

UNIVERSITY OF CALCUTTA



A PROJECT REPORT FOR SEMESTER VI ON **COMPETITIVE COMPARISON AMONG THE EUROPEAN LEAGUES**

Submitted in partial fulfilment of the requirements for the award of the degree of
BACHELOR OF SCIENCE (HONS.) IN STATISTICS

Submitted by:

Name: SAPTAJIT DUTTA
CU Roll No.: 213146-21-0043
CU Reg. No.: 146-1111-0508-21



Under the guidance of
Dr. Tuhin Subhra Bhattacharya
Assistant Professor
Department of Statistics
MAULANA AZAD COLLEGE

Declaration

I, **Saptajit Dutta**, a student of Semester 6, entitled to the programme B. Sc. Statistics Hons. at Maulana Azad College, **University of Calcutta** with CU Regn. No. **146-1111-0508-21** and CU Roll No. **213146-21-0043**, solemnly declare that the project entitled "**Study of competition level in European Football**" is a genuine and original work undertaken by me under the supervision of **Dr. Tuhinshra Bhattacharya** as a part of my Bachelor's programme.

- All the data collected for this project is authentic and sourced from reliable and verifiable sources.
- The analysis, interpretations, and conclusions presented in this project are the results of my own research and statistical analysis.
- This project has not been submitted for any other academic degree or assessment at any institution.
- I understand the academic integrity requirements of my university, and I have adhered to all guidelines and regulations regarding plagiarism and research ethics.
- In case of any unintentional oversight or error, I take full responsibility for the same and will cooperate with the necessary revisions or corrections as advised by my faculty or supervisor.
- I am aware of the consequences of academic misconduct, and I affirm that this project is entirely free from any act of dishonesty or deception.

I am sincerely committed to the successful completion of this project and the pursuit of knowledge in the realm of spiritual intelligence assessment.

Signature:

Date:

Place: Kolkata, West Bengal

Certificate

This is to certify that the project paper titled ***Competitive comparison among the European leagues*** has been submitted by Saptajit Dutta in partial fulfilment of the requirements for the B. Sc. Statistics Hons. degree.

The research work presented in this project has been conducted under my guidance and supervision. I can confirm that the results and findings presented in this project paper are the outcome of the investigator's efforts and the research work conducted during the specified period.

Furthermore, to the best of my knowledge, the results and findings reported in this project have not been submitted for the award of any other degree or diploma in any academic institution.

I recommend this project paper for evaluation and fulfilment of the academic requirements for the B. Sc. Statistics Hons. degree.

Should you have any further queries or require additional information, please do not hesitate to contact me.

Date:

Dr. Tuhinsubhra Bhattacharya,
Assistant Professor,
Department of Statistics,
Maulana Azad College

Contents

1	Objective	1
2	Introduction	2
3	Data Visualization	3
3.1	Theoretical visualization:	3
3.2	Graphical visualization:	4
4	Methodology	7
4.1	Logic behind the construction of the index	7
4.2	Principal Component Analysis	7
4.3	Test for normality	8
4.3.1	Shapiro-Wilk test	8
4.4	Bartlett's test(Homogeneity of variance)	8
4.5	Analysis of variance (ANOVA)	9
4.6	The Kruskal–Wallis test	9
5	Calculations	11
5.1	Analysis with the help of core index	11
5.1.1	Checking the homogeneity of variance	13
5.1.2	Checking the normality	14
5.1.3	Analysis with the help of Kruskal-Wallis test	15
5.2	Analysis with the help of logarithmic index	15
5.2.1	Checking the homogeneity of variance	17
5.2.2	Checking the normality	18
5.2.3	Analysis with the help of ANOVA	19
6	Conclusions	20
6.1	Interpretation of the results	20
6.2	Broader implications	20
6.3	Conclusive interpretation	20
7	Future Scopes	21
8	Bibliography	22
8.1	External links	22
9	Acknowledgement	23
10	Appendix	24
10.1	Useful links	24
10.2	Datasets used	24

1 Objective

This project assists in visualizing the broader landscape of European football, illustrating how the quality of football has progressed with the support of UEFA and FIFA, the governing bodies of European football. While it is evident that European football is significantly more advanced economically compared to other regions, the key question remains: has the quality been sustained?

UEFA has implemented various initiatives to enhance the sport, including comprehensive training for managers and staff, youth development programs, substantial financial investments, and the introduction of increasingly intricate tactical strategies. However, the crucial inquiry is whether these efforts have successfully elevated the quality of play. It is clear that an improvement in play quality would necessitate a corresponding increase in the level of competition. Player skills have progressed worldwide, with a specific focus on Europe in this context.

Therefore, analyzing the league performances of the five most competitive European leagues can provide insights into whether the level of competition is evenly distributed across the continent.

2 Introduction

The football landscape is constantly evolving, with official tournaments worldwide becoming increasingly competitive. This trend is particularly evident in European football leagues, where the level of competition is exceptionally high. Our objective is to assess and compare the competitiveness of the five prominent European football leagues by analyzing the performance of the participating clubs.

To accomplish this, it is essential to collect the standings data from a specific season, such as the 2022-23 season, for the leagues under consideration. These leagues include the Premier League (England), La Liga (Spain), Serie A (Italy), Bundesliga (Germany), and Eredivisie (Netherlands), each comprising 18 or 20 teams. To evaluate the leagues, we need to consider the overall performance, which is influenced by the collective performances of the clubs involved. Therefore, it is crucial to quantify the performance of each club using an appropriate performance index derived from available data in league tables (e.g., total matches played, total points earned in the season, goals scored, goals conceded, total number of wins, draws, and losses). By applying relevant statistical tests and analytical methods, we can draw conclusions regarding the level of competition across these selected European football leagues.

3 Data Visualization

3.1 Theoretical visualization:

In this context, we will be creating indices(based on their performance) for clubs in the top five European leagues for the 2022-23 season. This will involve analyzing league tables that include metrics such as the number of matches played, victories, draws, defeats, goal differentials, and total points earned.

Teams	P	W	D	L	F	A	Pts
Manchester City	38	28	5	5	94	33	89
Arsenal	38	26	6	6	88	43	84
Manchester United	38	23	6	9	58	43	75
Newcastle United	38	19	14	5	68	33	71
Liverpool	38	19	10	9	75	47	67
Brighton And Hove Albion	38	18	8	12	72	53	62
Aston Villa	38	18	7	13	51	46	61
Tottenham Hotspur	38	18	6	14	70	63	60
Brentford	38	15	14	9	58	46	59
Fulham	38	15	7	16	55	53	52
Crystal Palace	38	11	12	15	40	49	45
Chelsea	38	11	11	16	38	47	44
Wolverhampton Wanderers	38	11	8	19	31	58	41
West Ham United	38	11	7	20	42	55	40
AFC Bournemouth	38	11	6	21	37	71	39
Nottingham Forest	38	9	11	18	38	68	38
Everton	38	8	12	18	34	57	36
Leicester City	38	9	7	22	51	68	34
Leeds United	38	7	10	21	48	78	31
Southampton	38	6	7	25	36	73	25

Table 1: Premiere League 2022/23

The table provided displays the standings of the 2022/23 Premiere League season, which featured twenty English clubs. In the table, **P**, **W**, **D**, **L**, **F**, and **A** represent the total number of matches played, total matches won, total matches drawn, total matches lost, total goals scored, and total goals conceded by each club, respectively. The **+/-** column indicates the goal difference, while **Pts** represents the total points earned by the club during the 2022/23 season.

Now based on the available data we have to construct suitable performance indices to make comparisons among the leagues' competitiveness. So we have constructed two different indices.

