

Optimiser - Fantasy Football

- [Aim of the Project](#)
 - [Rules of Fantasy Premier League \(FPL\)](#)
 - [Squad Composition](#)
 - [Starting Lineup Constraints](#)
 - [Team Limits](#)
 - [Transfer Rules](#)
- [Data Collection](#)
 - [Player History Extraction](#)
- [Feature Engineering](#)
 - [Feature Correlation Insights](#)
- [XGBoost Model Training](#)
 - [Model Performance:](#)
 - [Overfitting Check](#)
- [Team Optimisation Methods](#)
 - [Predicting Next Gameweek Points](#)
 - [Captain Selection](#)
 - [Optimising Starting XI](#)
 - [Transfer Recommendations](#)
- [Example Findings](#)

Aim of the Project [↗](#)

[🔗](#) App link: [Streamlit App](#)

This project focuses on optimising Fantasy Football teams by maximising predicted points for upcoming gameweeks. The key objectives are:

- Predict total points a user's Fantasy Premier League (FPL) team will score.
- Select the optimal starting XI lineup.
- Choose the best captain for maximum points.
- Recommend player transfers to improve team performance.

Rules of Fantasy Premier League (FPL) [↗](#)

When building and optimizing a Fantasy Premier League squad, the following rules and constraints must be respected to ensure squad validity:

Squad Composition [↗](#)

- **Total squad size:** 15 players
- **Positions:**
 - 2 Goalkeepers
 - 5 Defenders
 - 5 Midfielders
 - 3 Forwards

Starting Lineup Constraints [↗](#)

- **Starting XI size:** 11 players

- **Formation:** The starting lineup must follow one of the allowed FPL formations:
 - 3-4-3, 3-5-2
 - 4-4-2, 4-3-3
 - 5-3-2, 5-4-1
- Exactly 1 goalkeeper must be in the starting XI.

Team Limits [↗](#)

- **Maximum players from one Premier League club:** 3 in the whole squad (including bench and starting XI).

Transfer Rules [↗](#)

- **Transfers per gameweek:** 1 free transfer per gameweek, if not used they carry over to next gameweek (max 3); additional transfers cost points (-4).
- **Squad validity:** All transfers must maintain a valid squad according to the above rules.
- **Transfer limits in optimization:** The model respects allowed number of transfers.

Data Collection [↗](#)

Data is sourced directly from the official Fantasy Premier League API:

```
1 url = "https://fantasy.premierleague.com/api/bootstrap-static/"
2 response = requests.get(url).json()
3
4 players_df = pd.DataFrame(response['elements'])      # All players
5 teams_df = pd.DataFrame(response['teams'])          # All teams
6 positions_df = pd.DataFrame(response['element_types']) # Position info
```

The API data is clean and structured, requiring minimal preprocessing.

Player History Extraction [↗](#)

A function fetches detailed historical stats for each player across all gameweeks:

```
1 def get_player_history(player_id):
2     url = f"https://fantasy.premierleague.com/api/element-summary/{player_id}/"
3     r = requests.get(url).json()
4     history_df = pd.DataFrame(r['history'])
5     return history_df
6
```

We loop through all players, gather their histories, and concatenate into a single DataFrame:

```
1 all_histories = []
2 for player_id in players_df['id']:
3     df = get_player_history(player_id)
4     df['player_id'] = player_id
5     all_histories.append(df)
6
7 full_history_df = pd.concat(all_histories, ignore_index=True)
```

Finally, player metadata is merged for easier identification:

```
1 full_history_df = full_history_df.merge(players_df[['id', 'first_name', 'second_name', 'team_name',  
2 'position']],  
                                         left_on='player_id', right_on='id')
```

Data size: 27,605 rows × 50 columns.

Feature Engineering [↗](#)

We use **XGBoost** to predict player points for the next gameweek.

