# Cricket Format Clustering using Single Linkage Hierarchical Clustering Technique

Project Number: CS4

Submitted by:

Deepak Mewada (20CS91P02) - [deepakmewada96@kgpian.iitkgp.ac.in](mailto:deepakmewada96@kgpian.iitkgp.ac.in)

## Objective:

Consider the batting data of a batsman in a cricket match. His score has features like runs scored, ball faced, dot balls faced, fours scored and sixes scored. You have to cluster these scores into formats of cricket matches say Test cricket, One day cricket or Twenty-twenty.

## Section 0 - Import Python Libraries

```
 1 import numpy as np
 2 import pandas as pd
 3 import csv
 4
 5 #Only For testing silhouette scores
 6 from sklearn.cluster import KMeans
 7 from sklearn.metrics import silhouette_score
 8
 9 #For measuring time taken by the code
10 import time
11 start = time.process_time()
```

```
 1 #checking data
 2
 3 cricket_df = pd.read_csv("cricket_4_unlabelled.csv")
 4 cricket_df
```

| | Unnamed: 0 | total_balls_faced | dot_balls_faced | number_of_fours | number_of_sixes |
|---|---|---|---|---|---|
| **0** | 0 | 157 | 70 | 22 | 17 |
| **1** | 1 | 236 | 67 | 33 | 13 |
| **2** | 2 | 113 | 69 | 16 | 14 |
| **3** | 3 | 76 | 11 | 14 | 4 |
| **4** | 4 | 238 | 116 | 33 | 13 |
| **...** | ... | ... | ... | ... | ... |

## 0.1 Data Preprocessing

| **397** | 397 | 210 | 125 | 17 | 8 |
|---|---|---|---|---|---|

```
1 cricket_df = pd.read_csv("cricket_4_unlabelled.csv")
2 X = np.array(cricket_df).astype(np.float) #Converting the dataset into a numpy array of
3
4 X = np.delete(X,0,1)
5 print("Performing Normalisation of data.... ")
6
7 row = X.shape[0]
8 col = X.shape[1]
9
10 for m in range(0,col):
11     mu = np.mean((X[:,m]))
12     sigma = np.std((X[:,m]))
13     for n in range(0,row):
14         X[n,m] = (X[n,m] - mu)/sigma
15 print("Normalization Done !!!")
```

```
    Performing Normalisation of data....
    Normalization Done !!!
```

## Task 1 : K-means Clustering

Write a program to perform k-means clustering on the given dataset. Consider k=3 clusters. .
Randomly initialize k cluster means as k distinct data points. Iterate for 20 iterations. After the
iterations are over, save the clustering information in a file. This file may be used in step 4 if the
value of k is the optimal number of clusters.

## 1.a K-mean algorihem

```
1 class K_mean_clustering:
2
3     def euclidean_distance(A, B):
4         d = np.linalg.norm(A-B)
```

```
 5        return d
 6
 7    #This function assigns the centroids to the datasets based on the minimum euclidean
 8    def centroids_assignment(centroids, dataSet):
 9        assigned_centroid = []
10        for i in dataSet:
11            d = []
12            for j in centroids:
13                d.append(K_mean_clustering.euclidean_distance(i,j))
14            assigned_centroid.append(np.argmin(d))
15        return assigned_centroid
16
17    #Based on the assignments, we define the new centroids by
18    # taking the mean of all the data points that are assigned
19    #to one centroid in the last iteration.
20    def set_centroids(clusters, X):
21        new_centroid = []
22        new_df = pd.concat([pd.DataFrame(X), pd.DataFrame(clusters, columns=['cluster']
23        for c in set(new_df['cluster']):
24            current_centroid = new_df[new_df['cluster'] == c][new_df.columns[:-1]]
25            centroid_mean = current_centroid.mean(axis=0)
26            new_centroid.append(centroid_mean)
27        return new_centroid
28
```

