

K means clustering

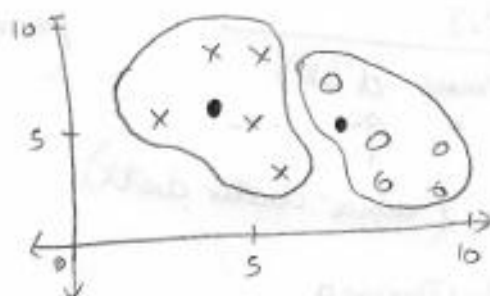
The output of K means is K clusters

Types of data in clustering analysis

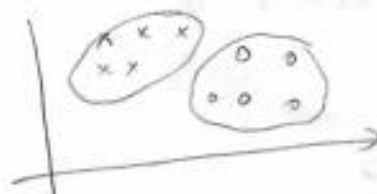
- 1.) P-dimensional points (vectors) in Euclidean space
- 2.) Nominal variables
- 3.) Strings, trees, graphs and other objects.
- 4.) Data of mixed types.

Algorithm

- 1.) Partition objects into k non empty subsets (clusters)
- 2.) Compute center for each cluster
- 3.) Assign each object to the cluster with the nearest center.
- 4.) Go back to step 2; terminate when there is no new assignment.



• → Center of each cluster.



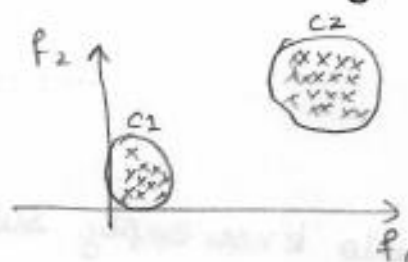
Metrics for clustering:

$$D = \{x_i\} ; \text{ no. } y_i \text{'s}$$

classification & regression y_i 's \leftarrow class labels
 \leftarrow regression values
{ground truth}

$$f(x) = y$$

geom: what is a clustering result (good)



Scatter plot

K-clusters

2-clusters

? Basis of how good our cluster is?

intra cluster \rightarrow means within a cluster \rightarrow this is small

inter cluster \rightarrow across between clusters \rightarrow this is kept large

✧ Dunn-index:

$$D = \frac{\min_{i,j} d(i,j)}{\max_k \min_{i,j} d'(k)}$$

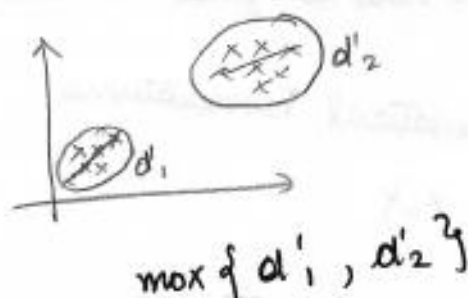
\rightarrow distance between C_i & C_j
(maximum inter cluster dist)
 \uparrow
(intra-cluster dist)

D is high \Rightarrow good clustering

$d(i,j)$ = dist between C_i & C_j 's farthest pts.



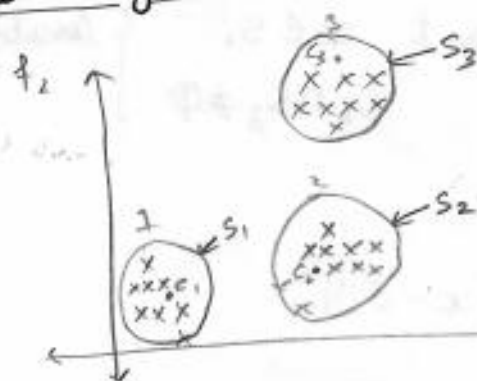
$\max d' \{K\}$



Many other measures / metrics for measuring clustering

ideal \rightarrow we should have high inter distance
we should have low intra distance

K-means: Geometric intuitions, Centroids



K-means
($K=3$)

c_1, c_2, c_3 : Centroids

$$S_1 \cup S_2 \cup S_3 = D$$

$$S_1 \cap S_2 = \emptyset$$

$$S_2 \cap S_3 = \emptyset$$

$$S_1 \cap S_3 = \emptyset$$

K : # Clusters

\hookrightarrow hyper parameter \rightarrow CV ideas like

K-Clusters \approx K-centroids $\rightarrow c_1, c_2, c_3, \dots, c_K$
K-Sets $\rightarrow S_1, S_2, S_3, \dots, S_K$

$$c_i = \frac{1}{n} \sum_{x_j \in S_i} x_j \leftarrow \text{mean-point to } S_i$$

K-means \rightarrow Centroid based clustering scheme.

Big Challenge :- How to find K-Centroid

K-means :- Mathematical Formulations

$$D = \{x_1, x_2, \dots, x_n\}$$

Task

K centroid :- c_1, c_2, \dots, c_k

Sets :- S_1, S_2, \dots, S_k

proximally \rightarrow (c_1, c_2, \dots, c_k)
 (S_1, S_2, \dots, S_k)

$$\text{argmin} \sum_{i=1}^k \sum_{x \in S_i} \|x - c_i\|^2 \quad \text{where cluster distance is minimized}$$

s.t $x \in S_i$
 $S_i \cap S_j \neq \emptyset$ } constraints

$$\text{argmin}_{c_1, c_2, \dots, c_k} \sum_{i=1}^k \sum_{x \in S_i} \|x - c_i\|^2$$

all clusters \uparrow

Sum of squared distance from centroid in cluster i

Very hard problem
 \downarrow
NP hard problem
exponential time complexity to solve problem

Approximation algorithm \rightarrow As hard to solve so introduced approximation algorithm

K-means: Lloyd's algorithm

-3-

① Initialization:

↙ randomly pick K points from D and call them our centroids c_1, c_2, \dots, c_K .

② Assignment:

for each point x_i in D

Loop: $\left\{ \begin{array}{l} \rightarrow \text{select the nearest } c_j \\ \text{dist}(x_i, c_j) \quad \forall j = 1, 2, \dots, K \\ \rightarrow \text{add } x_i \text{ to set } S_j. \end{array} \right.$

③ Recompute centroid / update

$x_i \rightarrow S_j \quad j = 1, 2, \dots, K$

\rightarrow recalculate / update centroids c_j 's as follows

$$c_j = \frac{1}{|S_j|} \sum_{x_i \in S_j} x_i$$

mean-pt.

