**TRƯỜNG ĐẠI HỌC BÁCH KHOA HÀ NỘI**
HANOI UNIVERSITY OF SCIENCE AND TECHNOLOGY

# Introduction to Data Science

Capstone Project Name: FIFA 2021 Clustering
Assistant Lecturer: PhD. Muriel Visani
Group number: 14

**Nguyen Nho Trung – 20204894**

**Nguyen Thi Thu Giang - 20194750**

**Tong Thi Thu Anh – 20194728**

# Introduction

Our subject for the Capstone is applying Machine Learning to cluster football players. In detail, our project aims to look at the class of young players and cluster them into groups based on their attributes. We will filter and consider the players who are promised based on their potential attributes.
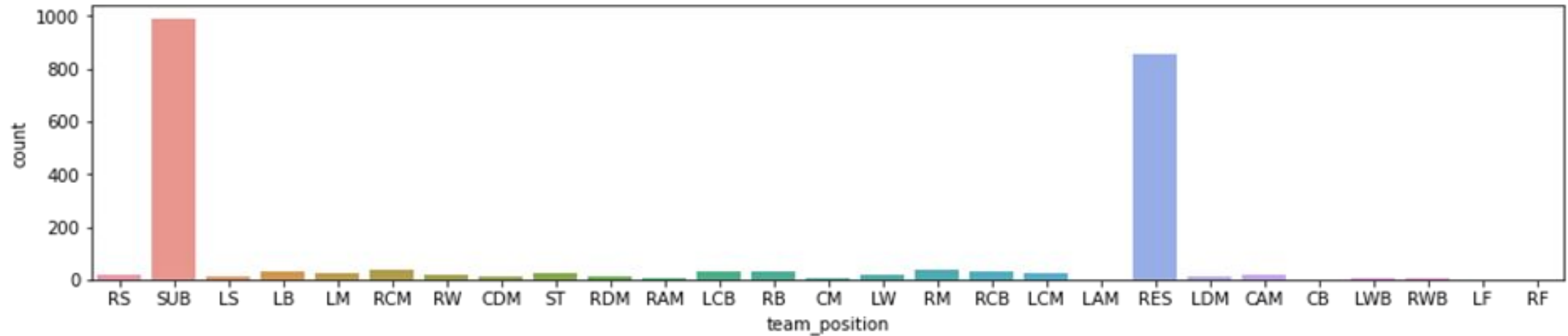
# DATASET

