

# Titolo del progetto

Domenico Plantamura, Eduardo David Lotto, Manuel D’Alterio Grazioli, Gabriele Fugagnoli

## Contents

<b>Introduction</b>	<b>1</b>
Objective of the project . . . . .	1
Steps followed . . . . .	1
Data collection . . . . .	1
Data integration and cleaning . . . . .	4
Data exploration . . . . .	4
Data analysis and interpretation . . . . .	6
Models . . . . .	6
Salaries analysis by position . . . . .	21
Conclusion . . . . .	31

## Introduction

### Objective of the project

Our goal is to investigate whether the salaries earned by the NBA players during the 2023-2024 season are fair in proportion to their performance during the current year’s Regular Season. To analyse performance, we selected several statistics: from the most common such as points, rebounds, assists to advanced metrics like Usage, Player Impact Estimated and Winning Shares. The idea is to explore the relationships between salaries and performance through different models in order to understand what kind of relationship there is and which model best fits the data. Then, we compared actual salaries with those predicted by our models to find out which players (according to the models) are the most overpaid or underpaid. In the end, we analysed the players by separating them by role (considering centers, forwards and guards separately) and constructed specific models for each position. The aim was to study similarities and differences of the specific models among themselves and with respect to the general ones implemented in the previous phase.

### Steps followed

To perform our analysis we followed these steps:

1. Data collection;

2. Data exploration;
3. Analysis;
4. Interpretation.

## Data collection

We performed a web scraping operation from the Official NBA Stats website, from which we collected most of the stats. Additionally, we downloaded data about the salaries from Hoopshype and other stats of interest from Basketball reference. All data concerns the 2023-2024 NBA Regular Season.

### Why consider only Regular Season data?

Considering only data about Regular Season without considering players performance during playoffs limits a bit the potential of our analysis. On one hand, it's reasonable to infer that player performance during playoffs should have an important weight in determining his salary. On the other hand, considering playoffs in the analysis carries different issues.

There are teams (and consequently players) that go further than others: 14 out of 30 teams can't qualify for the playoffs. For the teams which qualify, playoff stats are calculated on a number of games that could differ greatly between different teams (e.g. if a team loses in the first round, it plays from 4 to 7 games. If a team reaches the finals, it plays from 16 to 28 games). During Regular Season every team plays 82 games.

Additionally, coaches usually rotate players at their disposal in a different way during playoffs: for instance, during NBA Regular Season approximately 10-12 players for each team take part in the game; during playoffs it is not uncommon to observe only 7-8 players that come into play for each team. Furthermore, usually in a playoff game the best players are more involved compared to Regular Season games. It means that, first of all, they play several more minutes. Moreover, they have the ball in their hands for a lot of time and consequently their stats grow a lot; hence, it could happen that few players record a large part of the entire team's statistics. Considering this, including playoffs data in the analysis could lead to an overestimation of performance of 2-3 players and to an underestimation of the performance of the rest of the team.

All in all, it is undeniable that playoffs are a fundamental part of the season. It is also obvious that if a player has more responsibilities in that phase he probably deserves a higher salary. But we think that for the purposes of our analysis, the addition of statistics collected on a small sample of matches, different for practically every team, with highly polarized data between the various players may lead to biases if not handled properly.

We think that considering only the Regular Season, although leading to a limited analysis, may be sufficient to grasp the main relationships between salaries and performance.

## Glossary

- **PLAYER NAME**: name of a player;
- **SALARY**: salary earned by a player for 2023-2024 season (collected from Hoopshype);
- **AGE**: age of a player;
- **POS**: "Position", the playing position of a player.

### Traditional stats (collected from the NBA website)

- **GP**: “Games played”, the number of games played by a player during the 2023-2024 Regular Season;
- **FG\_PCT**: “Field Goal Percentage”, the percentage of field goal attempts that a player makes. Formula:  $(FGM)/(FGA)$ ;
- **FG3\_PCT**: “3 Points “Field Goal Percentage”, the percentage of 3pt field goal attempts that a player makes;
- **FT\_PCT**: “Free throws Percentage”, the percentage of free throws attempts that a player makes;
- **OREB**: “Offensive Rebounds”, the number of rebounds a player or team has collected while they were on offense;
- **DREB**: “Defensive Rebounds”, the number of rebounds a player or team has collected while they were on defense;
- **REB**: “Rebounds”, a rebound occurs when a player recovers the ball after a missed shot. This statistic is the number of total rebounds a player has collected on either offense or defense;
- **AST**: “Assists”, the number of assists (passes that lead directly to a made basket) by a player;
- **TOV**: “Turnovers”, a turnover occurs when a player on offense loses the ball to the defense;
- **STL**: “Steals”, number of times a defensive player takes the ball from a player on offense, causing a turnover;
- **BLK**: “Blocks”, a block occurs when an offensive player attempts a shot, and the defense player tips the ball, blocking their chance to score;
- **BLKA**: “Blocks Against”, The number of shots attempted by a player or team that are blocked by a defender
- **PF**: “Personal fouls”, the number of personal fouls a player or team committed;
- **PFD**: “Personal fouls drawn”, the number of personal fouls that are drawn by a player or team;
- **PTS**: “Points”, the number of points scored by a player;
- **MIN**: “Minutes played”, number of minutes played by a player during the 2023-2024 Regular Season;
- **MIN\_G**: “Minutes played per game”.

