

K-Means Clustering

Presented by ~

Abhishek Dukkar

(av20ms091@iitertkol.ac.in)

What is Clustering?

Clustering is an unsupervised learning technique that groups data points based on their inherent similarities, aiming to discover underlying patterns and structures within the data.

But What is Unsupervised Learning?

Types of Machine Learning:- 3 main types:-

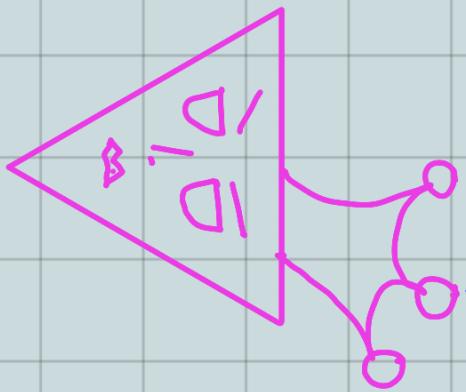
i) Supervised learning:- model is trained on labeled data, meaning each training example is paired with an output label. The model learns to map input to the correct outputs. Ex:- Classification

ii) Unsupervised Learning:- model deals with unlabeled data. Model tries to identify patterns or groupings within the data without any explicit labels. Ex:- Clustering

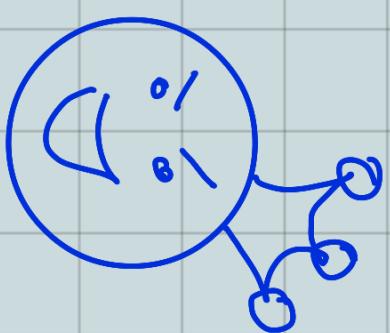
iii) Reinforcement Learning:- Agent learns to make decisions by taking actions in an environment to maximize a cumulative reward. The agent receives feedback in the form of rewards or penalties based on its actions. Ex:- Navigate a maze

What in natural grouping among them?