- Name
- Age
- Height
- Overall: Rated between 1-99
- Potential: Rated between 1-99
- Work Rate: Effort the player puts in the game
- Team Position
- Playing Skills: Rated between 1-99
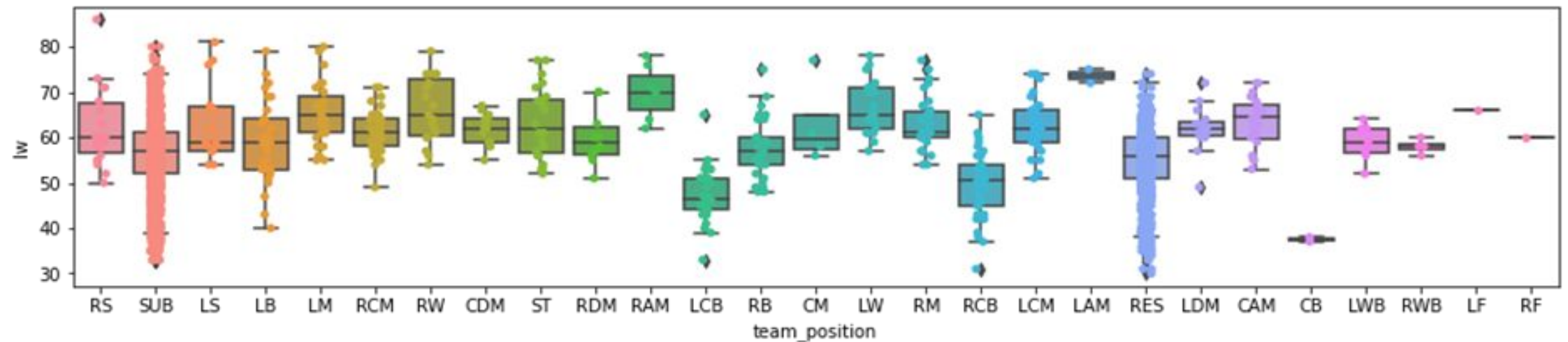- Positional