### Advanced stats (collected from the NBA website)

- **OFF\_RATING**: “Offensive Rating”, measures a team’s points scored per 100 possessions while a player is on the court. Formula:  $100*((Points)/(POSS))$ ;
- **DEF\_RATING**: “Defensive Rating”, the number of points per 100 possessions that the team allows while a player is on the court. Formula:  $100*((Opp\ Points)/(Opp\ POSS))$ ;
- **NET\_RATING**: “Net Rating”, Measures a team’s point differential per 100 possessions while a player is on the court. Formula:  $OFFRTG - DEFRTG$ ;
- **AST\_TO**: “Assist to Turnover Ratio”, the number of assists for a player compared to the number of turnovers committed;
- **TS\_PCT**: “True Shooting Percentage”, a shooting percentage that factors in the value of three-point field goals and free throws in addition to conventional two-point field goals. Formula:  $Points / [2(Field\ Goals\ Attempted + 0.44Free\ Throws\ Attempted)]$ ;
- **USG\_PCT**: “Usage Percentage”, the percentage of team plays used by a player when they are on the floor. Formula:  $(FGA + Possession\ Ending\ FTA + TO) / POSS$ ;
- **PIE**: “Player Impact Estimate”, measures a player’s overall statistical contribution against the total statistics in games they play in. PIE yields results which are comparable to other advanced statistics (e.g. PER) using a simple formula. Formula:  $(PTS + FGM + FTM - FGA - FTA + DREB + (.5 * OREB) + AST + STL + (.5 * BLK) - PF - TO) / (GmPTS + GmFGM + GmFTM - GmFGA - GmFTA + GmDREB + (.5 * GmOREB) + GmAST + GmSTL + (.5 * GmBLK) - GmPF - GmTO)$ .

The stats below are collected from Basketball Reference:

- **WS**: “Win Shares”, attempts to divvy up credit for team success to the individuals on the team. It is calculated using player, team and league-wide statistics and the sum of player win shares on a given

team will be roughly equal to that team's win total for the season (more details on the Basketball Reference page);

- **BPM**: “Box Plus/Minus”, a box score estimate of the points per 100 possessions that a player contributed above a league-average player, translated to an average team;
- **VORP**: “Value Over Replacement Player”, a box score estimate of the points per 100 TEAM possessions that a player contributed above a replacement-level (-2.0) player, translated to an average team and prorated to an 82-game season. Multiply by 2.70 to convert to wins over replacement.

BPM and VORP are calculated per 100 possessions; MIN and WS are calculated over the whole Regular Season, MIN\_G is calculated per game. The other stats are considered per 48 minutes.

### Why statistics per 48 minutes?

Considering most statistics projected over 48 minutes avoids overestimating performance for players who play, on average, more minutes in a game. In this way we think that the contribution of each player is fairly evaluated and not distorted by the minutes played.

## Data integration and cleaning

Once we had obtained the tables of interest, we selected from each table the statistics useful for analysis (those given in the glossary) and then merged the slices of the various datasets, removing all the players who played less than 480 minutes during the entire Regular Season.

```
data_traditional_tot <- data_traditional_tot[data_traditional_tot$MIN > 480, ]

final_dataset <- merge(data_salary, data_traditional_per48, by = "PLAYER_NAME", all = TRUE)
final_dataset <- merge(final_dataset, data_advanced, by = "PLAYER_NAME", all = TRUE)
final_dataset <- merge(final_dataset, data_miscellaneous, by = "PLAYER_NAME", all = TRUE)
final_dataset <- merge(final_dataset, data_traditional_tot, by = "PLAYER_NAME", all = TRUE)
final_dataset <- merge(final_dataset, data_vorp, by = "PLAYER_NAME", all = TRUE)
```

The reason why we selected players with at least 480 minutes played is that we wanted to avoid considering stats taken on a too small amount of minutes. After these operation, the final dataset consists of 360 rows and 31 columns.

At this stage, we cleaned the data following these other steps:

- NA removal;
- Matching players' names;
- Transforming the Salary column into a numeric one;
- Putting the players' name as row names for the dataset and thus removing the PLAYER\_NAME column.

## Data exploration

Before studying the data with formal models, we got an overview through an exploratory data analysis. For the first part of our analysis we used only numeric variables, so the categorical parameter Pos, which you can see on the table below, was removed from the dataset at this stage.