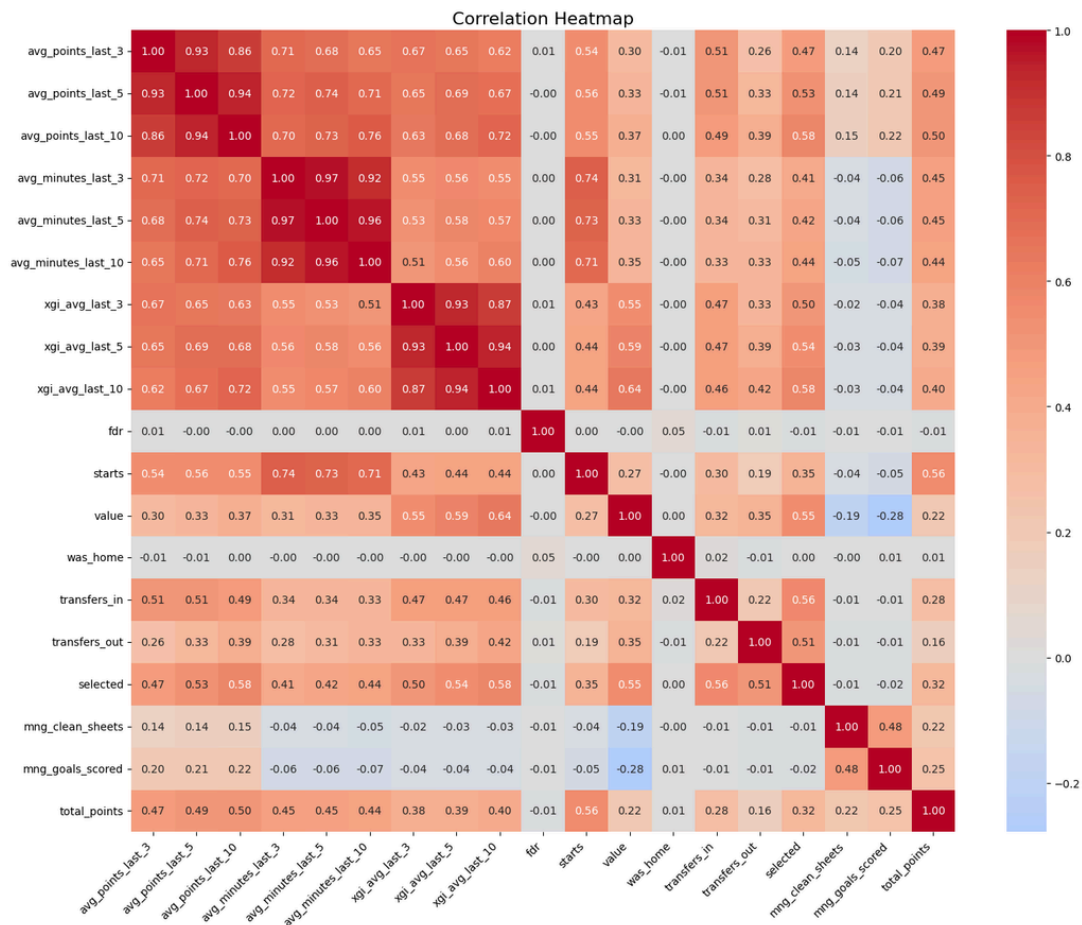
To enhance prediction, rolling averages for recent gameweeks are calculated for key features like points, minutes played, and expected goal involvements (xGI):

```
1 full_history_df = full_history_df.sort_values(by=["player_id", "round"])  
2 full_history_df['points_shifted'] = full_history_df.groupby('player_id')['total_points'].shift(1)  
3  
4 full_history_df['avg_points_last_3'] = full_history_df.groupby('player_id')['points_shifted'].rolling(window=3,  
min_periods=1).mean().reset_index(level=0, drop=True)  
5 full_history_df['avg_points_last_5'] = full_history_df.groupby('player_id')['points_shifted'].rolling(window=5,  
min_periods=1).mean().reset_index(level=0, drop=True)  
6 full_history_df['avg_points_last_10'] = full_history_df.groupby('player_id')  
['points_shifted'].rolling(window=10, min_periods=1).mean().reset_index(level=0, drop=True)
```

Similar rolling averages are created for minutes and xGI.

xGI (Expected Goal Involvement) is the sum of a player's **expected goals (xG)** and **expected assists (xA)**, estimating their overall contribution to potential goals.

