

Week 01

01] What is Clustering ?

- *) Clustering is an unsupervised learning technique used to group similar data points together based on some measure of similarity (usually distance)
- *) The goal is to discover patterns or structures in data without labels
- *) Each group is called a cluster, and points within the same cluster are more similar to each other than to points in other clusters

ex :- market segmentation, Grouping similar news

02] k - means Algorithm

- *) most simplest and popular clustering algorithms.
- *) It tries to partition the data into k clusters where each data point belongs to the cluster with the nearest centroid (mean points in that cluster)

K - means algorithm Steps

- 1) Choose the number of clusters (k)
 - Decide how many groups you want the algorithm to find.
- 2) Initialize centroids ($\mu_1, \mu_2, \dots, \mu_k$)
 - Randomly select k points (or use a smart method like k-Means fit) as initial Cluster Centers
- 3) Then assign points to the nearest centroid;
 - for $i = 1$ to m (training examples)
 - $c^{(i)} :=$ Index (from 1 to k) of cluster centroid closest to $x^{(i)}$

$$\min \|x^{(i)} - \mu_k\|^2$$

$(x^1, c^1 = ①) / x^{12}, c^{12} = ② / x^5, c^5 = ③)$
- 4) Update the centroids (move centroids)
 - for $k = 1$ to k
 - $\mu_k :=$ Average (mean) of points assigned to cluster k
- 5) Repeat 3, 4 until Convergence
 - When cluster assignments no longer change or the centroid position stop moving significantly

* If we found a cluster which has no training samples, we simply remove that cluster C_{k-1}

03] Optimization Objective (Cost Function)

*) K-means aims to minimize the sum of squared distance between each point and its cluster centroid

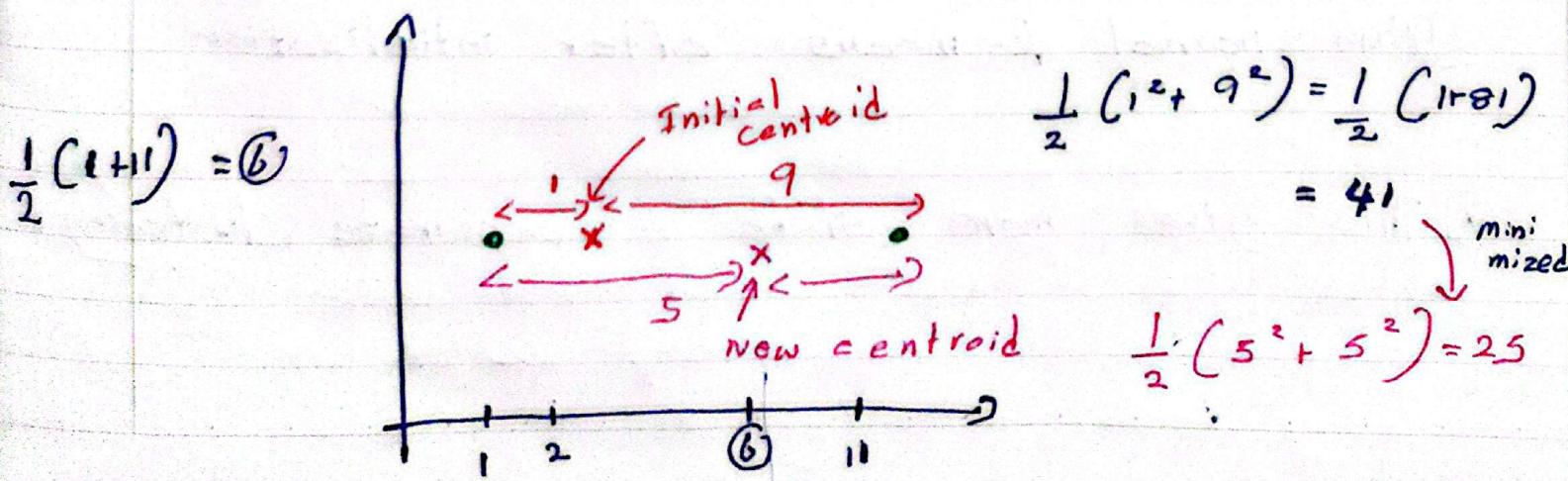
$$J(C^{(1)}, \dots, C^{(m)}, \mu_1, \dots, \mu_k) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{C(i)}\|^2$$

m = number of data points

$C^{(i)}$ = index of the cluster assigned to point i

$\mu_{C(i)}$ = Centroid of the cluster point i belongs to

*) The goal of k means is to minimize $J(C, \mu)$



05) Initializing K-Means

Random Initialization

- *]) Pick k random points from the dataset as initial Centroids
- *]) May lead to different results each time because k-means can converge to local minima

