_2' = \frac{10+11+12+20+25+30}{6} = 18$$

Step 2: consider the $\{D\}$ for new m_1' & m_2'

$$D_1 = \{2, 3, 4, 10\} \quad D_2 = \{11, 12, 20, 25, 30\}$$

$$m_1'' = \frac{2+3+4+10}{4} = 4.75 \approx 5 \quad m_2'' = \frac{11+12+20+25+30}{5} = 19.6 \approx 20$$

Step 3: consider $\{D\}$ for new m_1'' & m_2''

$$D_1 = \{2, 3, 4, 10, 11, 12\} \quad D_2 = \{20, 25, 30\}$$

$$m_1''' = \frac{2+3+4+10+11+12}{6} = 7 \quad m_2''' = 25$$

Step 4: consider $\{D\}$ for m_1''' & m_2'''

$$D_1 = \{2, 3, 4, 10, 11, 12\} \quad D_2 = \{20, 25, 30\}$$

$$m_1^4 = 7 \quad m_2^4 = 25$$

Since we are getting same mean, the grouping or clustering changes no further.
Hence KNN concludes.

Step 3 : continuing the same process for new centroids

c_1^2 & c_2^2

we get

A	0.3	3.8	— c_1
B	1.04	4.7	— c_1
C	1.22	3.5	— c_1
D	3.3	0.7	— c_2
E	4.6	0.7	— c_2

The clustering/grouping does not change & hence the centroids also do not change.

Hence ~~K++~~ K-Means converges

Supervised Learning after clustering

- * Clustering can be used for two purposes.
 - Used for data exploration, to understand structure of the data.
 - Used to find similarities between instances and thus group instances.
- * The mean of the cluster group formed gives the representative prototype of instances in the group.
Ex: Consider a cluster which is formed for the sales of a product in a particular region.
An instance from that group will tell the requirement of the people of the region (CRM).
- * Clustering is also used as a preprocessing stage.
- * The advantage of having a unsupervised learning clustering before supervised learning is that unsupervised learning does not need labelled data and labelled data is expensive.

Hierarchical clustering

- * Hierarchical clustering aims at finding groups such that instances in a group are more similar to each other than instances in different groups.

* Hierarchical clustering makes use of the Euclidean distance measure.

Euclidean distance $d(a, b) = \sqrt{(d_a - d_b)^2}$, which is

a special case of Minkowski distance with $p = 2$

$$d_m(x^r, x^s) = \left[\sum_{j=1}^d (x_j^r - x_j^s)^p \right]^{1/p}$$

other distance measure is city-block distance (L1-norm)

* Hierarchical clustering approach has various types viz agglomerative clustering algorithm, divisive clustering, single link clustering, complete link clustering.

→ Agglomerative clustering:

1. Starts with N groups, each initially containing one training instance, then merging similar groups to form larger groups, until there is a single one.
2. At each iteration, two closest groups are chosen to merge.
3. The result of Agglomerative clustering is a hierarchical structure called dendrogram.
4. Dendrogram represent a tree, where leaves corresponds to instances which are grouped in the order of their merge.

→ Dividing clustering algorithm:-

1. It starts with a large group, and dividing large groups into smaller groups, until each group contains a single instance.

→ single link clustering:

1. The distance between instance used for grouping is defined as the smallest distance between all possible pairs of elements of the two groups.

$$d(G_i, G_j) = \min_{x^i \in G_i, x^j \in G_j} d(x^i, x^j)$$

2. If we consider a weighted, completely connected graph with nodes being instances and the edges with weights being distance b/w instances. Then single link method corresponds to constructing minimal spanning tree.

→ complete-link clustering

1. The distance between two groups is taken as the largest distance between all possible pairs.

$$d(G_i, G_j) = \max_{x^i \in G_i, x^j \in G_j} d(x^i, x^j)$$

single link cluster fig 7.5 from text.