Skills: Rating of a player when play at a specified position rated between 1-99
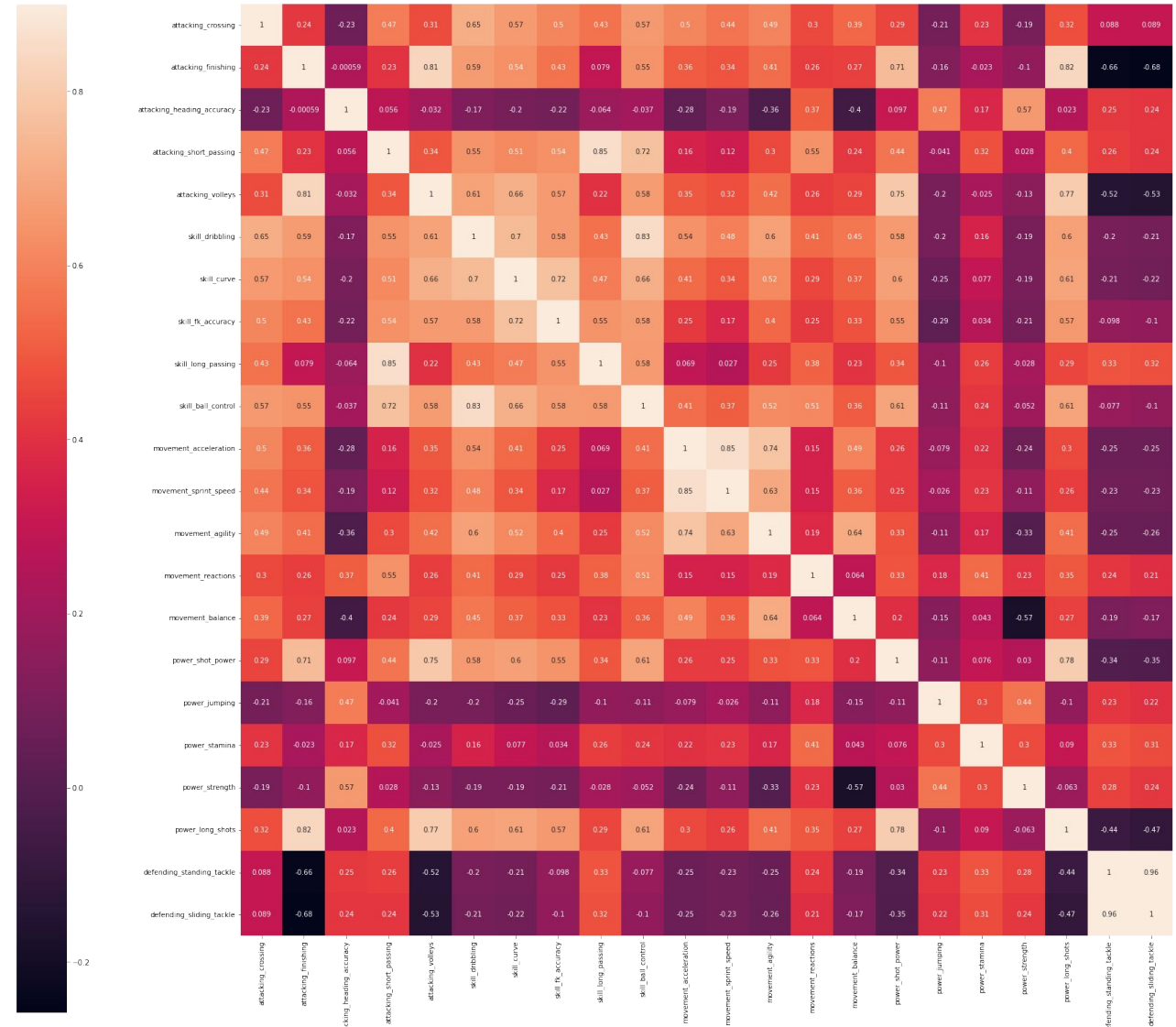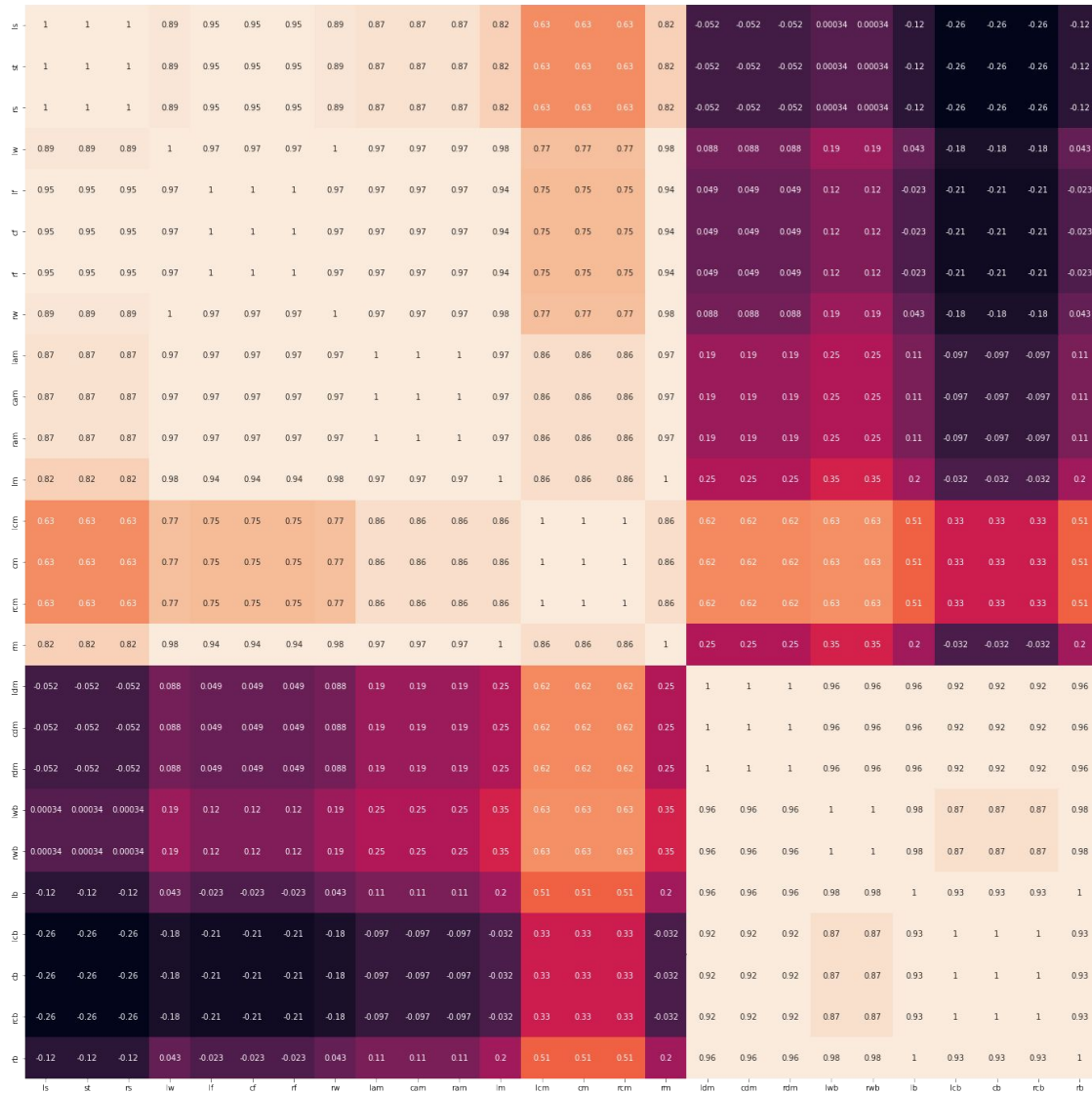
Plot the number of players in each position:



Plot the relationship between team positions and positional skills.

# Correlation between the position and playing skills features

Step 1: Drop unused attributes for clustering

Step 2: Reduce Dimension

Step 3: Standardized data.

```
Explained variance ratio for Striker PCA application:  [1.]
Explained variance ratio for Forward PCA application:  [1.]
Explained variance ratio for Attack Winger PCA application:  [1.]
Explained variance ratio for Midfielder PCA application:  [0.70388172]
Explained variance ratio for Defensive Winger PCA application:  [1.]
Explained variance ratio for Back PCA application:  [1.]
```

```
Explained variance ratio for defensive skills PCA application:  [0.98299723]
Explained variance ratio for movement PCA application:  [0.82825812]
Explained variance ratio for control skills PCA application:  [0.91608148]
Explained variance ratio for passing skills PCA application:  [0.92329027]
Explained variance ratio for defending skills PCA application:  [0.91549397]
```
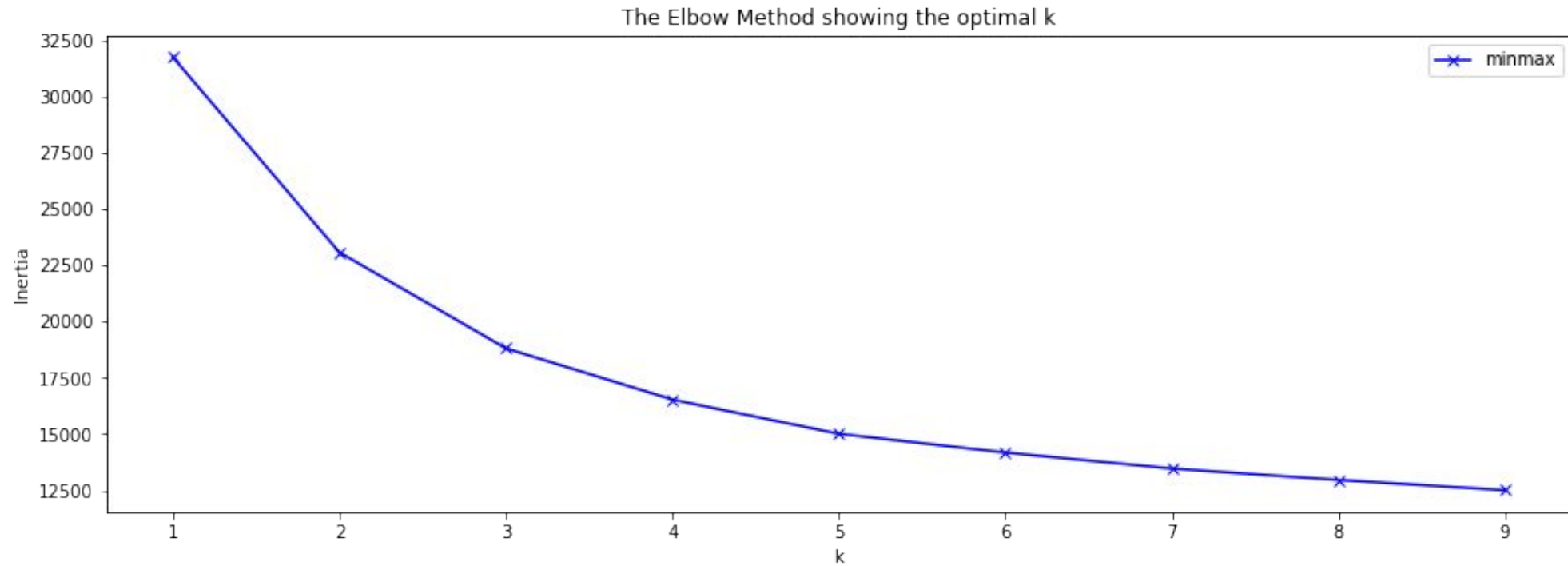
```
Explained variance ratio for Striker PCA application:  [1.]
Explained variance ratio for Forward PCA application:  [1.]
Explained variance ratio for Attack Winger PCA application:  [1.]
Explained variance ratio for Midfielder PCA application:  [0.70388172]
Explained variance ratio for Defensive Winger PCA application:  [1.]
Explained variance ratio for Back PCA application:  [1.]
```

```
Explained variance ratio for defensive skills PCA application:  [0.98299723]
Explained variance ratio for movement PCA application:  [0.82825812]
Explained variance ratio for control skills PCA application:  [0.91608148]
Explained variance ratio for passing skills PCA application:  [0.92329027]
Explained variance ratio for defending skills PCA application:  [0.91549397]
```