Core Performance Index(ζ):

We create our first performance indicator by convexly combining the goal ratio, wins/losses, and draws/losses.

$\zeta = A_1 \frac{W}{L} + A_2 \frac{D}{L} + (1 - A_1 - A_2) \frac{F}{A}$ [Symbols are defined above.]
Where all the coefficients belong to (0, 1).

Relative Performance Index(δ):

Our second performance index is constructed using the logarithm (base e) of the number of point-scoring games per defeat and the goal ratio(\mathbf{F}/\mathbf{A}).

$$\delta = B_1 \log_e \left(\frac{\frac{3}{4}W + \frac{1}{4}D}{L} \right) + (1 - B_1) \log_e \left(\frac{F}{A} \right) \quad [\text{Symbols are defined above.}]$$

Where all the coefficients belong to (0, 1).

The ratio of the points earned by winning and drawing is 3:1, so in this case, we used the convex combination of the total number of wins and draws with coefficients 3/4 and 1/4 respectively.

In addition to the Premier League and La Liga, we have included Serie A, Bundesliga, and Dutch Eredivisie in our comparison.

While there are twenty teams in the Bundesliga, La Liga, and Premier League, there are only eighteen in the Dutch Eredivisie. This is because each of their national football federations has certain financial rules. Since the two more clubs would be from lower levels, having 18 teams would, in theory, boost competition in the league. Being tiny in size means that it is more difficult for new clubs to advance to the upper divisions, which is a disadvantage.

We are computing the aforementioned indices for every team across all five leagues included in the comparison because we believe that the combined performances of the participating clubs accurately reflect the performance of the league. In order to carry out the action and ascertain the degree of competition among them, we will employ an appropriate statistical model based on the computed indices.

3.2 Graphical visualization:

Here in **Figure 1**, a dispersion plot, consisting the number of wins, draws, losses and total points scored of different teams from all the five leagues(which we are dealing with) is provided. The five different colors in each scatter plot are indicating the five different leagues. Some points really need to be observed.

Firstly, the number of wins has an almost perfect positive correlation with the total number of points scored, while the number of losses has an almost perfect negative correlation with it.

Secondly, the number of draws does not have any clear correlation with the rests(no. of wins, losses and total points scored).

Match statistics provide a comprehensive snapshot of the team's collective performance; however, it is imperative to concurrently emphasize the development of individual players. This scrutiny entails a meticulous assessment of both offensive and defensive units. Evaluating the offensive prowess involves scrutinizing metrics

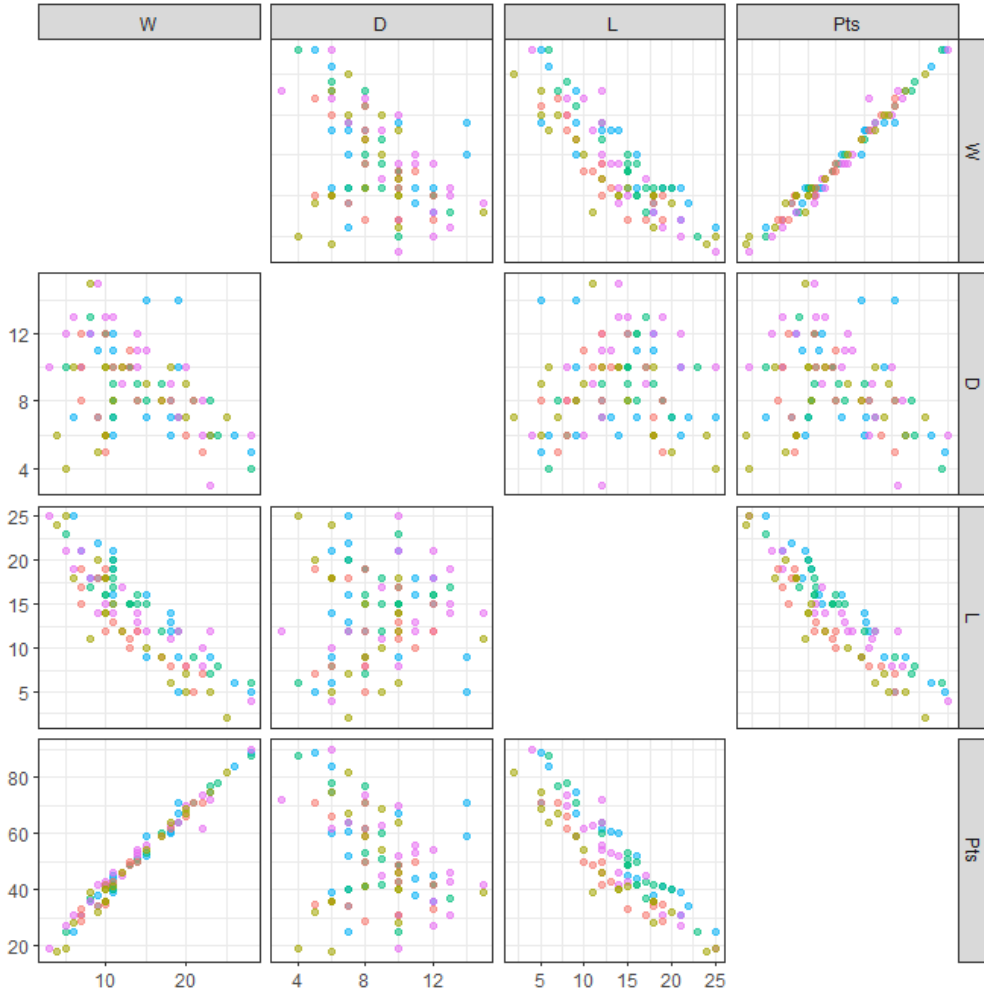


Figure 1: Dispersion plot

such as goal tally, scoring efficiency, and overall offensive efficacy. Conversely, the defensive aspect hinges on analyzing metrics like goals conceded and goal difference. A higher goal difference typically signifies superior structural integrity and cohesion within the team.

The **Figure 2** consists of three diagrams: the first exhibits a scatter plot where the Y-axis represents total points scored and the X-axis denotes goal ratio (top-left). The second diagram (top-right) mirrors the first but plots the number of goals scored against goals conceded on the X and Y axes, respectively. The third diagram (bottom-left) similarly features goal ratio and draw ratio plotted along their respective axes. Here, a palette of five distinct colors is employed to delineate points, each corresponding to a different league.