④ Repeat step 2 & 3 until convergence

↙ assignment
↘ update

Convergence \rightarrow means when centroids don't change much.
(old centroid \approx new centroid)

$$(c'_1 - c_1, c'_2 - c_2, c'_3 - c_3, \dots, c'_K - c_K)$$

Small values

How to initialize K-means:

Lloyd's algorithm : initialization stage



random \leftarrow random : pick K -pts randomly from D
int
↓
 C_1, C_2, \dots, C_K

Problem:-

initialization sensitivity

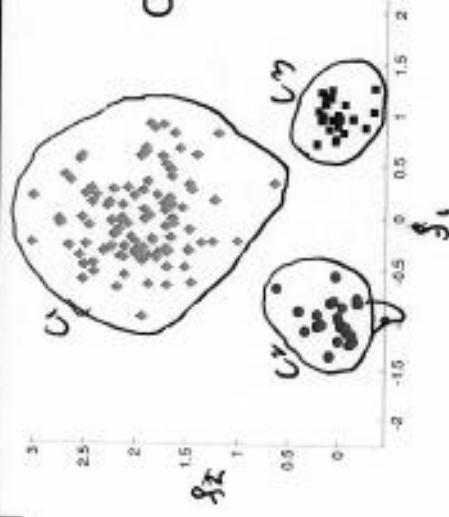
eg (s): Toy datasets.

final clusters & centroids
depends on how you
initialize

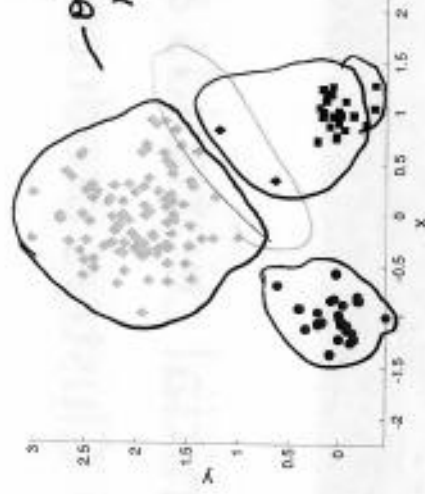
K-means Algorithm – Initialization

- Initial centroids are often chosen randomly.
- Clusters produced vary from one run to another.

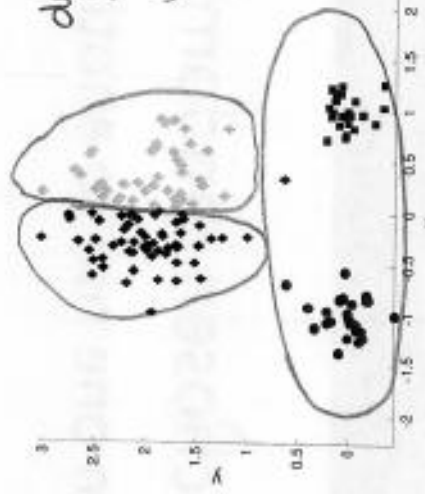
Two different K-means Clusterings



Original Points — human labels.