# Discover the relationship between Work Rate and other skills

|  | attacking_crossing | attacking_finishing | attacking_short_passing | attacking_volleys |
|---|---|---|---|---|
| **attack_work_rate** | | | | |
| **0.0** | 37 | 26 | 56 | 29 |
| **0.5** | 49 | 48 | 59 | 42 |
| **1.0** | 55 | 56 | 60 | 47 |

|  | defending_standing_tackle | defending_sliding_tackle |
|---|---|---|
| **defense_work_rate** | | |
| **0.0** | 28 | 26 |
| **0.5** | 52 | 49 |
| **1.0** | 62 | 59 |

The Elbow Method showing the optimal k

# Silhouette method



Silhouette analysis for KMeans clustering on sample data with n_clusters = 2
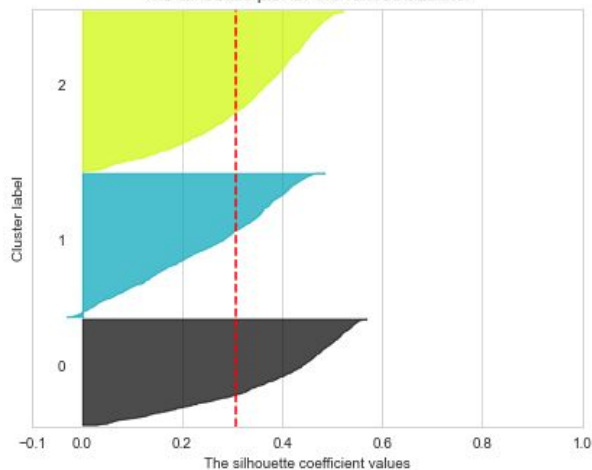
```
For n_clusters = 2 The average silhouette_score is : 0.28272602345704606
For n_clusters = 3 The average silhouette_score is : 0.3084938723509853
For n_clusters = 4 The average silhouette_score is : 0.2795438505249029
For n_clusters = 5 The average silhouette_score is : 0.22999870816740864
For n_clusters = 6 The average silhouette_score is : 0.21964100225379676
For n_clusters = 7 The average silhouette_score is : 0.20744387774940753
For n_clusters = 8 The average silhouette_score is : 0.19222045902464063
For n_clusters = 9 The average silhouette_score is : 0.1860086216190746
```
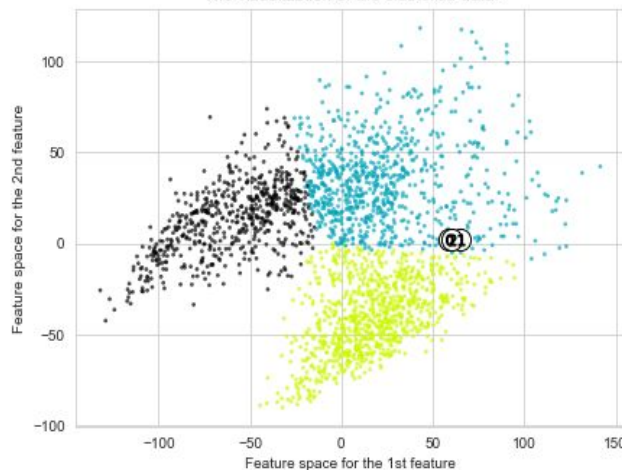
r of b(i) and a(i)

Silhouette analysis for KMeans clustering on sample data with n_clusters = 3
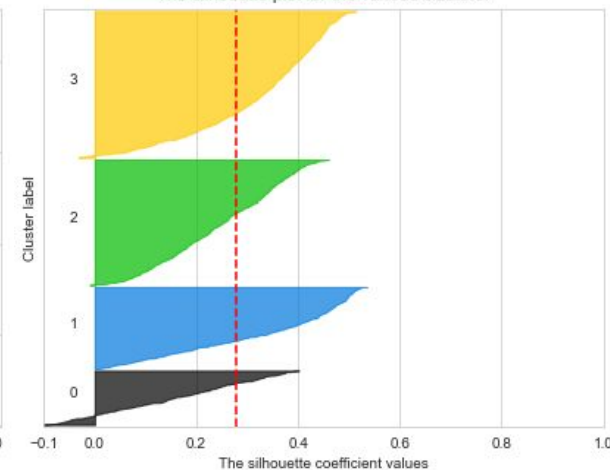
Silhouette analysis for KMeans clustering on sample data with n_clusters = 4

