

# Clustering & K-Means Clustering

Theory and Implementation with Programming Example



# Understanding the Basics: What is Clustering?

Clustering is an **unsupervised learning technique** in data science, used to group data points that are similar to each other. Unlike supervised methods, clustering does not rely on pre-labeled data. Its primary goal is to discover inherent or patterns within a dataset without prior knowledge of those groups.

## Why Unsupervised?

It learns directly from the data's intrinsic structure, making it ideal for exploratory data analysis where labels are unavailable or too costly to obtain.

## Core Idea

The aim is to maximize similarity **within** each group (intra-cluster similarity) while minimizing similarity **between** different groups (inter-cluster dissimilarity).



# Why Clustering Matters

Clustering helps us make sense of complex data by revealing natural groupings. Consider these real-world applications:



**Customer Behavior** Imagine a retailer analyzing purchase history. Clustering can group customers into segments like "frequent high spenders," "discount shoppers," or "occasional visitors," enabling tailored promotions and loyalty programs.



**Product Categorization** An e-commerce platform can cluster products by price, features, or user reviews. This helps in organizing inventory, optimizing pricing, and providing more accurate recommendations.



# Real-World Applications of Clustering

Clustering is a versatile technique with widespread applications across various industries, providing valuable insights from complex datasets.



## Customer Segmentation

Group customers based on purchasing behavior, demographics, or engagement to tailor marketing strategies.



## Product Grouping

Categorize products for inventory management, recommendation systems, or optimal store layouts.



## Fraud Detection

Identify unusual patterns in transactions that deviate from typical behavior, flagging potential fraud.



## Image Compression

Reduce the number of distinct colors in an image while maintaining visual quality, using color clusters.

# Diverse Approaches: Types of Clustering

While K-Means is a popular choice, several other clustering algorithms exist, each suited for different data structures and objectives.

1

## **Partitioning Methods (e.g., K-Means)**

Divides data into non-overlapping subgroups, where each data point belongs to exactly one cluster.  
Requires specifying the number of clusters in advance.

2

## **Hierarchical Methods**

Builds a hierarchy of clusters, either by starting with individual points and merging (agglomerative) or starting with one large cluster and splitting (divisive).

3

## **Density-Based Methods (e.g., DBSCAN)**

Discovers clusters of arbitrary shape based on data point density, distinguishing noisy outlier points from core clusters.



# K-Means Clustering: The Algorithm Explained

K-Means is an iterative algorithm that aims to partition 'n' observations into 'k' clusters, where each observation belongs to the cluster with the nearest mean (centroid).

The process repeats, refining cluster assignments and centroid positions until the clusters are stable and no data point significantly changes its cluster affiliation.

# Optimizing K: The Elbow Method

Choosing the optimal number of clusters ('k') is critical. The Elbow Method helps identify a suitable 'k' by examining the "inertia" or "within-cluster sum of squares" (WCSS).





# K-Means: Advantages & Limitations

Understanding the strengths and weaknesses of K-Means helps in deciding when and how to apply it effectively.

## Advantages

- **Simplicity:** Easy to understand and implement.
- **Scalability:** Efficient for large datasets with many observations.
- **Speed:** Relatively fast convergence, especially when data is well-separated.
- **Interpretability:** Clusters are defined by their centroids, offering clear interpretation.

## Limitations

- **Sensitive to Initialization:** Results can vary based on initial centroid placement.
- **Assumes Spherical Clusters:** Struggles with irregularly shaped or non-convex clusters.
- **Requires Predefined K:** The number of clusters must be known in advance, which is often not the case.
- **Sensitive to Outliers:** Outliers can significantly skew cluster centroids.



# Programming K-Means: My Implementation

Here's a glimpse into the code and libraries used to perform K-Means clustering on a dataset, showcasing my practical application.

```
# K-Means Clustering in Python
import numpy as np
import pandas as pd
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt

# Sample data (e.g., Mall customer spending dataset)
data = pd.read_excel("/content/customer_data_with_clusters_and_charts.xlsx")

# Drop non-numeric columns that are not needed for clustering
data_for_clustering = data.drop(['Customer ID', 'Gender', 'Cluster'], axis=1)

display(data.head())

# Fit K-Means
kmeans = KMeans(n_clusters=8, random_state=50)
data['Cluster'] = kmeans.fit_predict(data_for_clustering)

# Visualize clusters
plt.scatter(data['Income (₹)'], data['Spending Score'], c=data['Cluster'])
plt.xlabel('Income')
plt.ylabel('Spending Score')
plt.title('K-Means Clustering')
plt.show()
```

## Key Libraries Used

- **Scikit-learn (sklearn):** For the KMeans algorithm and related tools like StandardScaler.
- **Pandas:** For data manipulation and loading (e.g., CSV files).
- **Matplotlib:** For visualizing the clusters and the Elbow Method graph.
- **NumPy:** For numerical operations.

