

## PROJECT II: Player's performance in Premier League (FIFA22)

Usama ZAFAR

Bishwash NEUPANE

Kirubel WELDETENSAY

Nethaji RUTHIRAVEL

### Objective:

- Perform statistical analysis of a selected dataset using machine learning techniques
- Perform Kmeans clustering in the original space and reduced dimension following PCA
- Use linear model and develop insights/inferences from the analysis

### Dataset:

The dataset used for the analysis purpose is [FIFA 22](#) available on Kaggle. As the dataset is large (more than 19000 players with 110 features), data cleaning is done in the original dataset. Players playing in **English Premier League**, having an **overall rating of more than 80** are considered **excluding the goalkeepers**. The goalkeepers are removed from consideration since they are too different from the field players who are defensive by nature and they can be easily clustered in a separate group (**only 10% of goalkeepers** in the initial sample size).

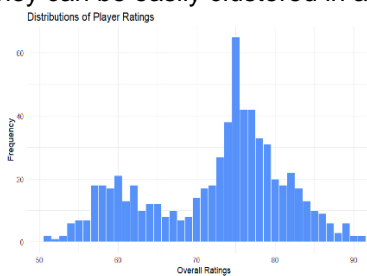


Figure 1(a) Distribution of overall player ratings

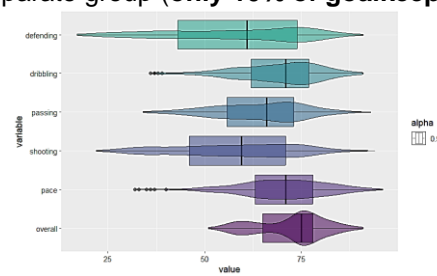


Figure 1(b) Violin plot for a few of the selected features

```
## [1] "pace" "shooting"
## [3] "passing" "dribbling"
## [5] "defending" "physic"
## [7] "attacking_crossing" "attacking_finishing"
## [9] "attacking_heading_accuracy" "attacking_short_passing"
## [11] "attacking_volleys" "skill_dribbling"
## [13] "skill_curve" "skill_fk_accuracy"
## [15] "skill_long_passing" "skill_ball_control"
## [17] "movement_acceleration" "movement_sprint_speed"
## [19] "movement_agility" "movement_reactions"
## [21] "movement_balance" "power_shot_power"
## [23] "power_jumping" "power_stamina"
## [25] "power_strength" "power_long_shots"
## [27] "mentality_aggression" "mentality_interceptions"
## [29] "mentality_positioning" "mentality_vision"
## [31] "mentality_penalties" "mentality_composure"
## [33] "defending_marking_awareness" "defending_standstill_tackle"
## [35] "defending_sliding_tackle"
```

Figure 1(c) Features based on which clustering is possible

### Feature Discussion and Clustering Idea:

The dataset contains characteristics that can be used to cluster players into different groups: attacking, defensive, and midfield players. Age, potential, salary, and a few other features are excluded as they only stratify the players making the clustering meaningless. Distribution of overall rating and the violin plots (**Figure 1(a) and 1(b)**) reveal that some characteristics have a **bimodal distribution**, indicating that clustering is feasible with the selected features. **Figure 1(c)** shows the list of features that are used to cluster the players. As there are a bunch of features, a **minimal analysis** is reasonable to extract more information from the dataset using less variables. For this purpose, two data frames have been created, the first with features used for analysis and the second one with the categorical values with relevant information about the players that is used to interpret the result later on.