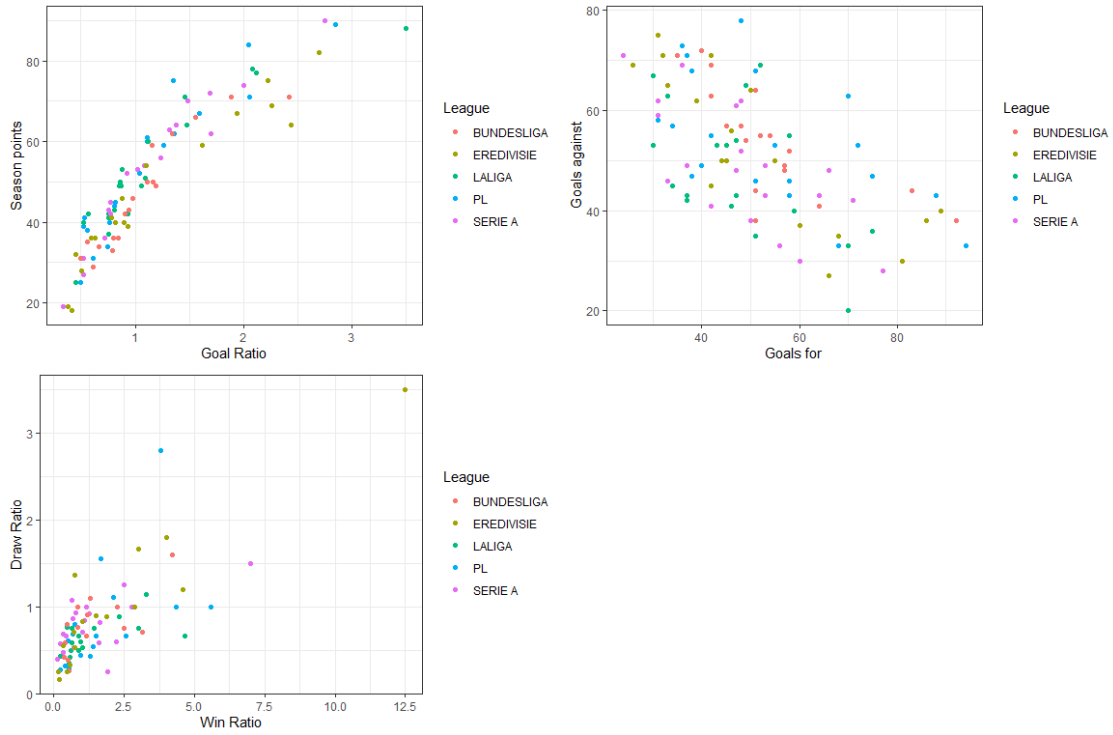


Figure 2: Scatterplots

From the first plot it is clear that there is positive correlation between the total points and the goal ratio indicates that the strength of the offensive unit of a team affects the performance of the team.

Also the goal ratio and the number of goals conceived are inversely proportional and so the team performance is also affected by the defensive unit's strength of the team. So it is obvious that the squad depth matters for the improvement of the team's game.

From the second plot, if we look carefully we can observe a negative correlation between the number of goals scored and the number of goals conceived by a team. It means that a good team should have the impact in the matches from both of their defensive and offensive units and so the teams in the top of the table have scored more goals and conceived less than the teams in the middle and the bottom of the table. The lack of perfect synchronization can make a negative impact in the team performance, whereas the balanced teams (teams with high quality players in both defensive and offensive unit) always dominate in their games.

4 Methodology

4.1 Logic behind the construction of the index

When contemplating comparisons among football leagues, the UEFA country coefficient system(*UEFA's country coefficient system ranks European football associations based on the performance of their clubs in UEFA competitions over the past five seasons.*) is the primary consideration. This metric proves beneficial when evaluating performance over a 5-6 season span, exclusively focusing on inter-European matches. While suitable for inter-European league comparisons, this approach prompts several inquiries.

Firstly, the challenge arises in comparing leagues within a single season, given that not all clubs participate in inter-European competitions. Consequently, evaluating a league based solely on the performances of select top-positioned clubs may be inadequate.

Secondly, even if this methodology is deemed appropriate, it's imperative to acknowledge the limited number of matches played in inter-European club competitions during a single season. Consequently, relying solely on these performances may not adequately capture the overall quality and consistency of the leagues.

Thirdly, the exclusion of certain European nations from UEFA club competitions does not inherently indicate a lower quality of their club football.

Given these considerations, the question arises: how do we effectively define the quality of European football in a single season? This necessitates the incorporation of various components into the assessment.

4.2 Principal Component Analysis

Suppose (X_1, X_2, \dots, X_p) is a random vector with mean vector $\mu = (E(X_1), E(X_2), \dots, E(X_p))$ and dispersion matrix $\Sigma = ((\sigma_{ij}))$ where $\sigma_{ii} = \text{Var}(X_i)$ and $\sigma_{ij} = \text{Cov}(X_i, X_j)$

Let us define,

$$Y_1 = a'_1 X = a_{11}X_1 + a_{12}X_2 + \dots + a_{1p}X_p$$

$$Y_2 = a'_2 X = a_{21}X_1 + a_{22}X_2 + \dots + a_{2p}X_p$$

.

.

.

$$Y_p = a'_p X = a_{p1}X_1 + a_{p2}X_2 + \dots + a_{pp}X_p$$

Now define $\Sigma_Y = A\Sigma A' = \text{diag}(\text{Var}(Y_1), \text{Var}(Y_2), \dots, \text{Var}(Y_p))$. Here we want to confine within Y_1, Y_2, \dots, Y_m for some $m < p$, where $\frac{\sum_{i=1,2,\dots,m} \text{Var}(Y_i)}{\sum_{i=1,2,\dots,p} \text{Var}(Y_i)}$ is large.

Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$ are p eigenvalues of A .

And (e_1, e_2, \dots, e_p) are corresponding eigenvectors. e_i, e_j are orthonormal.

So we can write $A = (e_1, e_2, \dots, e_p)$.
And $\Sigma_Y = A \Sigma A' = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$. So $\text{Var}(Y_i) = e_i' \Sigma e_i = \lambda_i (e_i' e_i) = \lambda_i$ and $\text{Cov}(Y_i, Y_j) = e_i' \Sigma e_j = \lambda_j (e_i' e_j) = 0$

4.3 Test for normality

To determine whether the population is normal, several tests are available. However, with a limited number of observations, we favor the Shapiro-Wilk test:

4.3.1 Shapiro-Wilk test

The Shapiro-Wilk test is a statistical test used to assess whether a sample of data comes from a normally distributed population. It's one of the most commonly used tests for normality and most preferable for small sample size. Suppose we have a sample x_1, x_2, \dots, x_n . We have to test:

H_0 : The given data are normally distributed. *against*,

H_1 : The given data are not normally distributed.

Test statistic: $W = \frac{(\sum_i a_i x_{(i)})^2}{\sum_i (x_i - \bar{x})^2}$,

$x_{(i)}$ is the i th order statistic in the sample.

\bar{x} is the sample mean.

The coefficients a_i are given by,

$(a_1, a_2, \dots, a_n) = \frac{m' V^{-1}}{(m' V^{-1} V^{-1} m)^{\frac{1}{2}}}$ and the vector $m = (m_1, m_2, \dots, m_n)'$ is made of

the expected values of the order statistics of independent and identically distributed random variables sampled from the standard normal distribution; finally, V is the covariance matrix of those normal order statistics.

We reject null hypothesis at α level of significance if $\text{obs}(W) < W_\alpha$, otherwise accept it.

4.4 Bartlett's test(Homogeneity of variance)

Bartlett's test (Snedecor and Cochran, 1983) is used to test if k samples have equal variances. Equal variances across samples is called homogeneity of variances. Some statistical tests, for example the analysis of variance, assume that variances are equal across groups or samples. The Bartlett test can be used to verify that assumption.

Bartlett's test is sensitive to departures from normality. That is, if your samples come from non-normal distributions, then Bartlett's test may simply be testing for non-normality. We have to test:

H_0 : $\sigma_1^2 = \sigma_2^2 = \sigma_3^2 = \dots = \sigma_k^2$ *against*,

H_1 : H_0 is not true.

Test statistic: $T = \frac{(N-k) \ln s_p^2 - \sum_{i=1}^k (N_i - 1) \ln s_i^2}{1 + (1/(3(k-1)))((\sum_{i=1}^k 1/(N_i - 1)) - 1/(N - k))}$,

In the above, s_i^2 is the variance of the i th group, N is the total sample size, N_i is the sample size of the i th group, k is the number of groups, and s_p^2 is the pooled

variance. The pooled variance is a weighted average of the group variances and is defined as:

$$s_p^2 = \sum_{i=1}^k (N_i - 1) s_i^2 / (N - k)$$

We reject null hypothesis at α level of significance if $obs(T) > \chi_{1-\alpha, k-1}^2$, otherwise accept it.

