

unsupervised machine learning

machine learning

supervised ml

unsupervised ml

clustering

updated

① k-means → k-means++

② hierarchical

③ DBSCAN

Regression

① linear → L1
→ L2

② logistic

③ SVR

④ DTR

⑤ RFR

⑥ GBR

⑦ XBR

⑧ KNNR

Classification

① logistic

② SVC

③ DTC

④ RFC

⑤ XBC

⑥ GBC / ABC

⑦ KNNC

target / dependent / supervised.
supervised ml
↓
supervisor
(body mass index)

Height	Weight	BMI
170	60	21
180	65	22
160	70	20
165	75	18
140	75	19

unsupervised ml

↓

Country

dataset	Height	Weight	BMI	Country
unsupervised ml	170	60	21	IND
↓	180	65	22	UK
clustering	160	70	20	USA
↓	165	75	18	IND
grouping	140	75	19	USA

2 cluster

3-groups

mathematically

- ① k means
- ② hierarchical
- ③ DBSCAN



① k - means

→ Data set

→ find similarity data

→ distance

→ deciding range

① centroid

② distance

③ mean

k - means

↑

no. of centroids

2, 3, 4, 5, ... we can't decide

~~cluster~~ → ~~same~~

→ ELBOW method → decided

→ wcss → inter cluster

based on distance we find

similarity

wcss = within cluster sum of square

Height

weight

→ clustering

1) c_1 {185} 72

2) c_2 {170} 56

3) {168} 60

4 179 68

5 182 72

6 188 77

7 180 71

8 160 70

9 183 84

10 180 88

11 180 67

12 167 76

① centroid (randomly)

$k=2$

initially take data point and build around cluster

No of centroid = Number of cluster

c_1
(185, 22)

③
(168, 60)

$$\Rightarrow \sqrt{(168-185)^2 + (60-22)^2}$$

$$\Rightarrow \sqrt{\quad}$$

$$D(c_1, 3) \Rightarrow 20.80$$

$$D(c_2, 3) \Rightarrow \begin{matrix} c_2 & ③ \\ (170, 56) & (168, 60) \end{matrix}$$

$$\Rightarrow \sqrt{(168-170)^2 + (60-56)^2}$$

$$\Rightarrow \sqrt{2^2 + 4^2}$$

$$D(c_2, 3) \Rightarrow \sqrt{20} = 4.4$$

if belongs to 2 cluster because distance minimum

* nearest distance we consider that cluster

Validation :

→ Dunn index

→ silhouette coeff

$k=2$

(185, 72)

distance 2 point c_1

$$E.D \Rightarrow d(c_1, 3) \rightarrow c_1 \rightarrow 3$$

$$\Rightarrow 5 \text{ dist}$$

$$\rightarrow d(c_2, 3) \rightarrow c_2 \rightarrow 3$$

$$\Rightarrow 8 \text{ dist}$$

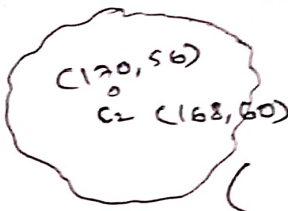
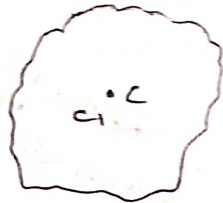
$[C_1 (170, 56), C_2 (168, 60)]$

↓

$(\frac{170+168}{2}, \frac{56+60}{2})$ update the centroid

$(169, 58)$

k-means $k=2$ → 2 centroid → around build cluster → dis → similarity → nearest point



$$D(C_1, 4) = \sqrt{(185-179)^2 + (72-68)^2} = 7.21 \checkmark$$

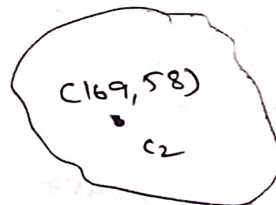
$$D(C_2, 4) = \sqrt{(169-179)^2 + (58-68)^2} = 14.14$$

why? low distance min distance

update centroid

$$\frac{179+185}{2} = 182$$

$$\frac{68+72}{2} = 70$$

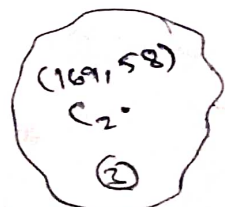


$$D(C_1, 5) = \sqrt{(182-182)^2 + (72-70)^2} = 2 \checkmark$$

$$D(C_2, 5) = \sqrt{(182-169)^2 + (72-58)^2} = 19.1$$

update

$$\frac{182+182}{2} = 182, \frac{70+72}{2} = 71$$



mean

- ① centroid
- ② distance (compare, min)
- ③ include point in cluster, update the cluster

