# Ranking Footballers with Multilevel Modeling

Gregor Grbec[1], Nino Bašić[1] and Marko Tkalčič[1,*]

[1]*University of Primorska, Faculty of Mathematics, Natural Sciences and Information Technologies, Glagoljaška 8, SI-6000 Koper, Slovenia*

**Abstract**

Despite football's collaborative nature, the inquiry into the identity of the best player is a frequent topic in the footballing realm. This discussion disproportionately highlights attacking players, creating an apparent bias, as every team role holds significance. Our study aimed to delineate player performance from team performance and ensure the inclusion of players from all positions in the ultimate ranking of the best players. We sourced data from FBref, encompassing every player in every match played by a top 20 European team in the current century's top 5 European leagues. Employing a multilevel linear mixed-effects model, we utilized team points as the response variable, accounting for both player and opponent team strength. The extraction of level-2 player residuals, averaged by player, facilitated the creation of a comprehensive ranking for the best players of this century. Surprisingly, two players widely regarded as among the best of all time, Messi and Ronaldo, secured relatively low positions on our list (Ronaldo at 12th, and Messi at 14th).

**Keywords**
ranking, football, multi-level modeling

## 1. Introduction

Despite football being a team sport, determining the best male football player remains a widely debated topic. In every era, a standout candidate emerges—Pelé in the 1960s, Diego Maradona in the 1980s, and more recently, Lionel Messi and Cristiano Ronaldo dominating the past 15 years. The ongoing debate centers on which of these four players is the all-time best. Our attempt to provide a data-driven solution was hindered by the historical match data's poor quality, leading us to focus on the 2000/2001 to 2022/2023 period due to its better availability and reliability.

To identify the premier football player, we sought to understand what defines greatness in football. Our contention was that a stellar player elevates the team through diverse contributions—be it goals, assists, defensive actions, or on-field leadership. The absence of such players correlates with a dip in team performance, underscoring their impact. To be deemed the best, a player must consistently make a difference on the grandest football stage, serving as a talisman for a top-tier team over several seasons. Our primary research goal was to gauge individual contributions to team performance.

All four candidates—Pelé, Maradona, Messi, and Ronaldo—share the role of attackers. While goal-scoring garners attention, players who score less frequently often go unnoticed. The current spotlight on attackers neglects the significance of midfielders, defenders, and goalkeepers in

CEUR Workshop Proceedings (CEUR-WS.org)

the best player debate. Goals scored, being the most coveted statistic, contributes to this bias. Defenders' performance is typically measured by goals conceded per game, yet this metric is assigned to the entire defense, creating an imbalance in player position evaluation. Our second research objective aimed to rectify this bias and provide a fair comparison among players irrespective of their position or style of play.

To fulfill our research goals—ranking players based on impact and ensuring equal opportunities for all positions—we employed a multilevel mixed-effects model. This model, utilizing achieved points in the game as a performance metric, underwent training on the last 23 seasons of every league match involving the top 20 European football teams across the top 5 European leagues.

## 2. Related Work

In this section, we explore prior research on ranking individual ability and multilevel modeling in team sports.

Brooks et al. [1] assessed players' offensive ability by analyzing completed passes leading to shots. They predicted pass quality by training a model on La Liga data, ranking players based on the quality of their passes. McHale and Scarf [2] also ranked players, correlating a team's and player's contributions to match outcomes. Their index awarded points for player contributions, validated in the Premier League, with a focus on eliminating player role bias.

Pappalardo et al. [3] used extensive event data from various leagues to rank players, employing weights for metrics. While successful in extracting player performance, these studies did not account for player role bias.

Mixed-effects modeling in football research includes Grund [4], who studied passing structures' impact on match outcomes. Beyond football, Casals and Martinez [5] analyzed basketball player performance, while Gerber and Craig [6] predicted baseball players' performance across leagues.

Inspired by Bell et al. [7], our study adopted multilevel modeling to extract player performance from team performance. Their F1 driver analysis, employing a cross-classified model, served as a valuable model for our approach. The model, controlling for team switches and opponent strength, allowed us to eliminate player role bias and extract meaningful player performance metrics.

## 3. Modeling

The linear mixed-effects model, also referred to as a mixed model, random effects model, multilevel model, or hierarchical model, serves as a statistical model tailored for hierarchical data [8].

This model empowers the control of variables at higher levels, effectively addressing the variation and correlation within the data structure to yield more precise outcomes. It encompasses fixed variables, representing coefficients with a consistent impact on the response variable across all groups, and at least one random variable, introducing a variable effect contingent on the group.

Building upon a fundamental linear regression model, the mixed-effects model enables the variation of intercepts and/or slopes of the regression line across different data groups for select variables. For instance, in a study by Bell et al. [7], team strength was controlled for, recognizing that the performance of drivers from superior teams, like Ferrari, may differ from those in weaker teams like Renault. This flexibility allowed for a nuanced assessment of driver quality, accounting for team effects on intercepts and slopes.

