

Premier League Match & Table Predictor

- [Aim of the Project](#)
- [Data Collection & Cleaning](#)
- [Feature Engineering](#)
 - [Correlation analysis](#)
 - [Most Predictive Features \(Good Correlation with Goals\)](#)
 - [Less Important Features \(Low or No Correlation\)](#)
- [Model Training \(TensorFlow\)](#)
 - [Input Preparation](#)
 - [Model architecture](#)
 - [Model Compilation & Training](#)
 - [Multiple Runs for Stability](#)
 - [Model Accuracy](#)
 - [Summary of Learning Behavior](#)
 - [Prediction Performance](#)
 - [Key Observations](#)
- [Example Findings](#)
 - [Example Predictions](#)

Aim of the Project [🔗](#)

[🔗](#) App link: [Streamlit App](#)

This project aims to:

1. Predict the number of **goals scored** by each team in a Premier League match.
2. Use the predicted scores to assign **match points**.
3. Aggregate all matches to generate the **final Premier League table** (20 teams ranked by points).

Data Collection & Cleaning [🔗](#)

Data was scraped from [fbref.com](#) using `pandas.read_html()`:

```
1 game_data = pd.read_html(  
2     'https://fbref.com/en/comps/9/schedule/Premier-League-Scores-and-Fixtures',  
3     attrs={'id': 'sched_2024-2025_9_1'})  
4 ][0]
```

Cleaned using:

```
1 game_data.columns = game_data.columns.str.lower().str.replace(' ', '', regex=False)  
2 game_data = game_data.dropna(how='all')  
3 game_data['wk'] = game_data['wk'].astype(int)  
4  
5  
6 game_data = game_data.rename(columns={  
7     'xg': 'xg_home',  
8     'xg.1': 'xg_away'  
9 })  
10
```

```

11
12 game_data = game_data.drop(columns=['venue', 'matchreport', 'notes'])
13
14 game_data

```

Final dataset shape: **(380, 11)** → matches expected for a full PL season.

Feature Engineering [🔗](#)

To capture team form, we created rolling averages of recent performance:

```

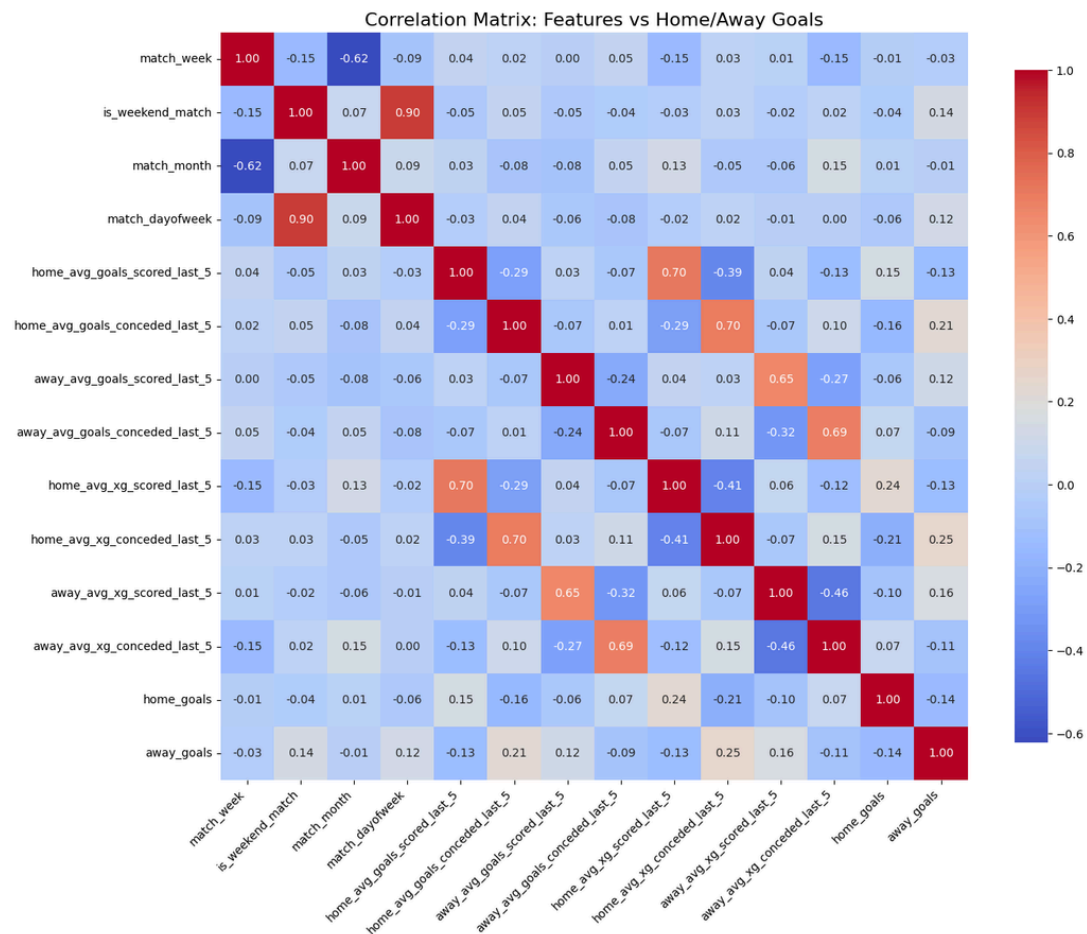
1 def compute_team_rolling_stats(df, team_col, goals_for_col, goals_against_col, prefix):
2     team_stats = []
3
4     for team in df[team_col].unique():
5         team_matches = df[(df[team_col] == team)].copy()
6         team_matches = team_matches.sort_values('date')
7
8         team_matches[f'{prefix}_avg_goals_scored_last_5'] = team_matches[goals_for_col].shift().rolling(5,
min_periods=1).mean()
9         team_matches[f'{prefix}_avg_goals_conceded_last_5'] =
team_matches[goals_against_col].shift().rolling(5, min_periods=1).mean()
10
11         team_stats.append(team_matches[[f'{prefix}_avg_goals_scored_last_5',
f'{prefix}_avg_goals_conceded_last_5']])
12
13     return pd.concat(team_stats).sort_index()