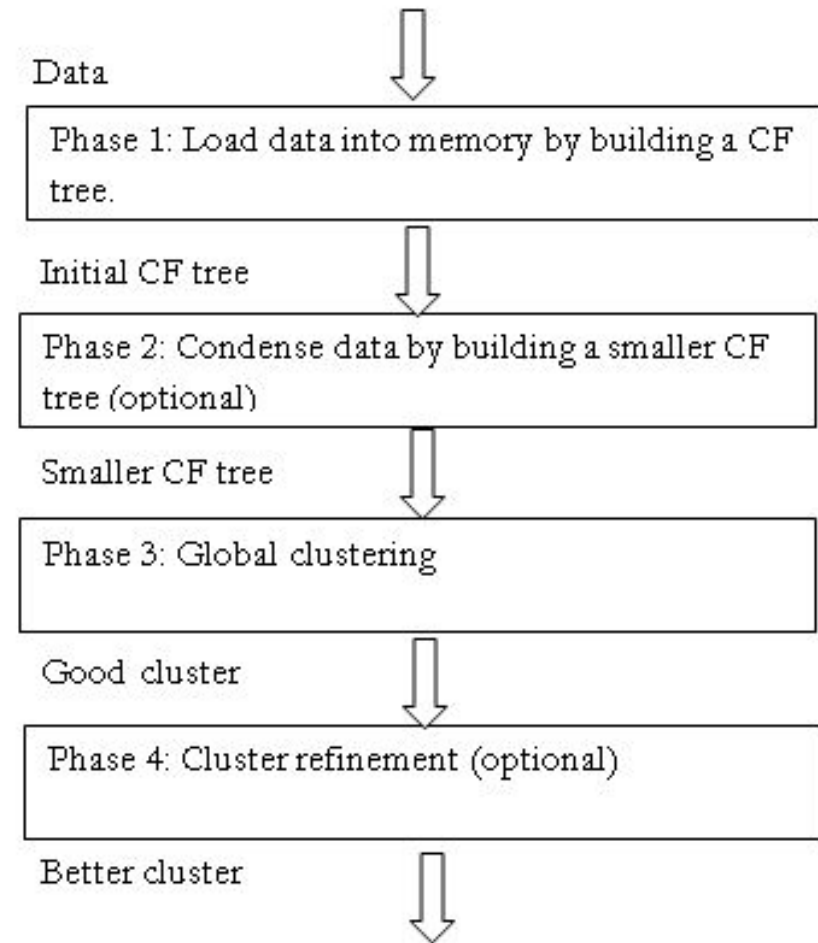$$L(r,s) = \max(D(x_{ri}, x_{sj}))$$



dendrogram

# Birch Algorithm

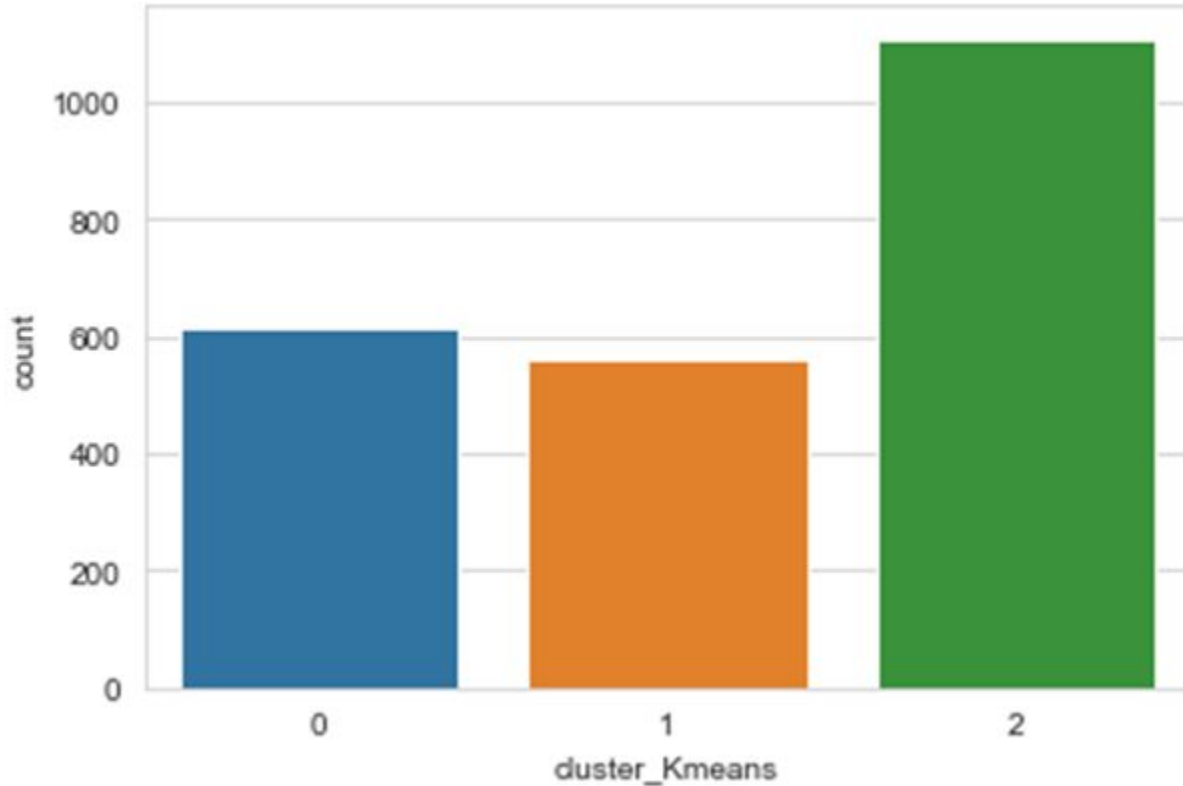Birch is a Hierarchical clustering algorithm using top-down

- Idea: Each successful iteration, a cluster is split into smaller clusters according to the value of some similarity measure until each object is a cluster or until the stopping condition is satisfied.
- This approach into using divide and conquer strategy in the clustering process

Data

Phase 1: Load data into memory by building a CF tree.

Initial CF tree

Phase 2: Condense data by building a smaller CF tree (optional)

Smaller CF tree

Phase 3: Global clustering

Good cluster

Phase 4: Cluster refinement (optional)

Better cluster

# Comparision three Algorithms

| | K cluster | kmean_silhouette_scores | birch_silhouette_scores | agg_silhouette_scores |
|---|---|---|---|---|
| 0 | 2 | 0.282726 | 0.238964 | 0.238964 |
| 1 | 3 | 0.308363 | 0.236685 | 0.236685 |
| 2 | 4 | 0.279544 | 0.212797 | 0.212797 |