4.5 Analysis of variance (ANOVA)

Suppose we have populations (or groups or classes) P_1, P_2, \dots, P_k , which are assumed to follow normal distribution with the same variance σ^2 , but different means $\mu_1, \mu_2, \dots, \mu_k$, i.e. the populations are homoscedastic. Let y_{ij} be the j -th observation of the i -th sample of size n_i , which is drawn from the i -th population P_i . These k populations are the only available populations in which we are interested, and that's why a fixed effects model is appropriate here.

Model: $y_{ij} = \mu_i + e_{ij}$; $j=1, 2, \dots, n_i$; $i=1, 2, \dots, k$

where, μ_i =fixed effect due to i -th population and e_{ij} =random errors and are i.i.d $N(0, \sigma^2)$ for all (i, j)

Hypothesis of interest: $H_0 : \mu_1 = \mu_2 = \dots = \mu_k$ against $H_1 : H_0$ is not true

SSE=SS due to error= $\sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i0})^2$ with d.f $(n-k)$, where $\bar{y}_{i0} = \frac{1}{n_i} \sum_j (y_{ij})$
SSB=SS between populations/groups= $\sum_{i=1}^k n_i (\bar{y}_{i0} - \bar{y}_{00})^2$ with d.f $(k-1)$, where $\bar{y}_{00} = \frac{1}{n} \sum_{i=1}^k n_i \bar{y}_{i0}$

Test Procedure: under H_0 , $\frac{SSB}{\sigma^2}$ and $\frac{SSE}{\sigma^2}$ follow χ_{k-1}^2 and χ_{n-k}^2 respectively and they are independently distributed. Hence,

$$F = \frac{MSB}{MSE} \sim F_{k-1, n-k} \text{ under } H_0$$

H_0 is rejected at level α iff $F > F_{\alpha; k-1, n-k}$, where $F_{\alpha; k-1, n-k}$ is the upper- α point of F-distribution with d.f $(k-1, n-k)$.

4.6 The Kruskal–Wallis test

The Kruskal-Wallis test is a non-parametric statistical test used to compare three or more independent groups of samples. It is the non-parametric alternative to one-way ANOVA and is used when the assumptions of ANOVA (normality and homogeneity of variances) are not met.

Hypothesis of interest: H_0 : The medians of all groups are equal against $H_1 : H_0$ is not true

Test statistic:

$$H = \frac{12}{n(n+1)} \sum_{i=1}^k \frac{R_i^2}{n_i} - 3(n+1)$$

Here k is the number of groups

n_i is the number of observations in the i -th group and $\sum_i n_i = n$

R_i is the sum of ranks in the i -th group (Here the rank is the ascending rank in the combined dataset)

If there are ties in the data, an adjustment factor T is needed, where:

$$T = 1 - \frac{\sum_{j=1}^G (t_j^3 - t_j)}{n^3 - n}$$

where t_j is the number of tied ranks in the j -th group of ties, and G is the number of groups of tied ranks. So adjusted H will be:

$$H_{\text{corrected}} = \frac{H}{T}$$

Test procedure: Under H_0 the distribution of the test statistic H approximates a chi-square distribution with $k-1$ degrees of freedom.

Let p be the p -value associated with the computed H value. If $p \leq \alpha$ we reject the null hypothesis otherwise we accept it.

5 Calculations

5.1 Analysis with the help of core index

We possess three variables ($\mathbf{W/L}$, $\mathbf{D/L}$, and $\mathbf{F/A}$) that we may use to model the core index for every club. Now, using the PCA, a distinct core index should be created for each league we are comparing. To create covex coefficients, we first find the dispersion matrix. Then, we utilize the eigenvector that corresponds to its greatest eigenvalue.

Premiere League:

Here the dispersion matrix is $\begin{pmatrix} 2.11 & 0.487 & 0.866 \\ 0.487 & 0.318 & 0.210 \\ 0.866 & 0.210 & 0.374 \end{pmatrix}$ and its eigenvalues and corresponding eigenvectors are given below:

Eigenvalue	Eigenvector	Variability(%)
0.0156	(-0.403, -0.046, 1)	0.56
0.195	(4.880, -20.977, 1)	6.96
2.59	(2.414, 0.609, 1)	92.48

So eigenvector corresponding to the largest eigenvalue is (2.414, 0.609, 1). So our required core index is:

$$\zeta_P = 0.5997 \frac{W}{L} + 0.1515 \frac{D}{L} + 0.2488 \frac{F}{A}$$

LaLiga:

Here the dispersion matrix is $\begin{pmatrix} 1.25 & 0.104 & 0.761 \\ 0.104 & 0.0306 & 0.0533 \\ 0.761 & 0.0533 & 0.487 \end{pmatrix}$ and its eigenvalues and corresponding eigenvectors are given below:

Eigenvalue	Eigenvector	Variability(%)
0.0106	(-0.689, 0.909, 1)	0.6
0.0295	(-0.497, -1.48, 1)	1.67
1.73	(1.62, 0.131, 1)	97.73

So eigenvector corresponding to the largest eigenvalue is (1.62, 0.131, 1). So our required core index is:

$$\zeta_L = 0.5889 \frac{W}{L} + 0.0476 \frac{D}{L} + 0.3635 \frac{F}{A}$$

Serie A:

Here the dispersion matrix is $\begin{pmatrix} 2.20 & 0.270 & 0.796 \\ 0.270 & 0.0848 & 0.0873 \\ 0.796 & 0.0873 & 0.333 \end{pmatrix}$ and its eigenvalues and corresponding eigenvectors are given below:

So eigenvector corresponding to the largest eigenvalue is (2.714, 0.336, 1). So our required core index is:

Eigenvalue	Eigenvector	Variability(%)
0.0323	(-0.447, 0.635, 1)	1.24
0.0585	(-0.159, -1.687, 1)	2.24
2.52	(2.714, 0.336, 1)	96.52

$$\zeta_S = 0.6698\frac{W}{L} + 0.083\frac{D}{L} + 0.2472\frac{F}{A}$$

Bundesliga:

Here the dispersion matrix is $\begin{pmatrix} 1.10 & 0.237 & 0.485 \\ 0.237 & 0.104 & 0.112 \\ 0.485 & 0.112 & 0.226 \end{pmatrix}$ and its eigenvalues and corresponding eigenvectors are given below:

Eigenvalue	Eigenvector	Variability(%)
0.009	(-0.414, -0.147, 1)	0.63
0.0511	(-5.562, 22.535, 1)	3.57
1.37	(2.237, 0.508, 1)	95.8

So eigenvector corresponding to the largest eigenvalue is (2.237, 0.508, 1). So our required core index is:

$$\zeta_B = 0.5977\frac{W}{L} + 0.1355\frac{D}{L} + 0.2668\frac{F}{A}$$

Eredivisie:

Here the dispersion matrix is $\begin{pmatrix} 8.14 & 2.05 & 1.75 \\ 2.05 & 0.600 & 0.493 \\ 1.75 & 0.493 & 0.574 \end{pmatrix}$ and its eigenvalues and corresponding eigenvectors are given below:

Eigenvalue	Eigenvector	Variability(%)
0.0615	(0.463, -2.678, 1)	0.66
0.203	(-0.303, 0.321, 1)	2.18
9.05	(4.525, 1.156, 1)	97.16

