

An overview of Football Analytics

Ioannis Ntzoufras

Athens University of Economics & Business



Introduction

Football/Soccer is the best sport for implementing Science/Statistics/Analytics

- Low number of events (so difficult to predict)
- High uncertainty (so difficult to predict)
- Very popular (because it is difficult to predict?)
- Very profitable (because it is difficult to predict?)
- High Financial Risk of investment (because passion becomes more important than numbers and science) – Professional Teams are usually acting as win-maximizers and not profit-maximizers



Main Topics Quantitative analysis of Football/Sports

- Prediction
- Player Evaluation & Performance analytics
- Physical Metrics of Players in training
- Inline game metrics with wearables
- Scheduling
- Sports Economics & Competitive Balance
- Other (Passing Network Analytics, Referee effects, Red card effect, Home effect, Corruption Analytics, Analysis of substitution times)



Prediction

- Offline (before the game)
- Inline (within the game)

Offline Prediction

Modeling of

- Game Scores
 - Poisson based models and extensions
 - Modeling the difference using the Skellam model
- Final outcome of a game (Win/Draw/Loss)
 - Multinomial regression model
 - Bradley Terry Model

Models for Scores

Models for Counts

- Simple Poisson Model (Maher, 1982; Lee, 1992; Dixon & Coles, 1997, Karlis and Ntzoufras, 2000)
- Bivariate Poisson Model (Karlis & Ntzoufras, 2003)
- Negative Binomial Model (see e.g. Ntzoufras 2009)
- Skellam Model for the goal difference (Karlis & Ntzoufras, 2009)
- Poisson-log-normal random effects model (not the best for football counts; see e.g. Ntzoufras 2009)

Models for Scores

Such models allow us not only to predict a single football game but also (simulation based results)

- Final League reproduction
- Estimate probabilities of winning a league, winning European tickets, or relegation.
- Estimate final rankings
- Estimate results under different scenarios/assumptions (by changing covariates i.e. conditions of the game)

Offline Prediction

Poisson Based models

- Vanilla model: home effect + teams attacking and defensive parameters
- Models with time evolved team parameters (time and form matters!)
- Additional covariates
 - Odds from betting teams (easily accessible – good covariates)
 - Team performance (ingame and before the game)
 - Information about events and formation (team strategy, formation, injuries etc.)
 - Economo-demographic variables (Stability, tradition, Budget, Player Value, Coach Value, Country of origin for European leagues)
 - Prior information (previous games between the teams)
 - Team form (e.g. performance in last 5 games)

Offline Prediction

The simple (vanilla) Poisson model

The model is expressed by

$$\begin{aligned} Y_{ij} &\sim \text{Poisson}(\lambda_{ik}) && \text{for } j = 1, 2 \\ \log(\lambda_{i1}) &= \mu + \text{home} + a_{\text{HT}_i} + d_{\text{AT}_i} \\ \log(\lambda_{i2}) &= \mu && + a_{\text{AT}_i} + d_{\text{HT}_i} \quad \text{for } i = 1, 2, \dots, n, \end{aligned}$$

where n = number of games, μ = constant parameter; home = home effect; HT_i and AT_i = home and away teams in i game; a_k and d_k = attacking and defensive effects–abilities of k team for $k = 1, 2, \dots, K$; and K = number of teams in the data (here $K = 20$).

In full (balanced) round-robin leagues, the parameters can be easily calculated by considering averaged of scored/conceded goals for each team

Offline Prediction

Data for the simple (vanilla) model

- **Observations**
 - $2 \times$ Number of games (N)
 - Each game will occupy two lines/observations (one for home team and one for away team)
- **Response Variable**: Goals scored by each team in each game
- **Covariates**
 - **Home effect**: Binary for home and away teams (1 for home teams and zero otherwise)
 - **Scoring team**: Categorical factor for the team scoring the number of goals (the corresponding coefficient will estimate the attacking ability of each team)
 - **Team accepting goals**: Categorical factor for the team receiving the number of goals (the corresponding coefficient will estimate the defensive ability of each team).

Offline Prediction

Important Assumptions

- Dependence/Independence of Goals of a game
- Time dependent attacking and defending parameters
- What about draw inflation?
- What about Over-dispersion?
- Shall we focus on modeling scores or outcomes (win/draw/loss)?

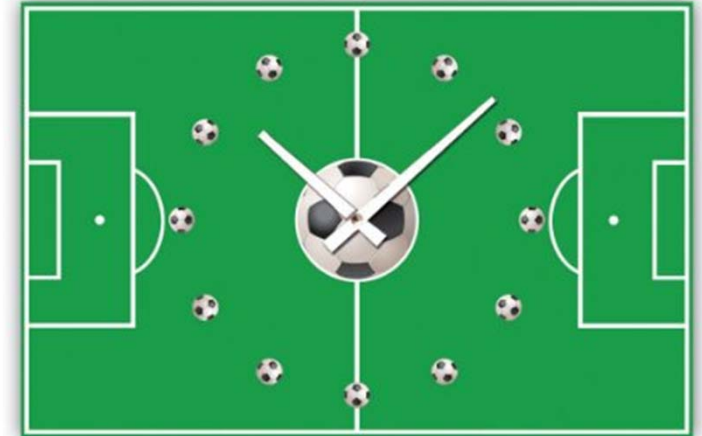
Checking the performance of the predictions

- Checking model fit and prediction using in-sample and out-of-sample measures

Prediction within the game

Modeling of

- Time to event (goal)
 - Survival analysis based models
 - Dixon & Robinson (1998, *RSSD*)
 - Nevo and Ritov (2013, *JQAS*)
 - Boshnakov, Kharrat, McHale (2017, *Int. J. Forecasting*)
 - Work in progress by our team
- Model the probability of event for short intervals (every 1 or 5 minutes)
 - Using Binomial mixed models for repeated measures



Player Evaluation

Aim

- Estimate the contribution of players in a team
- Rank, identify and reward best players
- Scouting – Early Identification of talents
- Estimate the future performance/value of a Player
- Help the manager to decide the best formation