# Data Interpretation

## ◆ Data Source:

- Secondary dataset of 100 customers with Annual Income & Spending Score.
- Commonly used in clustering demonstrations .

## ◆ What is Spending Score?

- A numerical value (0–100) reflecting customer spending behavior.
- Higher score = high-value customers; Lower score = low-value customers.

## ◆ How is it calculated?

- Based on past purchase history, frequency, and amount spent.
- In this dataset, pre-assigned as a scaled score (0–100).
- Used with income to form clusters.( For example: **90–100** → Very frequent/high-value spender; **50** → Moderate/average spend; **10–20** → Rare or low spender.)

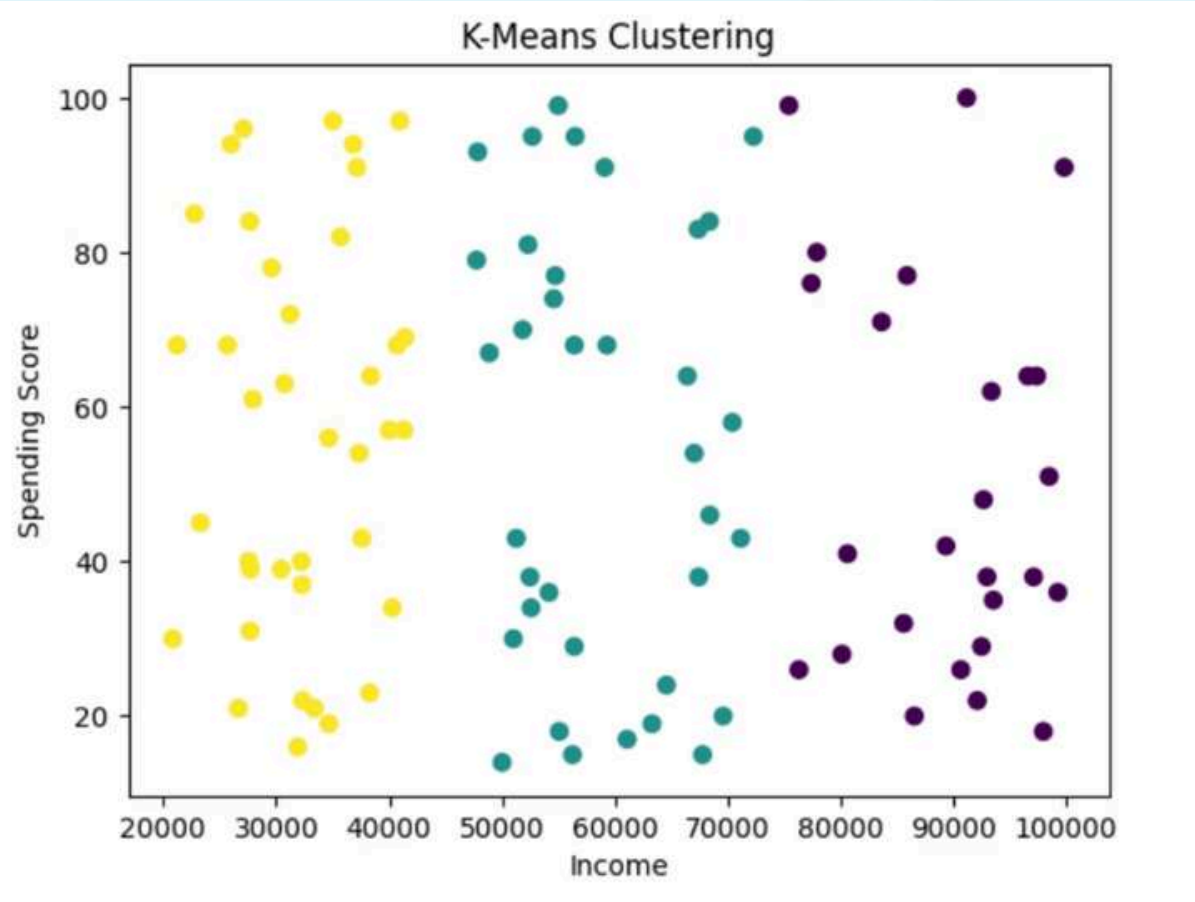


# K-Means Clustering: Interpreting Results

The output of K-Means provides grouped data points and cluster centroids, offering insights into underlying data structures.