So eigenvector corresponding to the largest eigenvalue is (4.525, 1.156, 1). So our required core index is:

$$\zeta_E = 0.6771\frac{W}{L} + 0.1734\frac{D}{L} + 0.1495\frac{F}{A}$$

The calculated league-wise indices for teams are shown below:

PL	La Liga	Serie A	Bundesliga	Eredivisie
4.21852303	4.052183334	5.4929	3.373076842	9.4743
3.259372093	2.559691667	2.41935	2.478553247	3.6553775
1.969157364	2.760417749	1.722419048	2.012343293	3.358862105
3.215738788	1.946082539	2.146175581	1.838398684	2.398428571
1.831387943	1.496354166	1.448833333	1.558678788	2.685744444
1.338542453	1.277804269	1.489208612	1.098009184	1.675532432
1.187774247	0.934512857	1.942850909	1.223644615	1.33616
1.112415873	0.975513953	1.218021705	1.146379545	0.95316
1.548871014	0.85186657	1.131812925	0.972565758	0.705816875
0.886688679	0.925440606	1.044783114	0.862224196	0.74205
0.764082041	0.850745409	0.95727033	0.87567963	0.730303571
0.717607447	0.858492963	0.845392857	0.601895906	0.868424242
0.543963521	0.703968036	0.70720675	0.589828472	0.528006989
0.572852727	0.72295	0.705127891	0.498459064	0.522403286
0.487070624	0.678406945	0.710894931	0.529411111	0.415425282
0.531468627	0.609517217	0.530361353	0.597958246	0.397933333
0.515940351	0.635007287	0.391905263	0.457338774	0.219497333
0.480136364	0.530959762	0.392675222	0.439657895	0.217993333
0.425150549	0.58747144	0.335878675		
0.30904389	0.311478585	0.197136563		

Table 2: ζ -index table

5.1.1 Checking the homogeneity of variance

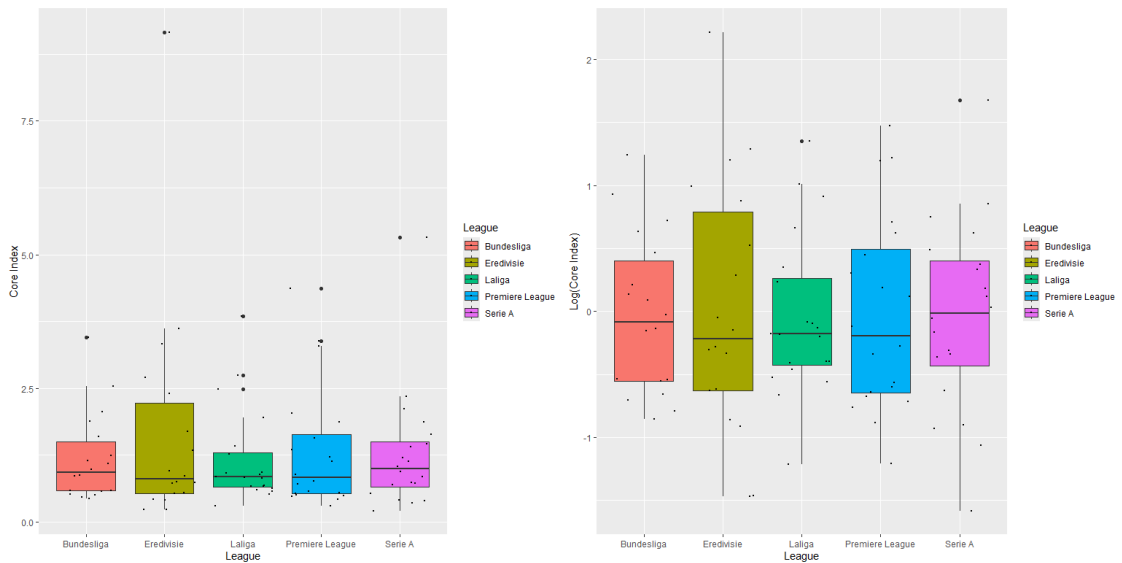


Figure 3: Boxplots

Here we have to test $H_0 : \sigma_P^2 = \sigma_L^2 = \sigma_S^2 = \sigma_B^2 = \sigma_E^2$ against $H_1 : H_0$ is not true
Now our observed test statistic $T=22.061$, which is greater than $\chi^2_{1-0.05,5-1}$.
So at 0.05 level of significance we reject the null hypothesis.

Now we are taking logarithmic transformation to check their homogeneity of variances.

The observed test statistic $T_1=5.2065$, which is less than $\chi^2_{1-0.05,5-1}$. So we conclude that the variances of the indices for the five different leagues are homogeneous.

5.1.2 Checking the normality

The purpose of this test is to determine whether the performance indices from the various leagues are normal or not.

Here the null hypothesis is that the population is normal.

The Shapiro-Wilk test is being used here, and it is suitable for a modest number of data. The test statistic value, p-value, and acceptance/rejection status(at 0.05 level of significance) are displayed in the table below.[Here we used R for the calculation].

League	Test statistic	p-value	Normality
PL	0.79144	0.0006458	No
LaLiga	0.75357	0.0001896	No
Serie A	0.73586	0.0001104	No
Bundesliga	0.83529	0.005007	No
Eredivisie	0.65051	0.0151	No

Here the decision is made on the basis of p-value. If the p-value is less than the level of significance(in this case 0.05) we reject the null hypothesis otherwise we accept it.

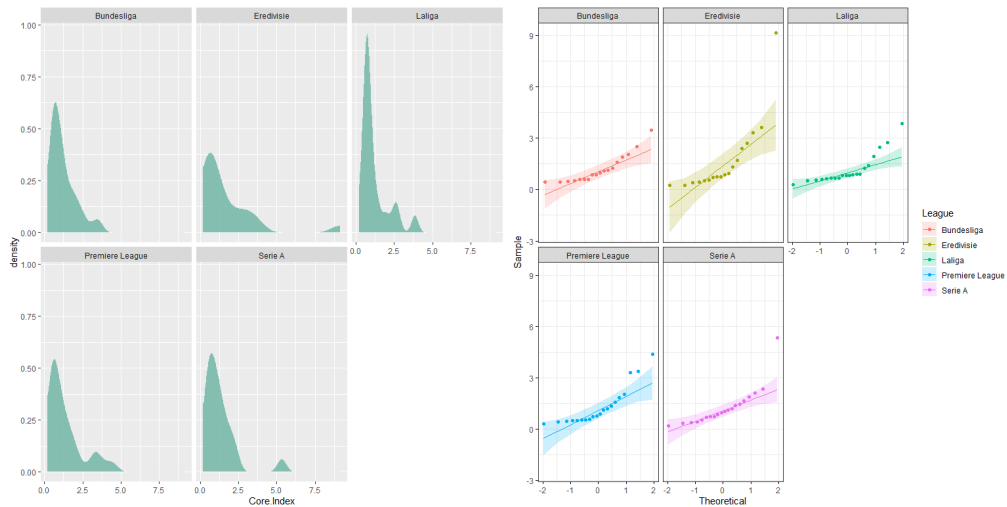


Figure 4: Density curves and Q-Q plots

5.1.3 Analysis with the help of Kruskal-Wallis test

Here we have to test H_0 : The median performance indices for all the leagues are equal against H_1 : H_0 is not true

Here the test statistic(Kruskal-Wallis chi-squared)=0.051704, df=4, p-value=0.9997
[Calculated using R-programming]

So we observed that the p-value>0.05 and so at 5% level of significance we accept the null hypothesis and conclude that the median performance indices for all the leagues under study are approximately equal or it can be clearly said that no significant difference exists among them.