Player Evaluation

Methods

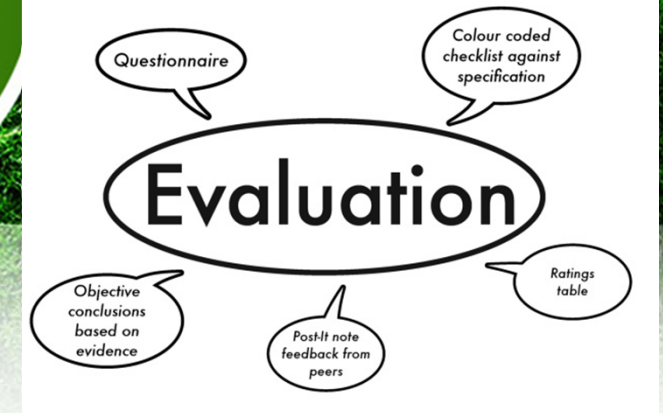
- Simple approach with binary indicators
- Random effects
- Analysis based on Game Performance Indicators
- Expected Goals (xG) and Expected Assists (xA)
- Player Economic/Marketing Value and performance



Player Evaluation

Methods (2)

- Simple approach with indicators
 - Build a model with indicators whether a player was in the field
 - Binary indicators for players
 - Difficult to build a dataset. Each game should be splitted in multiple lines according to substitution times
- Analysis based on Game Performance Indicators
 - Build a model to identify the importance of each event in the game (goals, shots, steals, passes, speed, stamina, area covered etc.)
 - Use model indicators to build an index of players
 - McHale, Scarf & Folker (2012, *Interfaces*) building different indexes based on different response measures



Player Evaluation

Methods (3)

- Random effects
 - Use random effects to identify individual contribution
 - Goal Scoring: McHale & Szczepanski (2014, *JRSSA*)
 - Passing Skills: Szczepanski & McHale (2016, *JRSSA*)
- Player Economic/Marketing Value and performance
 - Saebo & Hvattum (2018, *Journal of Sports Analytics*): *Modelling the financial contribution of soccer players to their clubs*
 - Evaluating the efficiency of the association football transfer market using regression based player ratings (pre-print only)



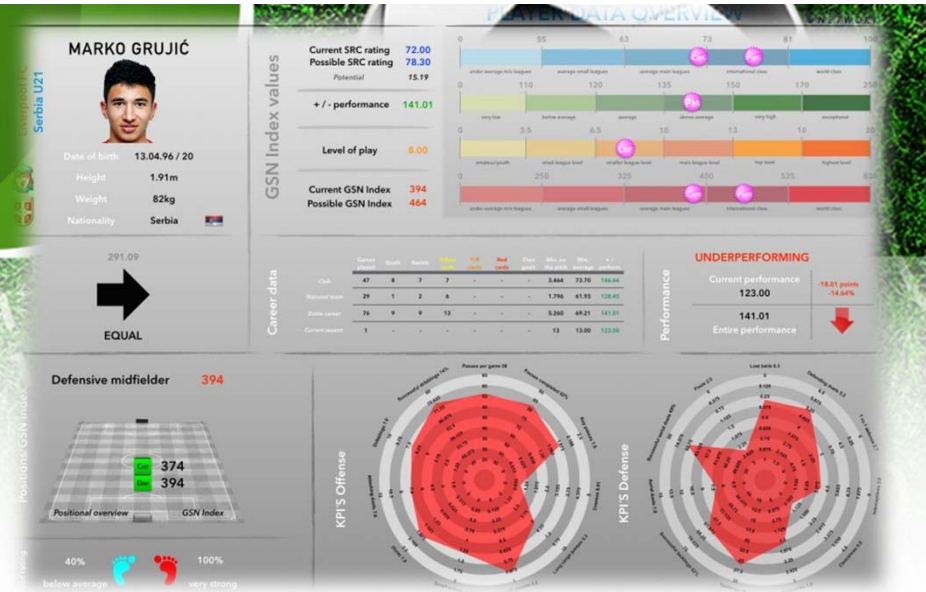
Player Evaluation

Methods (4)

McHale, Scarf & Folker (2012, *Interfaces*)
building different indexes based on different
response measures

Index ingredients:

- Subindex 1: Modelling Match Outcome (model based with outcome probability)
- Subindex 2: Points-Sharing Index (time played by each players and points)
- Subindex 3: Appearance Index (time played by each players)
- Subindex 4: Goal-Scoring Index
- Subindex 5: Assists Index
- Subindex 6: Clean-Sheets Index

