

Clustering with k-Means: Understanding the Role of Centroids and Distance Metrics

1. Introduction

What is k-Means Clustering?

Clustering is an essential technique in unsupervised machine learning, used to group data points into meaningful clusters based on shared characteristics or patterns. Among clustering methods, **k-Means** is one of the most widely used algorithms due to its simplicity and effectiveness. It partitions a dataset into k clusters by minimizing the variance within each cluster, making it highly efficient for exploratory data analysis and pattern recognition.

Why Use k-Means?

k-Means is particularly useful for:

1. **Grouping Data:** Identifies inherent structures in data by grouping similar points together.
2. **Versatility:** Applicable in diverse fields such as marketing, image compression, and bioinformatics.
3. **Scalability:** Handles large datasets effectively with relatively low computational complexity.
4. **Interpretable Outputs:** Provides intuitive cluster centroids and membership assignments for data points.

The Role of Centroids and Distance Metrics

At the heart of k-Means are two key concepts:

- **Centroids:** The central points of clusters, updated iteratively to best represent the data within their clusters.
- **Distance Metrics:** Measures like Euclidean or Manhattan distances, used to calculate the closeness of data points to centroids, determining cluster assignments.

Objectives of This Tutorial

This tutorial aims to:

1. Explain the theoretical foundations of k-Means, with a focus on centroids and distance metrics.
2. Demonstrate its implementation using Python, including visualization techniques.
3. Highlight its real-world applications and best practices for effective clustering.

2. Mathematical Foundations

The k-Means Algorithm

The k-Means algorithm partitions a dataset into k clusters by minimizing the **Within-Cluster Sum of Squares (WCSS)**. The objective function is defined as:

$$J = \sum_{i=1}^k \sum_{x \in C_i} ||x - \mu_i||^2$$

where:

- J : Total within-cluster variance.
- k : Number of clusters.
- C_i : Set of points in the i -th cluster.
- x : A data point.
- μ_i : Centroid of the i -th cluster.
- $||x - \mu_i||^2$: Squared distance between x and μ_i .

Steps of k-Means

1. **Initialization**: Choose k initial centroids randomly from the dataset to act as the starting points for clustering
2. **Assignment Step**: Assign each data point to the nearest centroid using the chosen distance metric. For **Euclidean distance**, this is:

$$d(x, \mu_i) = \sqrt{\sum_{j=1}^n (x_j - \mu_{ij})^2}$$

where x_j and μ_{ij} are the j -th feature values of the data point x and the centroid μ_i , respectively.

3. **Update Step**: Determine the new centroids by calculating the mean position of all data points within each cluster:

$$\mu_i = \frac{1}{|C_i|} \sum_{x \in C_i} x$$

where $|C_i|$ represents the number of data points in the i -th cluster.

4. **Repeat**: Alternate between the Assignment and Update steps until convergence is achieved, either when the centroids stabilize with minimal changes or the maximum number of iterations is reached.

The Role of Distance Metrics

The distance metric determines how the similarity between points is measured. Two commonly used metrics are:

- **Euclidean Distance:**

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

- **Manhattan Distance:**

$$d(x, y) = \sum_{i=1}^n |x_i - y_i|$$

Choosing the appropriate distance metric is crucial as it directly influences the clustering results. Euclidean distance is ideal for continuous numerical data, while Manhattan distance can be more robust to outliers.

3. Practical Implementation in Python

In this section, we demonstrate how to apply k-Means clustering using [mall-customers-dataset](#). The goal is to segment customers into distinct groups based on their annual income and spending scores, providing actionable insights for customer profiling.

Dataset Overview

The dataset contains 200 entries with the following key features:

1. **Annual Income (k\$):** Customers' yearly income in thousands of dollars.
2. **Spending Score (1-100):** A numerical rating provided by the mall, reflecting customer purchasing behavior and spending tendencies.

Steps to Implement k-Means

1. Load and Explore the Dataset

First, we load the dataset and select the relevant features for clustering.

```

# Import necessary libraries
import pandas as pd

# Load the dataset
data = pd.read_csv('Mall_Customers.csv')

# Select relevant features for clustering
features = data[['Annual Income (k$)', 'Spending Score (1-100)']]

# Preview the dataset
print(features.head())

```

Code Output:

	Annual Income (k\$)	Spending Score (1-100)
0	15	39
1	15	81
2	16	6
3	16	77
4	17	40

Figure 1: The preview shows a few initial rows, such as:

- Customer 0: Annual income of \$15k with a spending score of 39.
- Customer 3: Annual income of \$16k with a high spending score of 77.

