

```
1 # =====
2 # Foundations of Data Science Project
3 # Dataset: Cleaned_Matches_Dataset.csv
4 # Topic: Clustering (Distance Metrics, KMeans, Hierarchical)
5 # =====
6
7 import pandas as pd
8 import numpy as np
9 import matplotlib.pyplot as plt
10 import seaborn as sns
11 from sklearn.preprocessing import StandardScaler
12 from sklearn.cluster import KMeans, AgglomerativeClustering
13 from sklearn.decomposition import PCA
14 from scipy.spatial.distance import cdist, pdist, squareform
15
16 # =====
17 # 1. Load Dataset
18 # =====
19 df = pd.read_csv("Cleaned_Matches_Dataset.csv")
20
21 # Drop unwanted index column if present
22 if "Unnamed: 0" in df.columns:
23     df = df.drop(columns=["Unnamed: 0"])
24
25 print("Shape of dataset:", df.shape)
26
27 # Select numeric features for clustering
28 num_cols = ["result_margin", "target_runs", "target_overs"]
29 df_cluster = df[num_cols].fillna(0)
30
31 # Standardize
32 scaler = StandardScaler()
33 X_scaled = scaler.fit_transform(df_cluster)
```

```
35 # =====
36 # 2. Distance Metrics
37 # =====
38 print("\n==== Distance Metrics ===")
39
40 euclidean_dist = cdist(X_scaled, X_scaled, metric="euclidean")
41 manhattan_dist = cdist(X_scaled, X_scaled, metric="cityblock")
42 cosine_dist = cdist(X_scaled, X_scaled, metric="cosine")
43
44 print("Euclidean Distance (first 5x5 block):\n", euclidean_dist[:5,:5])
45 print("Manhattan Distance (first 5x5 block):\n", manhattan_dist[:5,:5])
46 print("Cosine Distance (first 5x5 block):\n", cosine_dist[:5,:5])
47
48 # =====
49 # 3. K-Means Clustering (Lloyd's Algorithm)
50 # =====
51 print("\n==== K-Means Clustering (Lloyd's Algorithm) ===")
52
53 k = 3
54 kmeans = KMeans(n_clusters=k, random_state=42, n_init=10)
55 labels_kmeans = kmeans.fit_predict(X_scaled)
56
57 print("K-Means Cluster Centers:\n", kmeans.cluster_centers_)
58
59 # PCA projection for visualization
60 pca = PCA(n_components=2, random_state=42)
61 proj = pca.fit_transform(X_scaled)
62
63 plt.scatter(proj[:,0], proj[:,1], c=labels_kmeans, cmap="viridis")
64 plt.title("K-Means Clustering (k=3)")
65 plt.xlabel("PC1")
66 plt.ylabel("PC2")
67 plt.show()
68
```

```

69 # =====
70 # 4. Hierarchical Agglomerative Clustering (HAC)
71 # =====
72 print("\n== Hierarchical Agglomerative Clustering (HAC) ==")
73
74 hac = AgglomerativeClustering(n_clusters=3, metric="euclidean", linkage="ward")
75 labels_hac = hac.fit_predict(X_scaled)
76
77 plt.scatter(proj[:,0], proj[:,1], c=labels_hac, cmap="plasma")
78 plt.title("Hierarchical Agglomerative Clustering (HAC)")
79 plt.xlabel("PC1")
80 plt.ylabel("PC2")
81 plt.show()
82
83 # =====
84 # 5. Dendrogram (for Hierarchical Clustering)
85 # =====
86 from scipy.cluster.hierarchy import dendrogram, linkage
87
88 Z = linkage(X_scaled, method="ward") # ward = variance minimization
89 plt.figure(figsize=(10,5))
90 dendrogram(Z, truncate_mode="lastp", p=10, leaf_rotation=45, leaf_font_size=10, show_contracted=True)
91 plt.title("Hierarchical Clustering Dendrogram")
92 plt.xlabel("Cluster Size")
93 plt.ylabel("Distance")
94 plt.show()
95

```