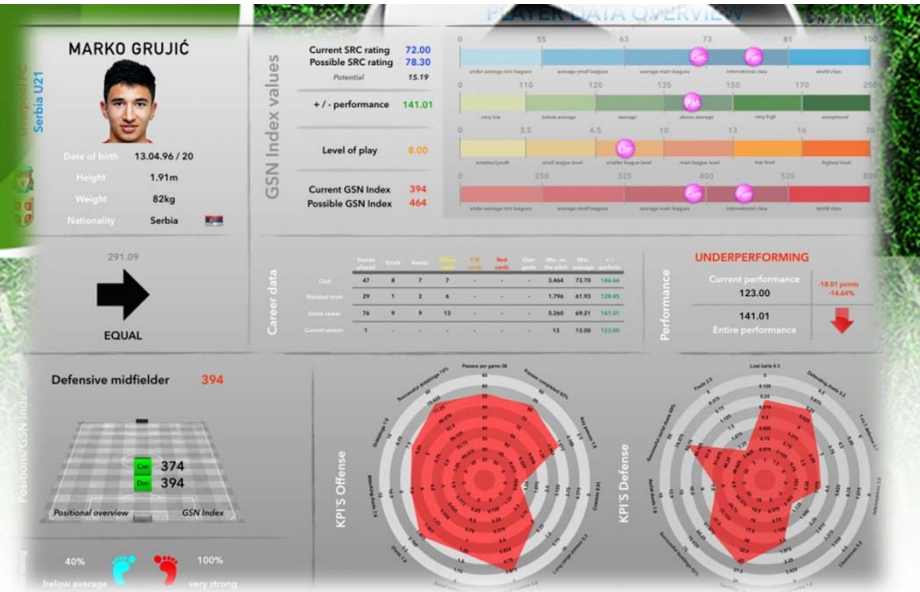


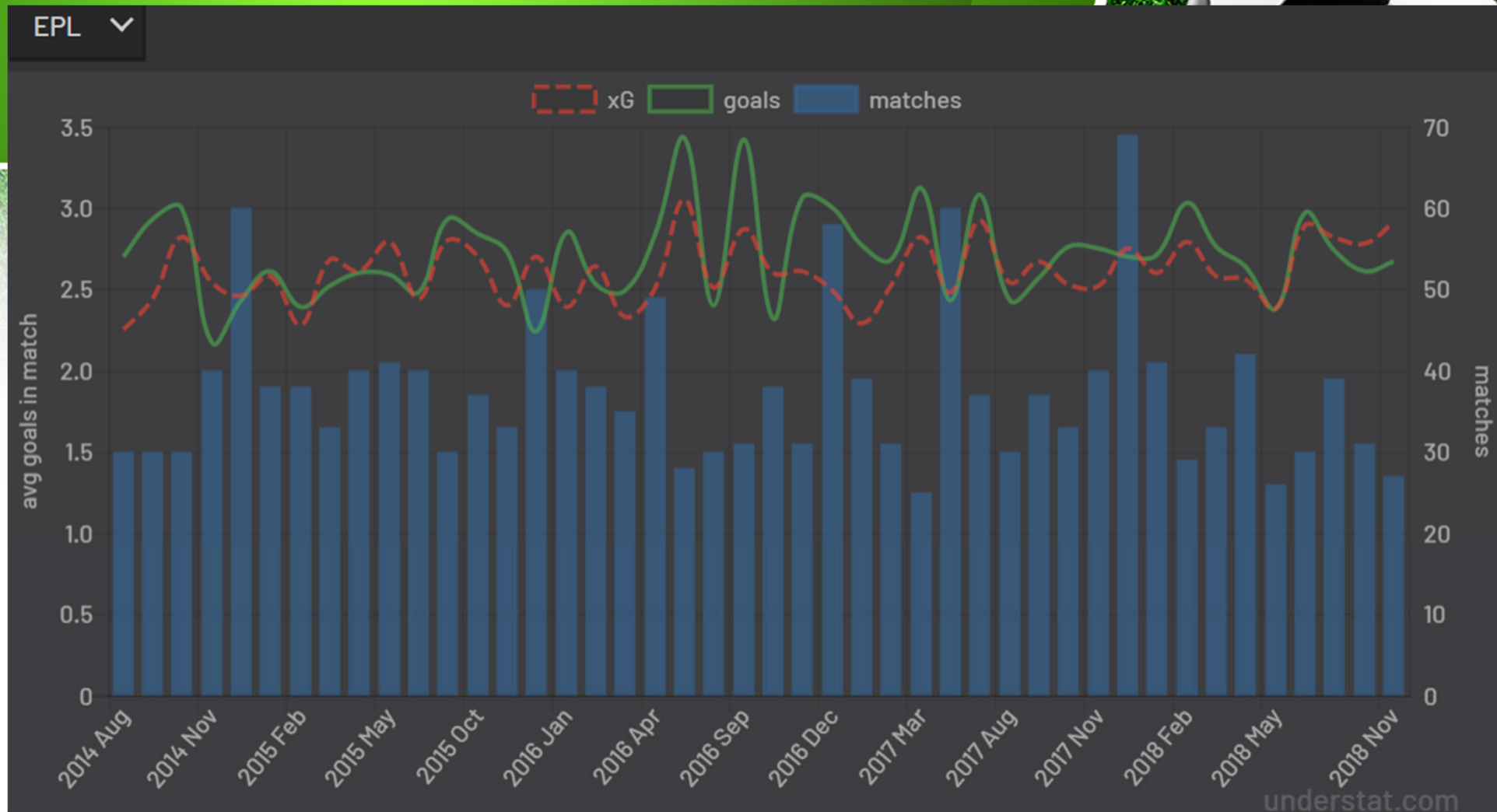
Player Evaluation

Methods (5)

Expected Goals (xG)

- We model every shot
- Response measure: is the probability of a shot resulting in a goal
- The sum of these probabilities will give the xG of a player and a team
- Similar for assists (xA)
- References:
 - <https://www.optasports.com/services/analytics/advanced-metrics/>
 - <https://understat.com/>





Expected Goals (xG): <https://understat.com/>

Player Evaluation

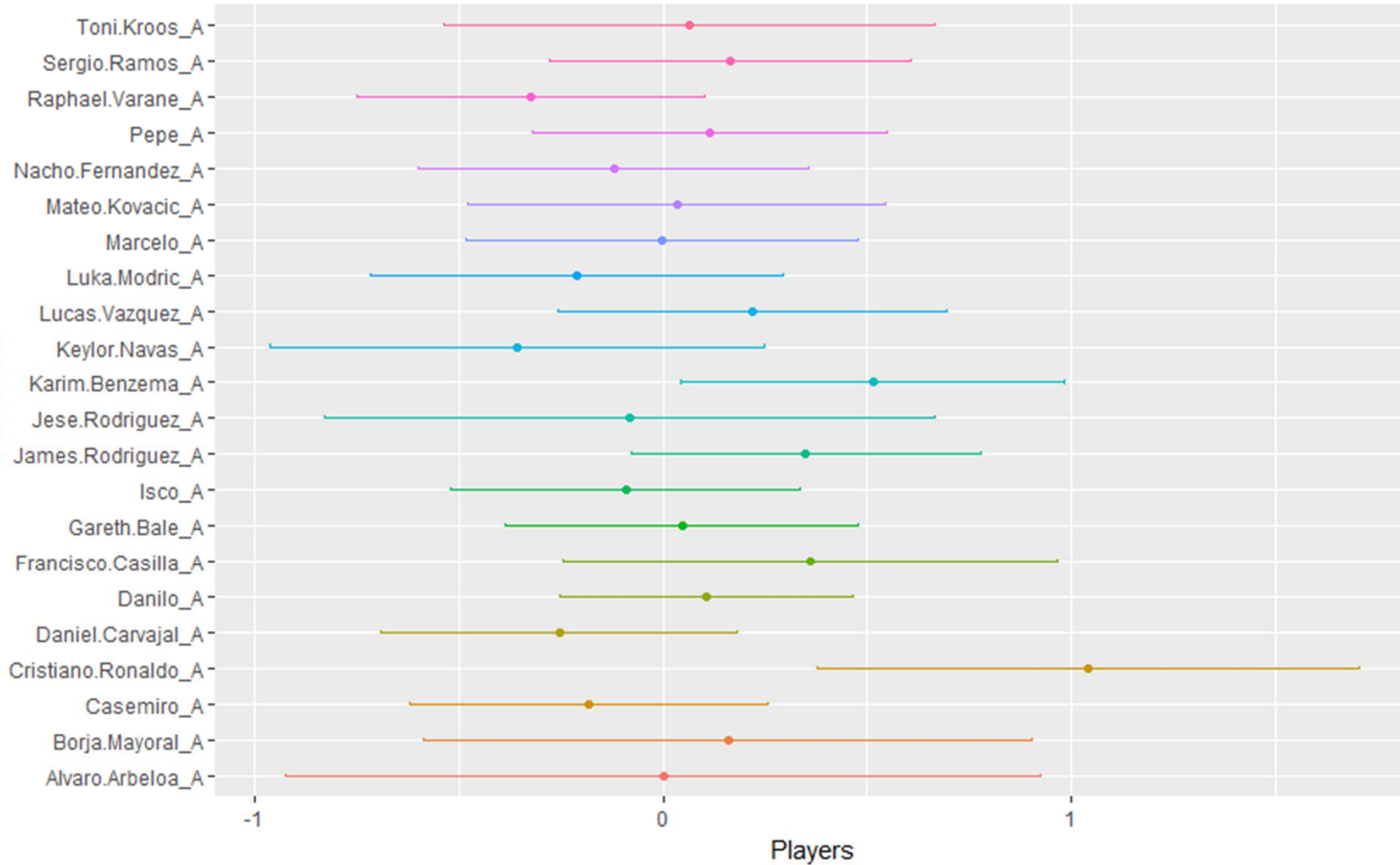
Example of the Simple approach with indicators