This sample reflects diverse customer behaviors, forming the basis for clustering analysis.

2. Apply k-Means Clustering

We use the **k-Means algorithm** to cluster customers into five groups. Before clustering, the features are standardized to ensure equal contribution from both attributes.

```

from sklearn.preprocessing import StandardScaler
from sklearn.cluster import KMeans

# Standardize the features
scaler = StandardScaler()
features_scaled = scaler.fit_transform(features)

# Apply k-Means clustering
kmeans = KMeans(n_clusters=5, random_state=42)
kmeans.fit(features_scaled)

# Add cluster labels to the dataset
data['Cluster'] = kmeans.labels_

# Display results
print("Inertia for k=5:", kmeans.inertia_)
print("Cluster Labels:\n", data['Cluster'].value_counts())
print("Centroids:\n", kmeans.cluster_centers_)

```

Code Output:

Inertia for k=5: 65.56840815571681

Cluster Labels:

Cluster

0 81

1 39

3 35

4 23

2 22

Name: count, dtype: int64

Centroids:

[[-0.20091257 -0.02645617]

[0.99158305 1.23950275]

[-1.32954532 1.13217788]

[1.05500302 -1.28443907]

[-1.30751869 -1.13696536]]

Figure2: Cluster Summary

- **Cluster Distribution:** Cluster 0 has the largest number of customers (81), while Cluster 2 represents a smaller, niche segment (22 customers).
- **Inertia:** At $k=5$, the inertia value is **65.57**, indicating compact and well-separated clusters.
- **Centroids:** The centroid coordinates reflect the average standardized annual income and spending score for each cluster. For example:
 - Cluster 0: $[-0.20, -0.03]$ $[-0.20, -0.03]$ $[-0.20, -0.03]$ (moderate spenders).
 - Cluster 1: $[0.99, 1.24]$ $[0.99, 1.24]$ $[0.99, 1.24]$ (high-income, high-spending customers).

This provides actionable insights into customer behavior and group profiles.

3. Visualize Clusters in 2D

To understand the clusters, we use a scatter plot where data points are color-coded based on their cluster assignments. The centroids are marked in yellow.

```
import matplotlib.pyplot as plt

# Scatter plot of clusters
plt.figure(figsize=(8, 6))
colors = ['red', 'blue', 'green', 'orange', 'purple']
for cluster in range(5):
    cluster_points = features_scaled[kmeans.labels_ == cluster]
    plt.scatter(
        cluster_points[:, 0], cluster_points[:, 1],
        s=50, c=colors[cluster], label=f'Cluster {cluster + 1}'
    )

# Plot centroids
plt.scatter(kmeans.cluster_centers_[:, 0], kmeans.cluster_centers_[:, 1],
            s=200, c='yellow', marker='X', label='Centroids')
plt.title("k-Means Clustering of Mall Customers")
plt.xlabel("Annual Income (scaled)")
plt.ylabel("Spending Score (scaled)")
plt.legend()
plt.grid()
plt.show()
```

