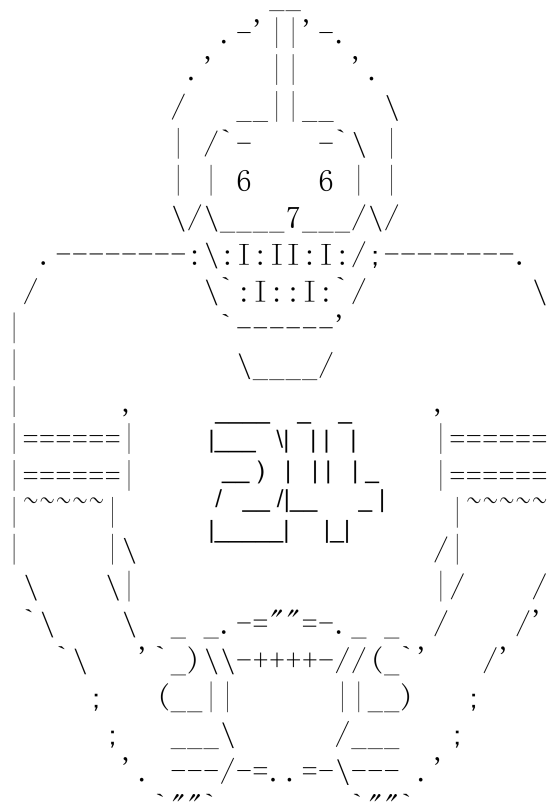


Introduction to Math for DS Group Mini-project

Analysis of factors affecting Premier League match results

IMDS Group 24

Zehao Qian, Chloe Mendez, Mohammad Jamshaid Iqbal



Introduction to Math for DS Group Mini-project

IMDS Group 24

Zehao Qian, Chloe Mendez, Mohammad Jamshaid Iqbal

December 7, 2023

Contents

1	Introduction	2
2	Model Assumptions	3
2.1	Performance Assumptions (Correlation Analysis)	3
2.2	Consistency Hypothesis (Principal Component Analysis)	3
2.3	Historical Performance Hypothesis (Entropy Weighting)	3
3	Data	4
4	Methods	5
4.1	Data Feature Extraction with Fourier Transform	5
4.2	Correlation Analysis	6
4.3	Principal Component Analysis	7
5	Conclusions	10
A	Appendix I: A breif proof for Fourier Transform	11
B	B	11

1 Introduction

The English Premier League is a ranking from one to twenty of the teams who won the most matches over the season. We are investigating the variables which factor into the ranking of the premier league; however, our dataset is fairly limited. We are examining historical data from previous seasons consisting of: the games won, lost and drawn by each team as well as the total number of goals for and against their team, and the goal difference. Ideally, we would analyse a dataset which featured variables not directly related to the ranking. In an attempt to enrich our dataset we performed the Fourier Transformation on the historical data from time domain to frequency domain. Time Domain feature extraction will allow us to examine the Time Domain evolution of the data. By analysing the time variability of metrics such as points, goals scored, goals conceded, etc., we hope to discover time-related patterns, such as seasonal changes or performance trends over a specific period. Of course, the most direct relationship is between the number of games won and the teams' placement on the league table, because this is how the table is curated. However, we did find some less obvious correlation; the teams' consistency in performance, as measured by a balanced distribution of wins, draws, and losses, is associated with a higher league position. To test this hypothesis, we have made use of Principal Component Analysis (PCA), to find the most significant components, and the Entropy method (Entropy Weighting), with the intention of finding the overall weighting of the team's performance. With the goal of establishing that teams with a consistent performance history over the seasons is likely to allow them to maintain their competitive positions. PCA allows us to establish which were the most significant factors in achieving first place on the premier league table. This is appropriate for our dataset because our dataset's multicollinearity is very high, so we are using PCA to eliminate this issue.

2 Model Assumptions

2.1 Performance Assumptions (Correlation Analysis)

- The number of wins positively correlates with the final league standing.
- Teams with a higher goal difference ($GF - GA$) tend to achieve higher league positions.
- Drawn matches have a minimal impact on final league standings.
- Teams with a higher number of goals scored (GF) are more likely to finish in the top positions.
- The defensive performance, measured by goals against (GA), influences the team's final standing.
- The number of points earned directly correlates with the team's final position in the league.

