

Analysis of Python Code for Estimating Premier League Player Transfer Values (Assignment 1, Part IV)

May 5, 2025

Abstract

This report provides a detailed analysis of a Python script designed to fulfill Part IV of Assignment 1, which involves collecting player transfer values for the 2024-2025 Premier League season and proposing a method to estimate these values. The script processes statistical data from `results.csv`, matches players with hard-coded transfer values due to inaccessible web scraping, and uses an XGBoost regression model to predict market values. Utilizing libraries such as pandas, NumPy, scikit-learn, XGBoost, and fuzzywuzzy, the script ensures robust data handling, fuzzy matching, and model evaluation. This report explains the scripts functionality, justifies the feature selection and modeling approach, and discusses the results and limitations, aligning with the assignment requirements.

1 Introduction

The Python script addresses Part IV of Assignment 1, focusing on collecting transfer values for Premier League players with over 900 minutes of playtime in the 2024-2025 season and proposing a method to estimate these values. The input data, stored in `results.csv` from Part I, contains detailed player statistics from [FBref.com](https://fbref.com). Due to the inability to access [FootballTransfers.com](https://footballtransfers.com), the script uses hard-coded transfer values for 50 players. It employs fuzzy matching to align player names and trains an XGBoost regression model to predict market values based on selected features. The script generates predictions, evaluates model performance, and visualizes results.

This report dissects the scripts structure, detailing its libraries, data processing, fuzzy matching, feature selection, modeling approach, and output generation. It justifies the chosen features and model, discusses limitations, and suggests improvements, ensuring compliance with the assignments requirements.

1.1 Objectives

The script aims to:

- Collect transfer values for players with over 900 minutes, using hard-coded data as a fallback.
- Match player names between `results.csv` and transfer data using fuzzy matching.

- Propose a method for estimating player values with XGBoost regression.
- Justify feature selection based on statistical relevance to market value.
- Evaluate model performance and visualize actual vs. predicted values and feature importance.
- Save results to specified files and generate a report summary.

2 Libraries Utilized

The script relies on the following Python libraries:

- **pandas**: For data manipulation, including reading CSV files and merging datasets.
- **numpy**: Supports numerical operations and array handling.
- **fuzzywuzzy**: Performs fuzzy string matching to align player names.
- **xgboost**: Implements the XGBoost regression model for value prediction.
- **sklearn.model_selection**: Provides tools for train-test splitting, grid search, and cross-validation.
- **sklearn.preprocessing**: Scales features for modeling.
- **sklearn.metrics**: Calculates mean absolute error (MAE) and R² for model evaluation.
- **matplotlib.pyplot**: Generates plots for actual vs. predicted values and feature importance.
- **os**: Manages file paths and directory permissions.
- **unicodedata**: Normalizes player names by removing diacritics.
- **stat**: Sets directory permissions.

```

1 import pandas as pd
2 import numpy as np
3 from fuzzywuzzy import fuzz
4 import xgboost as xgb
5 from sklearn.model_selection import train_test_split,
   GridSearchCV, cross_val_score
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.metrics import mean_absolute_error, r2_score
8 import matplotlib.pyplot as plt
9 import os
10 import unicodedata
11 import stat

```

Listing 1: Library Imports

3 Initial Setup

The script defines paths for input and output files and includes version checking for debugging:

```
1 SCRIPT_VERSION = "2025-05-09-v9"
2 results_csv_path = r'C:\Users\nguye\Downloads\results.csv'
3 transfer_values_csv_path = r'C:\Users\nguye\Downloads\
  transfer_values.csv'
4 model_results_csv_path = r'C:\Users\nguye\Downloads\
  player_value_predictions.csv'
5 plots_dir = r'C:\Users\nguye\Downloads\plots'
```

Listing 2: File Path and Version Definitions

The `plots` directory is created with appropriate permissions.