Cluster Visualization

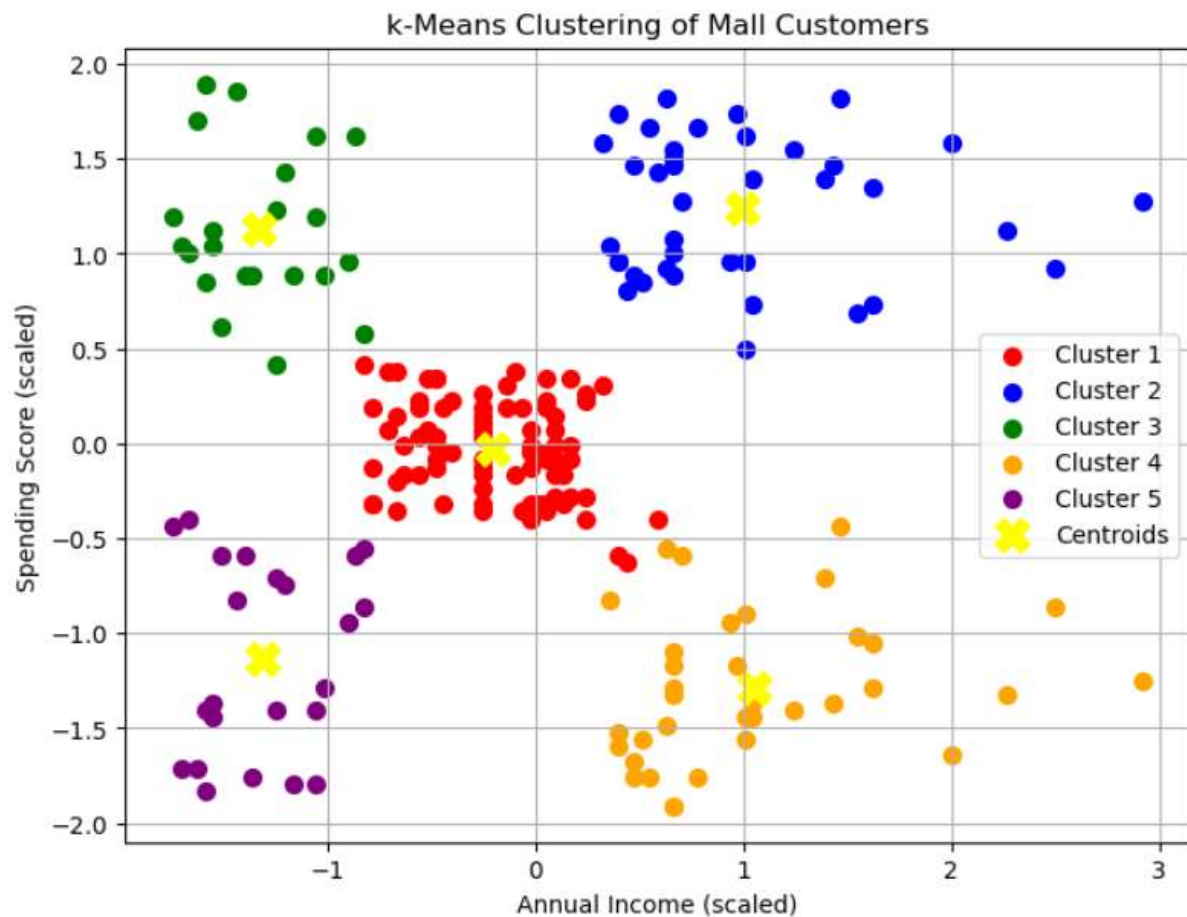


Figure 3: k-Means Clustering of Mall Customers :This scatter plot visualizes the clusters formed by the k-Means algorithm based on customers' annual income and spending score. The centroids (yellow markers) represent the average position of each cluster.

4. Determine Optimal Number of Clusters

The **elbow method** is used to validate the choice of five clusters by analyzing the inertia (sum of squared distances to the nearest centroid). A sharp bend in the elbow plot indicates the optimal number of clusters.

```

# Calculate inertia for different k values
inertia = []
for k in range(1, 11):
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(features_scaled)
    inertia.append(kmeans.inertia_)

# Plot the elbow curve
plt.figure(figsize=(8, 6))
plt.plot(range(1, 11), inertia, marker='o')
plt.title("Elbow Method for Optimal k")
plt.xlabel("Number of Clusters (k)")
plt.ylabel("Inertia")
plt.grid()
plt.show()

