# Investigating the validity of FPL data in determining player performance and the most impactful players in the English Premier League teams

MSc Research Project
MSc Data Analytics

# Nishant Meena
Student ID: x21221839

School of Computing
National College of Ireland

Supervisor:      Vitor Horta

| | |
|---|---|
| **Student Name:** | …Nishant Meena…………………………………………………… |
| **Student ID:** | …x21221839………………………………………………………..… |
| **Programme** | ……MSc Data Analytics………… **Year:** …2023……….. |
| **Module:** | …………MSc Research Project……………………………….……… |
| **Supervisor:** | ……………Vitor Horta……………………………………..……… |
| **Submission Due Date:** | …………14-08-2023……………………………………….……… |
| **Project Title:** | Investigating the validity of FPL data in determining player performance and the most impactful players in the English Premier League teams. |
| **Word Count:** | ……1002………………… **Page Count**………………11………..…….. |

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

<u>ALL</u> internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

**Signature:** …………Nishant Meena…………………………………………

**Date:** …………………………13/08/2023………………………………………

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST**

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies) | □ |
| **Attach a Moodle submission receipt of the online project submission,** to each project (including multiple copies). | □ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | □ |

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

| **Office Use Only** | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

<center>

**Configuration Manual**

**Nishant Meena**
**Student ID:x21221839**
**MSc Research Project in Data Analytics**
**14th August 2023**

</center>

## 1. Introduction

This manual's goal is to highlight the project's technical side, which includes system requirements and programming snippets that are not mentioned in the main report. The essential system requirements used are covered at the outset of this document, along with a discussion of how the methodology is put into practice.

### 1.1 System Requirements

- Hardware spec
  1. System Manufacturer: Apple Inc.
  2. Operating System: macOS
  3. Processor: Apple M1 chip, 8-core CPU with 4 performance cores and 4 efficiency cores
  4. Memory: 8GB unified memory

- Software spec
  1. Jupyter notebook
  2. Google colab
  3. Microsoft Excel

## 2. Project Development

Data preparation is achieved in multiple stages, between Excel, jupyter notebook and google colab. The screenshots have been included in the manual to make everything clear.

### 2.1 Data Preparation

This project focuses on two datasets: The first dataset was obtained from a github repository[1] which has all the FPL data starting from 2016 season to the latest season. All the game week data has been uploaded on this github

---

[1] https://github.com/vaastav/Fantasy-Premier-League/tree/master

repository. The second dataset was scraped from the website called understat.com, using the documentation available online [2]. The first dataset contains the stats of the players for each game week from game week 1 of 2016 season to game week 38 of 2021 season as seen in fig 1.



*fig 1 FPL dataset*

For the second dataset which includes secondary stats (xG and xA) was scraped with the help of the understat documentation. However, to scrape the players stats for each season player id was required which was obtained by running the below mentioned code in figure 2. Further, it is important to know that in English Premier League (EPL), three clubs are relegated (demoted) to lower league and three clubs from the lower leagues are promoted to EPL each season. For the same reason an array of list of teams was manually created to get all the player ids, as seen in figure 2.

```python
import pandas as pd
import asyncio
from understat import Understat

async def fetch_team_data(session, team_name):
    understat = Understat(session)
    players = await understat.get_league_players(
        "epl",
        2018,
        team_title=team_name
    )
    return players

async def main():
    list_teams = ['Manchester United','Chelsea','Tottenham Hotspur','Manchester City','Liverpool','Arsenal','Everton','Southampton','AFC Bournemouth',' West Bromwich Albion',
                  'West Ham United','Leicester City','Stoke City','Crystal Palace','Swansea City','Burnley','Watford','Hull City','Middlesbrough','Sunderland','Newcastle United',
                  'Brighton & Hove Albion','Huddersfield Town','Wolves','Cardiff City','Sheffield United','Aston Villa','Norwich City','Leeds United','Leeds United','Brentford']  # 1
    list_df = []

    async with aiohttp.ClientSession() as session:
        tasks = [fetch_team_data(session, name) for name in list_teams]
        team_data = await asyncio.gather(*tasks)

    for team_name, data in zip(list_teams, team_data):
        for player in data:
            player["team"] = team_name
            player_name = player["player_name"].replace(" ", "_").lower()
            player["player_name"] = player_name
            list_df.append(player)

    df = pd.DataFrame(list_df)

    return df

# Run the main coroutine using the `await` keyword
result_df = await main()

# Now you can work with the result_df DataFrame
print(result_df)
```