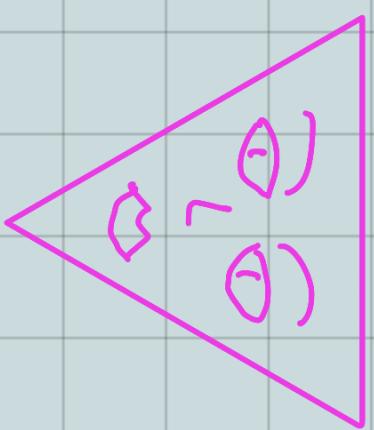
Queen



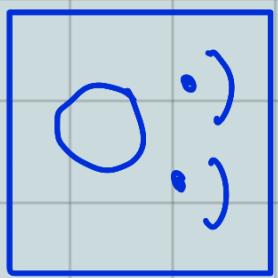
Prince



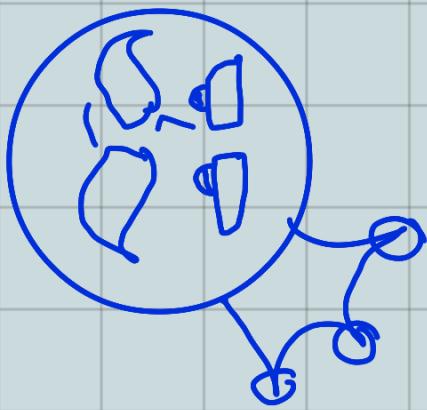
Woman



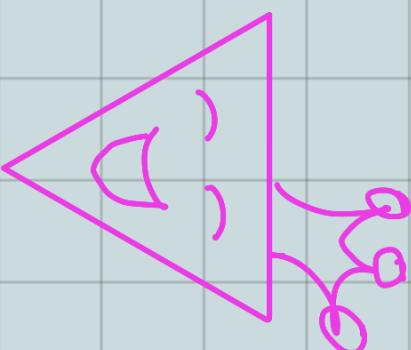
Boy



King



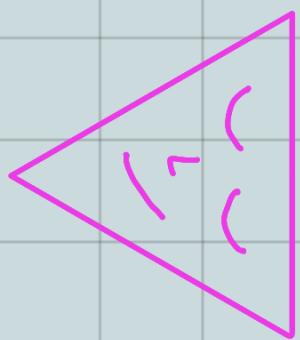
Princess



Man



Girl



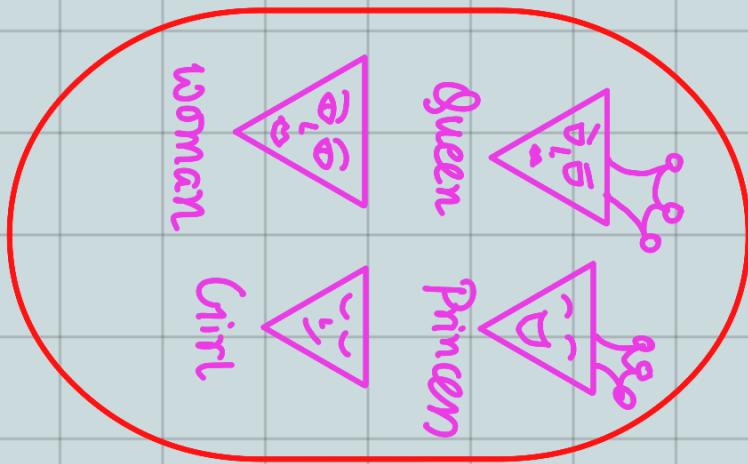
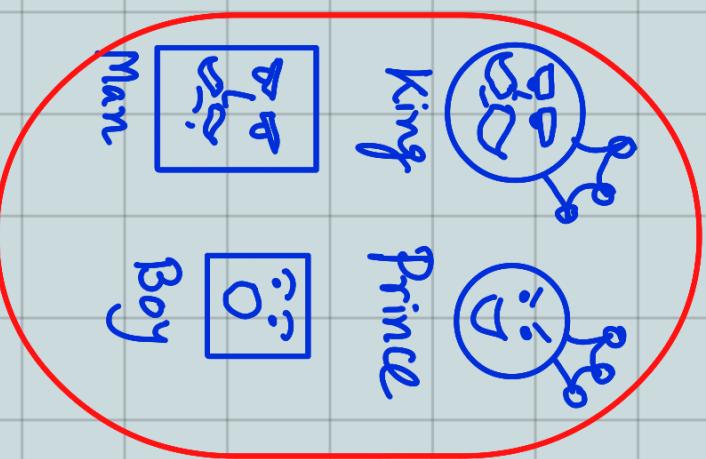
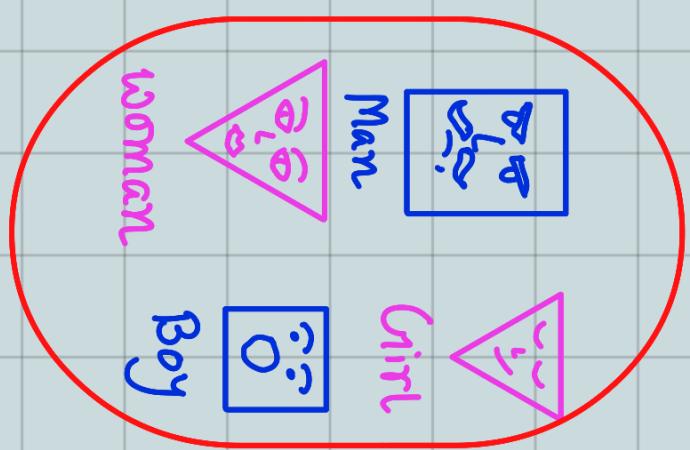
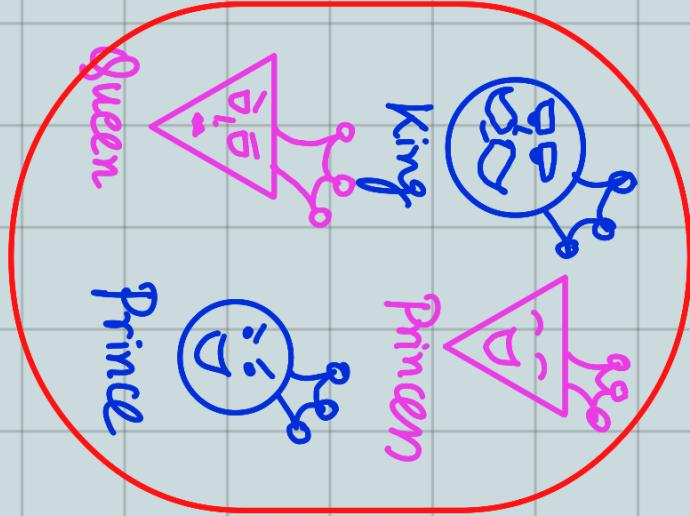
Clustering in subjective!

