

## Clustering

### Unsupervised Learning

In Supervised learning we are given a label associated with each data point, which helps in classification of data.

In unsupervised we don't have any labels associated with it and tasked to find some structure

One such approach could be to divide the given data into various distinct cluster based on some application.

- Use of to find pattern in stocks
- Dividing customer base into segment on some

### Extra

Can we use labeled data for clustering?

The thing is that we can use whatever we want to use, what matters is results. The results for labelled data is not as good as that of some algorithm for labelled data such as regression. They perform better.

Why?

Because Clustering tries to find relation b/w data points whereas in labelled data we want relation b/w data point and its label. So they perform poorly as they fail to do so.

for eg. if some person visits some website e-commerce label when he buy something

Classifier tries to predict whether he will buy this time or not (what he will buy)

whereas Clustering may try to cluster it type of user with its parameter based on the datapoints

### K-Means Clustering

This is an iterative algorithm which performs algorithm which has two steps in each

first: randomly takes k - data point as center

Second Then

for every iteration, colour the points based on its nearest centroid

Then move centroid to center of all colored points

then repeat until not a single ~~centroid~~ print change its color.

Can k-means only segment well separated data?

Even where it is not seen to segment data, k-means clustering separate our data points. This uses heavily on market segmentation

k-means clustering always give non overlapping k-segments

In Linear Regression we have cost function to reduce. In Clustering we take similar approach to minimize the avg. distance of all the cluster centroids with the datapoint in that cluster

Formally let

$c^{(i)}$  = index of cluster to which data point  $x^i$  is assigned

$\mu_k$  = cluster centroid

$\mu_i$  = centroid of cluster to which  $x^i$  is assigned

$$J = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

And our objective is to minimize this cost function

This cost function is also called distortion.

$$c_i = \operatorname{argmin}_{c \in K} \|x_i - c\|_2^2$$

$i \in \{1, 2, \dots, m\}$

$$J(c, u) = \frac{1}{m} \sum_{i=1}^m \|x_i - u_i\|_2^2$$

$$\delta(i, k) = \begin{cases} 1 & c^* = k \\ 0 & \text{otherwise} \end{cases}$$

Delta function

Delta assignment help us to write simpler formula

$$N_k = N_k + \delta(i, k)$$

$N_k$  will contain no. of data points in cluster

$$\hat{u}_k = \frac{1}{N_k} \sum_{i=1}^{N_k} x_i \delta(i, k)$$

$$\hat{u} = \{\hat{u}_1, \hat{u}_2, \dots\}$$

$$u = \hat{u}$$

Go back to cluster assignment until convergence

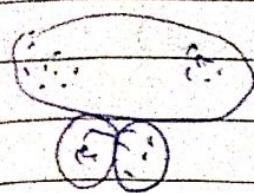
$$(J_K(c, u) - J(c, u)) \leq \epsilon_{ps}$$

↓  
convergence

by means of the given data point

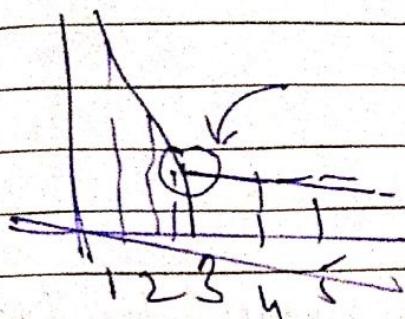


## Local optimum problem



Practical Solution : run it some n-times with random initialization.  
one P with least error is chosen then.

## Floyd method



$k=3$  is more accurate

After

After some "k" performance decreases.

P.

## Fuzzy C-means Clustering

### Fuzzy Set

Normal set  $A, x$

Either  $x \in A$ , or  $x \notin A$

Fuzzy set

i.e. membership function

$$0 \leq M_A(x) \leq 1$$

Not clear division in set.

### Fuzzy logic

A statement is true

with membership value

Neural Nets are considered to be part of soft computing, as it tells about "being something with some membership value".

## fuzzy c-partition

usual c-partition ( $c \geq 2, c \in I$ )

of a set  $S = \{x_1, \dots, x_n\}$   
represented

$$\text{Be } P(A_1, A_2, \dots, A_c)$$

Definition

$$A_i \neq \emptyset$$

$$A_i \cap A_j = \emptyset$$

$$\bigcup_{j=1}^c A_j = S$$

## fuzzy c-partition

Give membership value, represented by  $(U, S)$

where  $U_{n \times c}$  matrix

$U = ((U_{ij}))_{n \times c}$  where  $U_{ij}$  ! membership of  $i^{th}$  point to  $j^{th}$  fuzzy set

$$1 \leq i \leq n, 1 \leq j \leq c$$

Satisfying the property stated here.