- 351 matches of the La Liga Season 2015/2016
- 954 goals (555 goals were scored by home teams, 399 conceded)
- 110 scored by Real Madrid, 34 conceded
- M.Sc. Thesis at AUEB by A. Mourtopallas



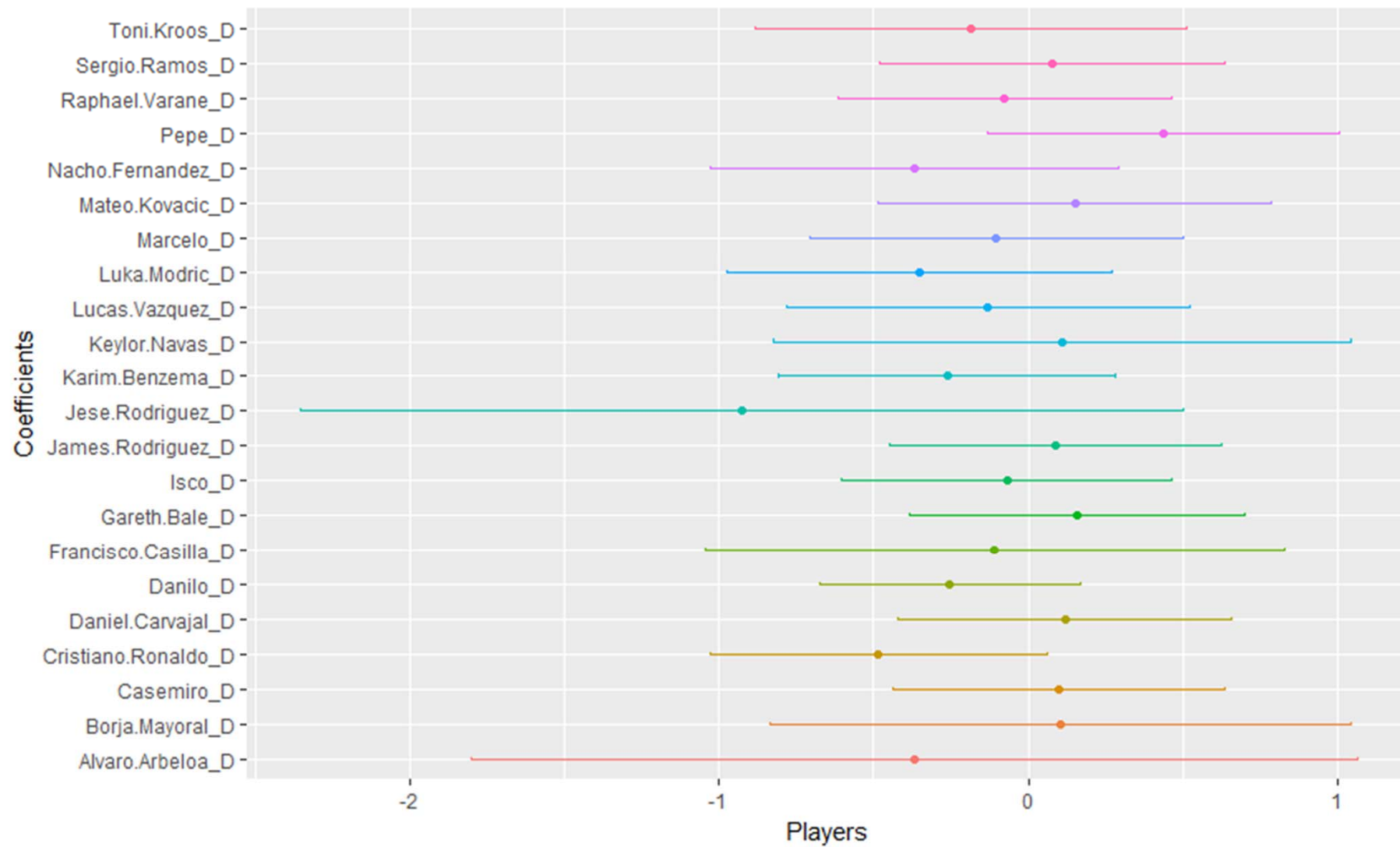
Realmadrid

Players errorbars for the attacking ability



Realmadrid

Players errorbars for the defensive ability