— one clustering

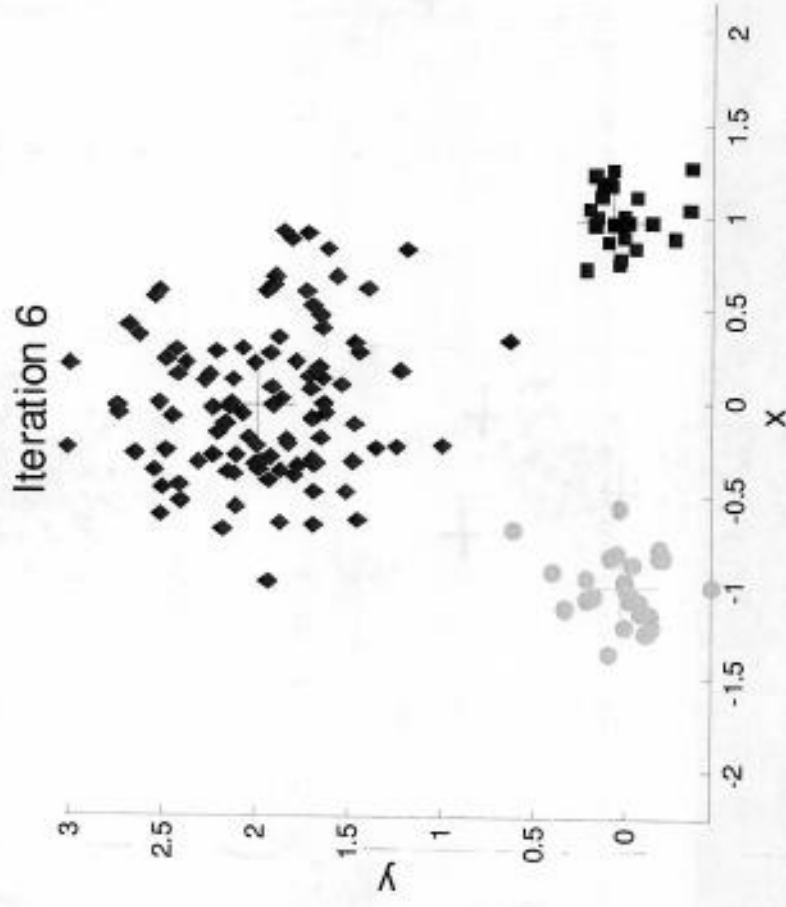


different initialization
under optimal
result

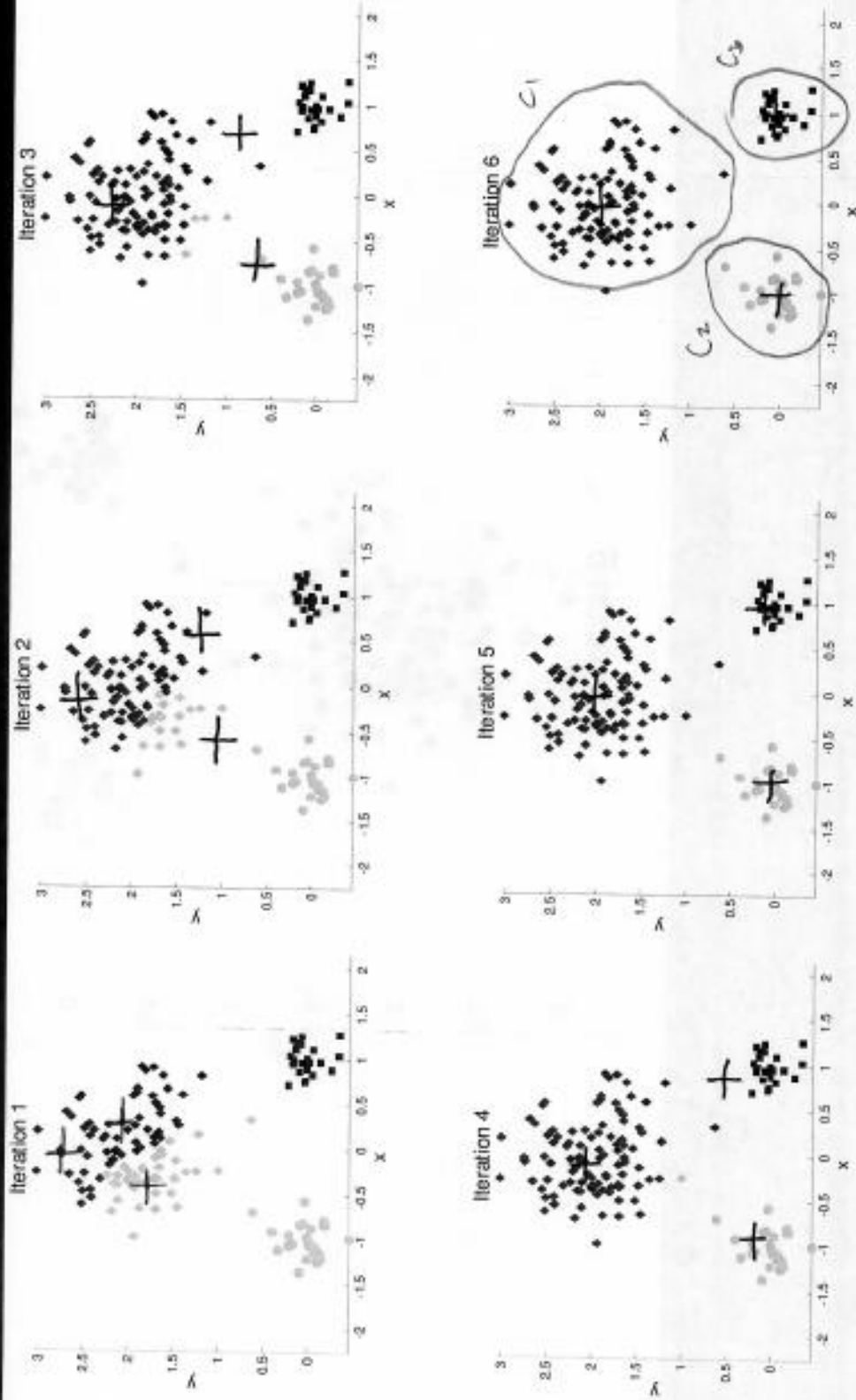
Optimal Clustering

Sub-optimal Clustering

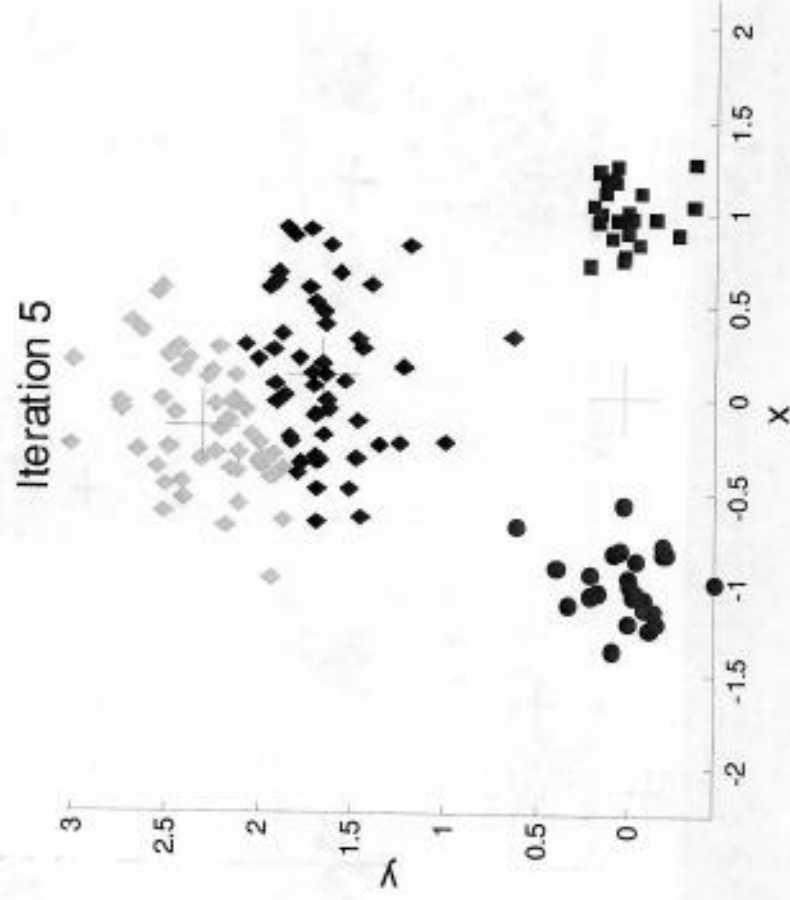
Importance of Choosing Initial Centroids