Feature Correlation Insights [↗](#)



- High correlation with points:
 - starts (0.56)
 - Rolling averages of points (last 3,5,10 gameweeks ~0.47–0.50)
 - Rolling averages of minutes (~0.44–0.45)
 - Rolling averages of xGI (~0.38–0.40)
 - transfers_in (0.28), selected (0.32)
- Low predictive value features:
 - fdr (-0.01), was_home (0.01), transfers_out (0.16)

Multicollinearity refers to a situation where two or more input features are highly correlated, making it hard to isolate their individual effects - but this is **not a concern for XGBoost**, as its tree-based structure selects the most informative features and is inherently robust to feature correlation.

XGBoost Model Training [↗](#)

The dataset is split by gameweek for training and testing:

```
1 max_round_adj = fpl_data['round_adjusted'].max()
2 train_data = fpl_data[fpl_data['round_adjusted'] < max_round_adj]
3 test_data = fpl_data[fpl_data['round_adjusted'] == max_round_adj]
4
5 X_train = train_data[features]
6 y_train = train_data['total_points']
7
```

```

8 X_test = test_data[features]
9 y_test = test_data['total_points']

```

Model hyperparameters and training:

```

1 # Define model with some common hyperparameters
2 xgb_model = XGBRegressor(
3     n_estimators=300,
4     max_depth=2,
5     learning_rate=0.03,
6     subsample=0.6,
7     colsample_bytree=0.6,
8     reg_alpha=2,
9     reg_lambda=3,
10    gamma=1,
11    random_state=42,
12    n_jobs=-1
13 )
14
15 # Train the model
16 xgb_model.fit(X_train, y_train)

```

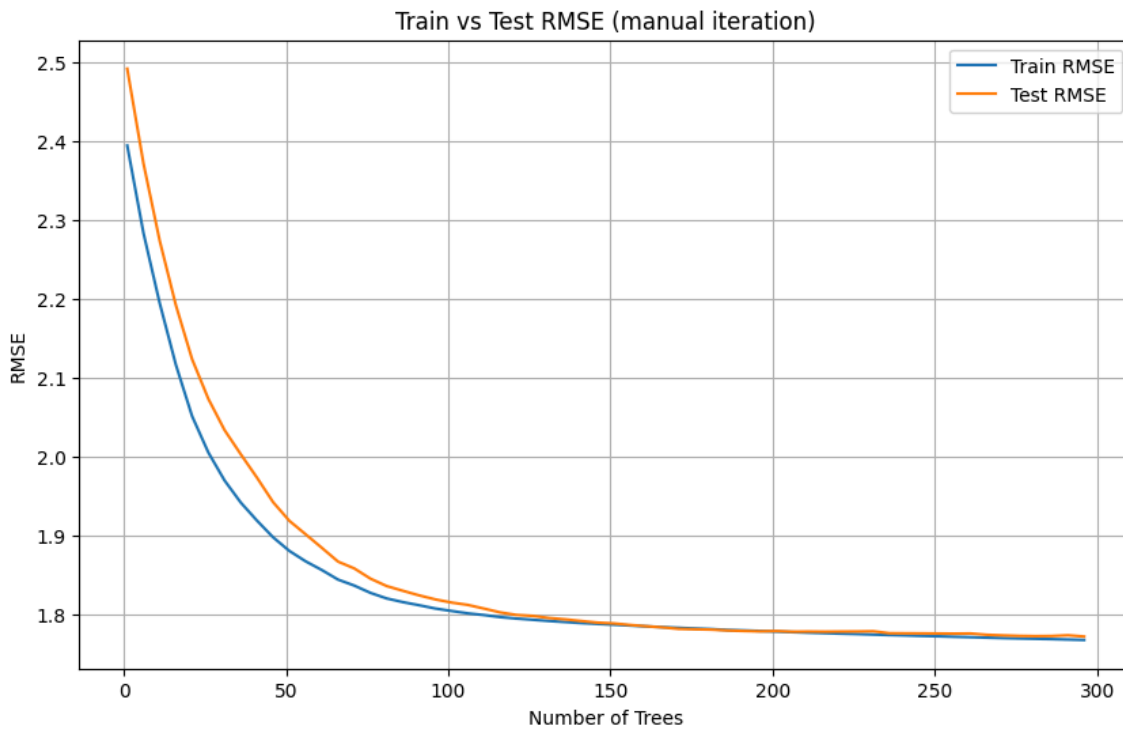
Model Performance: [🔗](#)