$$i) 0 \leq U_{ij} \leq 1$$

$$ii) \sum_{j=1}^c U_{ij} = 1 \quad \forall i, i = 1 \dots n$$

$$iii) 0 < \sum_{i=1}^n U_{ij} < n \quad \forall j, j = 1 \dots c$$

Let  $\sigma > 1$  let  $V_p = \left( \sum_{i=1}^n U_{ij}^{\sigma} \right) / \left( \sum_{i=1}^n U_{ij} \right)^{\sigma}$   
mean for  $j^{th}$  cluster

$$J_{\sigma}(U, S) = \sum_{j=1}^c \sum_{i=1}^n U_{ij}^{\sigma} \|x_i - v_j\|^2 / \left( \sum_{i=1}^n U_{ij} \right)^{\sigma}$$

minimize

$A$  is positive definite matrix ~~non-singular~~

Page No.  
Date

Given  $S, R, A$  need to find optimal point.  
So  $J$  is minimum.

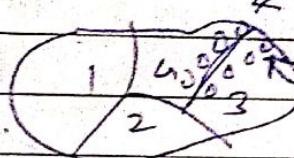
## FCM Algorithm

- 1  $\rightarrow$  we are given  $S, r, A$  and  $c$
- 2  $\rightarrow$  start with fuzzy  $c$ -partition  $V$  of  $S$
- 3  $\rightarrow$  compute fuzzy  $c$ -partition  $V_j^r \quad j=1 \dots c$
- 4  $\rightarrow$   $V_{ij} = \left( \sum_{k=1}^c \left( \frac{d_{ik}^2}{d_{ik}^2 + d_{jk}^2} \right)^{\frac{1}{2r}} \right)^{-1}$   $1 \leq i \leq n \quad 1 \leq j \leq c$

$$d_{ij}^2 = (x_i - v_j)^T A (x_i - v_j)$$

5  $\rightarrow$  ~~goto~~ if (convergence)  
break

else goto (3)



high membership func for  $k=1$

lower mem func for  $k=1, 2$

highest for 3, lower  
for 4

Data point near to a cluster mean will have higher membership function for that cluster.

### Advantage

worst case performance is same as k-means  
and where data is less fuzzy tends to give more stable cluster means.

### ② FCM

$\rightarrow$  Same initialization as of k-means.

FCM

K-means

Cluster assignment

$$c^i = \underset{1 \leq k \leq K}{\operatorname{argmin}} (x^i - \mu_k)^2$$

$$c^i = \frac{1}{1 + \sum_{k=1}^K (x^i - \mu_k)^2}$$

1 ≤ i ≤ M

$$\sum_{k=1}^K (x^i - \mu_k)^2$$

gives higher membership

to cluster which are

steps

$$\sum E(y, c) = \frac{1}{M} \sum_{i=1}^M (x^i - \mu_c)^2$$

$$\sum_{k=1}^K c_k = 1.0$$

$$J(c, y) = \frac{1}{M} \sum_{i=1}^M \sum_{k=1}^K (x^i - \mu_k)^2 c_k$$

New means

$$\hat{\mu}_k = \frac{1}{N_k} \sum_{i=1}^M x^i s(i|k)$$

$$\hat{\mu}_k = \frac{1}{N_k} \sum_{i=1}^M c_k x^i$$

Minimization

$$u = \hat{\mu}$$

Step ① → ④

$$u = \bar{x}$$

Step ① → ④

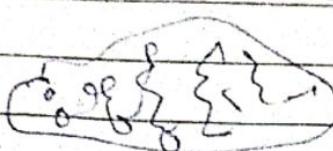
Model based Clustering

Principle:- Assume data points are from specific mathematical model

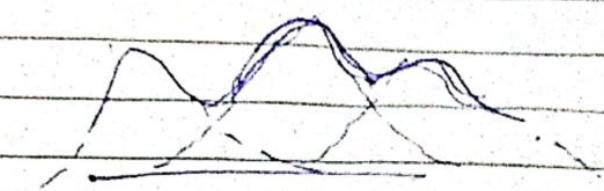
using training data point estimate model parameters.

## Gaussian Mixture Model

based on Gaussian method  
data points are distributed as Mixture of  
Many Gaussians models



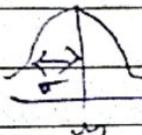
Top-view



Side view

can be shown with 2  
Gaussian distribution.

1) Gaussian Distribution  
or Univariate Gaussian distribution

$$g(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$


$$x = [x_1, x_2] \rightarrow 2D \text{ Multivariate GMM}$$

$\downarrow \quad \downarrow$

$$(\mu_1, \sigma_1) \quad (\mu_2, \sigma_2)$$

$$g(x_1) = \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{1}{2}\left(\frac{x_1-\mu_1}{\sigma_1}\right)^2}$$

$$g(x_2) = \frac{1}{\sqrt{2\pi}\sigma_2} e^{-\frac{1}{2}\left(\frac{x_2-\mu_2}{\sigma_2}\right)^2}$$

$$g(x) = g(x_1) \cdot g(x_2)$$

$$\approx \frac{1}{(2\pi)^{1/2}} \frac{e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}}{\sigma_1 \sigma_2} \quad |\Sigma|^{1/2}$$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$$

$$x = [x_1, x_2]$$

$$\mu = [\mu_1, \mu_2]$$

# probabilistic  
constraints.

Should be followed by EM

Page No.