5.2 Analysis with the help of logarithmic index

Here in this case, we possess two variables($\ln(0.75W/L+0.25D/L)$ and $\ln(F/A)$) that we may use to model the core index for every club. Now, using the PCA, a distinct core index should be created for each league we are comparing. To create covex coefficients, we first find the dispersion matrix. Then, we utilize the eigenvector that corresponds to its greatest eigenvalue.

Premiere League:

Here the dispersion matrix is $\begin{pmatrix} 0.651 & 0.387 \\ 0.387 & 0.250 \end{pmatrix}$ and its eigenvalues and corresponding eigenvectors are given below:

Eigenvalue	Eigenvector	Variability(%)
0.0144	(-0.608, 1)	1.6
0.886	(1.64, 1)	98.4

So eigenvector corresponding to the largest eigenvalue is (1.64, 1). So our required core index is:

$$\delta_P = 0.6212 \log_e \left(\frac{\frac{3}{4}W + \frac{1}{4}D}{L} \right) + 0.3788 \log_e \left(\frac{F}{A} \right)$$

LaLiga:

Here the dispersion matrix is $\begin{pmatrix} 0.397 & 0.293 \\ 0.293 & 0.237 \end{pmatrix}$ and its eigenvalues and corresponding eigenvectors are given below:

Eigenvalue	Eigenvector	Variability(%)
0.0137	(-0.764, 1)	2.16
0.621	(1.31, 1)	97.84

So eigenvector corresponding to the largest eigenvalue is (1.31, 1). So our required core index is:

$$\delta_L = 0.5671 \log_e \left(\frac{\frac{3}{4}W + \frac{1}{4}D}{L} \right) + 0.4329 \log_e \left(\frac{F}{A} \right)$$

Serie A:

Here the dispersion matrix is $\begin{pmatrix} 0.585 & 0.385 \\ 0.385 & 0.263 \end{pmatrix}$ and its eigenvalues and corresponding eigenvectors are given below:

Eigenvalue	Eigenvector	Variability(%)
0.00693	(-0.666, 1)	0.82
0.841	(1.5, 1)	99.18

So eigenvector corresponding to the largest eigenvalue is (1.5, 1). So our required core index is:

$$\delta_S = 0.6 \log_e \left(\frac{\frac{3}{4}W + \frac{1}{4}D}{L} \right) + 0.4 \log_e \left(\frac{F}{A} \right)$$

Bundesliga:

Here the dispersion matrix is $\begin{pmatrix} 0.430 & 0.254 \\ 0.254 & 0.164 \end{pmatrix}$ and its eigenvalues and corresponding eigenvectors are given below:

Eigenvalue	Eigenvector	Variability(%)
0.0106	(-0.606, 1)	1.78
0.584	(1.65, 1)	98.22

So eigenvector corresponding to the largest eigenvalue is (1.65, 1). So our required core index is:

$$\delta_B = 0.6226 \log_e \left(\frac{\frac{3}{4}W + \frac{1}{4}D}{L} \right) + 0.3774 \log_e \left(\frac{F}{A} \right)$$

Eredivisie:

Here the dispersion matrix is $\begin{pmatrix} 1.08 & 0.637 \\ 0.637 & 0.401 \end{pmatrix}$ and its eigenvalues and corresponding eigenvectors are given below:

Eigenvalue	Eigenvector	Variability(%)
0.0182	(-0.601, 1)	1.23
1.46	(1.66, 1)	98.77

So eigenvector corresponding to the largest eigenvalue is (1.66, 1). So our required core index is:

$$\delta_E = 0.6241 \log_e \left(\frac{\frac{3}{4}W + \frac{1}{4}D}{L} \right) + 0.3759 \log_e \left(\frac{F}{A} \right)$$

The calculated league-wise log indices for teams are shown below:

PL	La Liga	Serie A	Bundesliga	Eredivisie
1.323915024	1.279144469	1.440972934	1.12249957	1.825817358
1.049488938	0.823006001	0.780256386	0.818831241	1.125536444
0.56929486	0.899214344	0.453283168	0.618772247	1.079889869
1.060900306	0.548129309	0.628728791	0.522832672	0.794177439
0.562901676	0.331394547	0.299990736	0.363293747	0.948121676
0.275040361	0.176358221	0.325141794	0.082490342	0.490037027
0.138248435	-0.12522107	0.564209651	0.180140988	0.223122373
0.082768935	-0.053659118	0.176127127	0.131867138	-0.07461402
0.394690747	-0.205051516	0.102058468	-0.033422581	-0.330434818
-0.11495434	-0.091860267	0.021067938	-0.139489312	-0.249597341
-0.255582095	-0.185686688	-0.076481866	-0.119806244	-0.283935722
-0.313277845	-0.174953807	-0.12644754	-0.496409722	-0.101218348
-0.620682042	-0.350422073	-0.351997634	-0.517244916	-0.606850148
-0.532731588	-0.33553193	-0.312247622	-0.704585562	-0.629944512
-0.723506333	-0.360212877	-0.274982593	-0.620162795	-0.871428286
-0.617428352	-0.485231179	-0.548741842	-0.461426442	-0.844253472
-0.626305482	-0.4723154	-0.815306554	-0.756012091	-1.403344493
-0.69972044	-0.673008666	-0.855517852	-0.787191475	-1.376836684
-0.803140905	-0.467582571	-0.941222986		
-1.12895512	-1.08671491	-1.430289143		

Table 3: δ -index table

5.2.1 Checking the homogeneity of variance

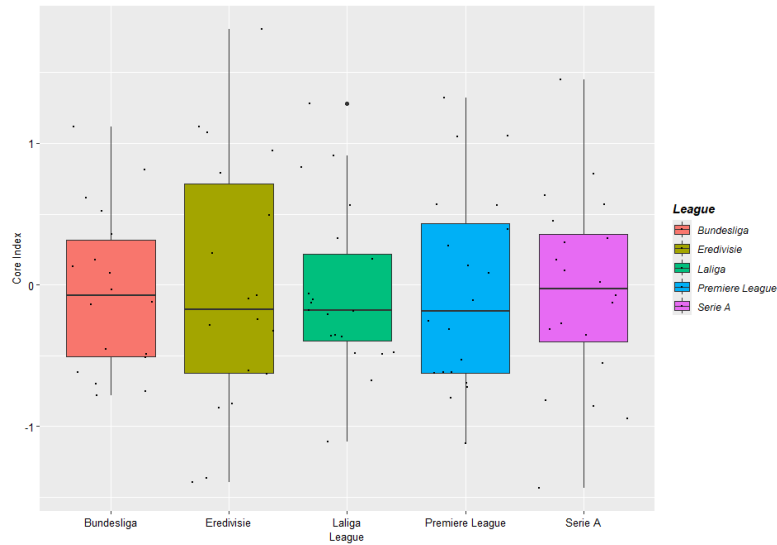


Figure 5: Boxplot

Here we have to test $H_0 : \sigma_{P_R}^2 = \sigma_{L_R}^2 = \sigma_{S_R}^2 = \sigma_{B_R}^2 = \sigma_{E_R}^2$ against $H_1 : H_0$ is not true

Now our observed test statistic $T=4.8588$, which is less than $\chi_{1-0.05,5-1}^2$.

So at 0.05 level of significance we accept the null hypothesis. So we conclude that the variances of the log indices for the five different leagues are homogeneous.

5.2.2 Checking the normality

The purpose of this test is to determine whether the log performance indices from the various leagues are normal or not.

Here the null hypothesis is that the population is normal.