## 1.b Run K-Means

To run K-Means, we first load the dataset and convert it into a numpy array for ease of operations.
We then take k = 3, as asked in the question and randomly identify 3 rows as the initial centroids.
After that, the functions *assign_centroids* and *update_centroids* are run iteratively 20 times to find the centroid for each dataset

```
 1 cricket_df = pd.read_csv("cricket_4_unlabelled.csv")
 2 print("loaded .CSV file : cricket_4_unlabelled.csv")
 3 X = np.array(cricket_df) #Converting the dataset into a numpy array
 4
 5 K = 3
 6 init_centroids = [0,3,41]
 7
 8 centroids = []
 9 for i in init_centroids:
10     centroids.append(X[i])
11 centroids = np.array(centroids)
12 print("K-means algo running for 20 iterations. . . ")
13 iters = 20
14 for i in range(iters):
15     get_centroids = K_mean_clustering.centroids_assignment(centroids, X)
16     centroids = K_mean_clustering.set_centroids(get_centroids, X)
17 print("K-means algo stopped. Cluster formed.")
```

```
loaded .CSV file : cricket_4_unlabelled.csv
K-means algo running for 20 iterations. . .
K-means algo stopped. Cluster formed.
```

## 1.c Write File

We now write down the clusters in a new file called *cricket_1_labelled_K*, where K = number of clusters (here, K=3).

```
1 cricket_df_cluster = cricket_df.copy() #initialise a new dataframe to store cluster inf
2 cricket_df_cluster['cluster'] = get_centroids
3 cricket_df_cluster.to_csv("cricket_4_labelled_K_is_3.csv")
4 print("File saved : cricket_4_labelled_K_is_3.csv")
```

```
File saved : cricket_4_labelled_K_is_3.csv
```

# Task 2 - Evaluation of the clustering algorithm

Evaluate the result of your clustering algorithm using the Silhouette coefficient metric and print the value of s. The Silhouette Coefficient is defined for each sample and is composed of two scores:

**a:** The mean distance between a sample and all other points in the same cluster.

**b:** The mean distance between a sample and all other points in the next nearest cluster.

The Silhouette Coefficient for a set of samples is given as the mean of the Silhouette Coefficient for each sample. The score is bounded between -1 for incorrect clustering and +1 for highly dense clustering. Scores around zero indicate overlapping clusters. Larger value of Silhouette Coefficient

## 2.a Function to calculate Silhouette Coefficient

Since we have to calculate the Silhouette Coefficient multiple times, we create a generic function to calculate Silhouette Coefficient. The function takes as inputs the dataframe along with cluster information and the value of the number of clusters (K).

We calculate a as the mean distance between a sample and all other points in the same cluster.

$$a = \frac{1}{|C_i| - 1} \sum_{j \in C_i,\ i \neq j} d(i, j)$$

We also calculate b as the mean distance between a sample and all other points in the next nearest cluster. Use of the minimum function is to check the next nearest cluster.

