

Department of Computer Science and Engineering



Semester: VIII

Academic Year: 2024-25

(-MEANS CLUSTERING K-means clustering is an unsupervised machine learning algorithm used to group similar data points into k clusters.

Subject: DAV AIFB

Role of K-means ductering in financial Data Analysis (1) Customer Segmentation:

- * Banks and financial institutions segment customers spending habits, credit scores and income levels.
- * Helps in targeted marketing, personalized loan offers and credit nick asserment.

A bank can group customers into high-income investors, Example:middle class salaried individuals, and low-income customers for customized financial services.

(a) Credit Risk Analysis:

- * Identifies groups of borrowers with risk profiles.
- * Segmenti clients based on loan repayment history, credit scores, and transaction behaviours.
- * Helps bank decide interest rates, loan eligibility and credit card limits

Example: Grouping customers into Low-Risk -> Regular payments, high credit score. Medium-risk -> Occasional late payments High-Risk -> Frequent defaults, poor credit scores.

Subject Incharge: Prof. Sarala Mary



Department of Computer Science and Engineering



Semester: VIII

Subject: DAV AIFB

Academic Year: 2024-25

(3) Fraud Detection:

* K-means clusters normal vs. fraudulent transactions based on spending patterns.

* Anamolies (outliers) can indicate potential fraud, identily theft or suspicious activities.

Example:

If a customer usually spends \$500 per month, but suddenly a \$10,000 withdrawal happens from another Country, k-means can flag it as fraud.

(4) Stock Market and Portfolio Clustering:

Groups stocks based on volatility, returns and risk levels.

*Helps investors create diversified portfolio.

Example: Slocks clustered into High Risk (Startups), Medium Risk (Tech stocks), and Low Risk (Blue-chip stocks).

(5) Anomaly Detection in Trading:

* Detect unusual trading patterns that may insider trading or market manipulation.

* Helps regulating bodies (SEC, RBI etc). ensure market integrity



PARCHICANATO EMARCARIA TRUST S A.P. SHAH INSTITUTE OF TECHNOLOGY

Department of Computer Science and Engineering



Semester: Vin

Subject: DAY AIFB

Academic Year: 2024-25

Step 1: Initialization:

Randomly choose K centroids.

Stepa: Assign Points

for each data point, calculate the distance to each centroid and assign it to the closest controid.

Steps: Updale Controids

*For each cluster, sum the feature values of all points assigned to that cluster.

* Calculate the new centroid as the average of those points.

Step4: Repeat

Reassign points and update certroids. Continue urtil centroids

no longer change significantly

Problem Statemen

A financial institute wants to segment 5 customers based on their Annual Income (\$k) and Spending Score (out of 100) using K-means clustering (K=2)

Given datuset:

Caslomer	Annual Income (\$K)	Spending Score		
A	15	39		
В	45	81		
c	25	55		
D	60	95		
E	30	60		



Department of Computer Science and Engineering
Data Science



Semester: VIII

Subject: DAY AIFB

Academic Year: 2024-25

Solution .

Step1: Choose Initial Centroids (Randomly)
We choose two initial centroids randomly from the dalaset

Centroid 1 (Ci) = (15,39) -> Based on Customer A

Centroid 2 ((2) = (60,95) -> Based on Customer D

Stepa: Compute Euclidean distance:

The Eudidean distance between two points (x1, y1) and

(x2, y2) is:

d = J(x2-x1)2+ (y2-y1)2

Cluster 1 -> A (15,89), C(25,55), E(30,60).

Cluster 2 -> B (45,81), D(60,95)

Compute the distance for each customer to both

centroids:

Customer		(60,95) (2-Distance	Assigned Cluster
A(15,39)	0	: tromstel?	Cluster
B (45,81)	$\sqrt{(45-15)^2+(81-39)^2}=48.79$	1(45-60)+(81-95)=2102	clusters
C(25,55)	1(25-15)2+(55-39)2 = 18.87	1(25-60)2+(55-95)=50	clusteri
D(60,95)	$\sqrt{(60-16)^2+(95-39)^2}=73.63$	0	clustera
E(80,60)	V(30-15)2+(60-39)2 = 27.90	V(30-60)2+ (60-95) = 47!	clusteri
Steps: Con	first iteration, the new du	isters are:	9

Subject Incharge: Prof. Sarala Mary



Department of Computer Science and Engineering



								ú
Se	m	ė	51	P	۲	4		5

Subject:

DAV

Academic Year: 2024-25

Compute New centroids:

New controld (Ci):

$$\left(\frac{15+25+30}{3}, \frac{39+55+60}{3}\right) = (23.33,51.33)$$

New centroid ((2):

$$\left(\frac{45+60}{2}, \frac{81+95}{2}\right) = (52.5, 88)$$

Stepa: Recompute Distances and Reassign Clusters:

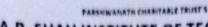
Recalculate distance using new centroids:

Recalcula		(52.5,88) Distance	New
A(15,29)	$\sqrt{(23.33-15)^2+(51.33-39)^2}$ =15.02	$\sqrt{(15-52.5)^2+(39-82)^2}$ = 62.42	Diges E
B(45,81)	$\sqrt{(45-23\cdot33)^2+(81-51\cdot33)^2}$ = 86.17	$\int (45-52.5)^2 + (81-88)^2$ = 10.60	Cluster 2.
C(25,56)	1(25-23.33)2+ (55-51.83)2 =3.99	$\sqrt{(a5-52.5)^2+(55-88)^2}$ = $4a.60$	Clustert
D(60,95)	(60-28-88) + (95-51-33) + = 57-01	10.60	Clustera
E(30,60)	\(\langle (80-28.83)^2 + (60-51.83)^2 \\ = 11.13	(80-52.5) ² +(60-88) ² = 35.15	· Cluster 1.

Since the clusters remain unchanged, the algorithm converges.

Subject Incharge: Prof. Sarala Mary

Department of CSE-Data Science | APSIT





Department of Computer Science and Engineering



Semester: Du

Subject: AIFB .

Academic Year: 2024-25

Final cluster:

Cluster 1: Low Income, Moderate Spending Customers

A (15,39)

C(25,55)

E (80,60)

Centroid: (23.33, 51.33)

Cluster 2: High Income, High - Spending Cuetomers

B(45,81)

D (60,95)

Certified: (58.5, 88)

SPARSITY AND CONNECTEDNESS OF UNDIRECTED GRAPH:

Graphs can be used to model various types of relationships, such as between assets, companies or individuals.

Understanding connectedness and sparsily of these graphe is crucial for analyzing network in finance, such as portfolio diresification, credit networks, marked interactions and fraud detection.

Francial Asset Network (Portfolio Diverification)

In portfolio diversification, we can model the relationship between different financial assets leg. stocks, bonds, commodifies) using a graph where:

* Vertices (nodes) represent different financial assets

* Edges represent relationships between assets, often based

Subject Incharge: Prof. Sarala Mary

Department of CSE-Data Science | APSIT