### k-means in the Original Space

After selecting relevant features from a dataset (35 variables), the k-means algorithm was applied to the original space without scaling as the selected features are on the same scale. The **Within-Cluster Sum of Squares (WSS)** graph is plotted to determine the optimal number of clusters, which is found to be **4**. The algorithm resulted in the formation of 4 clusters: **cluster 1 for attacking midfielders** (such as Kevin De Bruyn, Bruno Fernandes), **cluster 2 for defensive midfielders** (like Thomas Partey, Fernandino), **cluster 3 for defensive players** (such as Raphael Varane, Ruben Dias), and **cluster 4 for attackers/strikers** (such as Cristiano Ronaldo, Son Heung Min, Harry Kane).

overall	shooting	passing	dribbling	defending	physic
Min. :80.00	Min. :30.0	Min. :52.00	Min. :57.00	Min. :33.00	Min. :57.00
1st Qu.:81.00	1st Qu.:59.5	1st Qu.:72.00	1st Qu.:76.00	1st Qu.:47.75	1st Qu.:69.75
Median :83.00	Median :72.5	Median :76.00	Median :80.50	Median :72.50	Median :75.00
Mean :83.34	Mean :69.0	Mean :76.03	Mean :79.34	Mean :65.62	Mean :74.48
3rd Qu.:85.00	3rd Qu.:79.0	3rd Qu.:81.00	3rd Qu.:84.00	3rd Qu.:81.00	3rd Qu.:80.00
Max. :91.00	Max. :94.0	Max. :93.00	Max. :91.00	Max. :91.00	Max. :88.00

Figure 2 Summary of the few features

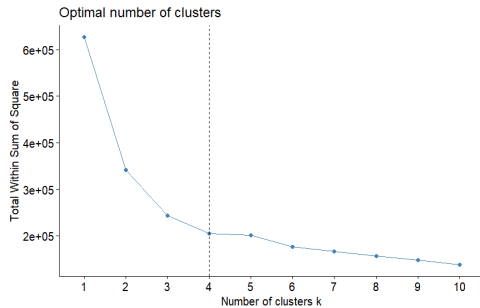


Figure 3(a) Optimal Number of Clusters

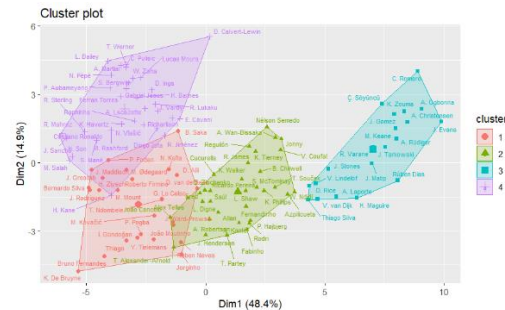


Figure 3(b) Clustering on original dimension

### Principal Component Analysis and Dimensionality Reduction

The first two dimensions **represent 65%** of the total inertia (the inertia is the total variance of the dataset i.e., the trace of the correlation matrix).

The variables "defending, defending-marking-awareness" are not correlated to the variables "attack-finishing, movement agility, movement-sprint-speed, attacking-shooting" because the skillset required for attackers is not relevant for defenders. For example, M. Salah lies on the right side because he has the skillsets of the attacker, meanwhile, Harry Maguire is on the left side because he has the skill sets of the defender. These skill sets have opposing features because they are used against each other in a match, as defenders are going against attackers. The first axis opposes players who are good at attacking like Mo. Salah between those who are good at defending like Harry Maguire.

The second axis opposes players who are good at both defending and attacking (possess the skillsets of both categories) against those who are not good at both. In other words, high level midfielders (top half) against low level midfielders (bottom half). So, attacking players who lie on the bottom right of **Fig 4b** (Anthony Martial, T. Werner) are not outstanding attackers, and also, they possess very low defending skills. We can divide the factorial plan into four parts:

<b>Top Left</b> <ul style="list-style-type: none"> <li>High Defending skill</li> <li>Medium Attacking Skill</li> <li>Fabian, Henderson</li> </ul>	<b>Top Right</b> <ul style="list-style-type: none"> <li>High attacking skill</li> <li>Medium Defending skill</li> <li>De Bryune, Thiago</li> </ul>
<b>Bottom Left</b> <ul style="list-style-type: none"> <li>High Defending skill</li> <li>Low Attacking Skill</li> <li>J Gomez, Raphael Varane</li> </ul>	<b>Bottom Right</b> <ul style="list-style-type: none"> <li>High Attacking Skill</li> <li>Low defending skill</li> <li>Diego Jota, Bukayo Saka</li> </ul>

The players at the top right and top left are mostly categorized as midfielders because they possess good level of both skills. The players mentioned are highlighted in **figure 4b** with boxes.