K-Means ++ Initialization

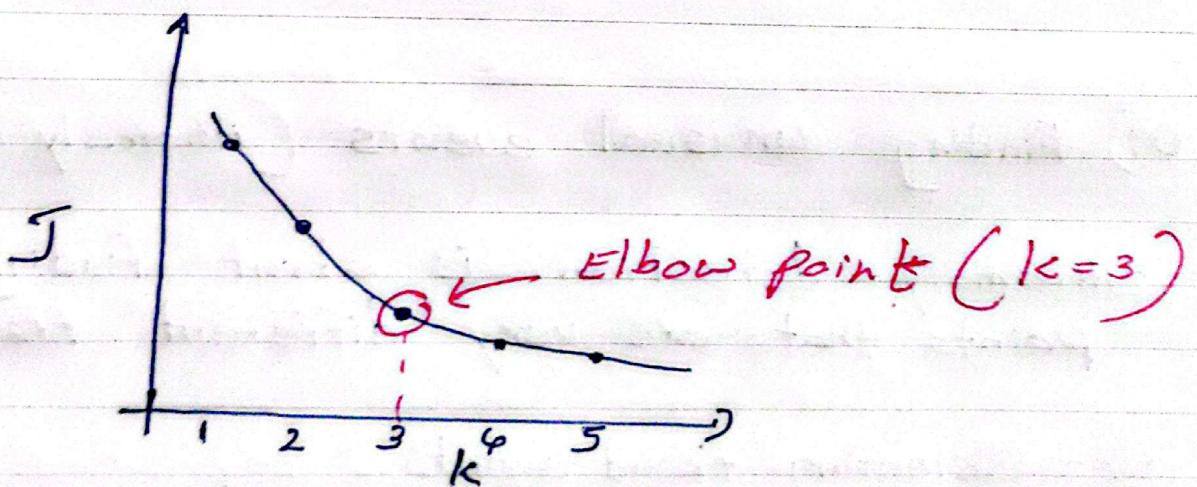
- *) A smarter method that spreads out the initial centroids to improve performance
 - 1] Pick the first centroid randomly
 - 2] For each remaining centroid, pick a point with probability proportional to its squared distance from the nearest existing centroid.
 - 3] Run normal k-means after initialization
- *]) This gives more stable & accurate clustering.

Q6] Chosing the number of Clusters (k)

- * Selecting the right k is crucial and often depends on experimentation.

Elbow method

- * Plot the cost J (distortion) vs number of clusters k .
- * As k increases, cost decreases (clusters fit data better)
- * Choose k at the "elbow" point, where improvement slows down.



Silhouette Score

- * Measures how similar a point is to its own cluster compared to others.

* Ranges -1 to +1

$\Rightarrow +1 =$ well clustered

$-1 =$ misclassified

$0 =$ on the boundary

* Choose k that gives highest silhouette score.

Gap statistic

* Compares the total within-cluster variation for different k values with reference random data

* Choose k where the gap is largest.

Q) Finding unusual events (Anomaly Detection)

* Anomaly detection is about finding data points that are very different from the majority

Ex:- Detecting fraud events

Detecting abnormal server behaviour

Detecting defective products in manufacturing

* we model what 'normal' looks like, Then if a new example has a very low probability of being normal, we flag it "anomaly"

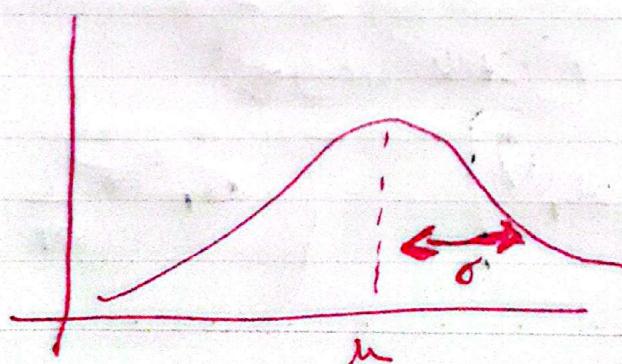
08] Gaussian Distribution (Normal Distribution)

$$P(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- μ = mean
- σ^2 = Variance
- x : data point
- $p(x)$: probability Density

- *) Most data points are near the mean (typical behaviour)
- *) Values far from the mean are rare (potential anomalies)
- *) If we have multiple features, we assume they are independent and compute

$$P(x) = \prod_{j=1}^n P(x_j; \mu_j, \sigma_j^2)$$



09) Anomaly Detection Steps

(Anomaly detection system for a server)

- * Each example (data point) describes the server's behaviour at a given time using features like
 - 1) CPU usage
 - 2) memory usage

1] Collect Training Data

- * we collect the many normal examples (times when the server was working fine)
- * we will learn what "normal" CPU and Memory values look like

