

```
1  # =====
2  # Foundations of Data Science Project
3  # Dataset: Cleaned_Matches_Dataset.csv
4  # Topic: DBSCAN, Clustering Tendency & Quality
5  # =====
6
7  import pandas as pd
8  import numpy as np
9  import matplotlib.pyplot as plt
10 from sklearn.preprocessing import StandardScaler
11 from sklearn.decomposition import PCA
12 from sklearn.cluster import KMeans, AgglomerativeClustering, DBSCAN
13 from sklearn.metrics import silhouette_score, davies_bouldin_score
14
15 # =====
16 # 1. Load Dataset
17 # =====
18 df = pd.read_csv("Cleaned_Matches_Dataset.csv")
19 if "Unnamed: 0" in df.columns:
20     df = df.drop(columns=["Unnamed: 0"])
21
22 # Use numeric features
23 num_cols = ["result_margin", "target_runs", "target_overs"]
24 df_cluster = df[num_cols].fillna(0)
25
26 # Standardize
27 scaler = StandardScaler()
28 X_scaled = scaler.fit_transform(df_cluster)
29
30 # PCA for visualization
31 pca = PCA(n_components=2, random_state=42)
32 proj = pca.fit_transform(X_scaled)
```

```

34 # =====
35 # 2. K-Means (baseline)
36 # =====
37 kmeans = KMeans(n_clusters=3, random_state=42, n_init=10)
38 labels_kmeans = kmeans.fit_predict(X_scaled)
39
40 # =====
41 # 3. Hierarchical Agglomerative Clustering (HAC)
42 # =====
43 hac = AgglomerativeClustering(n_clusters=3, metric="euclidean", linkage="ward")
44 labels_hac = hac.fit_predict(X_scaled)
45
46 # =====
47 # 4. DBSCAN
48 # =====
49 print("\n==== DBSCAN Clustering ===")
50 dbSCAN = DBSCAN(eps=1.2, min_samples=2) # tune eps based on dataset
51 labels_dbSCAN = dbSCAN.fit_predict(X_scaled)
52
53 print("DBSCAN Labels:", np.unique(labels_dbSCAN))
54
55 # Visualize DBSCAN clusters
56 plt.scatter(proj[:,0], proj[:,1], c=labels_dbSCAN, cmap="rainbow")
57 plt.title("DBSCAN Clustering")
58 plt.xlabel("PC1")
59 plt.ylabel("PC2")
60 plt.show()

```

```

62 # =====
63 # 5. Clustering Quality Metrics
64 # =====
65 print("\n==== Clustering Quality ===")
66
67 def evaluate_clustering(X, labels, name):
68     if len(set(labels)) > 1 and -1 not in set(labels): # at least 2 clusters, no noise on
69         sil = silhouette_score(X, labels)
70         dbi = davies_bouldin_score(X, labels)
71         print(f"{name}: Silhouette={sil:.3f}, DBI={dbi:.3f}")
72     else:
73         print(f"{name}: Not enough valid clusters for metrics.")
74
75 evaluate_clustering(X_scaled, labels_kmeans, "K-Means")
76 evaluate_clustering(X_scaled, labels_hac, "HAC")
77 evaluate_clustering(X_scaled, labels_dbSCAN, "DBSCAN")
78

```

Figure 1

DBSCAN Clustering