Royalty

Gender

Based on Royalty

Based on Gender



Aim of Clustering!

- high intra-class similarity
- low inter-class similarity

But what does determine the similarity?

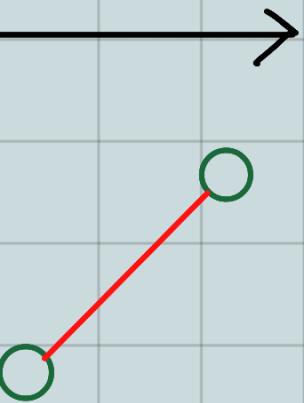
The answer in Distance Measure.

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

Now we will discuss some distance measure.



1. The Euclidean distance; (also called 2-norm distance)



measures the shortest distance between two real valued vectors.

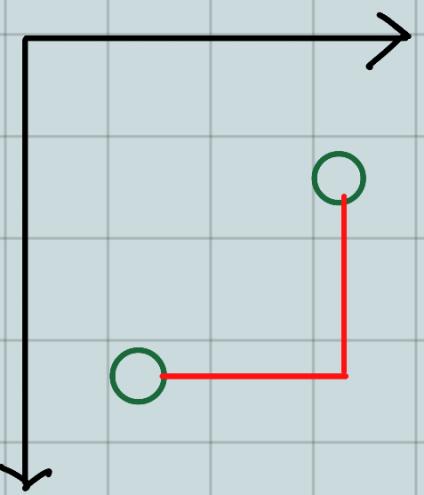
$$d = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

Code in
from scipy.spatial import distance
distance.euclidean(vector_1, vector_2)

- Does not work well on higher dimensional data than 2D & 3D

- If we don't normalize the features, the distance might be skewed due to different units.