Realmadrid

Impact of players

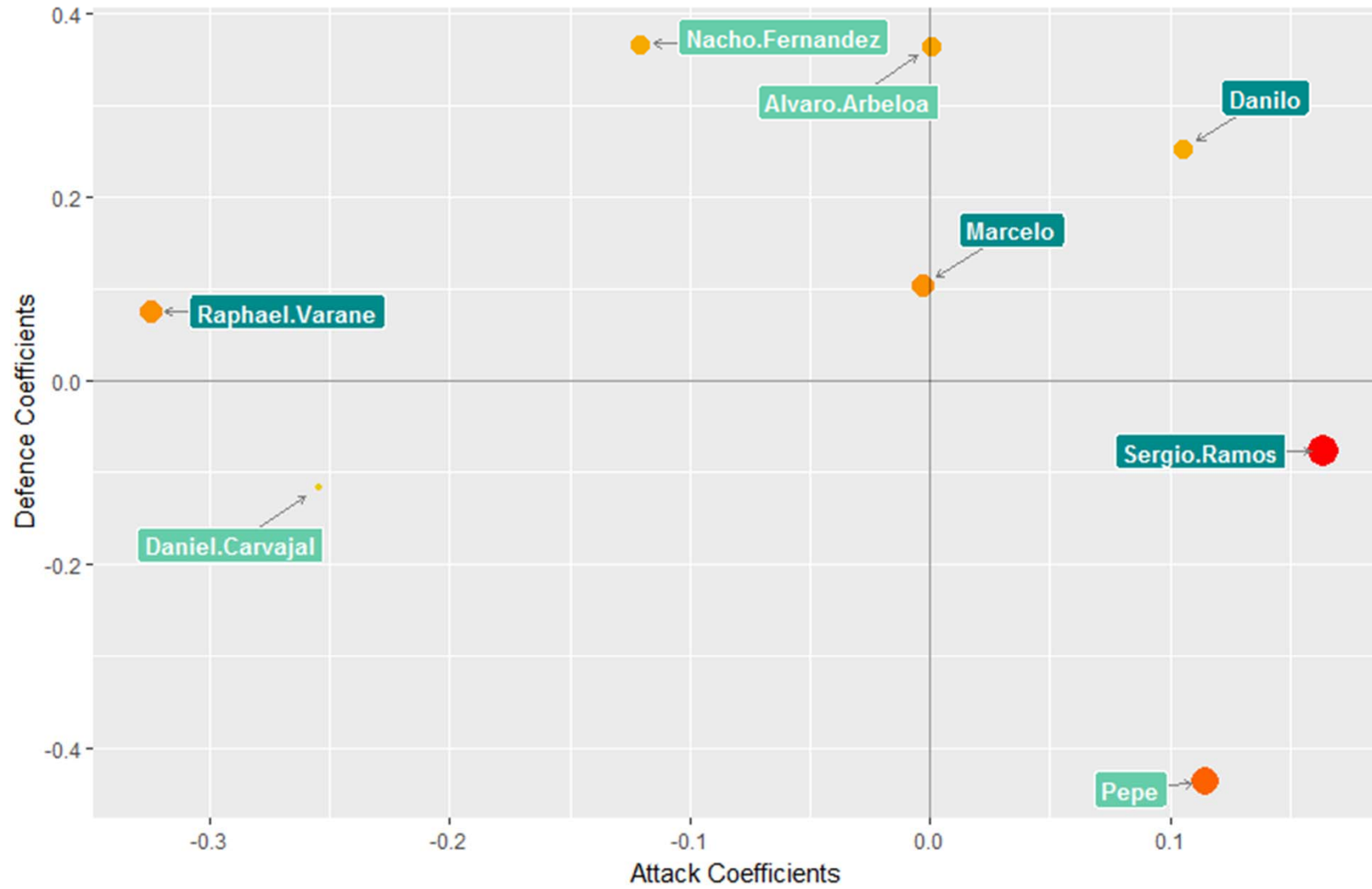


Category

- a Bad at both
- a Better in attack
- a Better in defence
- a Good at both



Impact of defenders

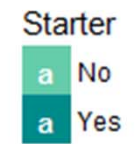
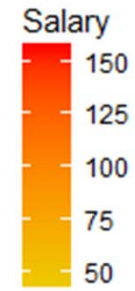
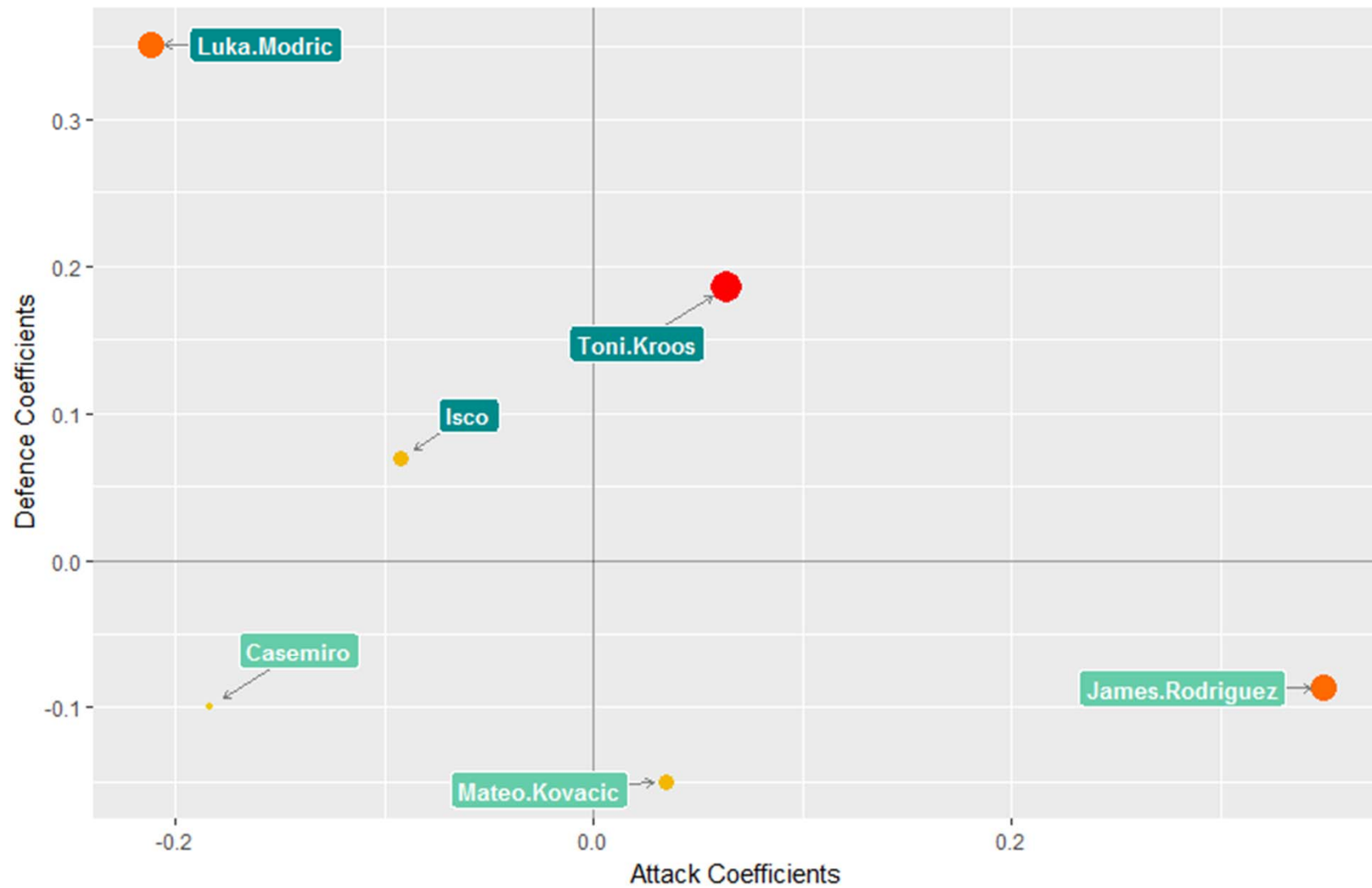