*fig 2 Code to obtain the players ids*

The next step in scraping the data was to obtain the required stats for the all the players. The focus was to get the 'xG' and 'xA' stats, but the date and seasons were also the fetched as these columns are required to merge with the original FPL dataset, the snippets of the code can be seen in figure 3 below. After which the dataframe is saved to a csv file which is used to merge with the original dataset for second experiment. To merge the scrapped dataset, few changes were made to the name and date column of the FPL dataset to merge the datasets, figure 4. This code removes any numerical characters from the name of the players, in the second line of code, any spaces, dashes, are replaced with underscore, all the non-English characters are converted to their closest English equivalent and finally all the names were converted into lower case. Another change which was made was to extract the date from the kickoff time available in the FPL dataset, as the kickoff time includes date and time of the fixture. The date was extracted and was converted into YYYY-MM-DD to match it with the scraped stats from understat.com, the snippet can be seen below.

```
import json
import pandas as pd
from bs4 import BeautifulSoup
from urllib.request import urlopen
import unidecode
def scrape(id,name1):
  scrape_url = f"https://understat.com/player/{id}"
  page_connect = urlopen(scrape_url)
  page_html = BeautifulSoup(page_connect, "html.parser")
  page_html.findAll(name="script")
  json_raw_string = page_html.findAll(name="script")[4].string
  print(json_raw_string)
  start_ind = json_raw_string.index("\\")
  stop_ind = json_raw_string.index("')")
  data = json_raw_string[start_ind:stop_ind]
  data = data.encode("utf8").decode("unicode_escape")
  data = json.loads(data)
  df = pd.DataFrame(data)
  dk = df[['xG','xA','date','season']].sort_values(by='date', ascending=False)
  dk['team_id'] = id
  dk['name'] = name1
  return dk
```

```
combined_df = pd.DataFrame()
for index, row in result_df.iterrows():
    try :
      first_name = row['player_name']
      id_value = row['id']
      print(f"Name: {first_name} , ID: {id_value}")
      combined_df = pd.concat([combined_df, scrape(id_value,first_name)], ignore_index=True)
    except Exception as e:
        print(f"Exception occurred: {e}")
        continue
```

*fig 3 Code to scrape the players stats using player IDs*

## Importing the dataset

```
In [92]: data = pd.read_csv("/Users/nishant/Desktop/Semester 3/Research Project/FPL.csv")
```

```
In [93]: data['full_name'] = data.name.str.replace('_\d+','')
         data['full_name'] = data['full_name'].str.replace(" ", "_").str.replace("-", "_").str.replace('_\d+','')
         data['full_name'] = data['full_name'].apply(lambda x: unidecode.unidecode(x))
         data['full_name'] = data['full_name'].str.lower()
```

```
In [108]: data['date1'] = pd.to_datetime(data['kickoff_time']).dt.strftime('%Y-%m-%d')
```

```
In [109]: data.head()
```

Out[109]:

| nus | bps | clean_sheets | creativity | ... | total_points | transfers_balance | transfers_in | transfers_out | value | was_home | yellow_cards | GW | full_name | date1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0.0 | ... | 0 | 0 | 0 | 0 | 55 | False | 0 | 1 | aaron_cresswell | 2016-08-15 |
| 0 | 6 | 0 | 0.3 | ... | 1 | 0 | 0 | 0 | 60 | True | 0 | 1 | aaron_lennon | 2016-08-13 |
| 0 | 5 | 0 | 4.9 | ... | 2 | 0 | 0 | 0 | 80 | True | 0 | 1 | aaron_ramsey | 2016-08-14 |
| 0 | 0 | 0 | 0.0 | ... | 0 | 0 | 0 | 0 | 50 | False | 0 | 1 | abdoulaye_doucoure | 2016-08-13 |
| 0 | 3 | 0 | 1.3 | ... | 1 | 0 | 0 | 0 | 45 | True | 1 | 1 | adam_forshaw | 2016-08-13 |

*fig 4 Changes made to FPL dataset name and date column*

This csv file was saved to local hard disk, which was imported for the second experiment, after which it was merged with the original FPL dataset using the

merge function and was joined based on the player's name and the date on which the game was played between the clubs, the snippet of the code can be seen below in figure 5.



fig 5 Merging the secondary stats with FPL dataset.

The dataset contained some unnecessary columns which were removed like the column 'Unnamed: 0' as it same as the number of rows and the 'kickoff_time' as it was not required for our project. The code can be seen in figure 6.



fig 6 Cleaning of the dataset

Then the missing values were checked if there were any, it was found out that from season 2016 to 2018 the team names were not available in the dataset an the same were added with the help of Microsoft Excel as seen in figure 7.