## 3.1. Toy Example

Consider predicting students' performance on the fictional National Test of Mathematics based on their average percentage of points achieved in their Mathematics class. The data is nested on two levels: the school (School A and School B) and the student. Each data point represents a student's average percentage in class and the competition.

School A is known for its strict grading, while School B is more lenient. Predicting overall performance without accounting for school variations would be inaccurate due to the substantial difference in expected competition scores between the schools.

The mixed-effects multilevel model addresses this issue by allowing control for school, enhancing prediction accuracy. In our case, random slopes are not suitable; thus, random intercepts and fixed percentages of points in school are included in the formula:

$$PercInComp_i = (b_0 + b_{0,School}) + (b_1 + b_{1,School}) \cdot PercInSchool + \epsilon_i$$

Here, $b_0$ and $b_1$ represent the overall intercept and slope, $b_{0,School}$ and $b_{1,School}$ account for variations by school, and $\epsilon_i$ is the student's residual. The overall regression line is:

$$y = PercInSchool - 5.3$$

The school-specific lines are:

$$y = 0.81 \cdot PercInSchool + 24.5 \quad \text{(School A)}$$
$$y = 1.18 \cdot PercInSchool - 35.1 \quad \text{(School B)}$$

Intercept and slope values for the overall and school-specific cases are presented in Table **??**.

## 3.2. Footballer's Ranking Case

In our study, we employed a similar framework featuring two fixed effects (opponents's points per game and home or away indicator) and 3 random effects (team, team in a particular season, and player). These incorporate random intercepts and slopes, varying based on the predicting variables. The data is nested across four levels, comprising 20 teams, each spanning multiple seasons, players associated with one or multiple clubs across different seasons, and repeated measures for every match of every player.

For instance, Cristiano Ronaldo participated in 597 matches over 14 seasons for 3 different clubs. Teams varied in participation, with RB Leipzig, for instance, joining the Bundesliga from the 18/19 season onwards but achieving significant success in those four seasons, securing a spot in our top 20 teams list.

### 3.3. Good Player Definition

To identify the best football player, we established criteria defining a standout player as someone who consistently elevated top teams in the premier European leagues—English, Spanish, Italian, German, and French divisions—over an extended period.

### 3.4. Data Acquisition

Data from the top five European leagues was scraped from FBref [9], a division of Sports Reference [10]. We utilized Python libraries, including "requests" by Reitz [11] and "bs4" by Crummy [12], for web scraping. The dataset encompassed matches from the 2000/2001 to 2022/2023 seasons, including columns such as team points, player name, team name, season, opponent's points per game, minutes played, and home or away status.

We filtered players with a minimum of 340 matches for the top 20 teams across the leagues, setting the threshold close to a full season. Additionally, players with fewer than 15 minutes of playtime were excluded, ensuring impactful player contributions.

### 3.5. Model Building

For the multilevel mixed-effects model, we utilized the "lmer" function from lme4 [13] in R. Model building, inspired by Bell et al. [7], involved iterative development, comparing versions using AIC and BIC values. The final model includes fixed effects (opponent's points per game and home/away), and random effects for team, team in a season, and player, with random intercepts and slopes for differentiation.

### 3.6. Level-2 Residual Extraction

Player-specific intercepts and slopes were obtained using the "ranef" function from lme4 [13]. A custom function calculated player contributions to matches, extracting level-2 residuals. Team residuals were similarly extracted for the top 20 list.

## 4. Results

In Table 1, the ranking displays players with over 340 league games for the top 20 teams in the top 5 European leagues. Players are ordered by mean residual values, showcasing their impact on team performance. Giorgio Chiellini leads the ranking with a mean value of $4.091820 \times 10^{-9}$, signifying an average improvement in his team's performance when he played—he contributed to scoring more points. Conversely, Marcelo had a negative impact on his team, indicated by a mean value of $-4.280433 \times 10^{-9}$. In simpler terms, when Giorgio Chiellini played, he, on average, exceeded predicted team performance by $4.091820 \times 10^{-9}$ points.

Table 1: Ranking of players by average level-2 residual.

| Rank | Player | Average level-2 residual |
|---|---|---|
| 1 | Giorgio Chiellini | 4.09181982129302e-09 |
| 2 | Andrea Pirlo | 3.98918255737386e-09 |
| 3 | Petr Čech | 3.82879209703914e-09 |
| 4 | Gianluigi Buffon | 2.94709848811944e-09 |
| 5 | Thiago Silva | 2.91890508830269e-09 |
| 6 | Pepe Reina | 2.62715962216929e-09 |
| 7 | John Terry | 2.58695437819242e-09 |
| 8 | Xabi Alonso | 2.35152000228812e-09 |
| 9 | Wojciech Szczesny | 2.34987638849159e-09 |
| 10 | Steve Mandanda | 2.26510689022037e-09 |
| 11 | Ashley Cole | 2.22746719720238e-09 |
| 12 | Cristiano Ronaldo | 2.14548852805352e-09 |
| 13 | Jordan Henderson | 2.03429603881283e-09 |
| 14 | Lionel Messi | 1.90504027044852e-09 |
| 15 | Víctor Valdés | 1.78940910084554e-09 |
| 16 | Karim Benzema | 1.76636176608884e-09 |
| 17 | Marek Hamšík | 1.759971404514e-09 |
| 18 | Javier Zanetti | 1.66454124921608e-09 |
| 19 | Jamie Carragher | 1.57078817377543e-09 |
| 20 | Sergio Busquets | 1.53245448858787e-09 |