Starter
a No
a Yes

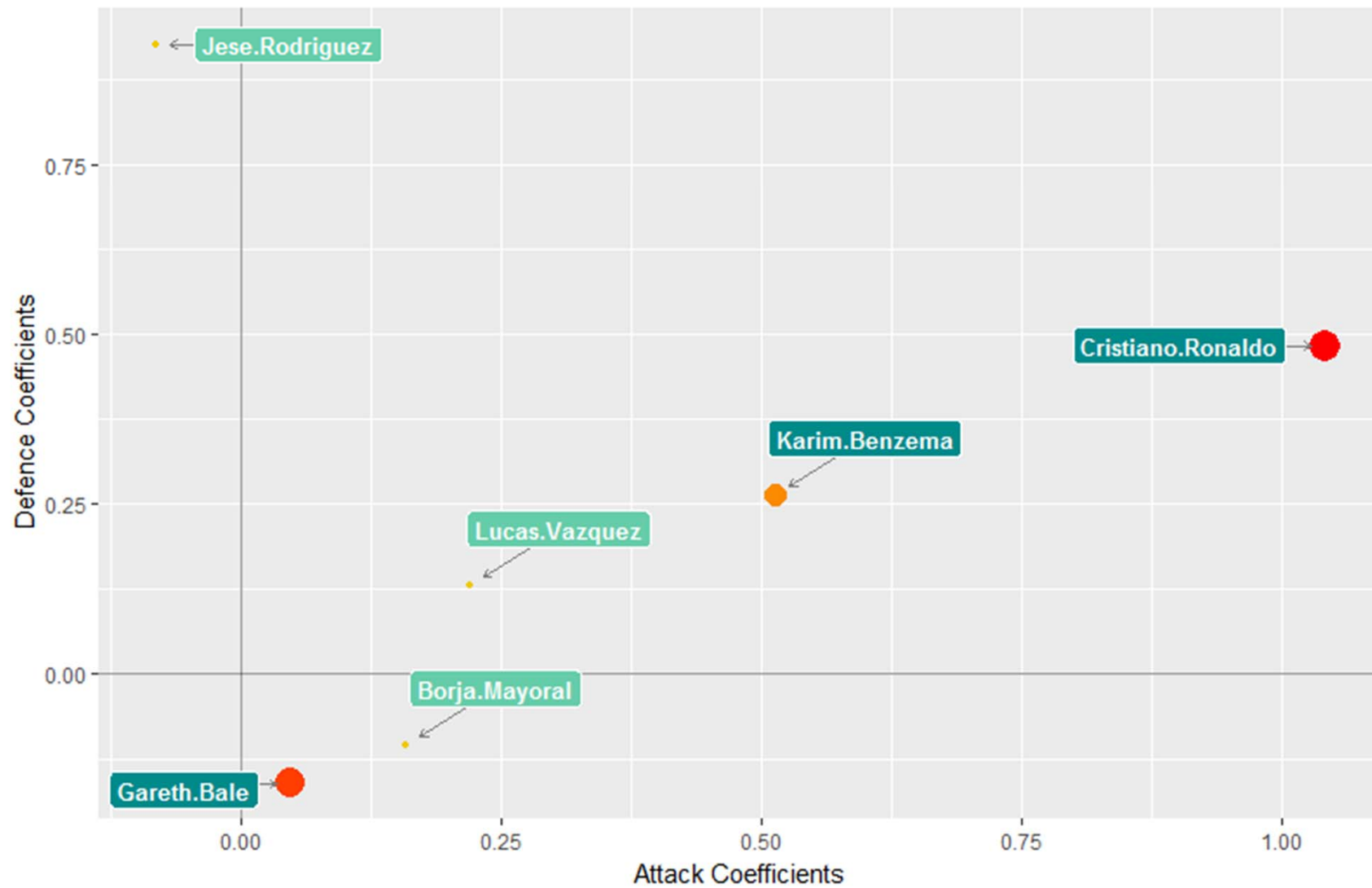
Salary
140
120
100
80
60
40



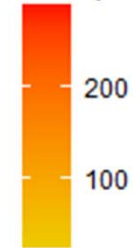
Impact of midfielders



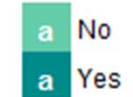
Impact of forwards



Salary



Starter



Realmadrid

Real Madrid 2015/16 Player Evaluation

Conclusions



Cristiano Ronaldo is the key player of the team

Tony Kroos' impact is higher than we may presume



Nacho Fernandez improved since last season (very high def contribution)

Lucas Vasquez is a very promising player (contributed positively in both attack and defensive dimensions with low salary)



Gareth Bale performed less than expected (overprized)

Pepe \Rightarrow low defensive contribution – high salary (overprized?)



Realmadrid

Metrics for physical improvement and training



Aim

- Improve the physical condition of athletes
- Focus on specific skills and measure them
- Avoid injuries
- Improves the team by optimizing allocated training time

Inline game metrics with wearables

The aim is to measure

- Movement of players in the game
- Speed and coverage
- Physical condition
- Physical and tactics performance

It helps

- Evaluate the performance of players and teams within a game
- The manager to decide formation and substitutions



League and Contest Scheduling

AIM

- Fair scheduling
- Eliminate bias due to the sequence of games
- Strengthen competitiveness (related with next slides)
- Incorporate constraints (incl. other sports, safety issues, other events, tv requirements etc.)

HOW?

- Using Operational Research and optimization methods
- Hybrid search methods
- Validate using simulation methods from Statistical models



Competitive Balance

- A balanced league increases the interest of the fans and improves the athletic product
- The notion of a balanced league is not yet very well defined
 - Equal Strength between all teams? or
 - Equal Strength between best teams (or the teams with the highest number of fans?)

Sports Economics & Competitive balance

What league do we want to see?