2. Manhattan distance := (also called taxiCab norm or 1-norm)



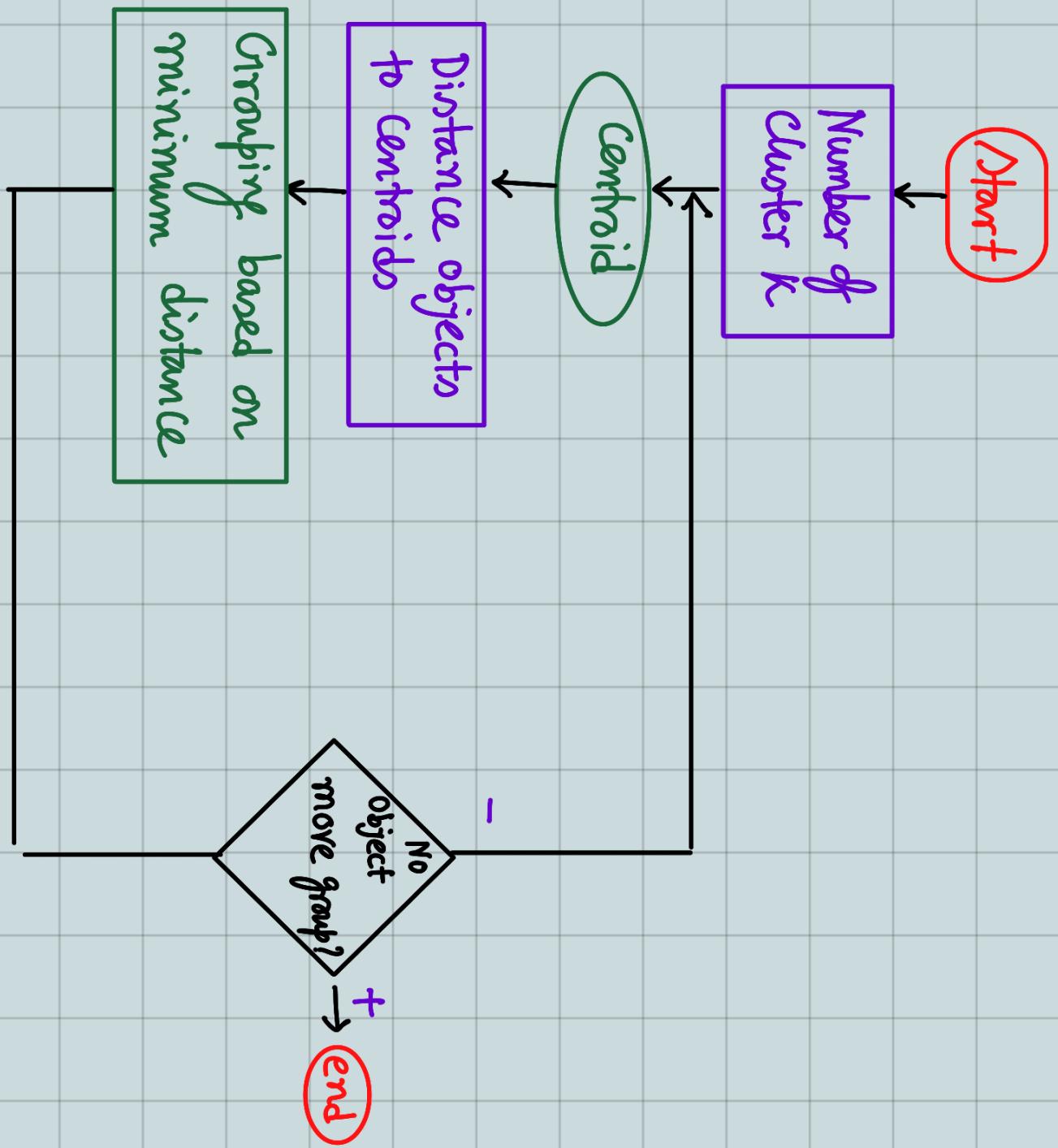
The distance betⁿ two real valued vectors is calculated as if one could only move at right angles.

Code:
from scipy.spatial import distance
distance.cityblock(vector_1, vector_2)

- It does not show the shortest path possible.
- less intuitive than Euclidean distance in high dimensional space.