	Salary	AGE	GP	FG_PCT	FG3_PCT	FT_PCT	OREB	DREB	REB	AST
Aaron Gordon	22266182	28	73	0.556	0.290	0.658	3.6	6.2	9.8	5.4

	Salary	AGE	GP	FG_PCT	FG3_PCT	FT_PCT	OREB	DREB	REB	AST
Aaron Holiday	2346614	27	78	0.446	0.387	0.921	0.9	3.8	4.7	5.3
Aaron Nesmith	5634257	24	72	0.496	0.419	0.781	1.5	5.1	6.6	2.6
Aaron Wiggins	1836096	25	78	0.562	0.492	0.789	2.3	4.9	7.3	3.4
Al Horford	10000000	37	65	0.511	0.419	0.867	2.3	9.1	11.4	4.6

TOV	STL	BLK	BLKA	PF	PTS	OFF_RATING	DEF_RATING	NET_RATING	AST_TO
2.2	1.2	0.9	1.2	3.0	21.2	119.8	111.1	8.7	2.47
2.0	1.6	0.2	0.8	4.7	19.4	110.5	107.6	2.9	2.64
1.5	1.6	1.2	1.2	5.8	21.1	119.3	115.0	4.3	1.69
2.2	2.2	0.7	1.3	3.6	21.2	115.6	110.0	5.7	1.54
1.3	1.0	1.7	0.3	2.6	15.5	120.9	109.5	11.4	3.50

TS_PCT	USG_PCT	PIE	PFD	MIN	MIN_G	Pos	WS	BPM	VORP
0.607	0.174	0.103	4.7	2296.810	31.46315	PF	7.1	1.3	1.9
0.578	0.158	0.078	2.5	1269.297	16.27303	PG	2.5	-1.5	0.2
0.631	0.158	0.071	3.5	1994.655	27.70354	SF	4.1	-0.5	0.8
0.664	0.163	0.096	2.3	1227.938	15.74280	SG	3.7	0.7	0.8
0.650	0.119	0.105	0.8	1739.797	26.76610	C	6.2	3.6	2.5

Firstly, we perform an analysis, which can be seen in the Figure 1, of the variable Salary, and its logarithmic transformation, that will be the dependent variable in the models.

The boxplot shows that the salary distribution is right skewed, with some outliers in the right side. We expected this kind of distribution, the outliers are the players earning the highest salaries. The histogram also highlights the right skewed distribution. It can be seen that Salary's log transformation reduces the skewness and makes the distribution of the variable closer to normal.

In order to study correlations between the predictors of the model, we used the `corrplot` function (Figure 2).

```
library(corrplot)
corrplot(cor(fd_numeric), method = 'color')
```

Different correlations between the variables emerge from the `corrplot`. With regard to the variable Salary, it is interesting to notice that Salary is positively correlated with PTS and advanced stats like USG\_PCT, BPM and VORP: all of these variables are related to players' shots and point contribution. For what concerns the other variables, there are some obvious correlations: for instance, between variables MIN (total minutes played during the Regular Season) and MIN\_G (minutes played per game) and between variables REB, OREB and DREB (all related to rebounds, with the relation  $REB = OREB + DREB$ ). Additionally, we expected the positive correlation between BPM and VORP because are both related to players point estimation. A strong positive correlation emerges between PTS and USG\_PCT. The usage percentage is "The percentage of team plays used by a player when they are on the floor. Formula:  $(FGA + Possession\ Ending\ FTA + TO) / POSS$ ". Thus, players with a high USG\_PCT often make the last play in an offensive possession (a shot, a free throw or a turnover): it is straightforward that if a player often ends the offensive possession of his team, he has more opportunities to score points. For what concerns the negative correlations, the most interesting are the ones between rebounds variables (OREB, DREB, REB), FT\_PCT and FG3\_PCT. Players that grab a lot of rebounds are usually the tallest ones and these players are not great free throws shooters or 3 point shooters (on average).