*fig 7 Team names added in Excel.*

The season's name in the FPL dataset were mentioned as 2016-17, 2017—18 which was changed to 2016, 2017 respectively to increase the readability and to ease of use the code snippet is shown in figure 8.

```python
In [421]: #Changes made to the season to increase readability and to make use case easier
import pandas as pd

# Sample DataFrame
sample = {'season': ['2016-17', '2017-18', '2018-19', '2019-20', '2020-21']}
df = pd.DataFrame(data)

# Extract the year from the 'season' column
data['season'] = data['season_x'].str.split('-').str[0]

# Display the DataFrame

data = data.drop(['season_x'],axis = 1)
data
```

*fig 8 Transforming the season column to increase readability and use case.*

## 2.2    Data Training and Feature Selection

After doing the introductory exploratory data analysis (EDA), the dataset was split into training and testing dataset, where the season from 2016 to 2020 was for training the dataset and season 2021 as testing dataset as seen in figure 9.

```python
In [447]: train_data = data[data.season != '2021']
          test_data = data[data.season == '2021']
```

*fig 9 Training and testing dataset*

Before training the models, a feature importance technique was applied called Permutation feature importance, which gives the important features for building our models. The snippet of the code and the graph can be seen in the figure 10.

```python
In [436]: from sklearn.ensemble import RandomForestClassifier
          import matplotlib.pyplot as plt

          # Initialize and fit a Random Forest model
          forest_model = RandomForestClassifier(n_estimators=100, random_state=42)
          forest_model.fit(X_train, y_train)

          # Get feature importances
          feature_importances = forest_model.feature_importances_

          # Sort features by importance
          sorted_indices = feature_importances.argsort()[::-1]

          # Plot feature importances
          plt.figure(figsize=(10, 6))
          plt.bar(range(X_train.shape[1]), feature_importances[sorted_indices])
          plt.xticks(range(X_train.shape[1]), X_train.columns[sorted_indices], rotation=90)
          plt.title("Feature Importance")
          plt.show()
```

*fig 10 Permutation feature importance*

Based on the result of feature importance, the important features were included and rest of the features were removed before training of the model.

### 3. Modeling

Both the experiment involves building three models: Linear Regression, Random Forest and XGBoost. At first snippets of first experiment model building can be seen then after merging the xG and xA stats second experiment snippets are shown below:

## Linear Regression ( First Experiment)

After selecting the features based on feature importance, the models were build based only on these features for both the experiments in all the three models.

```
from sklearn.linear_model import LinearRegression

# Initialize the Linear Regression model
LR_model = LinearRegression()


# Select relevant features
features = ['minutes', 'goals_scored', 'assists', 'clean_sheets','GW','was_home','value','threat','team_h_score'
            , 'saves','penalties_saved','opponent_team','influence','ict_index','fixture','creativity','season']
X_trainLR = train_data[features]  # Features for training
y_trainLR = train_data['total_points']  # Target variable for training
X_testLR = test_data[features]  # Features for testing
y_testLR = test_data['total_points']  # Target variable for testing

# Train the benchmark model
LR_model.fit(X_trainLR, y_trainLR)

# Predict using the benchmark model
LR_prediction = LR_model.predict(X_testLR)


# Calculate Mean Squared Error for the benchmark model
LR_mse = mean_squared_error(y_testLR, LR_prediction)

# Calculate Root Mean Squared Error (RMSE)
LR_rmse = np.sqrt(LR_mse)
#R square value for Linear Regression
from sklearn.metrics import mean_squared_error, r2_score
r2LR = r2_score(y_testLR, LR_prediction)

# Calculate Mean Absolute Error
maeLR = mean_absolute_error(y_testLR, LR_prediction)
print(f"Mean Absolute Error:", maeLR)
print("Linear Regression Mean Squared Error:", LR_mse)
print("Linear Regression Root Mean Squared Error:", LR_rmse)
print("Linear Regression R-Square value:", r2LR)
```