- All fans like the fact that a weaker team occasionally wins a game or a league

- May neutral fans follow the weakest team
e.g. Greece in Euro 2004

But

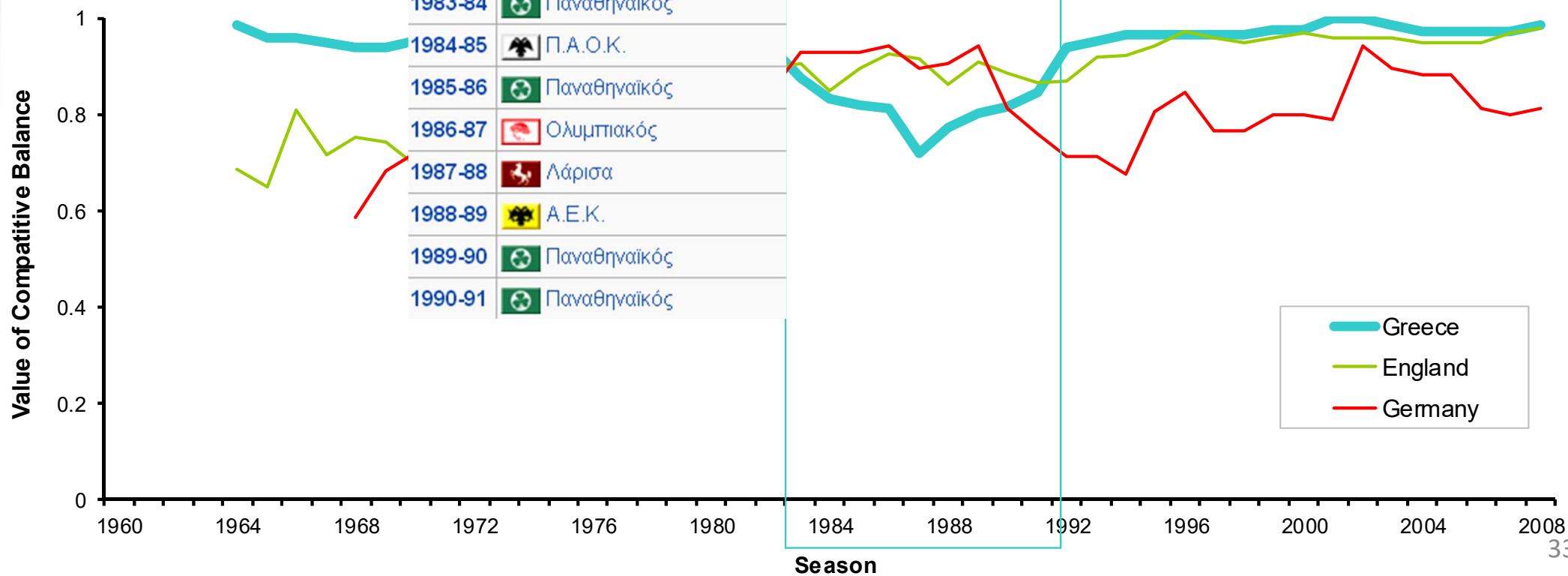
- They do not like their team to loose
- They like or they are willing to pay an expensive ticket to see a final with high ranked and expensive teams
e.g. Bayern-Barcelona



Moving Averages of lag five for DN_t (Champion) from 1959-20

1991-92	A.E.K.	2000-01	Ολυμπιακός
1992-93	A.E.K.	2001-02	Ολυμπιακός
1993-94	A.E.K.	2002-03	Ολυμπιακός
1994-95	Παναθηναϊκός	2003-04	Παναθηναϊκός
1995-96	Παναθηναϊκός	2004-05	Ολυμπιακός
1996-97	Ολυμπιακός	2005-06	Ολυμπιακός
1997-98	Ολυμπιακός	2006-07	Ολυμπιακός
1998-99	Ολυμπιακός	2007-08	Ολυμπιακός
1999-00	Ολυμπιακός	2008-09	Ολυμπιακός

1983-84	Παναθηναϊκός
1984-85	Π.Α.Ο.Κ.
1985-86	Παναθηναϊκός
1986-87	Ολυμπιακός
1987-88	Λάρισα
1988-89	A.E.K.
1989-90	Παναθηναϊκός
1990-91	Παναθηναϊκός



ManU won 13 out of 17 leagues for the period 1992-2009 and it was not ranked in lower position than 3rd.

3 cases in England ⇒ promoted team ⇒ won the championship:
Ipswich (1961) & Nottingham (1997) & Leicester (2015-16 – not in the Figure)

1992-93	Manchester United	+	W
1993-94	Manchester United	+	W
1994-95	Manchester United	+	RU
1995-96	Manchester United	+	W
1996-97	Manchester United	+	W
1997-98	Manchester United	+	RU
1998-99	Manchester United	+	W
1999-2000	Manchester United	+	W
2000-01	Manchester United	+	W
2001-02	Manchester United	+	3rd
2002-03	Manchester United	+	W
2003-04	Manchester United	+	3rd
2004-05	Manchester United	+	3rd
2005-06	Manchester United	+	RU
2006-07	Manchester United	+	W
2007-08	Manchester United	+	W
2008-09	Manchester United	+	W

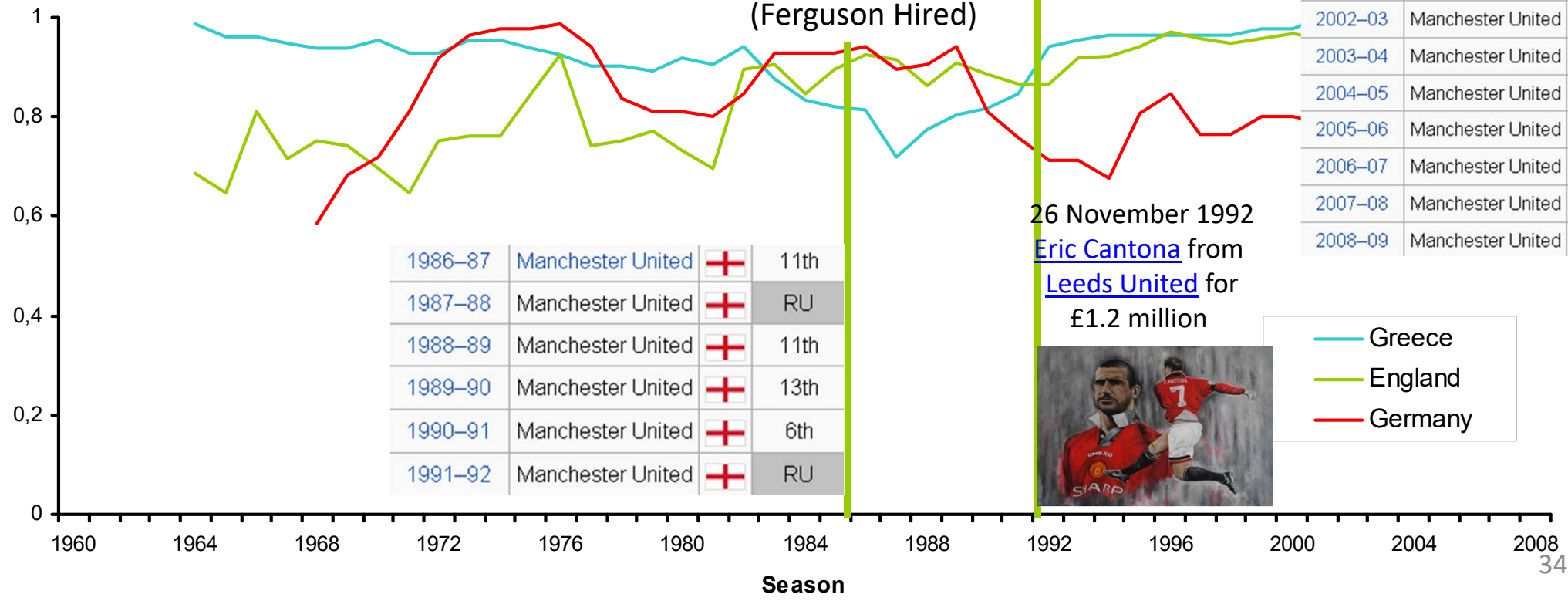
Moving Averages of lag five for DN_1 (Champion) from 1959-2008

