# Introduction to Statistic for Data Science
## Group Mini-Project Presentation: Happiness Ladder

ISDS Group 10

Christopher Barrow, Zehao Qian, Mengyuan Zhu, Hithein Augustine

Department of Natural Sciences
Durham University, England, UK

November 28, 2023

# Outline

# The Premier League

- Premier League: Top tier of English Football League System.
- 20 teams play 38 home and away matches.
- Globally renowned and challenging to predict outcomes.

**Background**

1. Outcome predictions involve expert analysis.
2. Factors include team performance, player form, and tactics.
3. Growing data, e.g., player touches, team running stats, manager experience.

# Outline

# Overview of Entropy Weight Method in Football

1. Introduction
   - The Entropy Weight Method is a powerful analytical technique used in football team evaluation. It goes beyond traditional methods by considering the inherent information entropy within various performance attributes.

2. Key Characteristics
   - Entropy: Reflects the degree of uncertainty or randomness within a dataset.
   - Weight Assignment: Assigns weights to attributes based on their information entropy.

3. Objective
   - The method aims to provide a nuanced evaluation, giving higher importance to attributes that contribute more to understanding a team's performance.

# Key Steps in Entropy Method

1. **Data Collection and Attribute Selection**
2. **Entropy Calculation:**
   - Utilize mathematical formulas to calculate the entropy of each selected attribute.
   - Entropy $= -\sum(p_i \cdot \log_2(p_i))$, where $p_i$ is the probability of each attribute value.
3. **Weight Assignment:**
   - Assign weights to attributes based on their calculated entropy.
   - Attributes with higher entropy receive lower weights, and vice versa.
   - The sum of weights equals 1 for normalization.
4. **Outcome:**
   - The result is a set of weights that reflect the relative importance of each attribute in evaluating a football team's performance.

# Entropy Method in Our Model

**Step 1: Get Data Set from FootyStats with web crawler**

| Football Stats | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Team | MP | Win | Draw | Loss | GF | GA | GD | Pts |
| MU | 38 | 25 | 10 | 3 | 80 | 22 | 58 | 87 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |

**Step 2: Attribute Selection**

**In order to reduce the complexity of data processing, the model input is simplified.**
**An Example:**

$$\left. \begin{array}{l} Loss = MP - Min - Draw \\ GA = GF - GD \end{array} \right\} \Rightarrow$$ They are negative and can be represented by other data

**Won't take these coloums in to consideration**

**Step 3: Normalization the Matrix**

Implement it with MinMaxScaler

$$x_{ij} = \frac{x_{ij} - x_{min}}{x_{max} - x_{min}}$$

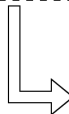**Step 4: Calulate the Information Entropy**

$$E_i = -\sum_{j=1}^{n} p_{ij} \log(p_{ij})$$

**Step 5: The entropy weight for criterion $i$**

$$W_i = \frac{1 - E_i}{n - \sum_{i=1}^{n} E_i}$$

**Step 6: Apply weights**

Multiply the weight by the value of the corresponding criterion to obtain a weighted sum.

**Team Rank for Reference**

Manchester City FC

Liverpool FC

...

# Outline

# Linear Transformation and Gradient Ascend

1. Linear transformation involves transforming input variables into a new space where they are more linear and easier to model

2. Our model will be taking in 5-dimension vector and transforming it to a 1-dimension vector

$$F_v : \mathbb{R}^5 \Rightarrow \mathbb{R}$$

3. The vector v basically consists of the following variables:
   - The number of games won by a team in $n^{th}$ week, $v_1$
   - The number of games lost by a team in $n^{th}$ week, $v_2$
   - The number of games drawn by a team in $n^{th}$ week, $v_3$
   - The goal difference (goals scored - goals conceded) of a team in $n^{th}$ week, $v_4$
   - The number of points scored by a team in $n^{th}$ week, $v_5$