Date

N-dimension

$$g(x) = \frac{1}{(2\pi)^{\frac{N}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(\bar{x} - \bar{y})^T \Sigma^{-1} (\bar{x} - \bar{y})}$$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \end{bmatrix}$$

$$\lambda(x) = \sum_{k=1}^K w_k g(x, y_k, \Sigma_k)$$

$$\sum_{k=1}^K w_k = 1.0 \leftarrow \text{weight factor.}$$

$k=3$  for provided eg.

Expectation Maximization Approach

- "x" + Initial Model ' $\lambda$ '

Train  $\uparrow$

unimportant

$\rightarrow$  Iterative procedure.

after 1st iteration through  $\lambda$

$$\lambda = \frac{\lambda^1}{\downarrow}$$

weights generated from ~~#~~  
features!

$\lambda$  should be more representative of  
data points

history is not carried away

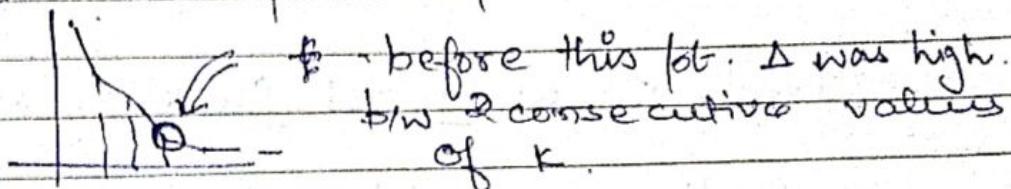
~~# f C~~

## Extras (Source: Medium and tds.com)

Q) Ways of finding  $k$  in k-means algorithm.

A) Elbow method

It is very common method start to tell when our accuracy increase at much slower rate than it was before we have reached some optimal fit.



B) Silhouette Curve

Silhouette coefficient calculates the density of the cluster the cluster by generating a score for each sample based on the difference between the average intra-cluster distance and mean nearest-cluster distance for that sample normalized by maximum value. We select best  $k$  using best score calculated.

To Do:

DBSCAN: Density Based clustering

Disadvantage of k-Means / Partition based Clustering.

Unable to fit arbitrary-shape cluster

- prone to outliers : outliers are assigned to clusters, after that they pull mean vector towards them which makes their cluster prone to far point than near points.

### DBSCAN can

DBSCAN Density based can find out any arbitrary shape cluster without being affected by noise.

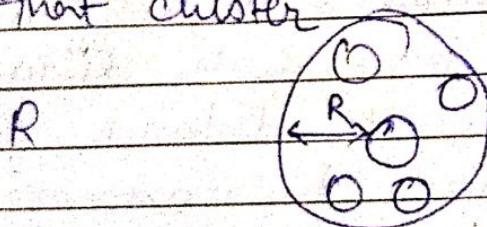
### Density Based Spatial Clustering of Application

R: Radius of Neighbour hood

M: Min number of neighbour

Institution / Idea.

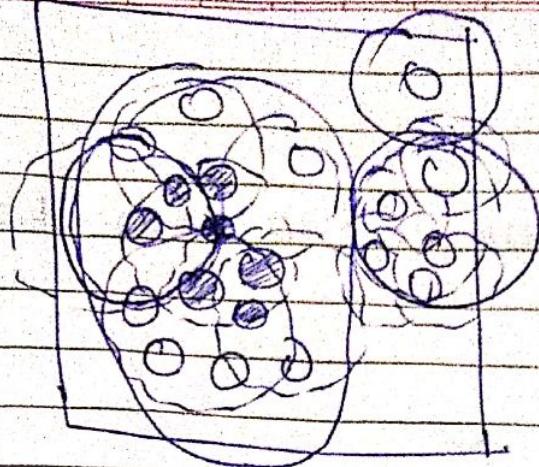
If a point is in a cluster of it should be near to lots of other points in that cluster.



If taking any point if in a given Radi 'R' If have more than M neighbours we call it Dense area or ~~area~~ cluster

3 types of Data point

- Core Points
- Border
- Outlier



$$R = 2$$

$$M = 6$$

Core : If there are  $M$  points or more in a given radius.

Border If there are points but less than  $M$  it is border

Outlier : if points in neighbourhood is zero

### Soft Clustering

Core points surrounded by border points generates arbitrary shape

### Disadvantage

- Need to select 2 parameter 'R' & 'M'

### Advantage

- Arbitrary shaped cluster
- Robust to classifier
- Doesn't require specification of cluster