$$b = \min_{k \neq i} \frac{1}{|C_k|} \sum_{j \in C_k} d(i, j)$$

We then calculate the silhouette score for each instance of the dataset as

$$s_i = \frac{b_i - a_i}{max(a_i, b_i)}$$

Finally, **Silhouette Coefficient** of the entire dataset is taken as the mean of the Silhouette Scores of each instance

```
1  #defining Sillhout here
2  def silhoutte_metric(cricket_df_cluster,K):
3  
4      df_with_cluster = cricket_df_cluster.copy()
5      dataset_size = len(df_with_cluster.index)
6      np_df_with_cluster = df_with_cluster.to_numpy()
7  
8      df_with_cluster["silhouette_value"]=""
9  
10     for i in range(0, dataset_size):
11         this_instance = np_df_with_cluster[i,:-1]
12         this_centroid = np_df_with_cluster[i,-1]
13  
14         a = None
15         total_distance = 0
16         c = 0
17  
18         for j in range(0, dataset_size):
19             if this_centroid == np_df_with_cluster[j,-1]:
20                 d = np.linalg.norm(this_instance - np_df_with_cluster[j,:-1])
21                 total_distance += d
22                 c += 1
23  
24         #a is the mean distance to all other instances in this cluster
25         if (c>0):
26             a = total_distance / c
27             #print("a", i, "=", a)
28  
29         b = float("inf") #initialising b as maximum
30         avg_dist = 0
31  
32         for centroid in range(0,K):
33             total_distance = 0
34             c = 0
35             if centroid != this_centroid: #for all centroids except the one of this ins
36                 for k in range(0, dataset_size):
37                     if np_df_with_cluster[k,-1] == centroid:
38                         d = np.linalg.norm(this_instance - np_df_with_cluster[k,:-1])
39                         total_distance += d
40                         c += 1
41                 #print("td", i, "=", total_distance)
42                 if c > 0:
43                     avg_dist = total_distance / c
44                 if avg_dist < b:
45                     b = avg_dist #updating with the minimum b to know the average dista
46         s = (b - a) / max(a,b)
```

```
46          s = (b - a) / max(a,b)
47
48          df_with_cluster.loc[i,"silhouette_value"] = s
49
50      silhouette_value = df_with_cluster["silhouette_value"].mean()
51      return silhouette_value
```

## ▾ 2.b Calculate Silhouette Coefficient for K = 3

We now run the function to find the Silhouette Coeffecient of the dataset with K = 3

```
1 s = silhoutte_metric(cricket_df_cluster, 3)
2 print("The Silhouette coefficient metric for K = 3 is", s)
3
4 #testing silhouette score using sklearn
5 km = KMeans(n_clusters=3, random_state=42)
6 km.fit_predict(X)
7 score = silhouette_score(X, km.labels_, metric='euclidean')
8 print("via library sklearn Silhouette coefficient metric for K = 3  is", score)
```

```
    The Silhouette coefficient metric for K = 3 is 0.2975215057651689
    via library sklearn Silhouette coefficient metric for K = 3  is 0.3087725312230132
```

This score indicates that there has been correct clustering. However, the clusters are not very dense, and there is a fair bit of overlapping.

Using the silhouette_score function of sklearn.metrics, we find a similar score using the library too, and hence can verify that our clustering is corrector not.

## ▾ Section 3 - Find optimal value of K

Repeat steps 1 and 2 for k = 4, 5 and 6 as well. Report the value of k for which you get the highest value of the Silhouette Coefficient. This will be the optimal number of clusters. You will be using this number in the next step.