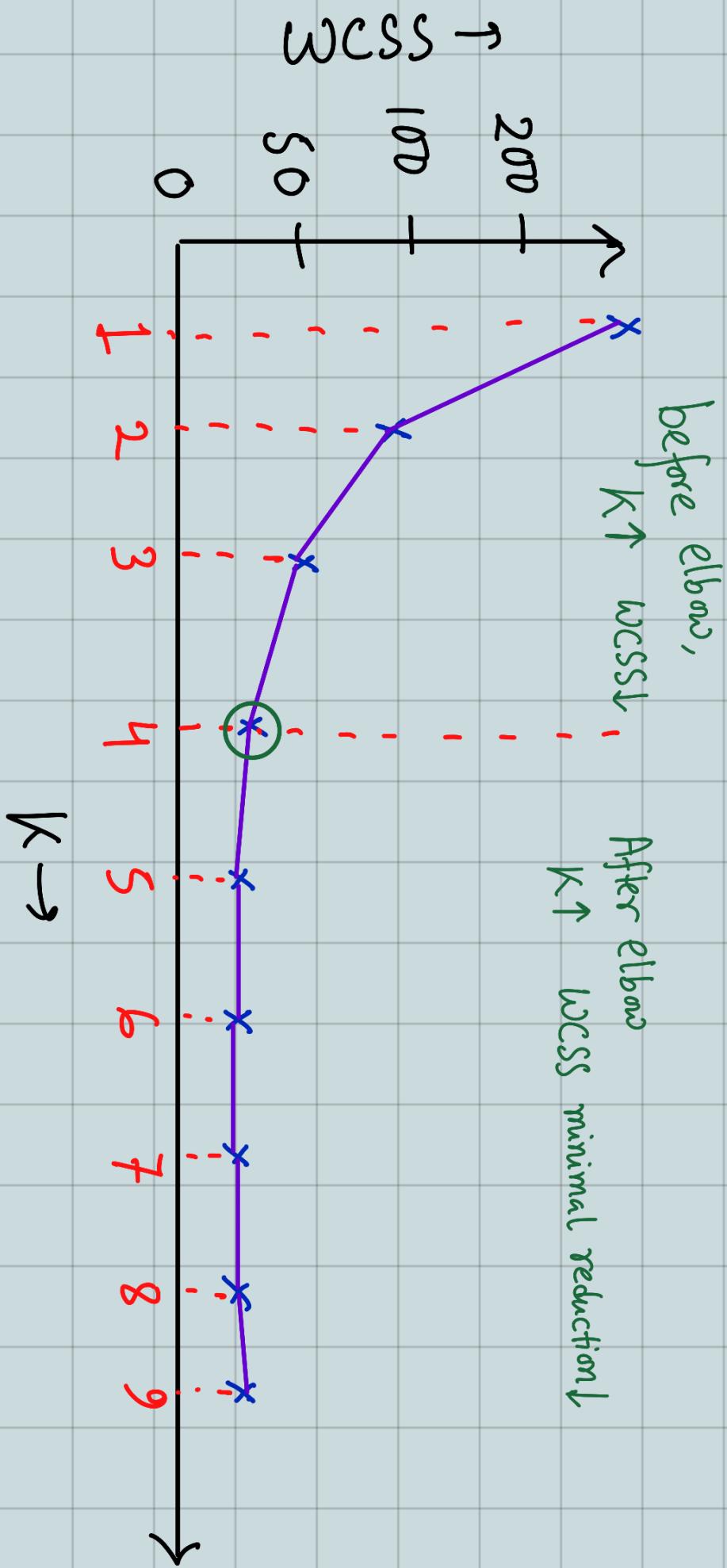
- K-means Clustering: The K-means clustering is an algorithm to cluster n objects based on attributes into K partitions where $K < n$
- K is a positive integer number
 - The grouping is done by minimizing the sum of squares of distances between data and the corresponding cluster centroid.

But how does it work?