Choosing the Number of clusters :-

- * The complexity of clustering depends on the no. of clusters ' k '.
- * There are various ways to fine tune k , such as.
 1. In case of colour quantization, k is defined by the application.
 2. Plotting the data in 2-D using PCA.
 3. An incremental approach, where maximum allowed distance is selected and made equivalent to maximum allowed reconstruction error.
 4. Manual check on clusters being meaningful groups of data.
- * Reconstruction error can be plotted as a function of ' k ' depending on the type of the clustering method used.
- * In hierarchical clustering, the difference between levels in the tree is considered to decide on a good split.

Expectation-Maximization Algorithm. (EM) Prof. Parag U

* EM works in a similar fashion to that of K-means, but EM yields a soft decision (elliptical curves) and K-means yields a hard decision (0/1, circles).

* K-means is applied on models of data which are independent.

EM is applied on a mixture model.

Note: Mixture model is a probabilistic model for representing the presence of subpopulations within an overall population, without subpopulation identity information. \Rightarrow Latent variables. (z_i)

- * Expectation-Maximization algorithm's approach is
1. To find maximum likelihood of parameters in statistical model (where the model depends on unobserved latent variables \Rightarrow Mixture model). the expectation (E) step which creates log likelihood function
 2. The Maximization (M) step, which computes parameters maximizing the expected log-likelihood found on the E-step.

* EM algorithm is an iterative method which alternates between E & M.

* The mixture density of a mixture model is given by

$$P(x) = \sum_{i=1}^K P(x|G_i) P(G_i) \quad \text{--- (1)}$$

eg $G_i \rightarrow$ mixture components

$P(x|G_i) \rightarrow$ component densities

Using (1), For the given sample $X = \{x^t\}_t$, the log likelihood is

$$\begin{aligned} L(\phi|X) &= \log \prod_t P(x^t|\phi) \\ &= \sum_t \log \sum_{i=1}^K P(x^t|G_i) P(G_i) \quad \text{--- (2)} \end{aligned}$$

* The goal of EM algorithm is to find vector ϕ that maximizes the likelihood of the observed values of X . Since z is also the parameter hidden in the model, likelihood L_c is

$$L_c(\phi|X, z) \quad \text{--- (3)}$$

* Since z values are not observed, L_c cannot be worked, hence working with its expectation

Q 2,

$$E \text{ step : } Q(\phi|\phi^l) = E[L_c(\phi|X, z)|X, \phi^l] \quad \text{--- (4)}$$

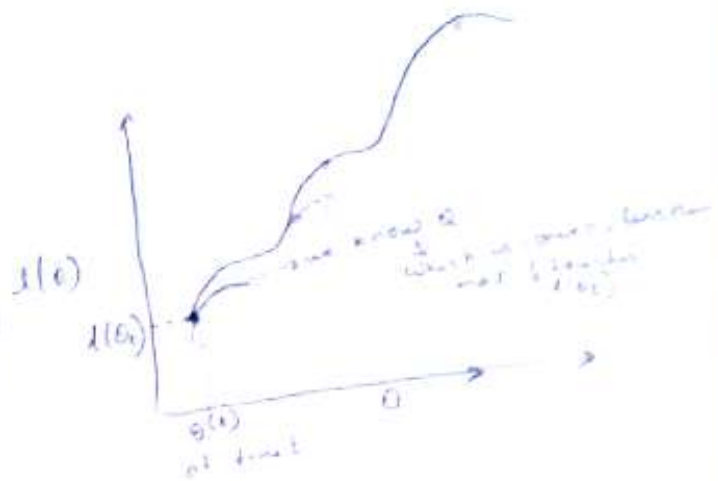
$\phi^l =$ current parameter value, l index iteration

In Maximization step (M), new parameter value ϕ^{t+1} is found, that maximizes ξ_q (4)

Thus,
M step: $\phi^{t+1} = \arg\max_{\phi} \xi(\phi | \phi^t)$ — (5)

Illustration

ϕ to be found where $\xi(\phi)$ is max
approach is gradient descent
if not possible
then it is EM



$$\xi(\phi^t) = \xi(\phi^t)$$

$\xi'(\phi) \leq \xi(\phi)$ ← bounded by concave function

Finding $\phi \rightarrow$ maximization

→ missing data
→ latent variables

while not converge
find ϕ using current ϕ
find next ϕ by maximizing ξ