We now run the above steps for K = 4, 5 and 6

```
 1 print("#######################################################################################
 2 print()
 3 #fining optimal value of K
 4 for K in range (3,7):
 5     init_centroids = np.random.choice(range(0, len(X)), K)
 6     centroids = []
 7     for i in init_centroids:
 8         centroids.append(X[i])
 9     centroids = np.array(centroids)
10
11     iters = 20
12     for i in range(iters):
```

```
13        get_centroids = K_mean_clustering.centroids_assignment(centroids, X)
14        centroids = K_mean_clustering.set_centroids(get_centroids, X)
15
16    cricket_df_cluster = cricket_df.copy() #initialise a new dataframe to store cluster
17    cricket_df_cluster['cluster'] = get_centroids
18    cricket_df_cluster.to_csv("cricket_4_labelled_K_is_"+ str(K) + ".csv") #Store the f
19    print("File saved : cricket_4_labelled_K_is_"+ str(K) + ".csv !")
20
21    s = silhoutte_metric(cricket_df_cluster, K)
22    print("Silhouette coefficient metric for K =", K, "is", s)
23
24    #testing silhouette score using sklearn
25    km = KMeans(n_clusters=K, random_state=42)
26    km.fit_predict(X)
27    score = silhouette_score(X, km.labels_, metric='euclidean')
28    print("via library sklearn Silhouette coefficient metric for K =", K, " is", score,
29    print("------------------------------------------------------------------------
30    print()
31
32 print("##################################################################################
```

```
    ##################################################################################

    File saved : cricket_4_labelled_K_is_3.csv !
    Silhouette coefficient metric for K = 3 is 0.3124821590171541
    via library sklearn Silhouette coefficient metric for K = 3  is 0.3087725312230132


    ------------------------------------------------------------------

    File saved : cricket_4_labelled_K_is_4.csv !
    Silhouette coefficient metric for K = 4 is 0.30939421470920025
    via library sklearn Silhouette coefficient metric for K = 4  is 0.302290513416582


    ------------------------------------------------------------------

    File saved : cricket_4_labelled_K_is_5.csv !
    Silhouette coefficient metric for K = 5 is 0.29381704722361957
    via library sklearn Silhouette coefficient metric for K = 5  is 0.28633439589405785


    ------------------------------------------------------------------

    File saved : cricket_4_labelled_K_is_6.csv !
    Silhouette coefficient metric for K = 6 is 0.2997417561012066
    via library sklearn Silhouette coefficient metric for K = 6  is 0.29049853037081896


    ------------------------------------------------------------------

    ##################################################################################
```

```
 1 K_optimal=3   #After analysing results in above cell
 2 print("Since for K = ",K_optimal,"; Silhouette coefficient is highest hence optimal val

    Since for K =  3 ; Silhouette coefficient is highest hence optimal value of K is :  3
```

We observe that the best Silhouette Coefficient comes from k = 3. Thus here we are submitting a copy of the file *cricket_4_labelled_K_is_3.csv* as **kmeans.txt** in the overall ZIPPED submission. The file **kmeans.txt** is created as follows:

```
1 print("loaded .CSV file : cricket_4_labelled_K_is_3.csv")
2
3 cricket_df_cluster_final = pd.read_csv("cricket_4_labelled_K_is_3.csv")
4 cricket_df_cluster_final = cricket_df_cluster_final[['Unnamed: 0','cluster']]
5
6 l = len(cricket_df_cluster_final)
7
8 #initialize the cluster arrays
9 cluster={}
10 for c in range (K_optimal):
11     cluster[c]=[]
12
13 for i in range(l):
14     for j in range(K_optimal):
15         if(cricket_df_cluster_final.at[i,'cluster']) == j:
16             cluster[j].append(cricket_df_cluster_final.iat[i,0])
17
18 C=[]
19 for p in range (K_optimal):
20     C.append(cluster[p])
21
22 df=pd.DataFrame(C)
23 df = df.sort_values(df.columns[0])
24 df = df.fillna('')
25 df.to_csv("kmeans.txt", header = False, sep = ",", index = False)
26 print("kmeans.txt is saved !!!")

    loaded .CSV file : cricket_4_labelled_K_is_3.csv
    kmeans.txt is saved !!!
```
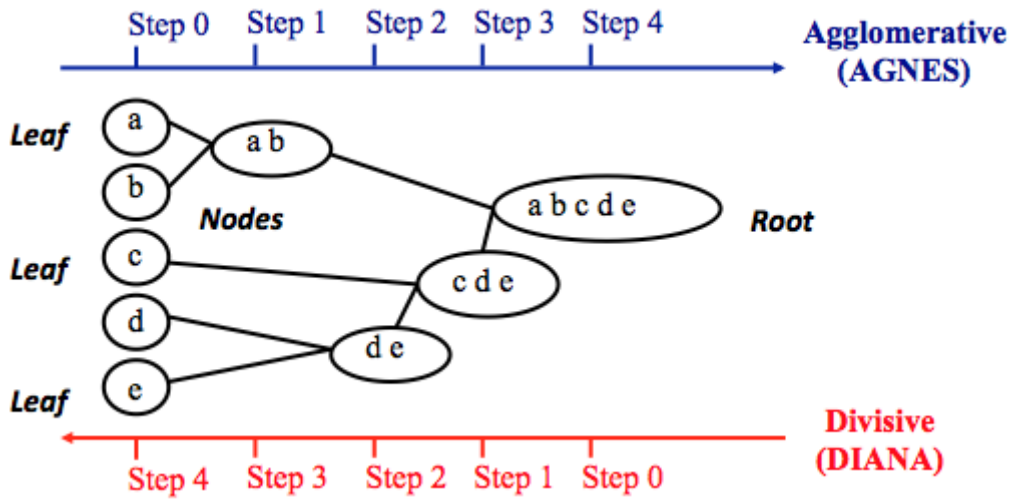
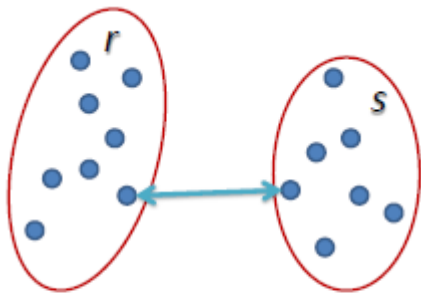# Task 4 - Bottom-Up Hierarchical Clustering with single linkage strategy