Importance of Choosing Initial Centroids

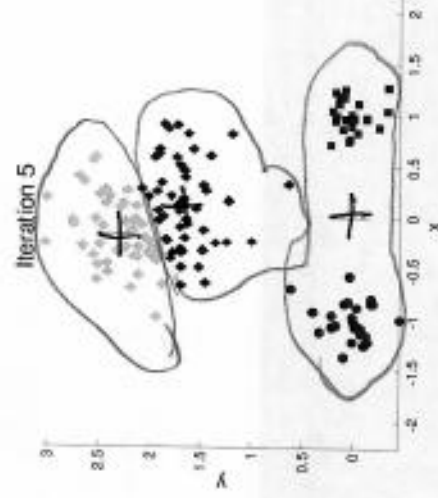
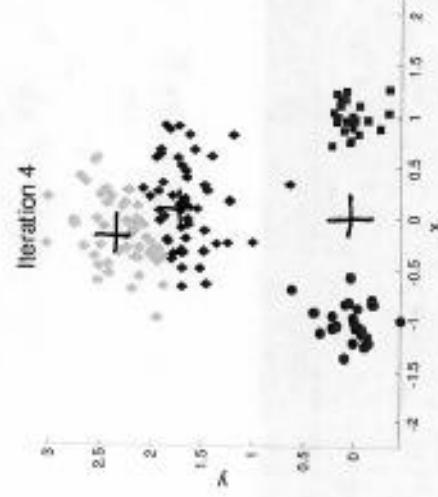
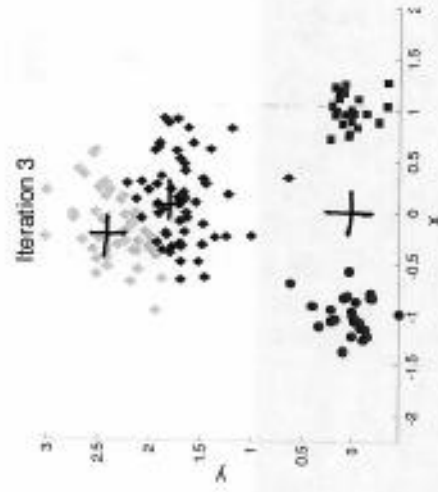
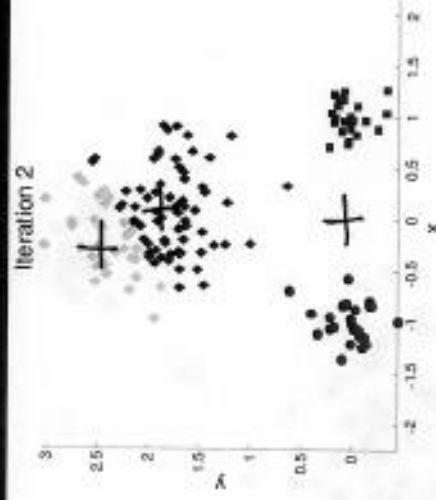
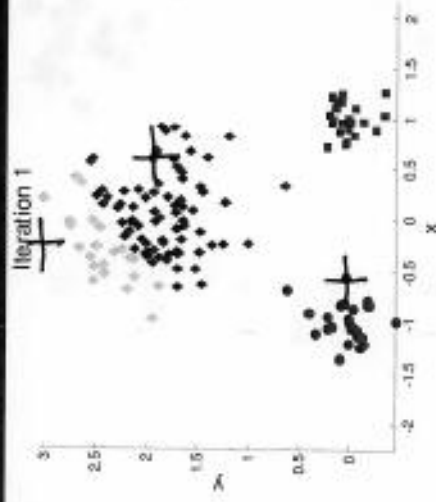


Importance of Choosing Initial Centroids



Importance of Choosing Initial Centroids ...

Choosing
different
cluster
centroid



Dealing with Initialization

- Do multiple runs and select the clustering with the smallest error
- Select original set of points by methods other than random . E.g., pick the most distant (from each other) points as cluster centers (K-means++ algorithm)

K-means Algorithm – Centroids

- The centroid depends on the distance function
 - The minimizer for the distance function
- 'Closeness' is measured by Euclidean distance (SSE), cosine similarity, correlation, etc.
- Centroid:
 - The mean of the points in the cluster for SSE, and cosine similarity
 - The median for Manhattan distance.
- Finding the centroid is not always easy
 - It can be an NP-hard problem for some distance functions
 - E.g., median for multiple dimensions

It can be improved

→ Repeat K-means multiple times with different initializations

↳ pick the clustering based on

{smaller intra-cluster}

{maximum inter-cluster}

② K-means ++

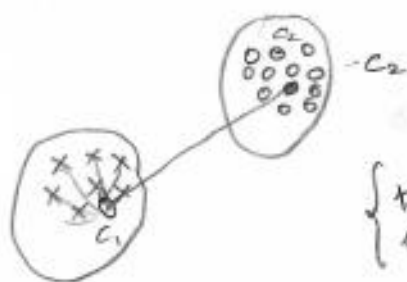
↳ random-init → smart init

initialization in K-means: (Task) pick c_1, c_2, \dots, c_k

① pick the first centroid randomly → c_1 from D .

② $\forall x_i \in D$ create a distribution:

$x_i \rightarrow \text{distance}^2(x_i, \text{nearest centroid})$



{ probabilistic
approach

x_1	d_1
x_2	d_2
x_3	d_3
\vdots	\vdots
x_n	d_n

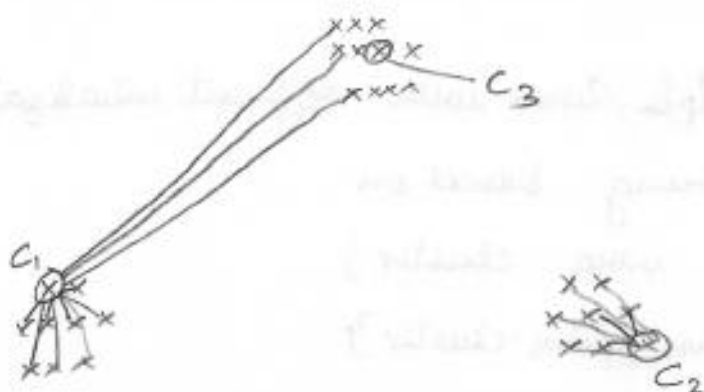
$$\|x_i - c_i\|^2$$

pick a point from

$D - \{c_1\}$

with a prob.

proportional to d_i



In initialization we are trying to pick points as (centroid) that are as far as possible from other centroids.

Q Why do ^{we} this probabilistically?

Ans pick a point which has the highest value of $d^2(x_i, \text{nearest centroid})$

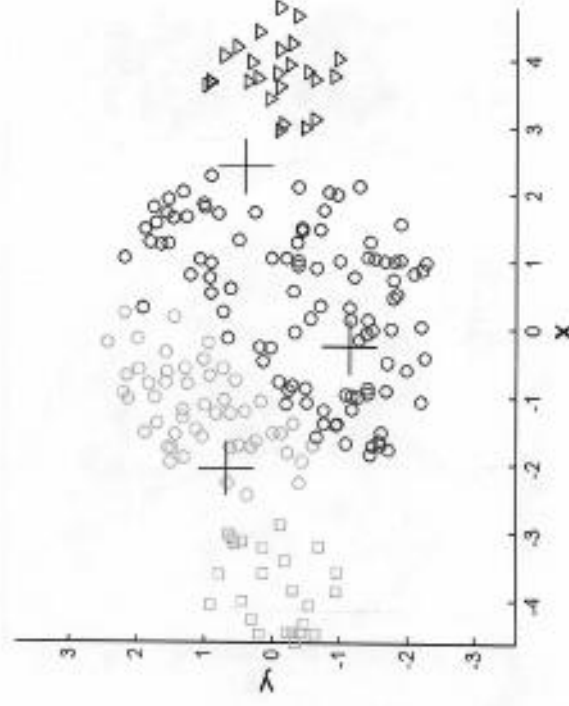
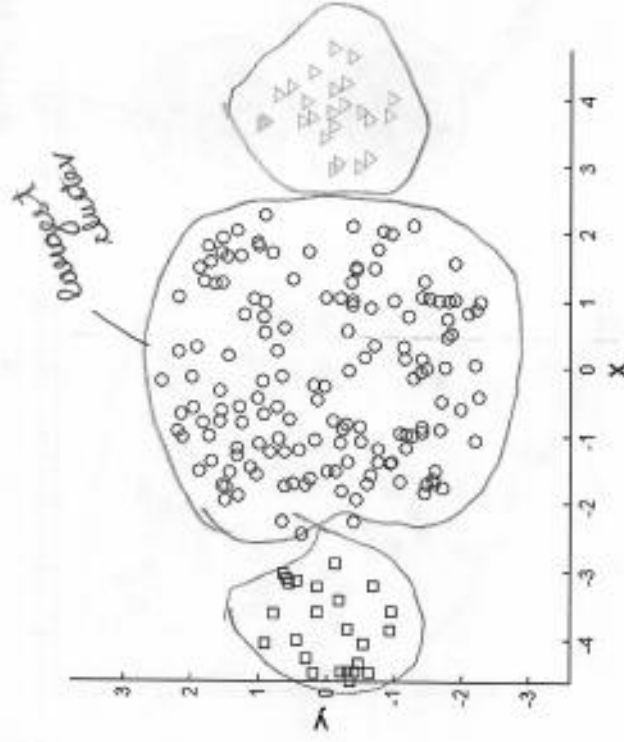
Because the point farther can be outlier. Which we don't want