*fig 11 Linear Regression (Exp 1)*

# Random Forest ( First Experiment)

```python
from sklearn.linear_model import LinearRegression

# Initialize the Linear Regression model
LR_model = LinearRegression()


# Select relevant features
features = ['minutes', 'goals_scored', 'assists', 'clean_sheets','GW','was_home','value','threat','team_h_score'
            , 'saves','penalties_saved','opponent_team','influence','ict_index','fixture','creativity','season']
X_trainLR = train_data[features]  # Features for training
y_trainLR = train_data['total_points']  # Target variable for training
X_testLR = test_data[features]  # Features for testing
y_testLR = test_data['total_points']  # Target variable for testing

# Train the benchmark model
LR_model.fit(X_trainLR, y_trainLR)

# Predict using the benchmark model
LR_prediction = LR_model.predict(X_testLR)


# Calculate Mean Squared Error for the benchmark model
LR_mse = mean_squared_error(y_testLR, LR_prediction)

# Calculate Root Mean Squared Error (RMSE)
LR_rmse = np.sqrt(LR_mse)
#R square value for Linear Regression
from sklearn.metrics import mean_squared_error, r2_score
r2LR = r2_score(y_testLR, LR_prediction)

# Calculate Mean Absolute Error
maeLR = mean_absolute_error(y_testLR, LR_prediction)
print(f"Mean Absolute Error:", maeLR)
print("Linear Regression Mean Squared Error:", LR_mse)
print("Linear Regression Root Mean Squared Error:", LR_rmse)
print("Linear Regression R-Square value:", r2LR)
```

*fig 12 Random Forest (Exp 1)*

# XGBoost ( First Experiment)

### XGBoost

```
In [465]: import pandas as pd
          from sklearn.model_selection import train_test_split
          from xgboost import XGBRegressor


          # Select relevant features
          features = ['minutes', 'goals_scored', 'assists', 'clean_sheets','GW','value','threat','team_h_score'
                      , 'saves','penalties_saved','opponent_team','influence','ict_index','fixture','creativity']
          X_train_XG = train_data[features]  # Features for training
          y_train_XG = train_data['total_points']  # Target variable for training
          X_test_XG = test_data[features]  # Features for testing
          y_test_XG = test_data['total_points']  # Target variable for testing
          # Convert categorical 'season' column to numerical using one-hot encoding
          # X = pd.get_dummies(X, columns=['season'], drop_first=True)

          # X = pd.get_dummies(X, columns=['was_home'], drop_first=True)


          # Create an XGBoost model
          model = XGBRegressor(n_estimators=100, learning_rate=0.1, max_depth=3, objective='reg:squarederror', random_state=42)
          # You can adjust parameters as needed

          # Train the model on the training data
          model.fit(X_train_XG, y_train_XG)

          # Predict on the test data
          y_predXG = model.predict(X_test_XG)
```

*fig 13 XGBoost (Exp 1)*

# Linear Regression (Second Experiment)

```
In [156]: from sklearn.linear_model import LinearRegression
          from sklearn.metrics import mean_squared_error
          from sklearn.metrics import mean_absolute_error
          # Initialize the Linear Regression model
          LR_model = LinearRegression()


          # Select relevant features
          features = ['minutes', 'goals_scored', 'assists', 'clean_sheets','GW','was_home','value','threat','team_h_score'
                      , 'saves','penalties_saved','opponent_team','influence','ict_index','fixture','creativity','season', 'xA',
                      'xG']
          X_trainLR = train_data[features]  # Features for training
          y_trainLR = train_data['total_points']  # Target variable for training
          X_testLR = test_data[features]  # Features for testing
          y_testLR = test_data['total_points']  # Target variable for testing

          # Train the benchmark model
          LR_model.fit(X_trainLR, y_trainLR)

          # Predict using the benchmark model
          LR_prediction = LR_model.predict(X_testLR)


          # Calculate Mean Squared Error for the benchmark model
          LR_mse = mean_squared_error(y_testLR, LR_prediction)

          # Calculate Root Mean Squared Error (RMSE)
          LR_rmse = np.sqrt(LR_mse)
          #R square value for Linear Regression
          from sklearn.metrics import mean_squared_error, r2_score
          r2LR = r2_score(y_testLR, LR_prediction)

          # Calculate Mean Absolute Error
          maeLR = mean_absolute_error(y_testLR, LR_prediction)
          print(f"Mean Absolute Error:", maeLR)
          print("Linear Regression Mean Squared Error:", LR_mse)
          print("Linear Regression Root Mean Squared Error:", LR_rmse)
          print("Linear Regression R-Square value:", r2LR)
```