2.2 Consistency Hypothesis (Principal Component Analysis)

- Consistency in performance, as measured by a balanced distribution of wins, draws, and losses, is associated with a higher league position.

PCA can help identify patterns and relationships among these variables, which can contribute to understanding the consistency in team performance.

2.3 Historical Performance Hypothesis (Entropy Weighting)

- Teams with a consistent performance history over the years are likely to maintain their competitive positions.

3 Data

4 Methods

4.1 Data Feature Extraction with Fourier Transform

Given that the Premier League data we crawled from the web is relatively limited, we are eager to enrich our feature set through feature extraction. To this end, we decided to treat historical match data as an information-rich signal and adopt time-domain and frequency-domain feature extraction methods to extract more key features from it, laying the foundation for a comprehensive analysis of the team's performance.

Therefore, in our research, we introduced Fourier transform, a powerful mathematical tool, to transform historical game data from time domain to frequency domain. Fourier transform is a method of converting time domain signals into frequency domain representation. Through this conversion, we are expected to reveal the underlying periodicity and frequency information in the data, providing strong support for deeper analysis and understanding.

We will give the proof of the Fourier transform in the appendix section (A). In its continuous form:

$$X(f) = \int_{-\infty}^{\infty} x(t) \cdot e^{-j2\pi ft} dt$$

Among them, X represents the frequency domain signal, and x represents the time domain signal. In the discrete form (because our data is discrete), we will use discrete Fourier transform methods such as fast Fourier transform (FFT) to calculate:

$$X(k) = \sum_{n=0}^{N-1} x(n) \cdot e^{-j2\pi kn/N}$$

Time Domain feature extraction will allow us to delve into the Time Domain evolution of the data, revealing overall trends in the team over the past few years. By analyzing the time variability of metrics such as points, goals scored, goals conceded, etc., we can hopefully discover time-related patterns, such as seasonal changes or performance trends over a specific period.

At the same time, frequency domain feature extraction will help us understand the periodicity and frequency present in the data. This approach will help identify recurring seasonal patterns, allowing us to gain a more complete understanding of how a team's performance changes throughout the season. We will also explore whether specific events have a frequency-domain impact on team performance.

By converting historical match data into time and frequency domain features, we expect to be able to mine deeper information and provide a richer and more accurate feature set for our data-driven models to better interpret and predict Premier League matches. The team's performance. Here are the Time and Frequency Domain Features we will use:

Table 1: Time Domain Features and Formulas

Time Domain Feature	Formula for Time Domain Feature
Mean	$\text{mean} = \frac{1}{N} \sum_{i=1}^N x_i$
Standard Deviation	$\text{std_dev} = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \text{mean})^2}$
Root Mean Square (RMS)	$\text{rms} = \sqrt{\frac{1}{N} \sum_{i=1}^N x_i^2}$
Skewness	$\text{skewness} = \frac{\frac{1}{N} \sum_{i=1}^N (x_i - \text{mean})^3}{\left(\frac{1}{N} \sum_{i=1}^N (x_i - \text{mean})^2\right)^{3/2}}$
Kurtosis	$\text{kurtosis} = \frac{\frac{1}{N} \sum_{i=1}^N (x_i - \text{mean})^4}{\left(\frac{1}{N} \sum_{i=1}^N (x_i - \text{mean})^2\right)^2} - 3$
Max Value	$\text{max_value} = \max(\text{signal_array})$
Min Value	$\text{min_value} = \min(\text{signal_array})$
Median	$\text{median} = \text{np.median}(\text{signal_array})$
Zero Crossing Rate	$\text{zero_crossing_rate} = \frac{\sum_{i=1}^{N-1} (\text{sign}(x_{i+1}) - \text{sign}(x_i)) \neq 0}{N}$

Table 2: Frequency Domain Features and Formulas

Frequency Domain Feature	Formula for Frequency Domain Feature
Dominant Frequency	$\text{dominant_frequency} = \arg \max(\text{magnitude_spectrum})$
Max Frequency Magnitude	$\text{max_frequency_magnitude} = \max(\text{magnitude_spectrum})$
Power Spectral Density	$\text{power_spectral_density} = \frac{1}{N} \sum_{i=1}^N \text{magnitude_spectrum}^2$
Spectral Entropy	$\text{spectral_entropy} = \text{entropy}(\text{magnitude_spectrum})$
Total Power	$\text{total_power} = \sum_{i=1}^N x_i^2$
Centroid Frequency	$\text{centroid_frequency} = \frac{\sum_{i=1}^N i \cdot \text{magnitude_spectrum}[i]}{\sum_{i=1}^N \text{magnitude_spectrum}[i]}$

We apply these formulas to our data (historical matches). 'Magnitude_spectrum' is the magnitude of the spectrum calculated by Fourier transform.