## Data analysis and interpretation

### Models

We started creating a linear regression model in order to predict salaries (Figure 3).

```
##
## Call:
## lm(formula = Salary ~ +., data = fd_numeric)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -18450260 -4028989   276645   4003025  20712902
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -12550803   28966256  -0.433   0.6651
## AGE          1057586     93293    11.336 <2e-16 ***
## GP           -21024     118116   -0.178   0.8588
## FG_PCT       35836089   19086013   1.878   0.0613 .
## FG3_PCT      -56045     4195524  -0.013   0.9893
## FT_PCT       975535     6135667   0.159   0.8738
## OREB         4054377     6865112   0.591   0.5552
## DREB         4473997     6846450   0.653   0.5139
## REB         -4315225     6838752  -0.631   0.5285
## AST          -98527      667680  -0.148   0.8828
## TOV          2003183     1516145   1.321   0.1873
## STL          -69046      985541  -0.070   0.9442
## BLK          601287      664109   0.905   0.3659
## BLKA        -2253383     1230646  -1.831   0.0680 .
## PF          -616626      640191  -0.963   0.3362
## PTS          1117890      623630   1.793   0.0740 .
## OFF_RATING   16646653     6955245   2.393   0.0172 *
## DEF_RATING  -16681974     6953008  -2.399   0.0170 *
## NET_RATING  -16610236     6957987  -2.387   0.0175 *
## AST_TO      -252115      978605  -0.258   0.7969
## TS_PCT      -63710004     30332753  -2.100   0.0365 *
## USG_PCT     -40539821     73391644  -0.552   0.5811
## PIE        -131534170    115933751  -1.135   0.2574
## PFD          101829      397847   0.256   0.7981
## MIN          -4774         4689  -1.018   0.3094
## MIN_G        696311      299872   2.322   0.0208 *
## WS          1845668      740418   2.493   0.0132 *
## BPM         -391702      764769  -0.512   0.6089
## VORP         663105      1436430   0.462   0.6446
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6712000 on 331 degrees of freedom
## Multiple R-squared:  0.7081, Adjusted R-squared:  0.6834
## F-statistic: 28.67 on 28 and 331 DF,  p-value: < 2.2e-16

## [1] "MSE of the complete linear model = 4.142347e+13"
```