The script loads `results.csv`, filters for players with over 900 minutes, and removes duplicates:

```
1 df = pd.read_csv(results_csv_path)
2 df['Min'] = pd.to_numeric(df['Min'], errors='coerce')
3 df = df[df['Min'] > 900]
4 df = df.drop_duplicates(subset=['Player'], keep='first')
```

Listing 3: Data Loading and Filtering

4 Name Normalization and Matching

Player names are normalized to handle variations (e.g., diacritics, aliases):

```
1 def normalize_name(name):
2     name = ''.join(c for c in unicodedata.normalize('NFD', name)
3         if unicodedata.category(c) != 'Mn')
4     parts = [part.strip('.').strip() for part in name.lower().
5         split()]
6     parts = [part for part in parts if part not in ['de', '
7         andrade']]
8     name = ' '.join(parts)
9     alias_map = {'rodri': 'rodrigo hernandez', 'j grealish': '
10         jack grealish', ...}
11     return alias_map.get(name, name)
```

Listing 4: Name Normalization

Fuzzy matching aligns hard-coded transfer values with players in `results.csv` using a two-pass approach (thresholds of 80 and 75):

```
1 for _, row in manual_df.iterrows():
2     normalized_name = row['Normalized_Player']
3     value = row['Transfer_Value']
4     best_match = None
5     best_score = 0
6     threshold = 80
```

```

7     for target_player, target_normalized in zip(players,
8         normalized_players):
9         score = fuzz.ratio(normalized_name, target_normalized)
10        if score > best_score and score >= threshold:
11            best_score = score
12            best_match = target_player
13        if not best_match:
14            threshold = 75
15        ...

```

Listing 5: Fuzzy Matching

Unmatched players are saved to `unmatched_players.txt` for debugging.

5 Transfer Value Collection

Due to the inability to access [FootballTransfers.com](https://www.footballtransfers.com), the script uses hard-coded transfer values for 50 players:

```

1 transfer_data = [
2     {"Player": "Aaron Wan-Bissaka", "Transfer_Value": 20},
3     {"Player": "Erling Haaland", "Transfer_Value": 150},
4     ...
5 ]
6 manual_df = pd.DataFrame(transfer_data)

```

Listing 6: Hard-Coded Transfer Data

The data is saved to `transfer_values.csv` after matching.

6 Player Value Estimation

6.1 Feature Selection

The script selects features relevant to market value:

- **Numeric Features:** Age, Min, Gls, Ast, xG, xAG, PrgC_standard, PrgP_standard, PrgR_standard, Tkl, Int, Blocks, Touches, Carries, Rec, Won%, Save%, CS%, Gls_Ast (derived).
- **Categorical Features:** Pos, Squad (one-hot encoded).

```

1 features = ['Age', 'Min', 'Gls', 'Ast', 'xG', 'xAG', ..., '
2   Gls_Ast']
3 categorical_features = ['Pos', 'Squad']

```

Listing 7: Feature Definition

Justification:

- **Age:** Younger players often have higher market values due to potential.
- **Min:** Indicates playing time and importance to the team.

- Gls, Ast, xG, xAG: Reflect offensive contributions and expected impact.
- PrgC, PrgP, PrgR: Capture progressive play, valued in modern football.
- Tkl, Int, Blocks: Indicate defensive contributions.
- Touches, Carries, Rec: Reflect involvement in play.
- Won%, Save%, CS%: Capture aerial and goalkeeping performance.
- Gls_Ast: Combines goal and assist contributions.
- Pos, Squad: Account for role-specific and team-specific value differences.

6.2 Model Selection

The script uses an XGBoost regression model due to its ability to handle non-linear relationships, feature interactions, and small datasets:

```

1 xgb_model = xgb.XGBRegressor(objective='reg:squarederror',
  random_state=42, early_stopping_rounds=10)
2 param_grid = {
3     'n_estimators': [100, 200, 300],
4     'max_depth': [3, 5, 7],
5     'learning_rate': [0.01, 0.05, 0.1],
6     'subsample': [0.8, 1.0]
7 }
8 grid_search = GridSearchCV(xgb_model, param_grid, cv=5, scoring='
  neg_mean_absolute_error')
```

Listing 8: XGBoost Model

Justification:

- XGBoost is robust to outliers and missing data, suitable for the small dataset.
- Grid search optimizes hyperparameters to improve performance.
- Early stopping prevents overfitting.
- Cross-validation ensures reliable performance estimates.