```

similar function was used for xG

Correlation analysis [🔗](#)

A heatmap of correlations between engineered features and target variables (`home_goals` , `away_goals`) revealed:



Most Predictive Features (Good Correlation with Goals) [↗](#)

Correlated with `home_goals`:

Feature	Correlation
<code>home_avg_goals_scored_last_5</code>	0.24
<code>home_avg_xg_scored_last_5</code>	0.23
<code>home_avg_xg_conceded_last_5</code>	-0.21 (inverse relation)

Correlated with `away_goals`:

Feature	Correlation
<code>home_avg_xg_conceded_last_5</code>	0.25
<code>away_avg_xg_scored_last_5</code>	0.16
<code>away_avg_goals_scored_last_5</code>	0.12

Less Important Features (Low or No Correlation) [↗](#)

Feature	Correlation with <code>home_goals</code>	Correlation with <code>away_goals</code>
---------	--	--

match_week	-0.01	-0.03
match_month	0.01	0.14
is_weekend_match	0.04	0.12
match_dayofweek	0.06	0.12

These were retained but weighted as weaker predictors.

Model Training (TensorFlow) [🔗](#)

Input Preparation [🔗](#)

- Categorical features (`home` , `away`) were one-hot encoded.
- Numerical features + encoded teams used as model input (`X_train`).
- Targets: `home_goals` , `away_goals` .

```
1     home_enc = pd.get_dummies(features_train['home'], prefix='home')
2     away_enc = pd.get_dummies(features_train['away'], prefix='away')
3     X_train = pd.concat([features_train[feature_cols].fillna(0), home_enc, away_enc], axis=1)
4     y_train = features_train[['home_goals', 'away_goals']].astype(float)
```

- **Input (X_train):** Combines your numerical features with one-hot encoded home and away teams.
- **Target (y_train):** The real number of home and away goals (used as regression targets).

Model architecture [🔗](#)

```
1     model = models.Sequential([
2         layers.Dense(90, activation='relu', input_shape=(X_train.shape[1],)),
3         layers.BatchNormalization(),
4         layers.Dropout(0.2),
5         layers.Dense(48, activation='relu'),
6         layers.BatchNormalization(),
7         layers.Dropout(0.2),
8         layers.Dense(30, activation='relu'),
9         layers.BatchNormalization(),
10        layers.Dropout(0.2),
11        layers.Dense(2)
12    ])
13    model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=learning_rate), loss='mse', metrics=
14    ['mae'])
14    history = model.fit(X_train, y_train, epochs=epochs, batch_size=batch_size, validation_split=0.2,
    verbose=int(verbose))
```

- 3 hidden layers: **90** → **48** → **30** units
- Each layer uses:
 - `ReLU` activation
 - `BatchNormalization` (improves stability)
 - `Dropout` (prevents overfitting)
- Final output layer has **2 nodes** to predict home and away goals