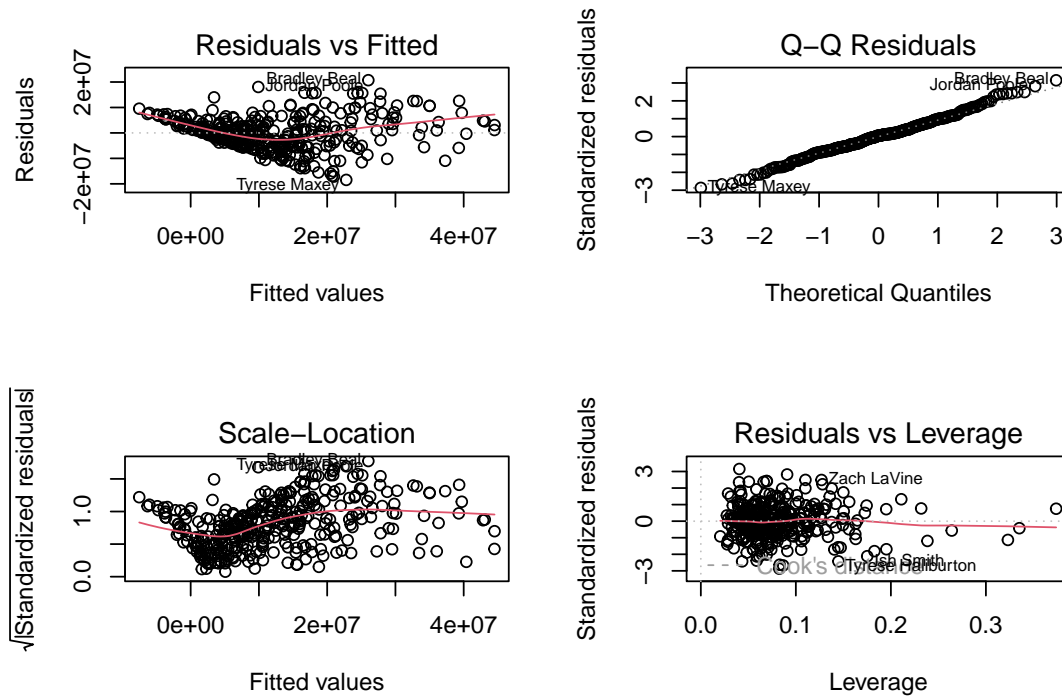


Figure 3: Residuals plot of the complete linear model

The complete model has a good adjusted R-squared of 0.68 and a MSE of  $4.14 \times 10^{13}$ . It emerges that many variables are not significant in determining the response. Through the residual analysis it is noticeable that the relationship between fitted values and residuals is not exactly linear (1st graph). Additionally, in the third graph the points are not included in a band of constant amplitude parallel to the x-axis, hence the homoscedasticity assumption can be doubted.

```
##
## Call:
## lm(formula = log(Salary) ~ +., data = fd_numeric)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.75872 -0.32207  0.02064  0.42623  1.62364
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.258e+01  2.744e+00   4.583 6.49e-06 ***
## AGE          9.669e-02  8.838e-03  10.940 < 2e-16 ***
## GP         -5.087e-04  1.119e-02  -0.045  0.9638
## FG_PCT       3.669e+00  1.808e+00   2.029  0.0432 *
## FG3_PCT     -7.625e-02  3.975e-01  -0.192  0.8480
## FT_PCT     -1.574e-01  5.813e-01  -0.271  0.7867
## OREB         3.180e-01  6.504e-01   0.489  0.6252
## DREB         4.105e-01  6.486e-01   0.633  0.5273
## REB        -3.435e-01  6.479e-01  -0.530  0.5963
## AST         3.954e-02  6.325e-02   0.625  0.5323
```



```

## TOV          1.536e-03  1.436e-01  0.011  0.9915
## STL          5.077e-02  9.337e-02  0.544  0.5870
## BLK          7.242e-02  6.291e-02  1.151  0.2505
## BLKA        -1.802e-01  1.166e-01 -1.545  0.1232
## PF          -5.294e-02  6.065e-02 -0.873  0.3834
## PTS          9.972e-02  5.908e-02  1.688  0.0924 .
## OFF_RATING   1.508e+00  6.589e-01  2.289  0.0227 *
## DEF_RATING  -1.499e+00  6.587e-01 -2.276  0.0235 *
## NET_RATING  -1.504e+00  6.592e-01 -2.282  0.0231 *
## AST_TO      -4.552e-02  9.271e-02 -0.491  0.6237
## TS_PCT      -6.498e+00  2.874e+00 -2.261  0.0244 *
## USG_PCT     -2.881e+00  6.953e+00 -0.414  0.6789
## PIE         -1.715e+01  1.098e+01 -1.562  0.1193
## PFD          3.123e-02  3.769e-02  0.829  0.4079
## MIN         -5.447e-05  4.442e-04 -0.123  0.9025
## MIN_G        5.720e-02  2.841e-02  2.014  0.0449 *
## WS           1.246e-01  7.014e-02  1.777  0.0765 .
## BPM          5.637e-02  7.245e-02  0.778  0.4371
## VORP        -1.639e-01  1.361e-01 -1.204  0.2293
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.6359 on 331 degrees of freedom
## Multiple R-squared:  0.6516, Adjusted R-squared:  0.6222
## F-statistic: 22.11 on 28 and 331 DF,  p-value: < 2.2e-16

```

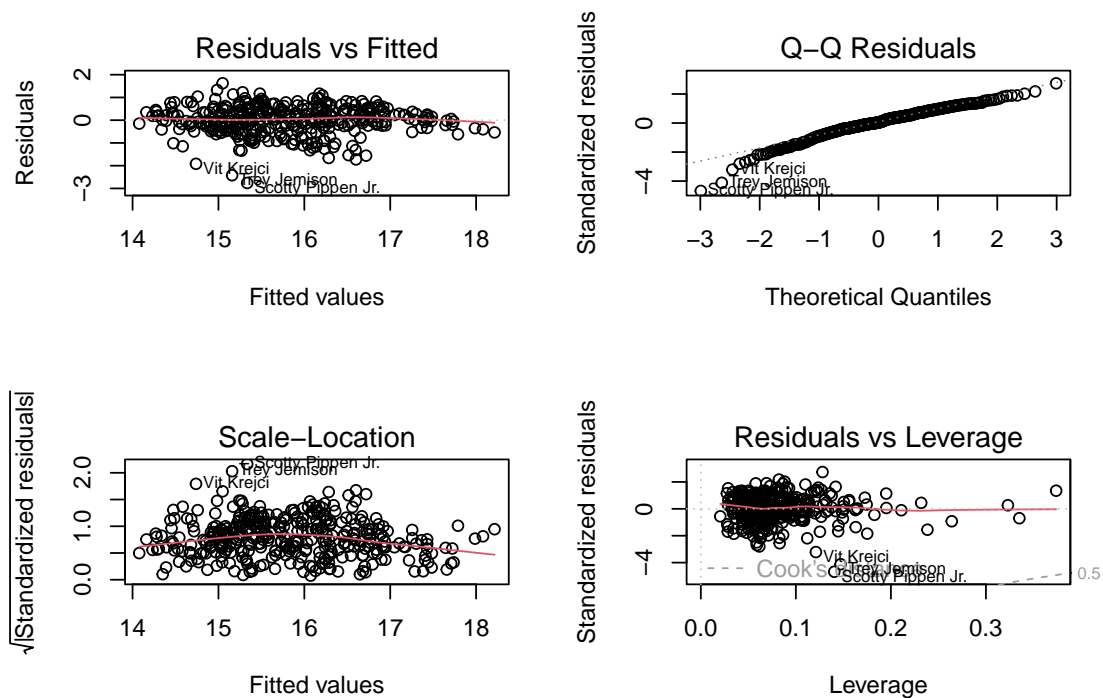


Figure 4: Residuals plot of the complete linear model with logarithmic Salary

```
## [1] "MSE of the complete linear model with logarithmic Salary = 4.703165e+13"
```

With a logarithmic transformation of the dependent variable, the model shows a slightly lower adjusted R-squared (0.62) and a slightly higher MSE ( $4.70e+13$ ). This can be seen in the Figure 4. Applying a logarithmic transformation to the dependent variable **Salary**, the first graph shows a more linear relationship and the third graph allows to infer a more constant variance in the error terms. In both models many variables are not significative in determining the response: for this reason, to avoid a model that is unnecessary complex, we performed a variable selection. A logarithmic transformation of the dependent variable **Salary** will be applied because, although it slightly worsens the performance of the model, it makes the salaries distribution closer to normal, it improves the linearity of the model and it reduces residuals eteroschedasticity.