| 3 | 5 | 0.229999 | 0.193872 | 0.193872 |
| 4 | 6 | 0.219155 | 0.185722 | 0.185722 |
| 5 | 7 | 0.206945 | 0.154607 | 0.154607 |
| 6 | 8 | 0.191965 | 0.155710 | 0.155710 |
| 7 | 9 | 0.185837 | 0.155961 | 0.155961 |

| cluster_Kmeans | potential | overall |
| --- | --- | --- |
| 2 | 79.738211 | 67.549593 |
| 1 | 74.443060 | 59.176157 |
| 0 | 74.352570 | 58.725879 |

```
2     1109
0      615
1      562
Name: cluster_Kmeans, dtype: int64
running time of Kmeans  70.96637034416199
```

# Result

| cluster_Kmeans | striker | forward | attack_winger | midfielder | defensive_winger | back |
|---|---|---|---|---|---|---|

| cluster_Kmeans | defensive | movement | control | passing | finishing |
|---|---|---|---|---|---|
| | | | | | 3 |
| | | | | | 3 |
| | | | | | 7 |

| cluster_Kmeans | attack_work_rate | defense_work_rate | defensive | finishing |
|---|---|---|---|---|
| 0 | 0.494662 | 0.588968 | 21.025102 | 23.755704 |
| | | | | 5.300393 |
| 1 | 0.638864 | 0.458070 | -13.352598 | 12.150520 |
| 2 | 0.704878 | 0.518699 | 4.864916 | |

# THANK YOU FOR LISTENING