*fig 14 Linear Regression (Exp 2)*

## Random Forest (Second Experiment)

```
: import pandas as pd
  from sklearn.ensemble import RandomForestRegressor  # For regression tasks



  # Select relevant features
  features = ['minutes', 'goals_scored', 'assists', 'clean_sheets','GW','was_home','value','threat','team_h_score'
            , 'saves','penalties_saved','opponent_team','influence','ict_index','fixture','creativity','season',
            'xG', 'xA']
  X_train_RF = train_data[features]  # Features for training
  y_train_RF = train_data['total_points']  # Target variable for training
  X_test_RF = test_data[features]  # Features for testing
  y_test_RF = test_data['total_points']  # Target variable for testing

  # Create a new RandomForest model
  model = RandomForestRegressor(n_estimators=100, random_state=42)  # You can adjust parameters as needed

  # Train the model on the training data
  model.fit(X_train_RF, y_train_RF)

  # Predict on the test data
  y_predRF = model.predict(X_test_RF)
```

*fig 15 Random Forest (Exp 2)*

## XGBoost (Second Experiment)

```
In [159]: import pandas as pd
          from sklearn.model_selection import train_test_split
          from xgboost import XGBRegressor


          # Select relevant features
          features = ['minutes', 'goals_scored', 'assists', 'clean_sheets','GW','value','threat','team_h_score'
                    , 'saves','penalties_saved','opponent_team','influence','ict_index','fixture','creativity', 'xG' , 'xA']
          X_train_XG = train_data[features]  # Features for training
          y_train_XG = train_data['total_points']  # Target variable for training
          X_test_XG = test_data[features]  # Features for testing
          y_test_XG = test_data['total_points']  # Target variable for testing
          # Convert categorical 'season' column to numerical using one-hot encoding
          # X = pd.get_dummies(X, columns=['season'], drop_first=True)

          # X = pd.get_dummies(X, columns=['was_home'], drop_first=True)


          # Create an XGBoost model
          model = XGBRegressor(n_estimators=100, learning_rate=0.1, max_depth=3, objective='reg:squarederror', random_state=42)
          # You can adjust parameters as needed

          # Train the model on the training data
          model.fit(X_train_XG, y_train_XG)

          # Predict on the test data
          y_predXG = model.predict(X_test_XG)
```

*fig 16 Random Forest (Exp 2)*

After building the model the best model was selected for prediction of points, which was Random Forest of first experiment. The below mentioned code snippet (figure 17) was used to show the best 11 players based on the predicted points. The user will be asked to enter the game week and the name of the English Premier League club for which they required the best performing players.

```python
# Prompt the user to input the football season
selected_GW = int(input("Enter the Game Week (between 1-38): "))

# Filter the DataFrame for the selected season (assuming you have a DataFrame named data_test)
selected_season_df = data_test[data_test['GW'] == selected_GW]

# Prompt the user to input the team name
selected_team = input("Enter the Team Name: ")

# Filter the DataFrame further for the selected team
selected_season_df = selected_season_df[selected_season_df['team'] == selected_team]

# Sort the DataFrame by predicted_total_points in descending order
selected_season_df = selected_season_df.sort_values(by='predictions', ascending=False)

# Initialize an empty dictionary to store selected players and their predicted points
selected_players = {}

# Loop through the sorted DataFrame to select unique players
for index, row in selected_season_df.iterrows():
    player_name = row['name']
    predicted_points = row['predictions']
    gw = selected_GW
    if player_name not in selected_players:
        selected_players[player_name] = {}
    selected_players[player_name][gw] = predicted_points
    if len(selected_players) == 11:  # Stop when 11 unique players are selected
        break

# Print the selected players along with their predicted points
print("Top 11 unique players from", selected_team, "with maximum predicted points for Game Week
for player, gw_points in selected_players.items():
    print(player)
    for gw, points in gw_points.items():
        print("- Predicted Points:", points, "- GW:", gw)
```

```
Enter the Game Week (between 1-38): 25
Enter the Team Name: Man City
Top 11 unique players from Man City with maximum predicted points for Game Week 25
Raheem Sterling
- Predicted Points: 18.75 - GW: 25
Ederson Santana de Moraes
- Predicted Points: 8.07 - GW: 25
Phil Foden
- Predicted Points: 8.01 - GW: 25
Kyle Walker
- Predicted Points: 6.89 - GW: 25
Rúben Santos Gato Alves Dias
- Predicted Points: 6.64 - GW: 25
Nathan Aké
- Predicted Points: 5.69 - GW: 25
Fernando Luiz Rosa
- Predicted Points: 4.04 - GW: 25
Ilkay Gündogan
- Predicted Points: 3.62 - GW: 25
Oleksandr Zinchenko
- Predicted Points: 3.34 - GW: 25
Bernardo Mota Veiga de Carvalho e Silva
- Predicted Points: 2.85 - GW: 25
Riyad Mahrez
```

*fig 17 code for best players in a team for a particular game week*