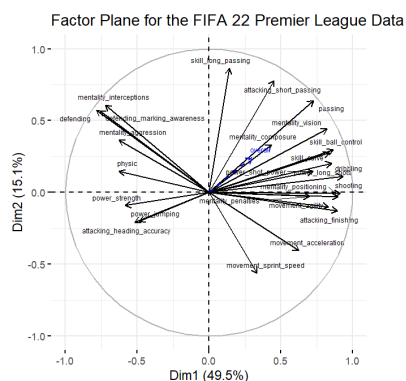


Figure 4(a) Factor Plane

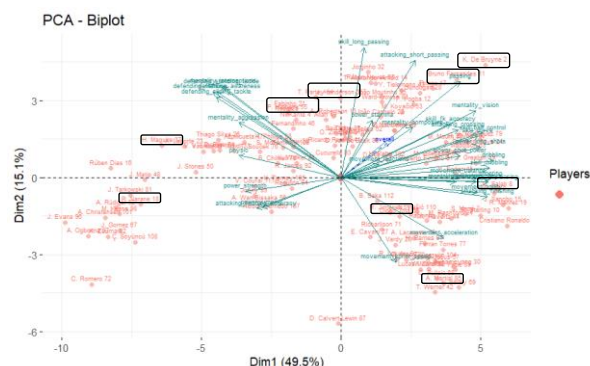


Figure 4(b) PCA Biplot

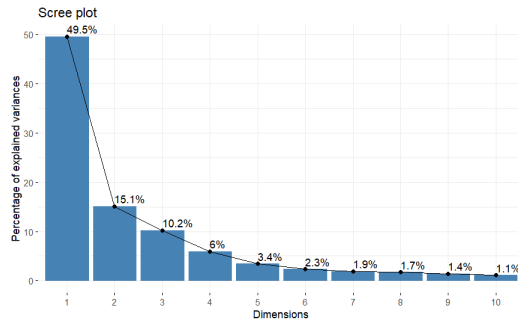


Figure 4(c) Scree Plot



Figure 4(d) PCA Coordinates

### k-means in Reduced Space

After the dimensionality reduction, k-means clustering is done by extracting the PCA coordinates of the reduced dimensions. Just like clustering done in the original space, the optimal number of clusters using the elbow method is identified (4), and clusters were visualized. Cluster 1 consists of all the defensive players, cluster 2 consists of strikers, cluster 3 consists of defensive midfielders, and attacking fullbacks and Cluster 4 consists of attacking midfielders.