The Shapiro-Wilk test is being used here, and it is suitable for a modest number of data. The test statistic value, p-value, and acceptance/rejection status(at 0.05 level of significance) are displayed in the table below.[Here we used R for the calculation].

League	Test statistic	p-value	Normality
PL	0.94114	0.252	Yes
LaLiga	0.93242	0.1719	Yes
Serie A	0.99145	0.9994	Yes
Bundesliga	0.94731	0.3844	Yes
Eredivisie	0.9649	0.6983	Yes

Here the decision is made on the basis of p-value. If the p-value is less than the level of significance(in this case 0.05) we reject the null hypothesis otherwise we accept it.

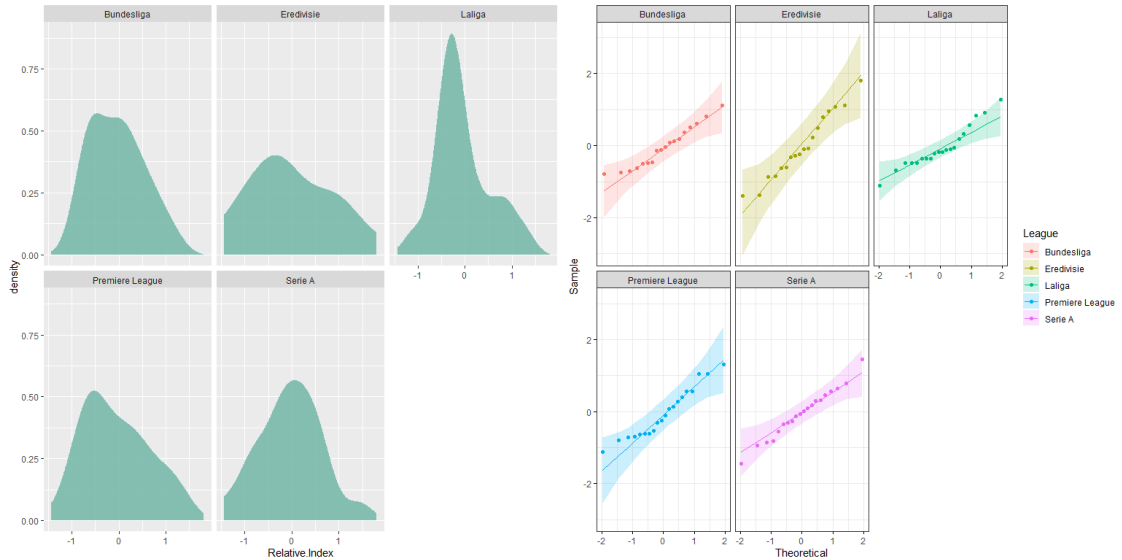


Figure 6: Density curves and Q-Q plots

5.2.3 Analysis with the help of ANOVA

The mean relative performance indices for all the leagues are equal against H_1 : H_0 is not true

Here the test statistic $F = 0.0087768$, num df = 4, denom df = 91, p-value = 0.9998
[*Calculated using R-programming*]

Here we observed that the p-value > 0.05 and so at 5% level of significance we accept the null hypothesis and conclude that there does not exist any significant difference among the mean relative indices of all the leagues.

6 Conclusions

6.1 Interpretation of the results

The comprehensive statistical analysis concluded that there is no significant difference in the competitive balance among the top European football leagues. This implies that, statistically, the leagues are equally competitive. The graphical visualizations provided intuitive insights into the performance dynamics, highlighting the relationships between wins, draws, losses, and points scored.

The use of PCA for constructing core indices helped in reducing the dimensionality of the data while preserving the essential competitive characteristics of each league. The normality and homogeneity of variance tests ensured that the data met the assumptions required for robust statistical testing, while the ANOVA and Kruskal-Wallis tests provided conclusive evidence of the competitive parity among the leagues. The project successfully demonstrated that the top European football leagues exhibit a comparable level of competition. The rigorous statistical approach and detailed visualizations provided a clear and comprehensive assessment of competitive balance, contributing valuable insights into the dynamics of European football.

6.2 Broader implications

The conclusion that top European football leagues have similar competitive balance is crucial. It indicates that, despite varying financial strengths and histories, the leagues are equally competitive, preserving the excitement and unpredictability that fans love. This balance highlights the success of governance and regulations in fostering fair competition, which is vital for the leagues' health and commercial success, maintaining fan interest, and ensuring a dynamic and unpredictable sporting environment.

6.3 Conclusive interpretation

In conclusion, the rigorous statistical approach and detailed visualizations employed in this project conclusively demonstrate that the Premier League, LaLiga, Serie A, Bundesliga, and Eredivisie are equally competitive. This equilibrium is vital for the overall health of the sport, as it sustains fan interest and ensures a dynamic, unpredictable competitive environment. The findings of this study contribute significantly to our understanding of competitive dynamics in elite football leagues, reaffirming the sport's intrinsic unpredictability and egalitarian spirit. Through meticulous analysis, this project highlights the robustness and fairness inherent in the competitive structures of Europe's top football leagues, supporting the notion that, despite various external differences, these leagues provide an even playing field in terms of competitive structure.

7 Future Scopes

This project can be integrated into the wider aspects of the sport, helping to evaluate various categories, stabilize the system, and make it more accessible to a larger group of people.

- **Comparative Studies with Other Continents:** Extending the comparative analysis to include football leagues from other continents (e.g., South America, Asia) would help understand the global football landscape and identify areas where European leagues excel or lag behind.
- **Inclusion of Additional Metrics:** Future research could incorporate more diverse performance metrics such as player statistics, financial data, and team investments, offering a more nuanced understanding of factors influencing league competitiveness.
- **Temporal Analysis:** Conducting a longitudinal study over multiple seasons would help identify trends and patterns in league competitiveness, assessing how changes in management, investment, and player development impact long-term performance.
- **Player Development Programs:** Examining the impact of youth academies and player development programs on league performance could provide actionable insights for clubs looking to enhance their competitive edge.
- **Fan Engagement and Influence:** Analyzing how fan engagement and support influence team performance and league competitiveness could offer valuable information for clubs and league organizers to improve fan experiences and loyalty.
- **Technological Integration:** - Exploring the role of emerging technologies such as VAR (Video Assistant Referee), wearable technology for player monitoring, and data analytics in shaping the future of football competition.
- **Policy and Regulation Impact:** Assessing how changes in UEFA and FIFA regulations impact the competitive balance within and between leagues, providing recommendations for policy adjustments to promote fair play and competitiveness.

By addressing these future scope points, the research can significantly contribute to the understanding and enhancement of competitive balance in European football, benefiting stakeholders across the sport.

8 Bibliography

- *The Elements of Statistical Learning*, Trevor Hastie, Robert Tibshirani, J. Friedman; *Springer*
- *Nonparametric Statistical Inference*, Jean Dickinson Gibbons, Subhabrata Chakraborti; *CRC Press*
- *Design and Analysis of Experiments*, Douglas C. Montgomery ; *Wiley*
- *Fundamentals of Applied Statistics*, Gupta SC, Kapoor V.K.; *Sultan Chand & Sons*

8.1 External links

- **Wikipedia:** <https://www.wikipedia.org/>
- **Comprehensive R Archive Network:** <https://cran.r-project.org/>
- **Posit:** <https://posit.co/>
- **Overleaf:** <https://www.overleaf.com/>

9 Acknowledgement

I would like to express my heartfelt gratitude and extend my sincerest appreciation to the individuals who have played a significant role in the successful completion of my project. Their guidance, support, and expertise have been invaluable throughout my journey. First and foremost, I am deeply indebted to my supervisor, **Prof. Tuhin Subhra Bhattacharya**, for his unwavering dedication, insightful feedback, and continuous encouragement. His profound knowledge in the field of statistics and his willingness to share his expertise have been instrumental in shaping the direction of my project. I am truly grateful for his mentorship and guidance, which have been invaluable assets in this endeavour. I would also like to extend my sincere thanks to the Head of the Statistics Department, **Prof. Partha Pal**, for his constant support and encouragement. His vision and leadership have provided me with a conducive environment to pursue my research interests. I am grateful for his valuable insights and the opportunities he has provided for my intellectual growth.