Implement a bottom-up hierarchical clustering algorithm considering the same notion of similarity as in step 1. Find k clusters (optimal number of clusters from step 3) using single linkage strategy.

## Bottom-Up Hierarchical Clustering is illustrated as follows :

Single linkage technique si explained in this figure



$$L(r,s) = \min(D(x_{ri}, x_{sj}))$$

```
1 class agglomerative_SingleLinkageStartegy:
2
3     def __init__(self, data, k):
4         self.k = k
5         self.data = data
6         self.fit()
7
8     def fit(self):
9         n = len(self.data)
10         self.clusters = {}
11
12         for i in range(n):
13             self.clusters[i] = []
14             self.clusters[i].append(i)
15
16         self.dist = np.sqrt((np.square(self.data[:,np.newaxis]-self.data).sum(axis=2)))
17
18         for i in range(n-self.k):
19             merge = self.merging()
20             self.clusters[merge[0]] = self.clusters[merge[0]] + self.clusters[merge[1]]
21             self.clusters.pop(merge[1])
22
23         for i in range(self.k):
```

```python
24              while not 1 in self.clusters:
25                  for j in [x for x in list(map(int, self.clusters.keys())) if x >= i+1]:
26                      self.clusters[j-1] = self.clusters.pop(j)
27
28          for i in self.clusters.keys():
29              self.clusters[i].sort()
30
31      def merging(self):
32          mini = 1e99 #taking big value
33          merge = (None, None)
34
35          for i in list(map(int, self.clusters.keys())):
36              for j in [x for x in list(map(int, self.clusters.keys())) if x >= i+1]:
37                  if self.dist[i][j] < mini:
38                      mini = self.dist[i][j]
39                      merge = (i, j)
40
41          return merge
```

```python
1 #Since our input matrix may have changed due to multiple operations, we re-initiate the
2 cricket_df1 = pd.read_csv("cricket_4_unlabelled.csv")
3 X1 = np.array(cricket_df1).astype(np.float) #Converting the dataset into a numpy array
4 X1 = np.delete(X1,0,1)
5 print("Performing Normalisation of data.... ")
6
7 row_count = X1.shape[0]
8 col_count = X1.shape[1]
9
10 for m in range(0,col_count):
11     mu = np.mean((X1[:,m]))
12     sigma = np.std((X1[:,m]))
13     for n in range(0,row_count):
14         X1[n,m] = (X1[n,m] - mu)/sigma
15 print("Normalisation Done!!!")
```

```
    Performing Normalisation of data....
    Normalisation Done!!!
```

```python
1 #calling agglomerative
2 print("Bottom Up Agglomerative single linkage cluster formation !!!")
3 hac = agglomerative_SingleLinkageStartegy(X1,K_optimal)
```

```
    Bottom Up Agglomerative single linkage cluster formation !!!
```

```python
1 #Saving the Clusters
2
3 all_cluster_agglo=np.empty(len(X1), dtype=object)
4 for i in range(3):
5     np.put(all_cluster_agglo,hac.clusters[i],i)
6
7 agglo_cluster={}
8 agglo_cluster[0] = hac.clusters[0]
9 agglo_cluster[1] = hac.clusters[1]
10 agglo_cluster[2] = hac.clusters[2]
```

```
11
12 all_cluster = [agglo_cluster[0], agglo_cluster[1], agglo_cluster[2]]
13 df=pd.DataFrame(all_cluster)
14 df = df.sort_values(df.columns[0])
15 df = df.fillna('')
16 df.to_csv("agglomerative.txt", header = False, sep = ",", index = False)
17 print("agglomerative.txt is saved !!!")
```

```
    agglomerative.txt is saved !!!
```