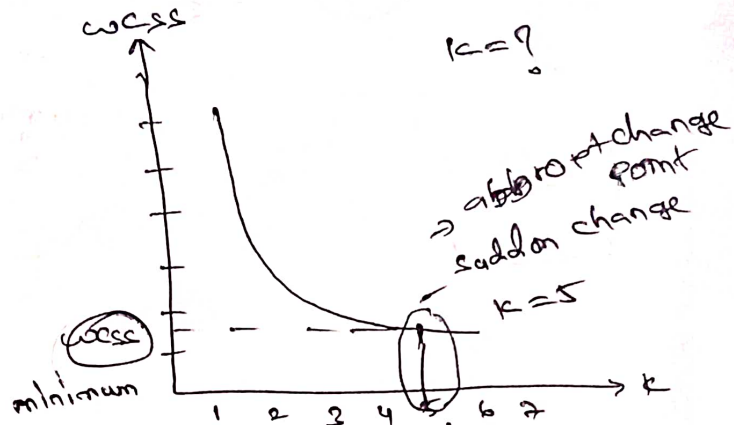
iterate all the data point

K-means

$K = 2, 3, 4, 5$

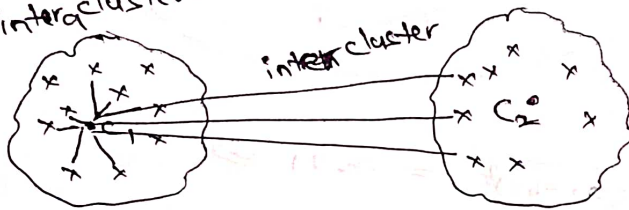
[2 cluster]

[ELBOW] method



WCSS \rightarrow within cluster sum of square optimized cluster point

intercluster distance



Point and centroid

inter = (C_1, C_2)

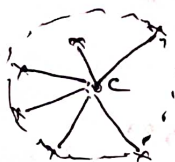
inter = within cluster

average of cluster

$K=1$

1 cluster

WCSS = 25



= WCSS₁

$K=2$

2 cluster

WCSS = 20



two groups



WCSS₂ =

WCSS₁ > WCSS₂

* distance minimized in two groups

* solid compare to $K=1$

wcss = within cluster sum of

$$\sum_{i=1}^n d(c, x_i)$$

$K=3$

wcss = 15



wcss₁ → wcss₂ → wcss₃

cluster is going to more solid
minimum wcss

difference between k-means / k-means++

Validate clustering ? $K=5$

① dunn index

(distance)

② silhouette score

Regression Coeff
 R^2 0 worst
Adj R^2 1 best

classification
→ for AUC
confusion
matrix

① dunn index =

$$\frac{\max \text{dis}(x_i, x_j)}{\max \text{dis}(y_i, y_j)}$$

→ c_1 to c_2

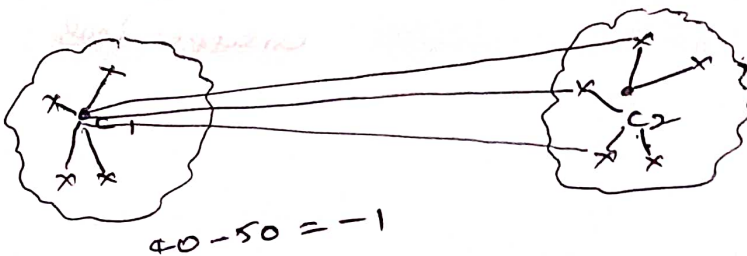
→ same cluster

worst best
→ -1 +1

② silhouette score =

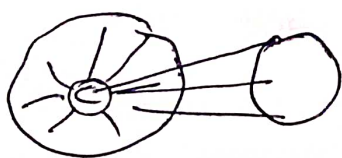
$$\frac{b_i - a_i}{\max(b_i - a_i, c_i - a_i)} \quad [-1 \text{ to } 1]$$

$$b_i = 50 - 40 = 10$$



Same cluster
 a_i

$$a_i > b_i$$



① Unsupervised

② k-means $\rightarrow k=9$

③ how to choose optimal $k=9$

Elbow Loss

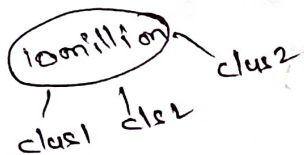
④ $k=5 \rightarrow$ how to validate

silhouette - $[-1, +1]$

how to make a best model or optimize sol

Custom learning or custom model = supervised + unsupervised
(semi supervised)