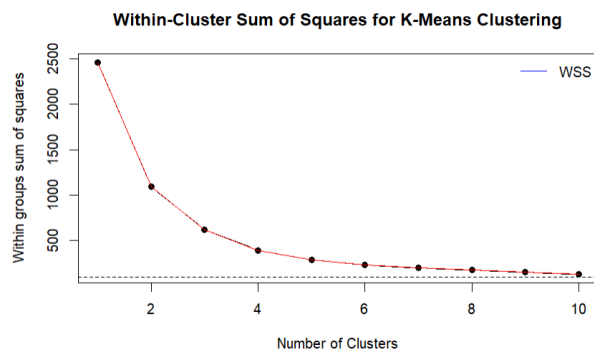


Figure 5(a) Optimal Number of Clusters

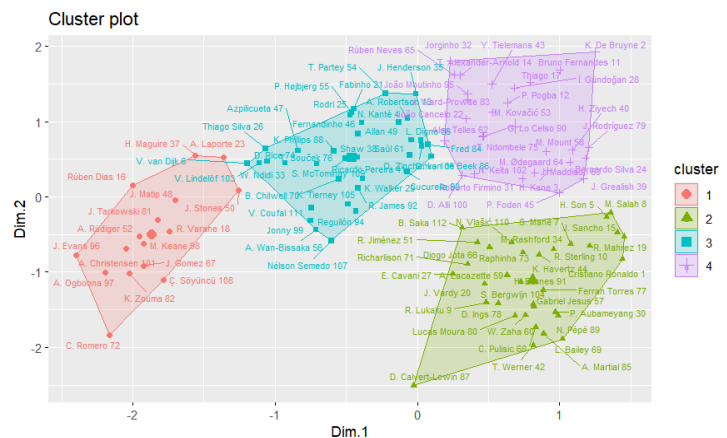


Figure 5(b) Clustering in Reduced Dimension

### Comparison between k-means clustering (original and reduced space)

To compare the performance of clustering models in original and reduced dimensions, inertia, clustering visualization, and the number of players on each cluster are taken into consideration for both models.

- **The inertia of k-mean clustering** is a measure of how well the data points in each cluster are grouped. Lower Inertia indicates better clustering. In this case, the inertia for the original dimension is much larger (**205447.8**) than the inertia for the reduced dimension after PCA (**394.6016**). This indicates that the clustering performance of the k-means algorithm is much better in the reduced dimension after PCA. This is because PCA has reduced the dimensionality of the data, removed redundant or noisy features, and focused on the most important ones.
- **Visualization:** It can be visualized from **Figure 3(b)** that there are four clusters but clusters 1,3 and 4 are **overlapped meaning** the clusters are **not well separated**. However, in **Figure 5(b)** after the dimension reduction, the clusters are **well separated** and there are no overlaps as well. This indicates that k-means after PCA is more accurate to identify the underlying structure of the data in the reduced dimension space and assign data points to their respective clusters. Also, PCA has been successful in capturing the most important and distinctive features of the data set that allow for a clear separation between clusters.
- **Cluster Numbers:**

k-means model	Defenders	Strikers	Attack. Midfield	Def. Midfield
Original Space	20	33	26	33
Reduced Space	35	33	28	16

After performing the dimensionality reduction, the clustering of players has significantly improved. In the original space, there were 33 defensive midfield players, but this was reduced to 16 after dimensionality

reduction. Additionally, the number of defenders increased from 20 to 35 as previously, the left back and full-backs were clustered as defensive midfield players, despite having more features of defensive players. This demonstrates the significance of dimensionality reduction in data analysis and machine learning.

### Linear Regression

A simple linear regression multimodal is taken into consideration for analyzing the relationship between a **response/dependent variable** ( $y = \text{value\_euro}$ ) and its interaction with one or more **independent variables** ( $x = \{\text{age, height, weight, overall, pace, shooting, passing, dribbling, defending, physic}\}$ ).

#### Residual vs Fitted Plot:

The residual vs fitted plot is a scatter plot of the residuals (the differences between the actual  $y$ -values and the predicted  $y$ -values) versus the fitted values (the predicted  $y$ -values). Looking at the plot, we see that the red line (which is just a scatterplot smoother, showing the average value of the residuals at each value of fitted value) is slightly curved but it has equally spread residual around the horizontal line without a distant pattern. This is good indication it is not a non-linear relationship. It indicates that the model's assumptions of linearity and equal variance are met, and the model fits the data well. We can also see the spread of the residual is more on the left side as compare to the right side. Finally, points 95, 44 and 1 seems to be far from the average data so they are outliers. These may indicate extreme observations that are not well explained by the model. We can remove these outliers or investigate them further.