K++ gets effected by outlier but so we use probability for this.

Limitations of K-means

- 1.) clusters of different sizes
- 2.) clusters of different densities
- 3.) cluster of Non-globular shape (non-convex shape)
- 4.) Kmeans has problems when the data contains outliers

Limitations of K-means: Differing Sizes

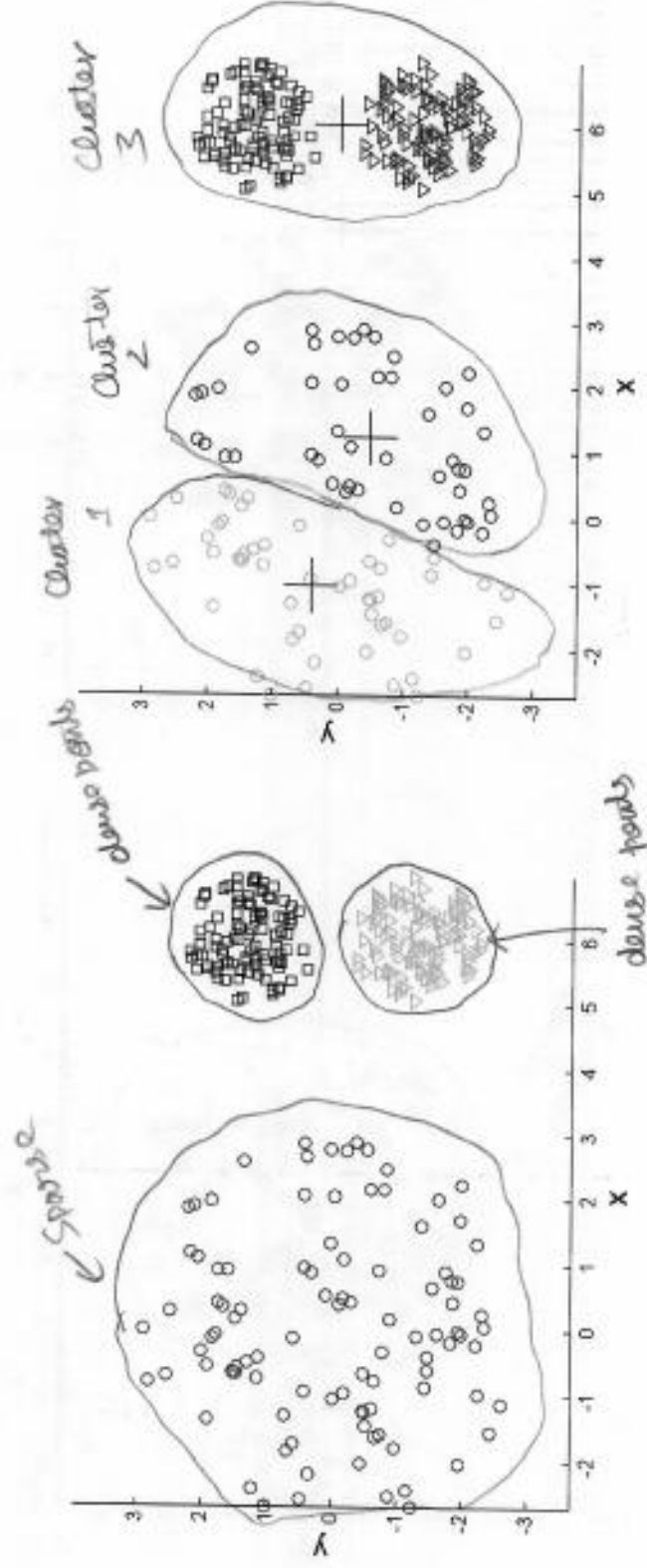


Original Points

K-means (3 Clusters)

clusters of same size always formed that is problem.