4.2 Correlation Analysis

After feature extraction, we are preparing to conduct a thorough analysis of the relationship between various team indicators (MP, Win, Draw, Loss, GF, GA, GD and other features generated by last part) and the final score (Points) through correlation analysis. To choose an appropriate correlation analysis method, we will compare the characteristics of the Pearson correlation coefficient and the Spearman rank correlation coefficient, ultimately selecting the Spearman rank correlation coefficient for correlation analysis in English Premier League football.

After comparing the above characteristics, we have decided to choose the Spearman rank correlation coefficient as our correlation analysis method. This is because in the context of English Premier League football matches, our data may not follow a normal distribution,

Table 3: Characteristics of Pearson and Spearman Correlation Coefficients

Characteristics	Pearson Correlation Coefficient	Spearman Rank Correlation Coefficient
Calculation	$r_{xy} = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2 \cdot \sum (Y_i - \bar{Y})^2}}$	$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$
Data Type	Continuous variables	Ordinal variables and non-linear relationships
Linear Assumption	Assumes a linear relationship between variables	Does not make a specific linear assumption about the relationship
Applicability	Works well when data is approximately normally distributed	Applicable to ordinal variables, non-linear relationships, or when data distribution does not follow a normal distribution
Outlier Sensitivity	Sensitive to outliers	Relatively insensitive to outliers as it is based on ranks
Interpretation	Emphasizes the strength and direction of linear relationships	Focuses more on the ordinal relationship between variables

and the Spearman rank correlation coefficient is more robust to non-linear relationships and ordinal variables, while being relatively insensitive to outliers. This makes it more suitable for our research purposes.

In the report, we will use the Spearman rank correlation coefficient to explore the relationship between various team indicators and the final score, providing a more comprehensive understanding of the factors influencing different aspects of English Premier League football matches.

4.3 Principal Component Analysis

Principal Component Analysis (PCA) is a technique for data dimensionality reduction and feature extraction. It achieves this by identifying the principal directions (principal components) in the data to reduce its dimensionality. Below are the detailed calculation steps for PCA, introducing an example dataset:

Assume we have the following dataset:

Table 4: Example Dataset

MP	Win	Draw	Loss	GF	GA	GD	Others
38	27	6	5	80	22	58	...
38	25	10	3	65	26	39	...
38	24	11	3	74	31	43	...
38	21	13	4	67	28	39	...
38	19	8	11	55	33	22	...

Steps:

1. Standardize Data:

Standardize each feature, making its mean 0 and standard deviation 1. The standardization formula is:

$$Z = \frac{(X - \bar{X})}{\sigma}$$

where X is the original data, \bar{X} is the mean, and σ is the standard deviation. Applying this to the dataset yields the standardized data.

2. Compute Covariance Matrix:

The covariance matrix is the covariance matrix of the standardized data. The covariance matrix's formula is:

$$\text{Cov}(X, Y) = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{n - 1}$$

where X and Y are two features, \bar{X} and \bar{Y} are their means. After calculating the covariance matrix, we get:

3. Compute Eigenvalues and Eigenvectors:

Perform eigenvalue decomposition on the covariance matrix to obtain eigenvalues and their corresponding eigenvectors. Eigenvalues represent variance in the data, and eigenvectors are the directions of principal components.

4. Select Principal Components:

Based on the magnitude of eigenvalues, choose the number of principal components to retain. Typically, you might select components that capture a certain percentage of variance, such as 90%.

5. Build Projection Matrix:

Compose a matrix with the selected eigenvectors as columns. This matrix serves as the projection matrix, mapping the original data into the new principal component space.

6. Project to New Principal Component Space:

Multiply the standardized data by the projection matrix to obtain the reduced-dimensional data.

5 Conclusions

A Appendix I: A breif proof for Fourier Transform

Equation

B B