Metric	Definition	Value	What the Value Means
R² Score	Proportion of variance in actual data explained by the model (Closer to 1 is ideal)	0.476	The model explains 47.6% of the variation in goal outcomes - moderate fit
Mean Squared Error	Average of squared prediction errors (penalizes big errors more closer to 0 is ideal)	3.335	On average, predictions are 3.335 goals² off - some larger mistakes present
Mean Absolute Error	Average of absolute prediction errors (same unit as target closer to 0 is ideal)	0.912	Predictions are off by ~0.91 goals on average - quite accurate for football data

The model shows moderate predictive power with an acceptable error margin, indicating it captures player point patterns reasonably well.

Overfitting Check [🔗](#)

Training and testing errors are close, indicating the model generalizes well without overfitting.



From this gplot we see train RMSE is very near Train MSE this means

- The model is **not overfitting** (train error isn't way lower than test error).
- It's **generalizing well** to unseen data.
- The fit is balanced and stable.

Team Optimisation Methods [🔗](#)

Predicting Next Gameweek Points [🔗](#)

The model predicts each player's points for the upcoming gameweek:

```
1 fpl_data['next_gw_pred'] = xgb_model.predict(fpl_data[features])
```

Captain Selection [🔗](#)

The captain is chosen as the player with the highest predicted points in the starting lineup:

```
1 agg_preds = starters_next_gw.groupby('player_id', as_index=False)['next_gw_pred'].sum()
2
3 # Find player with max aggregated predicted points
4 best_captain_row = agg_preds.loc[agg_preds['next_gw_pred'].idxmax()]
```

Optimising Starting XI [🔗](#)

Given a 15-player squad, players are first grouped by their playing position:

```
1 players_by_pos = {'Goalkeeper':[], 'Defender':[], 'Midfielder':[], 'Forward':[]}
2 for pid in squad_ids:
3     pos = player_info[pid]['position']
```

```
4         players_by_pos[pos].append(pid)
```

We then ensure there are two goalkeepers, iterating through each as a potential starter. The goalkeeper with the highest predicted points is selected as the starting keeper, and the other is assigned to the bench:

```
1         gk_candidates = players_by_pos['Goalkeeper']
2         if len(gk_candidates) < 2:
3             return None, None
4
5         for gk in gk_candidates:
6             bench_gk = [pid for pid in gk_candidates if pid != gk][0]
```

Next, the function iterates through all valid Fantasy Premier League formations to find the optimal lineup. The standard formations considered are:

- 3-4-3
- 3-5-2
- 4-4-2
- 4-3-3
- 5-3-2
- 5-4-1

For each formation, the function generates all possible combinations of defenders, midfielders, and forwards that meet the formation requirements:

```
1         for D_req, M_req, F_req in [
2             (3, 4, 3), (3, 5, 2),
3             (4, 4, 2), (4, 3, 3),
4             (5, 3, 2), (5, 4, 1)
5         ]:
6             if len(def_candidates) < D_req or len(mid_candidates) < M_req or len(fwd_candidates) < F_req:
7                 continue
8
9             def_combos = combinations(def_candidates, D_req)
10            mid_combos = combinations(mid_candidates, M_req)
11            fwd_combos = combinations(fwd_candidates, F_req)
```

For each position, all possible groups of players of the required size are generated to fill the formation slots. The function will iterate through these to find valid squads.

For each valid combination, the lineup is constructed by combining the goalkeeper, defenders, midfielders, and forwards:

```
1         mid_combos = list(mid_combos)
2         fwd_combos = list(fwd_combos)
3
4         for d_combo in def_combos:
5             for m_combo in mid_combos:
6                 for f_combo in fwd_combos:
7                     xi_ids = [gk] + list(d_combo) + list(m_combo) + list(f_combo)
8
9                     team_counts = {}
10                    valid = True
11                    for pid in xi_ids:
```

```

12         team = player_info[pid]['team_name']
13         team_counts[team] = team_counts.get(team, 0) + 1
14         if team_counts[team] > 3:
15             valid = False
16             break
17     if not valid:
18         continue

```

This ensures no more than three players from the same real-life team are included in the starting XI.