Limitations of K-means: Differing Density

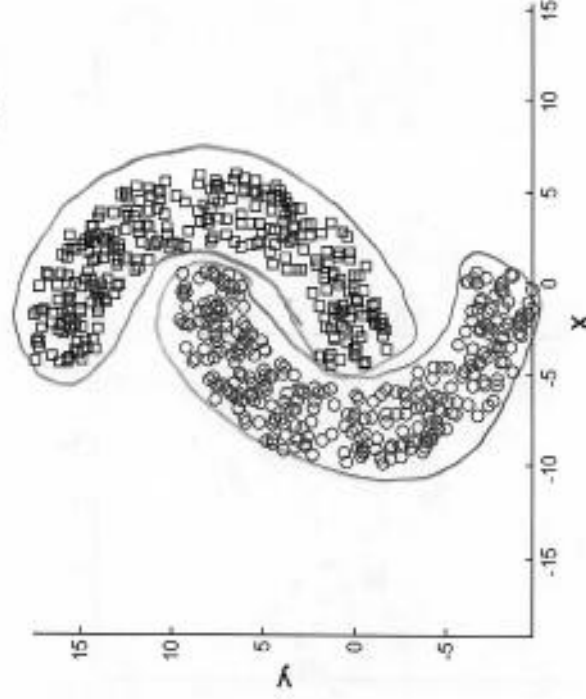


Original Points

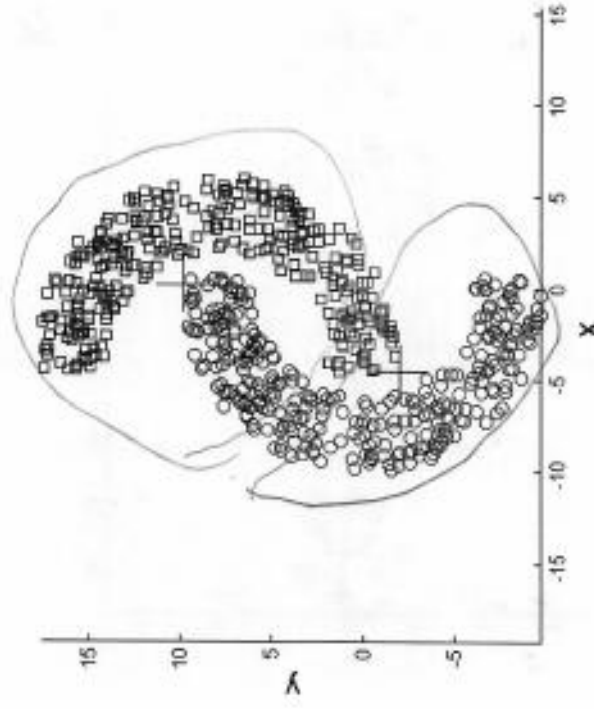
K-means (3 Clusters)

Limitations of K-means: Non-globular Shapes

non convex data / Non globular shapes

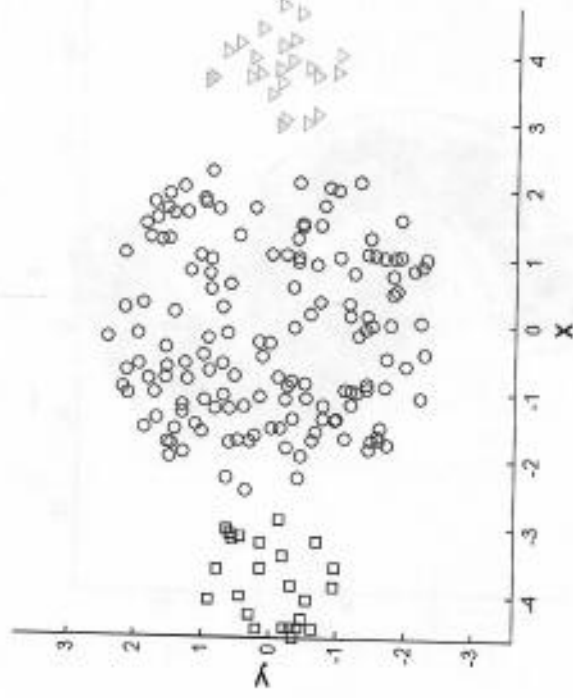


Original Points

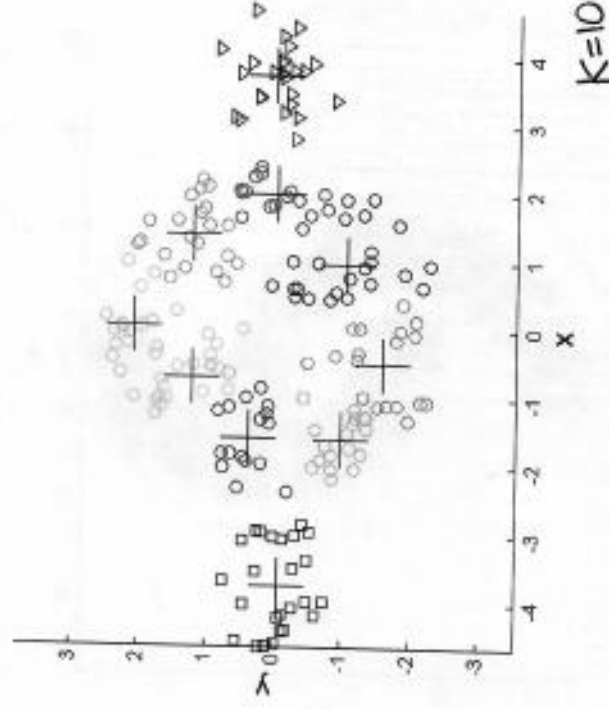


K-means (2 Clusters)

Overcoming K-means Limitations



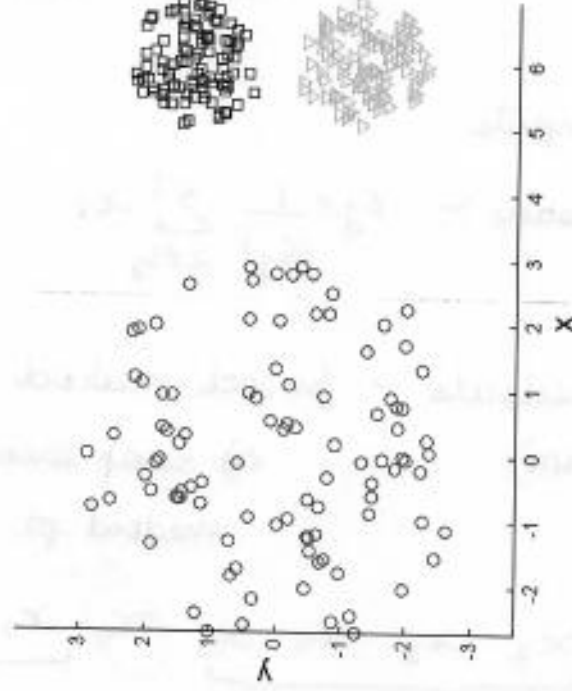
Original Points



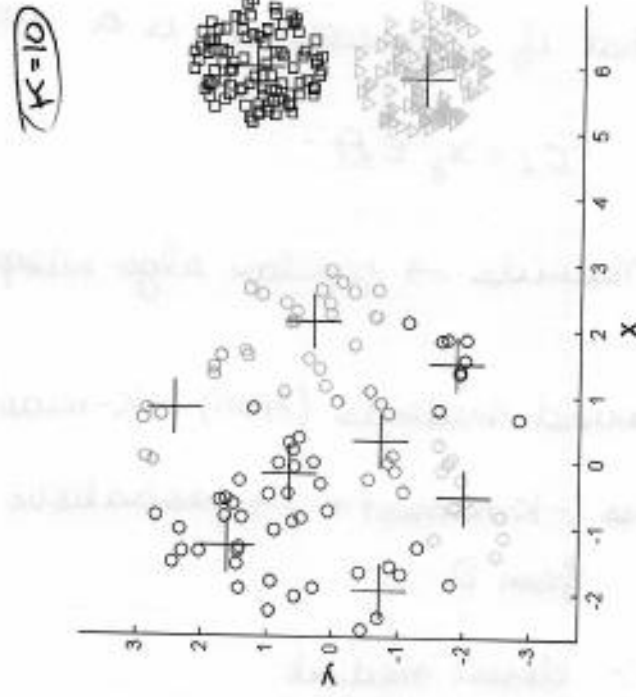
K-means Clusters

- One solution is to use many clusters.
Find parts of clusters, but need to put together.
We need to increase the value of K

Overcoming K-means Limitations



Original Points



K-means Clusters

problem K-medoids C_1, C_2, \dots, C_k \rightarrow centroid may not be interpretable

$$D = \{x_1, x_2, \dots, x_n\}$$

\rightarrow Big idea: what if each centroid is a datapoint in D .

$$C_i = x_j \in D$$

K-medoids \rightarrow popular algo interpret centroids

Partitioning around medoids (PAM): K-medoids

1) initialization :- K-means++ \rightarrow probabilistic methods pick K pts from D .