weight	height	Gender
170	55	m
180	60	p
165	70	3
180	80	p
155	50	p
160	100	p

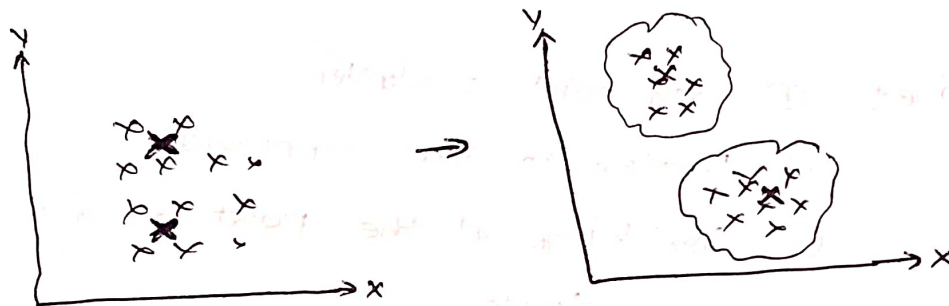


- ① k-means \rightarrow k-means++
- ② hierarchical
- ③ ~~linear~~ DBSCAN

① k-means

- ① random initialization of centroid
- ② find out the dist to all the point and make cluster (min dist)
- ③ update the centroid

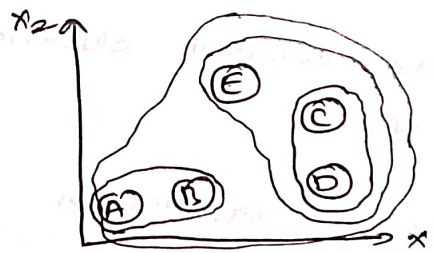
k-means algo is centroid based algo



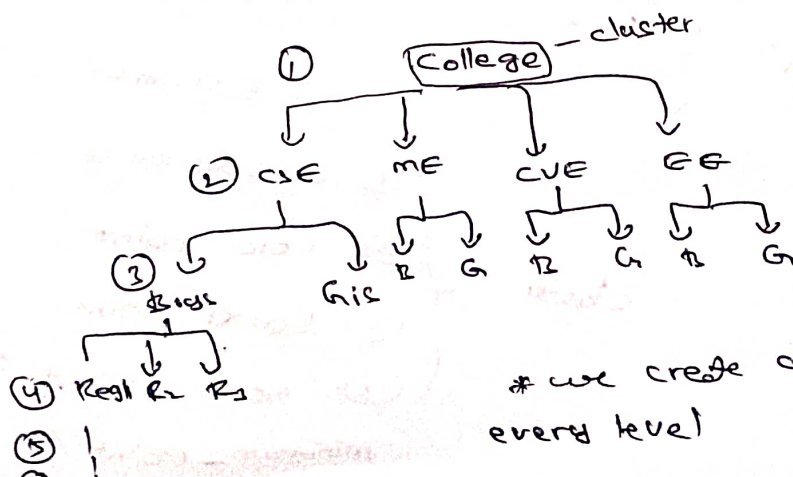
$k=2$
 $k=3$
 $k=4$

② hierarchical clustering:

(A) (B) (C) (D) (E)