## Variable selection

We selected a subset of relevant features starting from the predictors used in the complete model in order to have a simpler model that is easier to interpret, without redundant variables and less prone to overfitting. To do so, we used the `regsubsets` function which performs best subset selection by identifying the best model that contains a given number of predictors, according to the RSS metric. We set the function to return results up to the best 28-variables model.

To find the best balance between model simplicity and precision, we evaluated the number of parameters to be included in the model through Mallor's Cp, BIC and Adjusted R-squared.

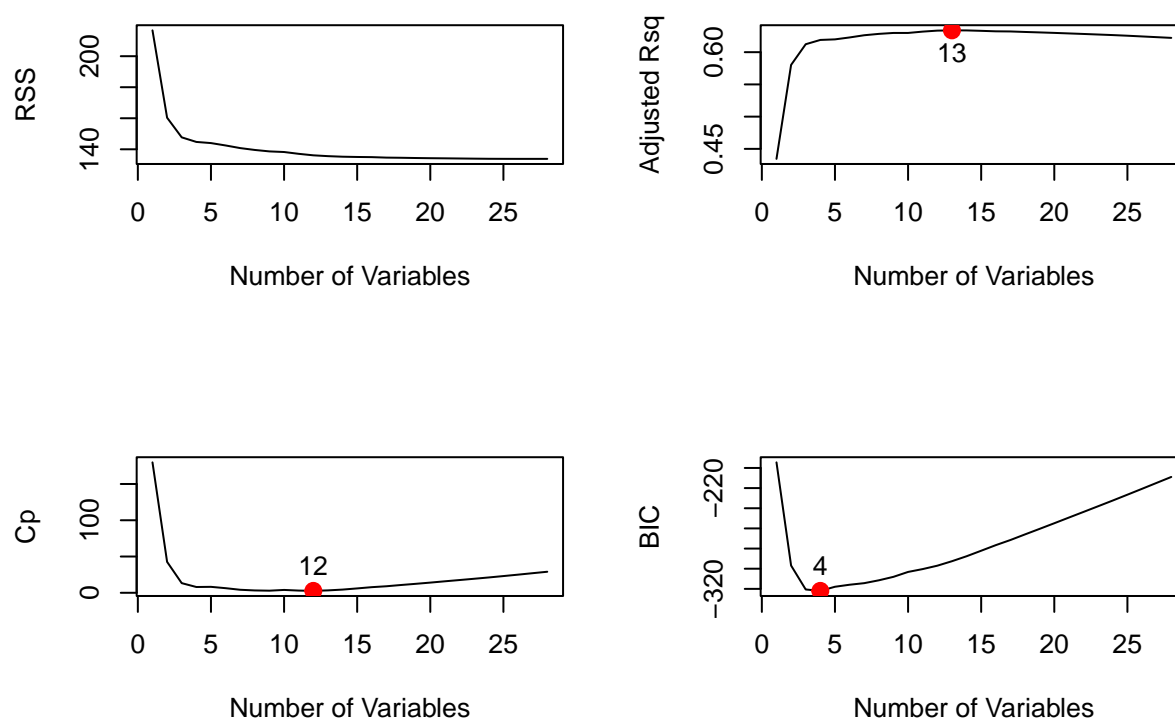


Figure 5: Evaluation of the number of parameters through RSS, Adjusted R-squared, Mallor's Cp and BIC

Considering Mallor's Cp, the best number of parameters for our model is 12. This result can be seen in the

Figure 5. We obtained the list of parameters from the `regsubset` function to get the best model with 12 parameters.

```
##
## Call:
## lm(formula = selected.formula, data = fd_numeric)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.6476 -0.3193  0.0482  0.4315  1.5814
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 10.852025   1.626347   6.673 9.95e-11 ***
## AGE          0.095227   0.008503  11.199 < 2e-16 ***
## FG_PCT       3.445692   1.155704   2.981 0.00307 **
## BLK          0.083159   0.049934   1.665 0.09674 .
## BLKA        -0.167470   0.106759  -1.569 0.11764
## PTS          0.058816   0.011460   5.132 4.78e-07 ***
## OFF_RATING   1.531106   0.629095   2.434 0.01544 *
## DEF_RATING  -1.513509   0.629188  -2.405 0.01667 *
## NET_RATING  -1.520988   0.629121  -2.418 0.01614 *
## TS_PCT      -5.839664   1.438599  -4.059 6.09e-05 ***
## PIE         -6.258145   2.750650  -2.275 0.02351 *
## MIN_G        0.059152   0.006890   8.585 3.10e-16 ***
## WS           0.053790   0.026074   2.063 0.03986 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.6263 on 347 degrees of freedom
## Multiple R-squared:  0.6457, Adjusted R-squared:  0.6334
## F-statistic: 52.7 on 12 and 347 DF,  p-value: < 2.2e-16
```