```

Elbow Curve

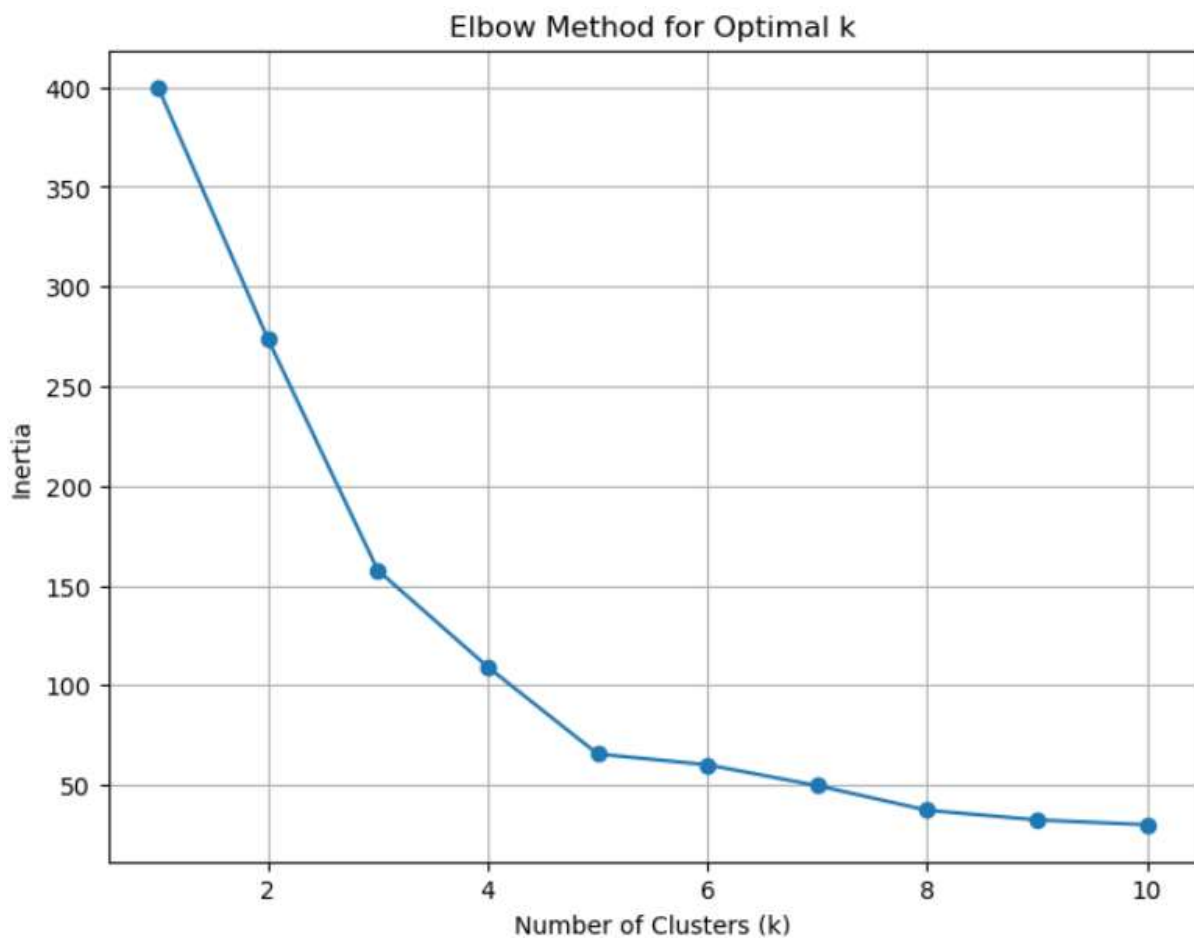


Figure 3: Elbow Method for Optimal k :This elbow plot shows the inertia (sum of squared distances to the nearest centroid) for different numbers of clusters (k). The sharp bend at $k=5$ indicates the optimal number of clusters.

Interpreting Results

1. **Cluster Visualization:** **Figure 1** shows distinct clusters based on customers' annual income and spending scores. Each cluster represents a unique group of customers with different behavioral patterns.
 - **Examples of clusters:**
 - **Cluster 2:** Customers with high income and elevated spending scores, often indicative of premium or high-value shoppers
 - **Cluster 5:** Low-income customers with low spending scores, indicating cost-sensitive or budget-conscious customers.
2. **Centroids:** The centroids, represented by yellow markers in **Figure 1**, indicate the average position of customers within each cluster. These centroids are useful for summarizing cluster characteristics and making data-driven decisions, such as tailored marketing strategies or personalized offers.
3. **Elbow Method:** The elbow plot in **Figure 2** demonstrates the optimal number of clusters, with a sharp bend at $k=5$. This confirms that five clusters provide the best trade-off between compactness (inertia) and interpretability, ensuring meaningful segmentation without overfitting.

4. Applications and Use Cases

k-Means clustering is a versatile algorithm with numerous real-world applications across various domains. By grouping data into clusters, it provides valuable insights for decision-making and resource optimization. Below are some practical use cases:

1. Customer Segmentation

k-Means is widely used in marketing to segment customers based on their behavior and demographics. For example:

- **Retail:** Grouping customers by spending habits and income levels to offer personalized promotions.
- **E-commerce:** Identifying high-value customers who frequently make large purchases to enhance loyalty programs.

In this tutorial's example, clustering customers in the Mall Customers dataset enables businesses to design targeted campaigns for different customer groups.

2. Image Compression

Clustering helps reduce the size of images by grouping similar colors and representing them with their centroids. This technique minimizes file size while preserving visual quality, making it useful for online platforms with storage constraints.

3. Document Categorization

In text analytics, k-Means is used to organize documents into topics based on word embeddings or feature vectors. This application is vital for tasks like news classification, customer feedback analysis, or organizing large repositories of textual data.

4. Anomaly Detection

k-Means can identify anomalies by highlighting data points that do not fit well into any cluster. This is crucial in fields like:

- **Fraud Detection:** Identifying unusual patterns in credit card transactions.
- **Healthcare:** Detecting irregular patient behaviors or outlier health metrics.

5. Geographical Mapping

k-Means clustering is applied to geographical data for resource planning:

- **Urban Planning:** Grouping neighborhoods based on socio-economic factors.
- **Logistics:** Optimizing warehouse locations by clustering delivery addresses.