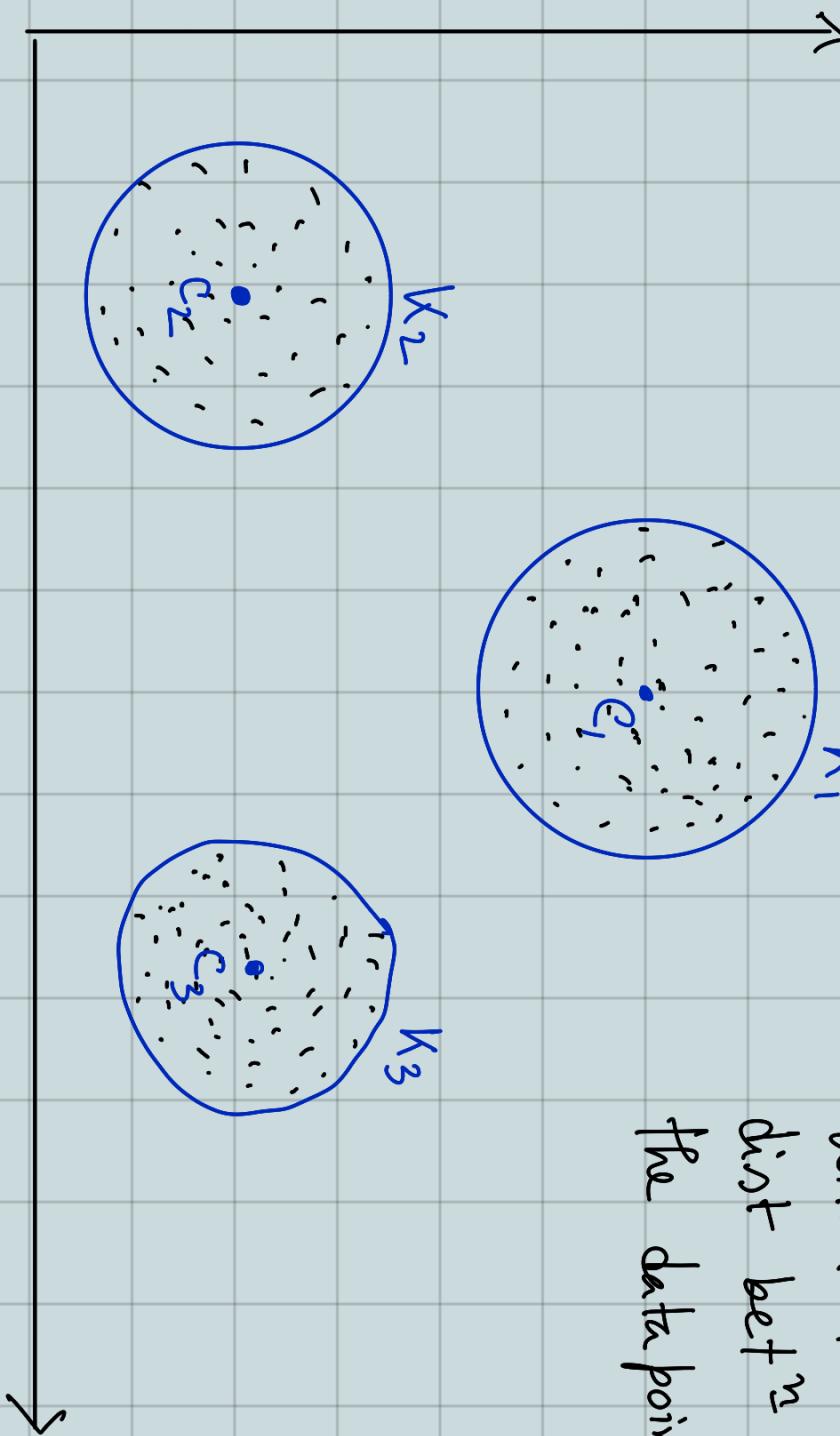


Elbow Method

Elbow Method - Choosing optimal number of clusters in a crucial step. Since we don't have predefined cluster counts in unsupervised learning, we need a systematic approach to determine the k -value.



WCSS (Within Cluster Sum of Squares)



$$WCSS = \sum_{k=1}^{C_m} \left(\sum_{i \in C_k} \text{distance}(d_i, c_k)^2 \right)$$

Within the cluster the dist betw' centroids and the data points.

- 
 - We calculate a distance measure called WCSS (Within-Cluster-Sum of Squares). This tells us how spread out the data points are.
 - We try different K -values and run KMeans and calculate the WCSS.
 - We plot a graph with K on x axis and WCSS on y axis.
 - We increase K , the WCSS typically decreases because we are creating more clusters, which tend to capture more data variations. However, there comes a point where adding more clusters results in only a marginal decrease in WCSS. This is where we observe "Elbow" shape in the graph.

 abhiseksarkar2001

CS2201_Spring2025

Introduction to Computation: Second Year Python Course of IISER Kolkata

☆ 0 stars ⚡ 0 forks

☆ STAR



Issues

0

Pull Requests

0

Actions

▼

More

Current branch
main

Code

Visit The GitHub Repo for all the Slides, Notes and Codes for this Course

https://github.com/abhiseksarkar2001/CS2201_Spring2025