The reduced model shows a slightly higher adjusted R-squared, 0.63, compared to the complete logarithmic model (0.62). It means that, despite the lower number of variables, this model fits better the data. Different variables are strongly significant:

- **AGE**: the positive coefficient associated to the variable shows that older players earn, on average, more than young ones. This makes sense since the youngest players in the league, rookies (first year in NBA) and sophomores (second year in NBA), usually earn less in the first years due to particular specifications in their contracts. However, there are also many veterans that sign for very low salaries in order to play with better teams and thus have a chance of winning the title.
- **PTS**: this is quite straightforward. Players who score more points, on average, have higher salaries.
- **TS\_PCT**: for what concerns true shooting percentage, the situation is peculiar. TS\_PCT weights a player's shooting percentages based on the shot type (3-pointer, 2 pointer or free throw). The negative coefficient seems counter intuitive: a better TS\_PCT reflects, on average, a lower salary. A possible explanation is that this metric is high for two players categories. The first one is composed by tall players who take most of their shots near the basket, thus getting a high percentage. The second category is composed by 3-point shooting specialists, because the weight for a 3 point shoot is higher for the metric. These players are crucial into a team, but we can say that they often have a limited role: the former have to score mostly near the basket, the latter from behind the 3-point line. Consequently, it makes sense if the model assigns a lower salary for players with a limited role. Additionally, shooting

percentages are also high for players that shoot only few shots in a game; it is reasonable to think that scoring only few shots it's not enough to earn a high salary. Moreover, it is important to highlight that the best players (the ones that should earn more) attract more attention from opposing defenders, so it is normal that they have more fluctuating shooting percentages than previous mentioned specialists.

- **MIN\_G**: players that play on average more minutes in a game earn, on average, a higher salary.

The variable FG\_PCT is less significative than TS\_PCT, but the coefficient here is positive. Both the stats measure shooting percentages, but FG\_PCT does not weight shots and does not consider free throws. In this way, the previous mentioned effect on 3 point shooting specialists is reduced. It is possible to infer that FG\_PCT represents better, within this model, the positive impact of good shooting percentages on wages.

The variables OFF\_RATING, DEF\_RATING, NET\_RATING, PIE and WS have a level of significance between 0.01 and 0.05. The positive sign of OFF\_RATING and WS coefficients and the negative sign of DEF\_RATING coefficient are in line with what we expected. OFF\_RATING (DEF\_RATING) represents the points scored (conceded) by the team when the player is playing, WS measures the player contribution to the team wins. We didn't expected a negative signs for NET\_RATING (OFF\_RATING - DEF\_RATING) and PIE, that measures the player impact in the game.

For what concerns PIE, the negative sign has different possible explanations: projecting PIE per 48 minutes inflates the metric for players who have a high impact on the game but few minutes played. It considers a lot of stats, even stats that seem to be not significative in determining salary; PIE difference between high salary players and low salary ones is not proportional to the differences in salaries. It is always difficult consider defensive contribution with this kind of metric and it is reasonable to think that defensive contribution plays an important role in determining a players salary. Furthermore, PIE does not consider aspects like leadership and IQ that, as defensive contribution, will certainly have an impact on the salaries. Anyway, beyond all the possible explanations, these unexpected negative signs likely depend from other variables not included in the model.