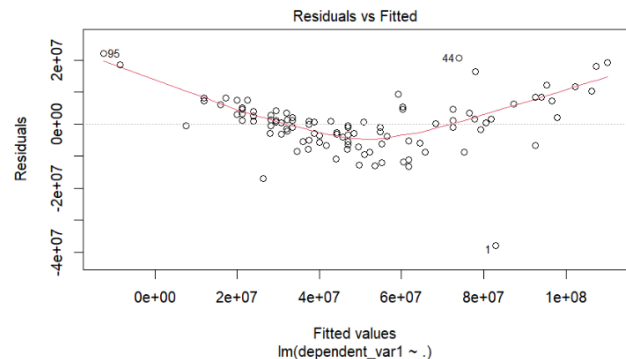


Figure 6(a) Residual vs Fitted Plot

#### Normal Q-Q Plot:

The Q-Q (quantile-quantile) plot is a graphical technique used to check the normality assumption of the linear regression model. It plots the standardized residuals (the residuals divided by their standard deviation) against the expected values of a normal distribution. If the residuals follow a straight line, the normality assumption is satisfied. For our model, the Q-Q plot shows pretty good alignment to the line with a few points at the top slightly offset. They are probably not significant so we have a reasonable alignment. It indicates that the normality assumption of the linear regression model is met. The model is appropriate for the data.

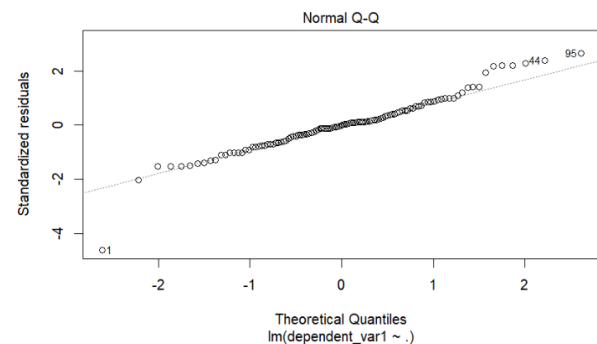


Figure 6(b) Normal Q-Q Plot

#### Scale-Location Plot

The scale-location plot is a scatter plot of the square root of the absolute standardized residuals (the absolute residuals divided by their standard deviation) versus the fitted values. This plot is used to check the assumption of constant variance (homoscedasticity) across the range of the predictor variable. Looking at the plot, we can see the residuals are reasonably well spread above and below a pretty horizontal line however the beginning of the line does have fewer points so slightly less variance there. It indicates that the model fits the data well and is appropriate for the data.

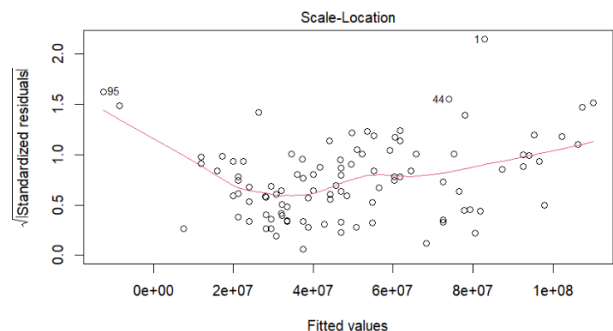


Figure 6(c) Scale-Location Plot

## Residuals vs Leverage:

The residuals vs leverage plot is a scatter plot of the standardized residuals against the leverage (a measure of how extreme a given value of the predictor variable is). This plot is used to check for influential points, which are points that have high leverage and/or high residuals. Influential points can have a large effect on the regression line, so they should be examined carefully. An influential case will appear in the top right or bottom left of the chart inside a dotted line which marks Cook's Distance. Looking at the plot, we can see point 1 is outside the dash line with high leverage and low residual. We need to investigate it further to determine if it is influential point that is distorting the regression line. We can remove this influential point or consider a different model that account for it.