# Linear Transformation and Gradient Ascend

- The linear transformation will transform the 5-dimensional vector into a one-dimensional vector which is the total number of points in the final week $F(v)$
- To calculate we will use the following formula:

$$F(v) = wv + b$$

- $w$ is a 5-dimensional vector, weights of each variable in vector v
- $b$ is the scalar bias - systematic error

## Using Gradient Ascend

- To calculate w and b, we will use gradient ascend on historical data.
- Normalize $v$ and $F(v)$ by centering and standardizing it.
- Use normalized data and take out predicted points using any initial value of weight.
- See the error by finding the difference between normalized $F(v)$ and predicted $F(v)$
- Find the gradient of the weight using the formula - transposed-v.error
- Find the gradient of scalar bias using the formula - sum of all elements in error vector
- Then to find the updated gradient the following formula is:

$$w = w + a * w_{gradient}, \ b = b + a * b_{gradient}$$

- Then denormalize w and b to find the actual values

# Outline

# Grabbing Data from Web



**Data Source: Footystats Website**

**Read and Analyse it!**

**Write Them in CSV files**

**Using Python Selenium Framework to automatically get data**

**Parse the data from HTML**

# Outline

# Evaluating the Teams's Tier with Entropy Method

**Attribute Weights:**
**MP: 0.2**
**Win: 0.16056986450138694**
**Draw: 0.15832737217198417**
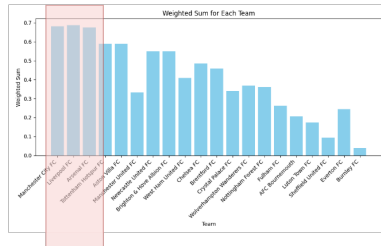**GF: 0.1613777266748812**
**GD: 0.1594934964983485**
**Points: 0.16023154015339913**

**BarChart of weight sum**



Weighted Sum for Each Team

| Table of Weight Sum | | |
|---|---|---|
| | Team | Weighted Sum |
| 0 | Manchester City FC | 0.681254 |
| 1 | Liverpool FC | 0.686918 |
| 2 | Arsenal FC | 0.676012 |
| 3 | Tottenham Hotspur FC | 0.588490 |
| 4 | Aston Villa FC | 0.588984 |
| 5 | Manchester United FC | 0.332012 |
| ... | ... | ... |
| ... | ... | ... |
| ... | ... | ... |

**After the entropy weight method calculation, the three teams in the first echelon after summing are also the current top scorers.**

# Outline

**Updated weights: [ 0.01791479 -0.01297361 -0.00399248  0.04477793  0.07853813]**
**Updated bias: 90.0**

## Prediction of 2023/24's final point:
## 93.66113213617407

| Increase the Accurency of Our Prediction | | |
|---|---|---|
| **Dataset** | | |
| More Variables | | |
| More Years | | |
| **Gradient Ascend Algorithm (without fix LR)** | | |
| Random Gradient Ascend | | |
| Newton's Method | | |
| Adagrad | | |
| ... | | |

Games Played

Games Won

Games Drawn

Goal Difference
(Goals Scored - Goals Conceded)

No. of Matches won (Last 5)

Points

# Outline

## Future plans

- Once our models are all complete, we will compare them with respect to efficiency and accuracy in predicting the winning team with this year's data as input. We could make use of Mallows' Cp and BIC statistic to establish which model has the best predictive performance.

- We will then critically analyse the most successful model and input a test dataset from earlier years to ensure we have not overfitted to our training dataset.

- To improve this model, it would be wise to consider variables other than the teams' performance and establish whether they are confounding and extraneous; for example, the about of money invested in each team.

## What is this model useful for

1. Once the most significant variables have been established, a researcher could potentially use the model to advise a team's manager of where best to focus their efforts, or indeed they could consult with a bookie to help them predict the odds of the game.

2. Of course, this model overlooks the human element of the game. Factors that a manager might be aware of such as player motivation, injuries, and team dynamics can have a significant impact on the outcome of a game, but they are challenging to factor into our model.

3. Ultimately though, it is impossible to control for every variable and, even if we did the result could surprise us. But this model is beneficial in finding which variables give the teams the best chance of success.