These applications illustrate the power of k-Means clustering in simplifying complex datasets and deriving actionable insights. Its adaptability across industries makes it a cornerstone of modern data analytics and decision-making.

5. Best Practices for Effective Clustering

To achieve meaningful and reliable results with k-Means clustering, consider the following best practices:

1. Standardize Your Data

Standardizing features is essential when working with data that has varying scales (e.g., income in thousands vs. spending scores). Without standardization, features with larger magnitudes may dominate the clustering process.

Tip: Use techniques like z-score normalization to ensure all features contribute equally.

2. Determine the Optimal Number of Clusters

Selecting the right value for k is critical for accurate clustering. The **elbow method** is a commonly used technique that identifies k by analyzing the inertia curve. However, other methods like the **silhouette score** or domain-specific knowledge can also guide this decision.

Tip: Always validate the chosen k with multiple techniques for robustness.

3. Handle Outliers Carefully

Outliers can distort the clustering process by pulling centroids away from meaningful groupings.

Tip: Remove or transform extreme values before applying k-Means, or use robust clustering methods if outliers are expected.

4. Use Domain Knowledge

Incorporate domain expertise when selecting features for clustering. Irrelevant features can dilute the meaningfulness of clusters.

Tip: Perform exploratory data analysis (EDA) to identify relevant features that align with your objectives.

5. Interpret Results Thoughtfully

Clusters are statistical groupings and may not always align perfectly with real-world categories. Validate the clusters using additional metrics or domain insights.

Tip: Use visualization techniques, like scatter plots or heatmaps, to assess the quality and separability of clusters.

By following these best practices, you can enhance the accuracy and interpretability of k-Means clustering, ensuring actionable insights that align with your objectives.

6. Conclusion

k-Means clustering is a foundational unsupervised learning technique that provides actionable insights by grouping data into meaningful clusters. It is highly effective in identifying patterns, segmenting data, and aiding decision-making processes across various domains such as customer segmentation, anomaly detection, and resource optimization.

This tutorial explored the theoretical foundations of k-Means, emphasizing the importance of centroids and distance metrics in clustering. Through practical implementation using the Mall Customers dataset, we demonstrated how to segment customers based on their spending behavior and income levels. The use of visualization and the elbow method further highlighted how to validate and interpret clustering results effectively.

By adhering to best practices like data standardization, optimal cluster selection, and thoughtful interpretation, k-Means becomes a powerful tool for simplifying complex datasets. With its versatility and scalability, k-Means clustering remains a cornerstone of modern data analytics, enabling data-driven decision-making in real-world scenarios.

7. References and Resources

Resources

1. **GitHub Repository:** Explore the complete implementation of the k-Means clustering tutorial, including Python scripts, dataset, and detailed instructions. The repository provides all the necessary resources to replicate the tutorial, making it easy for learners

to follow and apply the concepts. Access the repository here: [k-Means Clustering Tutorial on GitHub](#).

2. **Google Colab Notebook:** An interactive notebook where you can execute the clustering tutorial step-by-step in a cloud-based environment. It eliminates the need for local installations and provides a user-friendly interface for testing and learning. Access the notebook here: [Interactive Notebook Link](#).
3. **Dataset:**
The tutorial utilizes the Mall Customers Dataset, which includes features such as annual income and spending scores for customer segmentation. The dataset is freely available on Kaggle and serves as an excellent resource for practicing clustering techniques. Download the dataset here: [Mall Customers Dataset on Kaggle](#).

References

Md Raiesh (2022). *Mall Customers Dataset*. [online] Kaggle.com. Available at: <https://www.kaggle.com/datasets/mdraiesh9648/mall-customers-dataset/>.

Scikit-learn (n.d.). *scikit-learn: Machine Learning in Python*. [online] Scikit-learn.org. Available at: <https://scikit-learn.org/stable/>.

Matplotlib (2012). *Matplotlib: Python plotting — Matplotlib 3.1.1 documentation*. [online] Matplotlib.org. Available at: <https://matplotlib.org/>.

Hastie, T., Tibshirani, R. and Friedman, J. (2009). *The Elements of Statistical Learning*. [online] *Springer Series in Statistics*. New York, NY: Springer New York. doi:<https://doi.org/10.1007/978-0-387-84858-7>.

Jain, A.K., Murty, M.N. and Flynn, P.J. (1999). Data clustering: a review. *ACM Computing Surveys*, [online] 31(3), pp.264–323. doi:<https://doi.org/10.1145/331499.331504>.