6.3 Model Training and Evaluation

The data is split into training (80%) and test (20%) sets, with features scaled:

```

1 X_train, X_test, y_train, y_test = train_test_split(X_scaled, y,
  test_size=0.2, random_state=42)
2 scaler = StandardScaler()
3 X_scaled = scaler.fit_transform(X)
```

Listing 9: Train-Test Split and Scaling

The model is evaluated using MAE, R², and cross-validation:

```

1 train_mae = mean_absolute_error(y_train, y_train_pred)
2 test_mae = mean_absolute_error(y_test, y_test_pred)
3 test_r2 = r2_score(y_test, y_test_pred)
4 cv_scores = cross_val_score(xgb_model_no_early, X_scaled, y, cv
    =5, scoring='neg_mean_absolute_error')

```

Listing 10: Model Evaluation

6.4 Visualization

The script generates two plots:

- **Actual vs. Predicted:** A scatter plot comparing actual and predicted market values (`actual_vs_predicted.png`).
- **Feature Importance:** A bar plot of the top 10 features (`feature_importance.png`).

```

1 plt.barh(importance_df['Feature'][:10], importance_df['Importance'][:10])
2 plt.title('Top 10 Features Influencing Market Value')
3 plt.savefig(os.path.join(plots_dir, 'feature_importance.png'))

```

Listing 11: Feature Importance Plot

7 Outputs

The script produces:

- **transfer_values.csv:** Player statistics with matched transfer values.
- **player_value_predictions.csv:** Player data with actual and predicted market values.
- **unmatched_players.txt:** List of unmatched players (e.g., Rodri, Jack Grealish).
- **actual_vs_predicted.png:** Scatter plot of actual vs. predicted values.
- **feature_importance.png:** Bar plot of top feature importances.
- **report_summary.md:** Summary of data collection, feature selection, model performance, and limitations.
- **Console Output:** Debugging information, model metrics, and top features.

8 Limitations

The script notes several limitations:

- **Small Dataset:** Only 50 players were matched (out of 307), limiting model reliability (low R^2 , high MAE).

- **Hard-Coded Data:** Reliance on manual transfer values due to inaccessible web scraping.
- **Unmatched Players:** Five players (e.g., Rodri, Jack Grealish) were unmatched, possibly due to low minutes or name mismatches.
- **Model Performance:** Low R^2 and high MAE indicate poor predictive power, likely due to the small dataset.

9 Conclusion and Potential Improvements

The script fulfills Part IV of Assignment 1 by collecting transfer values (via hard-coded data), matching players with fuzzy logic, and estimating market values using XGBoost. The feature selection captures key performance metrics, and the model provides insights despite limitations. The outputs align with the assignments requirements, with clear visualizations and a detailed report.

9.1 Potential Improvements

- **Web Scraping:** Implement robust scraping for [FootballTransfers.com](https://www.footballtransfers.com) to expand the dataset.
- **Enhanced Matching:** Use additional identifiers (e.g., team, position) to improve fuzzy matching accuracy.
- **Feature Engineering:** Add features like contract length, international caps, or recent form.
- **Model Alternatives:** Test other models (e.g., Random Forest, Neural Networks) for comparison.
- **Data Augmentation:** Incorporate historical transfer data to increase sample size.

10 References

- [FBref.com About Page](https://fbref.com)
- [XGBoost Documentation](https://xgboost.ai)
- [scikit-learn Model Evaluation Documentation](https://scikit-learn.org/stable/tutorial/model_evaluation.html)
- [Matplotlib Scatter Plot Documentation](https://matplotlib.org/3.1.1/tutorials/plotting/scatter.html)