Model Compilation & Training [🔗](#)

```
1 model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=learning_rate), loss='mse', metrics=
  ['mae'])
2 history = model.fit(X_train, y_train, epochs=epochs, batch_size=batch_size, validation_split=0.2,
  verbose=int(verbose))
```

- **Loss function:** Mean Squared Error (MSE) - standard for regression
- **Metric:** Mean Absolute Error (MAE) - helps track training progress
- Trained for **120 epochs** with **batch size 128**
- **20% of data** used for validation during training

Multiple Runs for Stability [🔗](#)

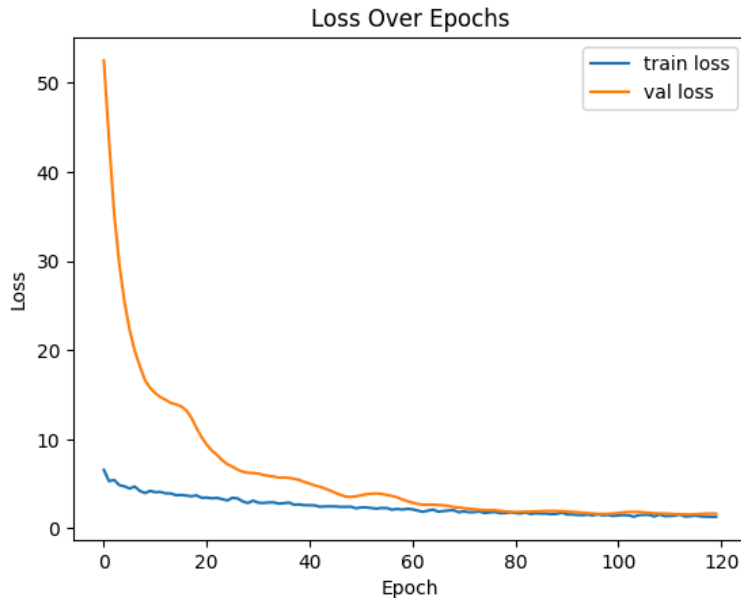
The model was trained **multiple times** (`n_runs = 10, 50, 100, 500`) and results were averaged to:

- Reduce randomness from model initialization and dropout
- Provide a **stable and fair estimate** of match outcomes and final table

Model Accuracy [🔗](#)

To ensure robustness, the model was trained **multiple times** and the predictions were **averaged**. This helps reduce the effect of randomness from model initialization and dropout layers.

The plot of **training vs validation loss** over 120 epochs gives a clear picture of how well the model learned:



Early Epochs (0–20): [🔗](#)

- **Validation loss drops rapidly** → The model is quickly learning meaningful patterns.
- **Training loss decreases steadily** → Learning is effective.

Middle Epochs (20–60): [🔗](#)

- Both losses continue to **decline smoothly**.
- No sign of **overfitting** - validation performance keeps improving.

Later Epochs (60–120): 🔗

- Loss curves **flatten out**, indicating convergence.
- **No sharp divergence** between training and validation → Model generalizes well.

Summary of Learning Behavior 🔗

- **Good convergence**: Model effectively learns from training data.
- **No overfitting**: Validation loss remains low and stable.
- **Final MSE is very low**, suggesting that the model is accurately predicting goal outcomes on average.

Prediction Performance 🔗

The final predictions were evaluated against actual match outcomes:

Metric	Accuracy
Match Outcome (W/D/L)	56.25%
Exact Home Goals Predicted	36.96%
Exact Away Goals Predicted	43.24%

Key Observations 🔗

- The model does **well at predicting match results** (win/draw/loss), with an outcome accuracy over 56%.
- **Away goals** were predicted **more accurately** than home goals.
- Predicting the **exact number of goals** is inherently difficult due to football’s unpredictability - but the model still provides **reliable estimations**.

Example Findings 🔗

Example Predictions 🔗

Match	Predicted Score	Predicted Outcome	Actual Score	Actual Outcome	Result Type
Arsenal (Home) vs Tottenham (Away)	2 - 1 (2.09 - 1.03)	Home Win	2 - 1	Home Win	Accurate
Arsenal (Home) vs Man City (Away)	2 - 1 (2.25 - 1.18)	Home Win	5 - 1	Home Win	Partial Miss

- The model often **predicts the outcome correctly**, even when the **exact scoreline varies**.
- This consistency in **outcome prediction** makes it well-suited for **league table forecasting**.