## Task 4.b : Jaccard similarity between corresponding sets of both the cases

Now you have k clusters from the k-means algorithm and k clusters from hierarchical clustering on the same dataset. Or in other words, the dataset is divided into k sets of data points as a result of the k-means algorithm (case A). Similar is the case for the hierarchical clustering algorithm (case B). You need to compute the Jaccard similarity between corresponding sets of both the cases.
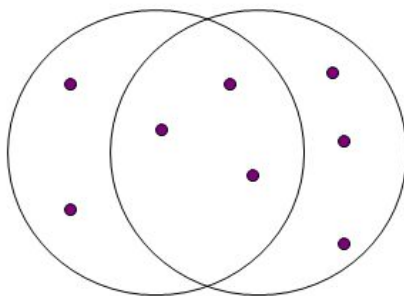
## Jaccard similarity coefficient is calculated as follows

```
1 #defining Jaccard function
2 def jaccard(list1, list2):
3     s1 = set(list1)
4     s2 = set(list2)
5
6     return float(len(s1.intersection(s2)) / len(s1.union(s2)))
```

```
 1 #calling Jacckard function
 2 jc={}
 3 for c in range (K_optimal):
 4     jc[c]=[]
 5
 6 print("Calculating Jaccard Similarity scores for all the mappings.....")
 7 print("------------------------------------------------------------------------")
 8 print()
 9 for i in range(K_optimal):
10   for j in range(K_optimal):
11     jc[i].append(jaccard(cluster[i],agglo_cluster[j]))
12 jc[i][j]*=10
13 print("Jaccard similarity score bet cluseter{0,1,2} of K-mean & cluster_agg{0,1,2} of A
14 print()
15 print()
16
17 print(jc[0])
18 print(jc[1])
19 print(jc[2])
20
21
22 print()
23 print("-----------------------------------------------------------------------------")
```

```
    Calculating Jaccard Similarity scores for all the mappings.....
    ------------------------------------------------------------------------

    Jaccard similarity score bet cluseter{0,1,2} of K-mean & cluster_agg{0,1,2} of Agglom


    [0.41904761904761906, 0.0, 0.013157894736842105]
    [0.13577023498694518, 0.45384615384615384, 0.08396946564885496]
    [0.3932926829268293, 0.03940886699507389, 0.4458598726114649]


    ------------------------------------------------------------------------
```

Here the Jaccard value shows the similarity between the clusters. It is dynamic;depends upon the formed clusters of K-mean & Agglomerative respectively. Higher the value in a row means both the corresponding cluster have high similarity.

```
1 print("Time taken for overall program single run is", time.process_time() - start, "sec

    Time taken for overall program single run is 63.66798628700002 seconds !!!
```

```
Thank you for reading so far!!!
  This is the end of file.
```