2) Estimate Parameters

- * we find the mean and variance for each feature (CPU & Memory)

\Rightarrow For Feature 1 (CPU usage)

$$\mu_1 = \frac{1}{m} \sum_{i=1}^m x_1^{(i)}$$

$$\sigma_1^2 = \frac{1}{m} \sum_{i=1}^m (x_1^{(i)} - \mu_1)^2$$

\Rightarrow For feature 2 (Memory Usage)

$$\mu_2 = \frac{1}{m} \sum_{i=1}^m x_2^{(i)}$$

$$\sigma_2^2 = \frac{1}{m} \sum_{i=1}^m (x_2^{(i)} - \mu_2)^2$$

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

i = example number (row in the dataset)

j = feature number (Column in the dataset)

* ex: $x_2^{(5)}$ = Value of feature 2 (memory) for example #5

* We compute μ and σ^2 for each column of data

3) Compute probability of Each example

* we assume each feature allows a Gaussian (normal) Distribution

For one feature :

$$P(x_j) = \frac{1}{\sqrt{2\pi\sigma_j^2}} e^{-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}}$$

If we have multiple features:

$$P(x^{(i)}) = P(x_1^{(i)}) \times P(x_2^{(i)})$$

- * This gives probability score for each example.

High $P(x)$ \Rightarrow looks normal
 Low $P(x)$ \Rightarrow looks strange (maybe an anomaly)

4) Decision phase

- * Now let's test a new data point:

CPU - 90t., Memory - 95t.

we compute : $P(x_{\text{test}}) = P(90; \mu_1, \sigma_1^2) \times P(95, \mu_2, \sigma_2^2)$

- * If this probability is very small,

$$P(x_{\text{test}}) < \epsilon$$

We classify its as a anomaly

5) Decision Rule

$$P(\alpha) \geq \epsilon \Rightarrow \text{Normal } (y=0)$$

$$P(\alpha) < \epsilon \Rightarrow \text{Anomaly } (y=1)$$

10] Developing and Evaluating an Anomaly Detection System

1] Choose a Dataset

We split our data into 3 parts

- 1] Training set → Mostly normal examples
 - Learn what normal behaviour looks like
- 2] Cross-validation → Normal + Some anomalies set
 - Tune (Chose) threshold (ϵ)
- 3] Test set → Mix of normal + anomalies
 - Final evaluation of system performance

Ex: 10,000 Normal / 50 Anomaly

Training - 8000 normal

CV set - 1000 normal + 25 anomalies

Test - 1000 normal + 25 anomalies

2) Fit the model

* Now we use the training set (normal data) to model normal behaviour

For each feature j ,

Compute mean (μ_j)

Compute variance (σ_j^2)

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

* Then we calculate, for any new data point x , the probability that it came from this distribution

$$P(x) = \prod_{j=1}^n P(x_j; \mu_j, \sigma_j^2)$$

* This gives each point "normality score"

* Low $P(x) \rightarrow$ unusual \rightarrow maybe an anomaly

3) Choose Threshold ϵ

- * We use CV set for this
- * We try different possible thresholds $\epsilon \in [0.1, 0.001, 0.01, 0.0001, \text{etc.}]$
- * For each threshold,
 - \Rightarrow Predict anomaly if $(p(x)) < \epsilon$
 - \Rightarrow Predict normal if $(p(x)) \geq \epsilon$
- * Then we check how many predictions are correct
(Semi-Supervised learning)

4) Evaluate using metrics

- * Because anomalies are rare, normal "accuracy" isn't enough (99% normal data, still get 99%)
- * Instead we use precision, recall, F1 score
- * Usually we calculate the F1 Score and select the ϵ with the highest F1

5) Test the system

- * Finally, we test the model on the test set using the chosen ϵ
- * Then gets a fair estimate (by, f1, recall, precision)

1) Anomaly Detection vs Supervised learning

A . D

S. L

- *) Very small no. of positive examples ($y=1$). ($0-20$)
Large number of negative examples ($y=0$)

- *) Many different 'types' of anomalies. Hard for any algorithm to learn from positive examples what the anomalies look like:

future anomalies may look nothing like any of the anomalous examples we've seen so far

Fraud Detection

- *) Large no. of positive and negative examples

- *) Enough positive examples for algorithm to get a sense of what positive examples are like, future positive examples likely to be similar to ones in training set

- *) Email spam detection

1²] Choosing what features to use

- *) The success of anomaly detection heavily depends on feature selection
- *) Normalize features to similar scales
(z-score normalization)
- *) Use domain knowledge - features that correlate with failures or anomalies
- *) Sometimes apply transformations (log, ratio)
- *) Remove irrelevant / noisy features