6 November 1986
(Ferguson Hired)

UND Wins the title

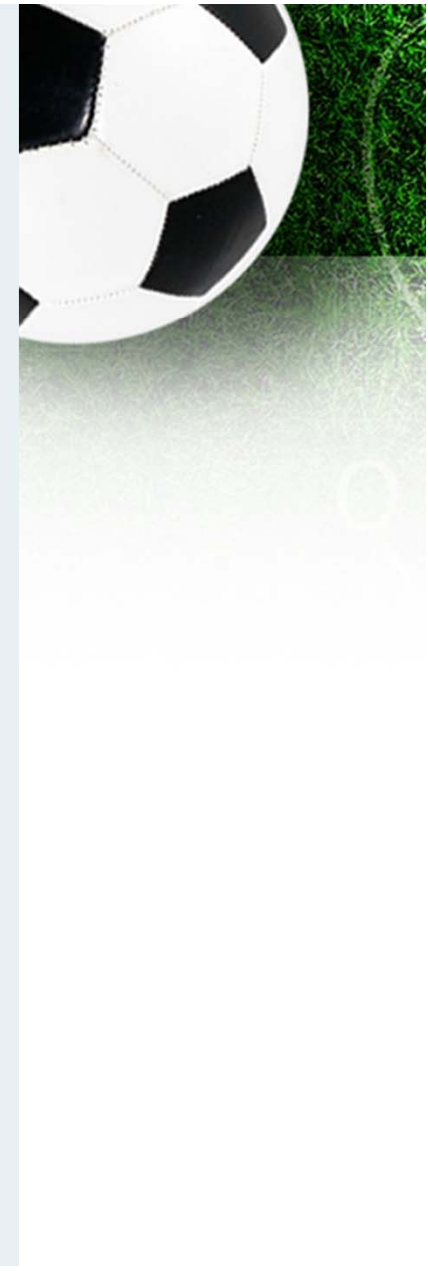
26 November 1992
[Eric Cantona](#) from
[Leeds United](#) for
£1.2 million

Value of Competitive Balance



— Greece
— England
— Germany

Premier League after 13 games of the 2015/16 season (when Leicester won)



Premier League after 13 games
of the 2015/16 season
(when Leicester won)



1968/1969		FC Bayern München
1967/1968		1. FC Nürnberg
1966/1967		Eintracht Braunschweig
1965/1966		TSV 1860 München
1964/1965		SV Werder Bremen
1963/1964		1. FC Köln

1984/1985		FC Bayern München
1983/1984		VfB Stuttgart
1982/1983		Hamburger SV
1981/1982		Hamburger SV
1980/1981		FC Bayern München
1979/1980		FC Bayern München
1978/1979		Hamburger SV

1996/1997		FC Bayern München
1995/1996		Borussia Dortmund
1994/1995		Borussia Dortmund
1993/1994		FC Bayern München
1992/1993		SV Werder Bremen
1991/1992		VfB Stuttgart
1990/1991		1. FC Kaiserslautern
1989/1990		FC Bayern München

2008/2009		VfL Wolfsburg
2007/2008		FC Bayern München
2006/2007		VfB Stuttgart
2005/2006		FC Bayern München
2004/2005		FC Bayern München
2003/2004		SV Werder Bremen

Moving Averages of lag five for DN_1 (Champion) 2008



One case => promoted team => won the championship: Kaiserslautern in 1998

Sports Economics & Competitive balance

How to design Knockout Tournaments?

- Do we support the stronger or the weakest teams?

We do not wish to see

- **many** strong teams to be disqualified early
- Two weak or not popular teams in the final

We do wish to see

- **Some** strong teams to be disqualified early
- Some weak teams to qualify further against all odds

Sports Economics & Competitive balance



In round-robin contests (National leagues)?

- Do we support the stronger or weakest teams?
- Small or large leagues?
- Playoffs?
- Give more money to strong teams (reward) or to weak teams (balance)?
- What about promotion/relegation rules (refreshes the interest or just recycles bad teams?)

We do wish to see

- A large enough group of teams to be close and compete for the championship
- A large enough group of teams to be close and compete for European tickets

We do not wish to see

- A team having big margin of points from all the rest (so the champion is known early)
- Teams with low number of points so they are not competitive (early relegation)
- Teams with economic problems

Sports Economics & Competitive balance

For UEFA Champions League

- Does it need improvement?
- Not metrics to measure balance
- Big discussion of how to reward teams and share income
- Closed or Open League?
- How many teams from each National League/Country
- The current income share and reward system destroys the balance in National teams in second ranked leagues like Greece.

Concluding remarks

To conclude with

- **Prediction** is important for fans (in terms of betting) \Rightarrow increases profits of bet companies and interest for the sport product (in macro perspective)
- **Inline prediction** is important for fans (in terms of betting) \Rightarrow increases profits of bet companies and interest for the sport product (Media – TV, Radio, Internet).



Concluding remarks

- **Player performance and evaluation** \Rightarrow Of main interest for: the fans (Player Ranking), Teams (Scouting, Future Performance and Value), Companies (Sponsoring), Players (A lot of money from all previous), Coaches/Managers (Selection of better players)
- **Physical Measurements** (Training and Games): It is related with player evaluation. Main value to help managers/coaches to improve their teams. In macro perspective also the teams financial position is also improving.
- **Scheduling and Competitive Balance**: More Fair and Balanced contests lead to better product and more profit.





That's all Folks!

**THANK
YOU**



**NO Matter
How Many
Goals
You
Save
People Always
Remember
The One You
Miss.**