2) Assignment :- closest medoid

$$\begin{cases} x_i \in C_j & \text{if medoid } j \text{ is the closest} \\ & \text{medoid to } x_i \end{cases}$$

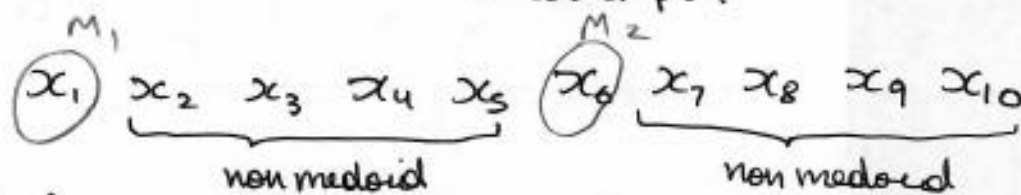
3) Update / Recompute :

$$\rightarrow \text{K-means :- } C_j = \frac{1}{|S_j|} \sum_{x \in S_j} x_i \quad \times$$

\rightarrow K-medoids :- for each medoid

(PAM)

a) swap each medoid with a non-medoid pt.



b) if loss decreases keep the swap else undo the ~~loss~~ swap

loss in K-means
↓
min

$$\sum_{i=1}^K \sum_{x \in S_i} \|x - m_j\|^2$$

↑
medoid_j

Determining the right 'K'

↑ hyperparameter

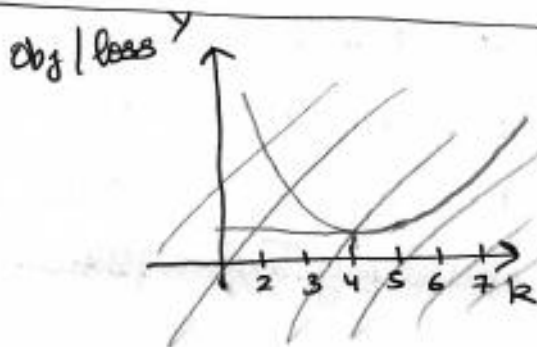
① domain - knowledge: Food reviews

↳ +ve & -ve

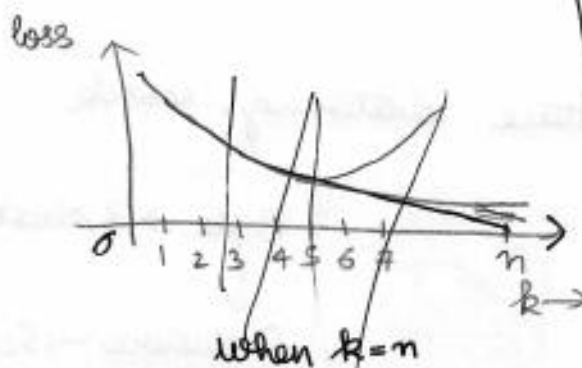
K=2

• We use elbow method or knee method

loss / Objective function: $\sum_{i=1}^K \sum_{x \in S_i} \|x - c_i\|^2 \rightarrow \text{minimize}$



best $K=4$
minimum values of obj is $K=4$



NOTE

$K = n = \# \text{ point}$

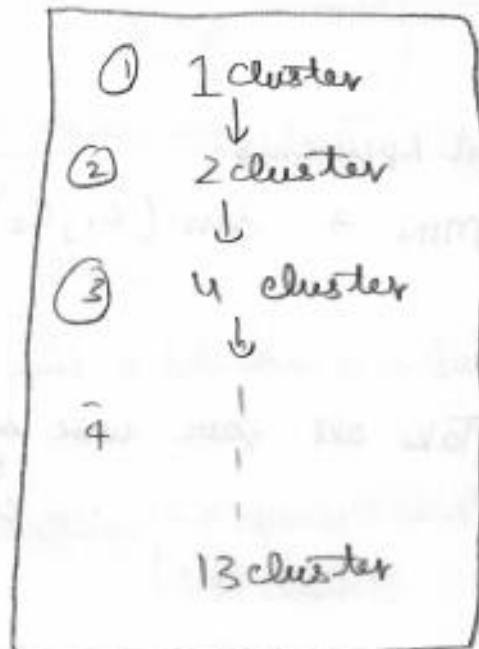
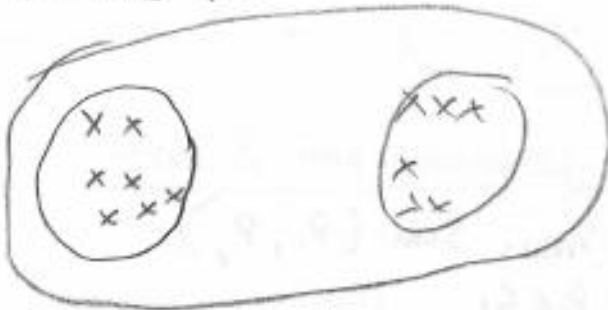
↳ 1000 cluster

↓

every point is a cluster

Agglomerative cluster → combine clusters close to each other.

2) Divisive :-



Agglomerative is more used in comparison to divisive

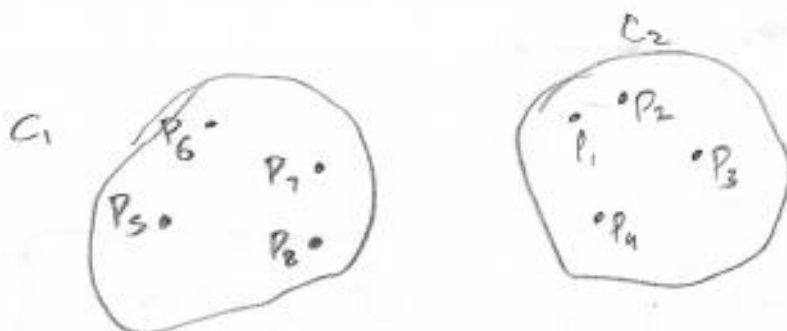
Key ingredient → (similarity or distance between clusters)

NOTE - 1 It is called Hierarchical because it forms trees which group datapoints into clusters. This tree is called dendrogram.