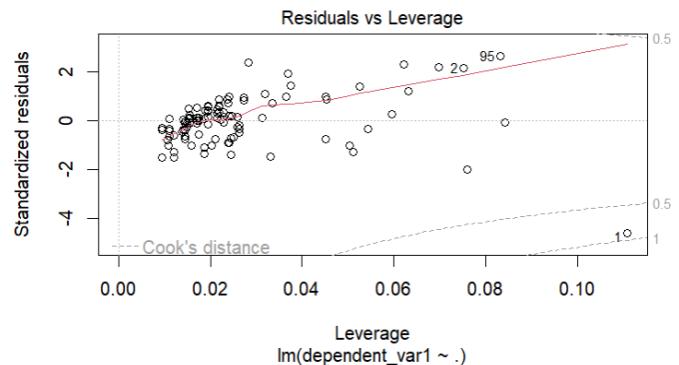


Figure 6(d) Residual vs Leverage Plot

## Linear Regression Summary:

Based on the summary of the linear regression model, it can be concluded that the market value (in euros) of players is significantly **related to age and overall rating**. This means that **younger players** and those with **higher overall ratings** tend to have higher market values. Additionally, the pace and defending ratings are also statistically significant, but to a lesser extent than age and overall rating. However, features such as weight, shooting, passing, dribbling, and physics do not appear to have a significant relationship with the market value of players in the Premier League. The residual standard error gives a difference between the observed value and the predicted values. In this case, **the residual standard error is 8.46 million euros**, which means that the model's prediction is off by around 8.46 million euros and it seems agreeable as the young players with a rating at around 90 have a market value of more than 100 million.

```
Call:
lm(formula = dependent_var ~ ., data = data_lm2)

Residuals:
    Min       1Q   Median       3Q      Max
-36062391 -4953981  1162101  4345130 20911816

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -571016647  45359407  -12.589  < 2e-16 ***
age          -4331004   279510  -15.495  < 2e-16 ***
height_cm    -93104    228751   -0.407  0.68486
weight_kg    -278217   178054   -1.563  0.12129
overall      9941471   396431   25.077  < 2e-16 ***
pace        -312591   113825   -2.746  0.00714 **
shooting     -104981   137293   -0.765  0.44626
passing       118025   238209    0.495  0.62135
dribbling    -259775   292282   -0.889  0.37623
defending    -235420   107907   -2.182  0.03145 *
physic        30799    202180    0.152  0.87923
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 8463000 on 101 degrees of freedom
Multiple R-squared:  0.9141,    Adjusted R-squared:  0.9056
F-statistic: 107.5 on 10 and 101 DF,  p-value: < 2.2e-16
```

Figure 6(e) Linear Model Summary

**Multiple R squared** value measures how well the independent variables in the model explain the variation in the dependent variable. In this case, the **multiple R-squared is 0.9141**, which means the independent variables in the **model explain around 91.41%**

of the variation in the dependent variables. The **adjusted R-squared value** penalizes the inclusion of unnecessary independent variables that do not improve the model's performance. The adjusted R-squared value is **0.9056**, which is slightly lower than the multiple R-squared, indicating that some of the independent variables do not add much to the model's performance.

**A high F-statistic and a low value** (less than 0.05) indicate the model's fit is statistically significant. In this case, the F-statistics is 107.5, and the p-value is less than 2.2e-16, which means the **overall fit of the model is highly significant**.

## Conclusion

All the objectives defined in the first section of this report have been achieved and the report is made more concise and to the point. Logistics regression is not considered here as it is optional and also, the report seemed to exceed the limit. The learning experience for this course has been great so far. Starting from the basics of inferential statistics and hypothesis testing, statistical modeling, and ending up with the machine learning frameworks: mathematical foundation and practical application through the project/lab sessions in R has been a good exploration curve. Even more, when applying mathematical concepts through a project on the selected dataset for analysis and logical reasoning has helped to clear doubts and enhanced the learning ability.