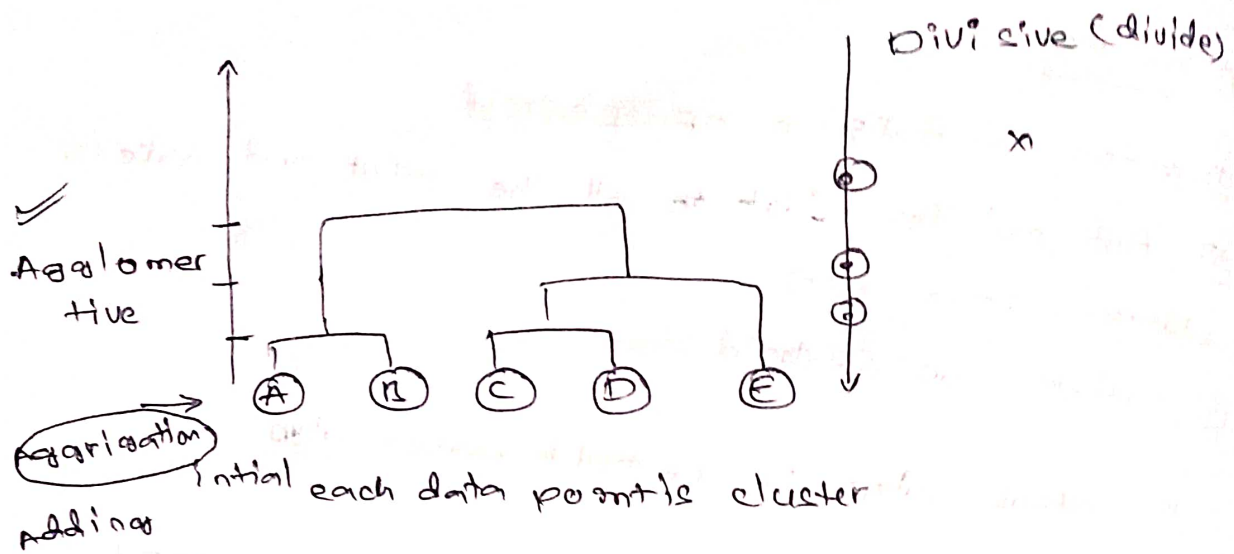
make cluster near each other
 * each data point is cluster



* we create cluster each
 every level

how to represent hierarchical cluster?

Dendrogram → just a representation of hierarchical clustering



Agglomerative: ① each point is cluster

② Bottom to top approach

③ Combining all the point as a single cluster

② DBSCAN

Density based spatial clustering with application to noise

⇒ Density based approach



k-means x

hierarchical x

DBSCAN ✓

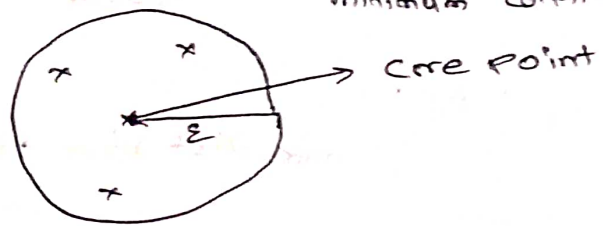
Cluster

- ① eps distance (unit distance)
- ② Core point
- ③ Border point
- ④ Noise point (outlier)

Minimum Point

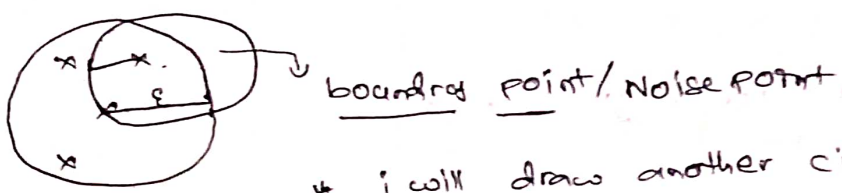
EPS distance = ϵ \Rightarrow draw circle

minimum point $\geq \epsilon$ this much point should be inside circle
 * if inside circle
 minimum condition satisfy is core point



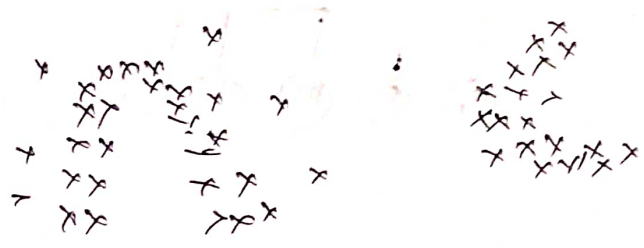
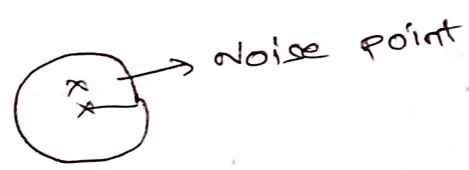
hyper parameters ① EPS distance

② min point to be consider in a circle



* i will draw another circle if
 inside any ^{core} point come in the circle it is
 called boundary point
 * it should be neighbour of the core point