Among all valid lineups, the one with the highest predicted points is selected as the optimal starting XI.

```

1         pts = total_pred_points(xi_ids)
2         if pts > best_pts:
3             best_pts = pts
4             best_xi = xi_ids

```

Transfer Recommendations [🔗](#)

To make informed transfer suggestions, we employ a **Top-K heuristic** to manage computational complexity effectively:

1. **Identify the bottom performers:** We select the bottom **k=6** players from the current squad based on predicted points.
2. **Identify potential top players:** We find the top **k=20** players outside the squad with predicted points greater than 1.

Select the bottom performing players in the current squad based on their predicted points for the next gameweek. These are the prime candidates to be transferred out

```

1     # Select the bottom topk_out players by predicted points from the current squad (worst performers)
2     bottom_current = current_players_df.nsmallest(topk_out, 'next_gw_pred')['player_id'].tolist()

```

From the pool of all available players outside the current squad, select those with a predicted points value greater than 1 for the upcoming gameweek. From these, take the top performers as potential transfer-in options:

```

1     # Potential players to bring in (outside current squad)
2     potential_ins_all = set(unique_players['player_id']) - current_player_ids
3     potential_ins_df = pred_next_gw[pred_next_gw['player_id'].isin(potential_ins_all)].copy()
4
5     # Filter to include only those with predicted points > 1
6     potential_ins_df = potential_ins_df[potential_ins_df['next_gw_pred'] > 1]
7
8     # Select the top topk_in players by predicted points
9     top_potential_ins = potential_ins_df.nlargest(topk_in, 'next_gw_pred')['player_id'].tolist()

```

For every number of transfers **k** from 0 up to the allowed number, all combinations of transferring out **k** players from the bottom candidates and transferring in **k** players from the top outside candidates are tested.

- When **k=0**, simply evaluate the current squad as is.
- For **k > 0**, generate all combinations of possible transfers and evaluate each candidate squad:

```

1     for k in range(0, transfers_allowed + 1):
2         # When k=0: no transfers, just check the current squad
3         if k == 0:

```

```

4         if valid_squad(current_player_ids):
5             xi_ids, bench_ids = pick_best_starting_xi(current_player_ids)
6             if xi_ids is not None:
7                 total_pts = total_pred_points(xi_ids)
8                 if total_pts > best_score:
9                     best_score = total_pts
10                    best_squad = current_player_ids
11                    best_starting_xi = xi_ids
12                    best_bench = bench_ids
13                    best_out = []
14                    best_in = []
15            continue
16
17        # For k > 0, consider combinations
18        outs_combos = combinations(bottom_current, k)
19        ins_combos = list(combinations(top_potential_ins, k))
20
21        for out_ids in outs_combos:
22            out_set = set(out_ids)
23            remaining_players = set(current_squad_list) - out_set
24
25            for in_ids in ins_combos:
26                in_set = set(in_ids)
27
28                candidate_squad = remaining_players | in_set
29
30                if not valid_squad(candidate_squad):
31                    continue
32
33                xi_ids, bench_ids = pick_best_starting_xi(candidate_squad)
34                if xi_ids is None:
35                    continue
36
37                total_pts = total_pred_points(xi_ids) # Transfer penalty
38
39                if total_pts > best_score:
40                    best_score = total_pts
41                    best_squad = candidate_squad
42                    best_starting_xi = xi_ids
43                    best_bench = bench_ids
44                    best_out = out_ids
45                    best_in = in_ids

```

This approach balances thoroughness and computational feasibility by focusing only on the worst-performing players in the current squad and the best available players outside it. By limiting transfer candidates to these Top-K groups, the algorithm efficiently searches for the optimal squad that maximizes predicted points for the next gameweek.

Example Findings [🔗](#)

For Gameweek 25, with 3 transfers allowed and Matz Sels as captain:

- Initial predicted points: 38
- After optimisation:
 - Captain changed to Salah
 - Transfers in Minteh, Diaz, Marmoush
 - New predicted points: 87 (an increase of 48 points)

This demonstrates the effectiveness of the model and optimiser in improving team potential.