Proximity methods:

Advantages and Limitations

☆ How to Define Inter Cluster Similarity



Different Approaches:

1) MIN $\rightarrow \text{Sim}(C_1, C_2) = \min_{\substack{P_i \in C_1 \\ P_j \in C_2}} \text{Sim}(P_i, P_j)$

Take all pair wise similarity and pair wise similarity of cluster.

2) MAX \rightarrow Farthest or dissimilar point

MAX: $\text{Sim}(C_1, C_2) = \max_{\substack{P_i \in C_1 \\ P_j \in C_2}} \text{Sim}(P_i, P_j)$

3) Group Average \Rightarrow

AVG: $\text{Sim}(C_1, C_2) = \frac{\sum_{\substack{P_i \in C_1 \\ P_j \in C_2}} \text{Sim}(P_i, P_j)}{|C_1| * |C_2|}$

Size of C_1 \leftarrow $|C_1| * |C_2|$ \rightarrow Size of C_2

NOTE

Kernel Trick can be applied for all 3 approaches
 $\Rightarrow \text{Ker}(\text{Sim}(P_i, P_j))$

4) Distance between centroids:

Define centroid and compute distance

Space & Time complexity for hierarchical clustering

Space: $O(n^2)$ \longrightarrow Sim matrix

n # data points \longrightarrow a lot when n is large

Time: $O(n^3)$:- atmost n iteration \longrightarrow group 2 cluster
 \downarrow
 \hookrightarrow 1 cluster.
 update similarity matrix $O(n^2)$

Drawback: not very useful with large data points.

Limitations of Hierarchical clustering

① No objective (math) function that we are directly solving :- (K means - clear math obj)

\hookrightarrow next algo soln.

② { MIN \rightarrow problem with outliers
 MAX \rightarrow breaks large clusters can't accommodate different sized clusters
 \hookrightarrow group average have problem too as it is combination of min and max

3.) MOST IMPORTANT LIMITATION:-

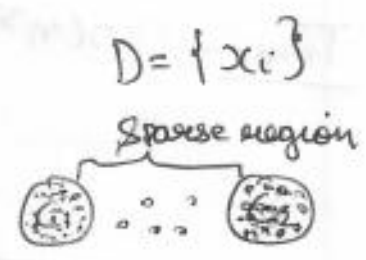
Space & time Complexity
 \downarrow
 $O(n^2)$ $\rightarrow O(n^2 \lg n)$
 $O(n^3)$

K-means have nice complexity compared to it.

DBSCAN: Density Based clustering Technique

- Centroid-based: - K-means
- Hierarchical based: - Agglomerative

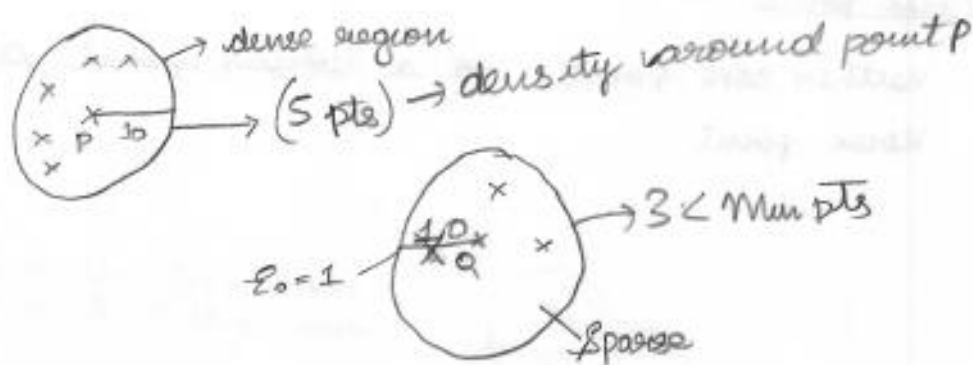
density - region \rightarrow clusters
 Sparse regions \rightarrow noise



Measuring density: Min Pts, EPS
 \rightarrow Hyperparameter of DBSCAN

① density at a P: # pts within a hypersphere of radius ϵ_0 around P.
 $\epsilon_0 = 1.0$
 Circle 2D
 Sphere 3D

② dense region: a hypersphere/circle of radius ϵ_0 that contains at least Min Pts points
 \downarrow
 4




If number of points in radius is greater than threshold then it is dense region

★ Core, Border and Noise points

$$D = \{x_i\} \quad \text{Min Pts} \cdot E_0$$

① Core point (P): if P has \geq Min Pts in an E_0 radius around it.

 Every core pt. \in belongs to dense region.

② Border point (P)

① P is not a core point & has $<$ Min pt points E_0 radius

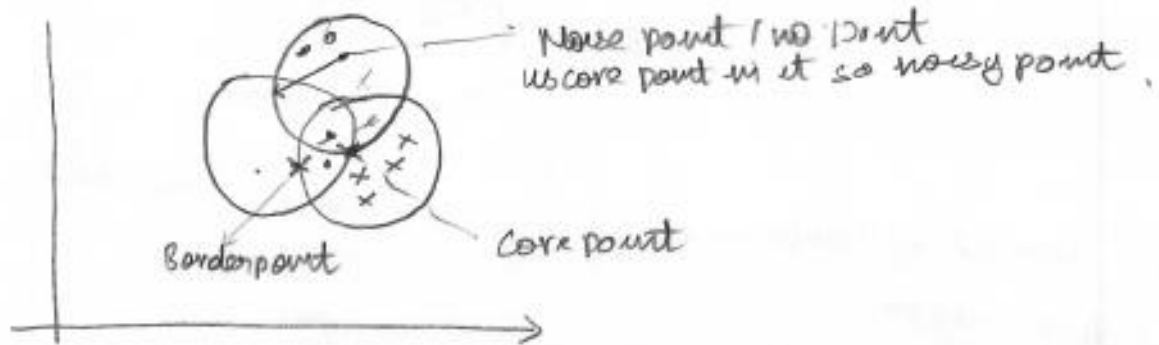
② $P \in \text{Neighbourhood}(Q)$

Q : core point

$$\text{dist}(P, Q) \in E_0$$

③ Noise point :-

neither core point nor a border point are called noise point.



☆ Density edge & density connected points

① density edge → connection

↳ P, Q : core point

② density connected points : $P \& Q$: core points