Noise point which not satisfy both conditions



k-mean is unable to solve this type of data

so we use DBSCAN

k-mean
 hierarchical
 DBSCAN } point matrix

How many point \rightarrow 5 point

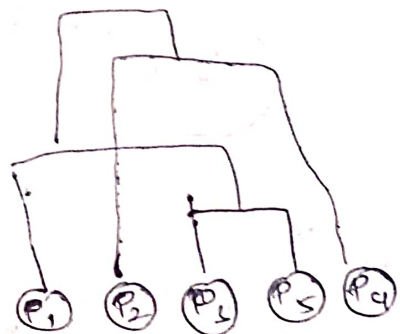
P_1 P_2 P_3 P_4 P_5

Find dist between the point and create cluster

	P_1	P_2	P_3	P_4	P_5
P_1	0				
P_2	9	0			
P_3	3	7	0		
P_4	6	5	9	0	
P_5	11	10	2		0

min-dist b/w point

	P_1	P_2	$(P_3 P_5)$	P_4	
P_1	0				10
P_2	9	0			9
$(P_3 P_5)$	3	7	0		8
P_4	6	5	8	0	7
					6
					5
					4
					3
					2
					1
					0



$$d(P_1, \{P_3, P_5\})$$

$$\min [d(P_1, P_3), d(P_1, P_5)]$$

$$\min [3, 11]$$

$$= 3$$

$$d(P_2, \{P_3, P_5\})$$

$$[d(P_2, P_3), d(P_2, P_5)]$$

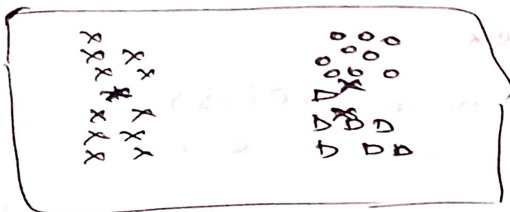
$$\min [7, 10]$$

$$\Rightarrow 7$$

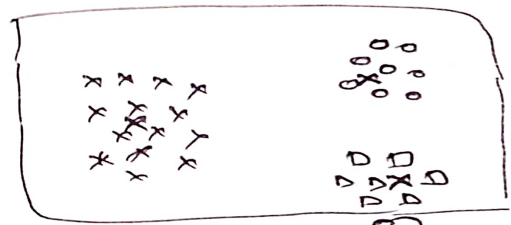
	P_1	P_2	P_5	P_2	P_5
P_1, P_2, P_5			0		
P_2				0	
P_5				5	0

$$D(P_2, [P_1, P_2, P_5])$$

$$\min D([P_2, P_1], [P_2, P_2], [P_2, P_5])$$



①



good

②

① it's all about centroid initialization

k-means - Randomly x C in many cases we are find ① doing

↓ updated

k-means ++

↓ initialized 1 centroid

Probabilistic

$C \xrightarrow{d} x_i$

(farthest)

→ make another centroid

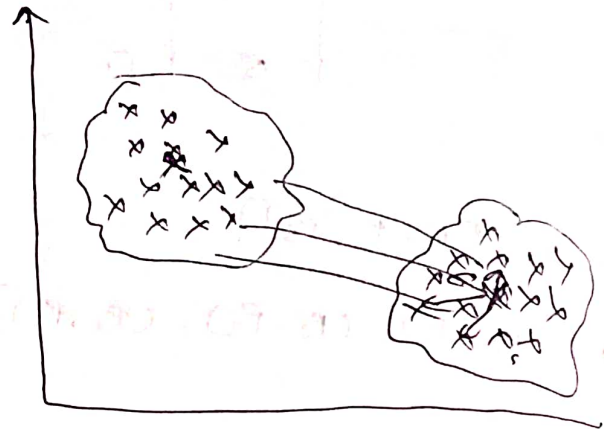
* initialize centroid as far as possible

known

↓
x x x x
x x x x
x x x x
x x

x x x
x x x
x x x
x x x
x x x

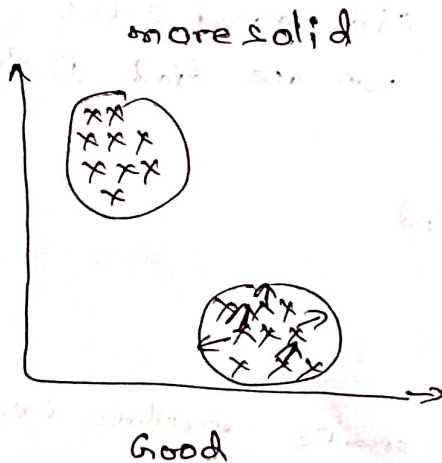
- cluster
- ① intra cluster dist (inside same cluster)
 - ② inter cluster dist (between 2 clusters)



Dunn index

$$D = \frac{\max d(i, j)}{\max d(k)}$$

\Rightarrow maximum inter cluster dist
max intra cluster dist



\uparrow as much as possible
 \downarrow as low as possible

$[0, 1]$ best dunn index $= \frac{\max(\text{inter})}{\max(\text{intra})}$