Additionally, I would like to express my gratitude to the other esteemed professors of the Statistics Department, **Prof. Anup Kumar Giri**, and **Prof. Nilanjan Chakraborty**. Their expertise and willingness to share their knowledge have been of immense help in refining my research ideas and broadening my understanding of statistical concepts. Their constructive criticism and valuable suggestions have played a pivotal role in shaping my project. I would also like to acknowledge my fellow students and friends who have supported me throughout this journey. Their camaraderie and encouragement have been a constant source of motivation. Lastly, I would like to thank **Maulana Azad College** for providing me with a nurturing academic environment and the necessary resources to pursue my project. I am grateful for the opportunities and platform that the college has offered me to explore my research interests and enhance my statistical knowledge.

10 Appendix

Here in this project we used league tables of 2022/23 season for all the leagues understudy. For calculation purposes, these tables are downloaded as Comma-Separated Values(.csv) format from an external website but are then verified from the league tables uploaded in the official websites of the respective leagues along with the one of the most trusted football network Skysports provided below:

10.1 Useful links

External Source: <https://www.footballwebpages.co.uk/>

Premiere League: <https://www.premierleague.com/>

Laliga: <https://www.laliga.com/en-GB/laliga-easports>

Serie A: <https://www.legaseriea.it/en>

Bundesliga: <https://www.bundesliga.com/en/bundesliga>

Eredivisie: <https://eredivisie.eu/home/>

Skysports: <https://www.skysports.com/football/tables>

10.2 Datasets used

The datasets(league tables) are provided below:

Teams	P	W	D	L	F	A	Pts
Manchester City	38	28	5	5	94	33	89
Arsenal	38	26	6	6	88	43	84
Manchester United	38	23	6	9	58	43	75
Newcastle United	38	19	14	5	68	33	71
Liverpool	38	19	10	9	75	47	67
Brighton And Hove Albion	38	18	8	12	72	53	62
Aston Villa	38	18	7	13	51	46	61
Tottenham Hotspur	38	18	6	14	70	63	60
Brentford	38	15	14	9	58	46	59
Fulham	38	15	7	16	55	53	52
Crystal Palace	38	11	12	15	40	49	45
Chelsea	38	11	11	16	38	47	44
Wolverhampton Wanderers	38	11	8	19	31	58	41
West Ham United	38	11	7	20	42	55	40
AFC Bournemouth	38	11	6	21	37	71	39
Nottingham Forest	38	9	11	18	38	68	38
Everton	38	8	12	18	34	57	36
Leicester City	38	9	7	22	51	68	34
Leeds United	38	7	10	21	48	78	31
Southampton	38	6	7	25	36	73	25

Table 4: Premiere League 2022/23

Teams	P	W	D	L	F	A	Pts
Barcelona	38	28	4	6	70	20	88
Real Madrid	38	24	6	8	75	36	78
Atlético Madrid	38	23	8	7	70	33	77
Real Sociedad	38	21	8	9	51	35	71
Villarreal	38	19	7	12	59	40	64
Real Betis	38	17	9	12	46	41	60
Osasuna	38	15	8	15	37	42	53
Athletic Bilbao	38	14	9	15	47	43	51
Mallorca	38	14	8	16	37	43	50
Girona	38	13	10	15	58	55	49
Rayo Vallecano	38	13	10	15	45	53	49
Sevilla	38	13	10	15	47	54	49
Celta Vigo	38	11	10	17	43	53	43
Valencia	38	11	9	18	42	45	42
Getafe	38	10	12	16	34	45	42
Cádiz	38	10	12	16	30	53	42
Almería	38	11	8	19	49	65	41
Real Valladolid	38	11	7	20	33	63	40
Espanyol	38	8	13	17	52	69	37
Elche	38	5	10	23	30	67	25

Table 5: LaLiga 2022/23

Teams	P	W	D	L	F	A	Pts
Napoli	38	28	6	4	77	28	90
Lazio	38	22	8	8	60	30	74
Inter Milan	38	23	3	12	71	42	72
AC Milan	38	20	10	8	64	43	70
Atalanta	38	19	7	12	66	48	64
Roma	38	18	9	11	50	38	63
Juventus	38	22	6	10	56	33	62
Fiorentina	38	15	11	12	53	43	56
Bologna	38	14	12	12	53	49	54
Torino	38	14	11	13	42	41	53
Monza	38	14	10	14	48	52	52
Udinese	38	11	13	14	47	48	46
Sassuolo	38	12	9	17	47	61	45
Empoli	38	10	13	15	37	49	43
Salernitana	38	9	15	14	48	62	42
Lecce	38	8	12	18	33	46	36
Spezia	38	6	13	19	31	62	31
Hellas Verona	38	7	10	21	31	59	31
Cremonese	38	5	12	21	36	69	27
Sampdoria	38	3	10	25	24	71	19

Table 6: Serie A 2022/23

Teams	P	W	D	L	F	A	Pts
Bayern Munich	34	21	8	5	92	38	71
Borussia Dortmund	34	22	5	7	83	44	71
RB Leipzig	34	20	6	8	64	41	66
1 FC Union Berlin	34	18	8	8	51	38	62
SC Freiburg	34	17	8	9	51	44	59
Bayer 04 Leverkusen	34	14	8	12	57	49	50
Eintracht Frankfurt	34	13	11	10	58	52	50
VfL Wolfsburg	34	13	10	11	57	48	49
Mainz 05	34	12	10	12	54	55	46
Borussia Mönchengladbach	34	11	10	13	52	55	43
1 FC Köln	34	10	12	12	49	54	42
1899 Hoffenheim	34	10	6	18	48	57	36
Werder Bremen	34	10	6	18	51	64	36
VfL Bochum	34	10	5	19	40	72	35
FC Augsburg	34	9	7	18	42	63	34
VfB Stuttgart	34	7	12	15	45	57	33
FC Schalke 04	34	7	10	17	35	71	31
Hertha Berlin	34	7	8	19	42	69	29

Table 7: Bundesliga 2022/23

Teams	P	W	D	L	F	A	Pts
Feyenoord	34	25	7	2	81	30	82
PSV Eindhoven	34	23	6	5	89	40	75
Ajax	34	20	9	5	86	38	69
AZ Alkmaar	34	20	7	7	68	35	67
FC Twente	34	18	10	6	66	27	64
Sparta Rotterdam	34	17	8	9	60	37	59
FC Utrecht	34	15	9	10	55	50	54
sc Heerenveen	34	12	10	12	44	50	46
RKC Waalwijk	34	11	8	15	50	64	41
Vitesse	34	10	10	14	45	50	40
Go Ahead Eagles	34	10	10	14	46	56	40
NEC Nijmegen	34	8	15	11	42	45	39
Fortuna Sittard	34	10	6	18	39	62	36
Volendam	34	10	6	18	42	71	36
Excelsior Rotterdam	34	9	5	20	32	71	32
FC Emmen	34	6	10	18	33	65	28
SC Cambuur	34	5	4	25	26	69	19
FC Groningen	34	4	6	24	31	75	18

Table 8: Eredivisie 2022/23