# K-Means Clustering

**K-Means Clustering** is an unsupervised machine learning algorithm used to partition a dataset into distinct, non-overlapping groups or clusters. Each cluster is characterized by its centroid, which represents the mean position of all points within that cluster. The algorithm iteratively assigns data points to the nearest centroid and updates the centroids based on the mean of the assigned points until convergence.

**Key Features:**

1. **Data Loading and Preparation:**

   - Utilize **pandas** to load and manage the dataset efficiently.
   - Select relevant numerical features for clustering.

2. **K-Means Clustering:**

   - Apply the **K-Means algorithm** to partition the data into a predefined number of clusters.
   - Each data point is assigned to the nearest centroid.

3. **Visualization:**

   - Use **Matplotlib** to visualize the clusters, especially when dealing with two-dimensional data.
   - Plot data points with colors representing different clusters and centroids marked distinctly.

**How It Works:**

1. **Data Loading and Feature Selection:**

   - Load the dataset using **pandas** and select the numerical features relevant for clustering.

2. **K-Means Clustering:**

   - Initialize the **K-Means** algorithm with the desired number of clusters.
   - Fit the model to the selected features and predict the cluster labels for each data point.

3. **Visualization:**

   - For two-dimensional data, plot the data points with colors representing different clusters.
   - Mark the centroids of the clusters distinctly on the plot.

**Code Walkthrough:**

1. **Data Loading and Feature Selection:**

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans

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

# Select only numerical features for clustering
X = data.select_dtypes(include=[np.number])

# Display the first few rows
print(X.head())
```

2. **K-Means Clustering:**

```
 # Initialize k-Means with the number of clusters
kmeans = KMeans(n_clusters=3, random_state=42)

# Fit the model and predict cluster labels
clusters = kmeans.fit_predict(X)
```

3. **Visualization:**

```
 # Visualize the clusters (for 2D data only)
plt.scatter(X.iloc[:, 0], X.iloc[:, 1], c=clusters, cmap='viridis', s=50)
plt.scatter(kmeans.cluster_centers_[:, 0], kmeans.cluster_centers_[:, 1], color='re
plt.title('k-Means Clustering')
plt.xlabel('Feature 1')
plt.ylabel('Feature 2')
plt.legend()
plt.show()
```

**Advantages:**

- **Simplicity and Efficiency:** K-Means is straightforward to implement and computationally efficient, making it suitable for large datasets.
- **Scalability:** The algorithm scales well with the number of data points, especially when combined with dimensionality reduction techniques.
- **Versatility:** K-Means can be applied to various types of data and is widely used in customer segmentation, image compression, and pattern recognition.

**Future Work:**

- **Optimal Number of Clusters:** Explore methods like the Elbow Method or Silhouette Score to determine the optimal number of clusters.
- **Dimensionality Reduction:** Implement techniques such as Principal Component Analysis (PCA) to reduce the dimensionality of the data before clustering.
- **Comparison with Other Clustering Algorithms:** Evaluate the performance of K-Means against other clustering algorithms like DBSCAN or Hierarchical Clustering to determine the most suitable approach for specific datasets.

**References:**

- K-Means Clustering in Python: Step-by-Step Example - Statology
- Python Machine Learning - K-means - W3Schools
- K-Means Clustering with Python - Kaggle