## Visualizing Clusters

The plot on the left visually represents the formed clusters. Each color signifies a distinct cluster, and the large markers indicate their respective centroids.



We observed **3 distinct clusters** from our dataset, chosen using the Elbow Method.

## Understanding Centroids

The table on the right displays the final coordinates of each cluster's centroid. These values represent the characteristics of the data points within that cluster.

We observed **3 distinct clusters** from our dataset, chosen using the Elbow Method.

Customer ID	Age	Gender	Income (₹)	Spending Score	Cluster
C1	58	M	23278	45	Cluster 3: Middle Income, Moderate Spending
C2	33	M	38289	23	Cluster 3: Middle Income, Moderate Spending
C3	52	M	97397	64	Cluster 3: Middle Income, Moderate Spending
C4	20	M	32280	37	Cluster 3: Middle Income, Moderate Spending
C5	32	M	93563	35	Cluster 1: High Income, Low Spending
C6	59	F	48893	67	Cluster 3: Middle Income, Moderate Spending
C7	55	F	20851	30	Cluster 3: Middle Income, Moderate Spending
C8	45	F	56421	29	Cluster 3: Middle Income, Moderate Spending
C9	31	F	33396	21	Cluster 3: Middle Income, Moderate Spending
C10	42	M	67052	54	Cluster 3: Middle Income, Moderate Spending
C11	56	F	25695	68	Cluster 2: Low Income, High Spending
C12	52	M	69615	20	Cluster 1: High Income, Low Spending
C13	53	F	67400	83	Cluster 3: Middle Income, Moderate Spending
C14	30	M	26006	94	Cluster 2: Low Income, High Spending
C15	32	F	30458	39	Cluster 3: Middle Income, Moderate Spending
C16	24	F	56434	68	Cluster 3: Middle Income, Moderate Spending
C17	58	F	41319	57	Cluster 3: Middle Income, Moderate Spending
C18	40	M	54993	99	Cluster 3: Middle Income, Moderate Spending
C19	59	M	99840	91	Cluster 3: Middle Income, Moderate Spending
C20	28	M	41417	69	Cluster 3: Middle Income, Moderate Spending
C21	42	F	93000	38	Cluster 1: High Income, Low Spending
C22	38	M	50021	14	Cluster 3: Middle Income, Moderate Spending
C23	38	F	55093	18	Cluster 3: Middle Income, Moderate Spending
C24	31	F	47869	93	Cluster 3: Middle Income, Moderate Spending
C25	49	F	80142	28	Cluster 1: High Income, Low Spending
C26	34	M	52325	81	Cluster 3: Middle Income, Moderate Spending
C27	52	F	96622	64	Cluster 3: Middle Income, Moderate Spending
C28	55	F	67447	38	Cluster 1: High Income, Low Spending
C29	26	F	31915	16	Cluster 3: Middle Income, Moderate Spending
C30	25	M	40969	97	Cluster 3: Middle Income, Moderate Spending
C31	45	M	70432	58	Cluster 3: Middle Income, Moderate Spending
C32	56	F	89352	42	Cluster 1: High Income, Low Spending
C33	53	M	35014	97	Cluster 2: Low Income, High Spending
C34	52	F	64587	24	Cluster 1: High Income, Low Spending
C35	36	F	40730	68	Cluster 3: Middle Income, Moderate Spending
C36	18	F	85612	32	Cluster 1: High Income, Low Spending

- **Cluster 1** (yellow): Customers with **low to medium income and a wide range of spending habits**. They form the largest group.
- **Cluster 2** (green): Customers with **moderate to high income and higher spending scores**. These could be premium or valuable customers for businesses.
- **Cluster 3** (purple): Customers with **high income but relatively lower spending scores**. These are potential customers that can be targeted with promotions or strategies to increase spending.

# Key Takeaways

1

## Clustering Unveiled

Clustering is an unsupervised learning technique to discover hidden groupings and patterns in data, fundamental for insights.

2

## K-Means Explained

K-Means is a popular, iterative algorithm that groups data points into K clusters based on proximity to centroids.

3

## Practical Impact

Used across industries for tasks like customer segmentation, their spending habits and so on.

4

## Dive Deeper

Explore other clustering algorithms (DBSCAN, Hierarchical), evaluate different metrics, and practice with real datasets.