## 5. Conclusion

Football, despite being a team sport, perpetually raises the question of the best player, creating endless debates and discussions. Opinions, often subjective, vary based on personal criteria. Notably, offensive players dominate discussions, overshadowing the defensive aspect, crucial but overlooked. This study aims to objectively extract player performances from team data, offering an equitable assessment of all roles.

Our definition of a good player hinges on their team's reliance—a player missed when absent, impacting team performance. To ensure an accurate evaluation, we employed a multilevel mixed-effects model, controlling for team strength. Data from FBref encompassed player and team details, match points, home/away status, season, and opposition's average points per game. A linear mixed-effects model allowed us to control for team strength, with extracted level-2 player residuals forming the final rankings.

The list featured impactful players in this century, with Giorgio Chiellini topping, followed by Andrea Pirlo and Petr Čech. Surprisingly, iconic players like Ronaldo and Messi ranked 12th and 14th. The top 30 showed balance across positions—8 goalkeepers, defenders, and midfielders, and 6 attackers.

Player level-2 residuals were small due to players' extensive playing time, emphasizing team and team season effects. Future exploration could widen the timeframe, create league-

specific rankings, and incorporate diverse metrics for player contribution, potentially examining managerial impact and expanding related studies to extensive periods.

## Acknowledgments

## References

[1] J. Brooks, M. Kerr, J. Guttag, Developing a Data-Driven Player Ranking in Soccer Using Predictive Model Weights, in: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ACM, San Francisco California USA, 2016, pp. 49–55. URL: https://dl.acm.org/doi/10.1145/2939672.2939695. doi:10.1145/2939672.2939695.

[2] I. McHale, P. Scarf, Ranking Football Players, Significance 2 (2005) 54–57. URL: https://doi.org/10.1111/j.1740-9713.2005.00091.x. doi:10.1111/j.1740-9713.2005.00091.x, _eprint: https://academic.oup.com/jrssig/article-pdf/2/2/54/49108761/sign_2_2_54.pdf.

[3] L. Pappalardo, P. Cintia, P. Ferragina, E. Massucco, D. Pedreschi, F. Giannotti, PlayeRank: Data-driven Performance Evaluation and Player Ranking in Soccer via a Machine Learning Approach, ACM Transactions on Intelligent Systems and Technology 10 (2019) 1–27. URL: https://dl.acm.org/doi/10.1145/3343172. doi:10.1145/3343172.

[4] T. U. Grund, Network structure and team performance: The case of English Premier League soccer teams, Social Networks 34 (2012) 682–690. URL: https://linkinghub.elsevier.com/retrieve/pii/S0378873312000500. doi:10.1016/j.socnet.2012.08.004.

[5] M. Casals, A. J. Martinez, Modelling player performance in basketball through mixed models, International Journal of Performance Analysis in Sport 13 (2013) 64–82. URL: https://www.tandfonline.com/doi/full/10.1080/24748668.2013.11868632. doi:10.1080/24748668.2013.11868632.

[6] E. A. E. Gerber, B. A. Craig, A mixed effects multinomial logistic-normal model for forecasting baseball performance, Journal of Quantitative Analysis in Sports 17 (2021) 221–239. URL: https://www.degruyter.com/document/doi/10.1515/jqas-2020-0007/html. doi:10.1515/jqas-2020-0007.

[7] A. Bell, J. Smith, C. E. Sabel, K. Jones, Formula for success: Multilevel modelling of Formula One Driver and Constructor performance, 1950–2014, Journal of Quantitative Analysis in Sports 12 (2016) 99–112. URL: https://www.degruyter.com/document/doi/10.1515/jqas-2015-0050/html. doi:10.1515/jqas-2015-0050.

[8] A. Bryk, S. Raudenbush, Hierarchical Linear Models: Applications and Data Analysis Methods, Advanced Quantitative Techniques in the Social Sciences, SAGE Publications, 1992. URL: https://books.google.si/books?id=eE-CAAAAIAAJ.

[9] FBref, Football Statistics and History, 2023. URL: https://fbref.com/en/.

[10] Sports Reference, Sports Reference | Sports Stats, fast, easy, and up-to-date, 2023. URL: https://www.sports-reference.com/.

[11] K. Reitz, requests: Python HTTP for Humans., 2023. URL: https://requests.readthedocs.io.

[12] Crummy, Beautiful Soup 4.12.0 documentation, 2023. URL: https://www.crummy.com/software/BeautifulSoup/bs4/doc/.

[13] lme4, lme4 package, 2023. URL: https://github.com/lme4/lme4.