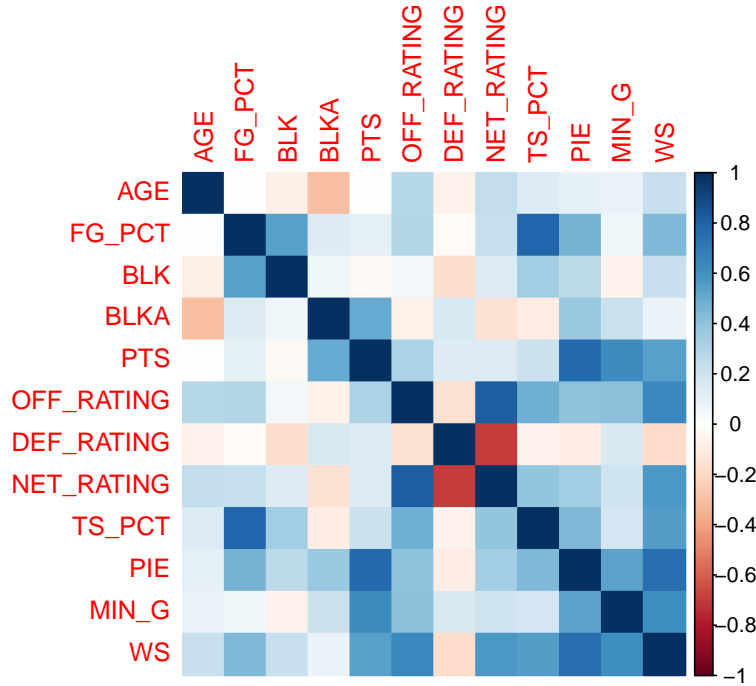


Figure 6: Correlation between dependent variables of the reduced model

**Correlation between selected variables** It can be seen in the Figure 6 that there are, also in this case, different correlations between the dependent variables.

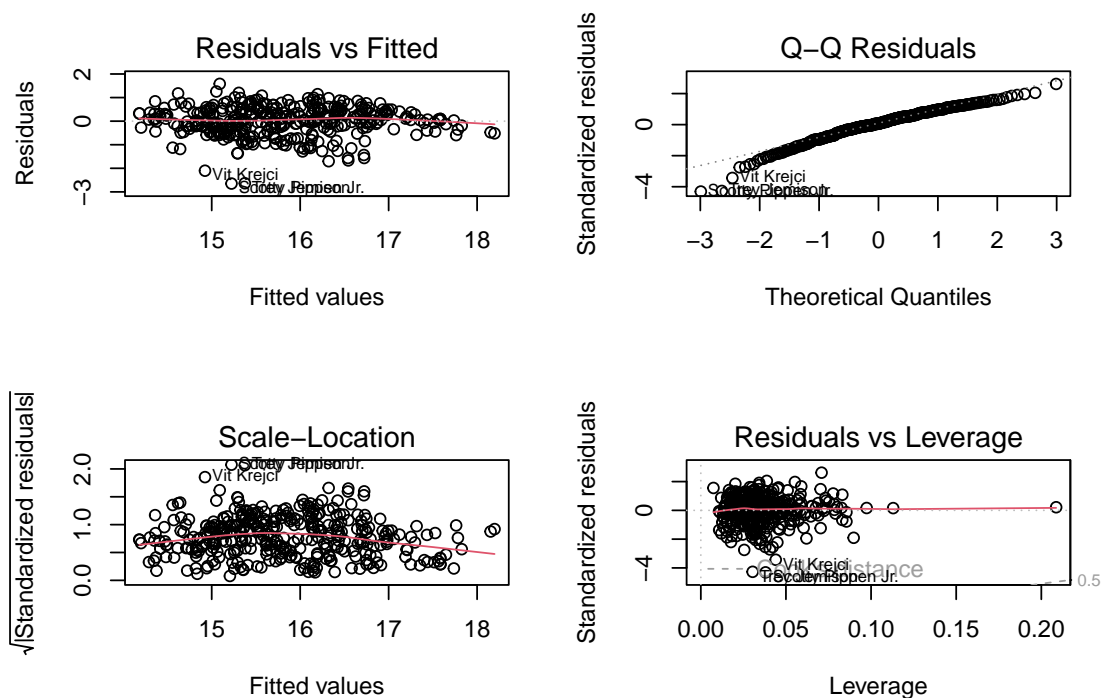


Figure 7: Residual plot of the reduced model with 12 covariates

**Residual analysis** From what can be seen in the Figure 7, the assumptions of the linear model seem to be fulfilled. There are some players who are outliers in each graph: they probably have special contracts (two-way contracts). This means that they usually play in the team's second team (in a so called development league) and occasionally in the first team, so they have really low salaries compared to the league average.

## MSE

```
## [1] "MSE of the redacted linear model with logarithmic Salary = 4.775116e+13"
```

The Mean Squared Error of the reduced model is really close to the one of the complete linear model with the logarithmic transformation of the Salary,  $4.77e+13$  against  $4.70e+13$ . Considering that the complete model has 28 variables and the reduced one 12, the latter model represents quite an improvement.

## Real salaries vs salaries prediction

Table 4: Ten most overpaid players according to the reduced model

	Salary	Predicted salary	Difference
Bradley Beal	\$46,741,590	\$21,945,158	\$24,796,432
Darius Garland	\$34,005,250	\$10,821,355	\$23,183,895

	Salary	Predicted salary	Difference
Zach LaVine	\$40,064,220	\$17,030,435	\$23,033,785
Trae Young	\$40,064,220	\$18,771,385	\$21,292,835
Deandre Ayton	\$32,459,438	\$12,251,585	\$20,207,853
Michael Porter Jr.	\$33,386,850	\$13,910,889	\$19,475,961
Zion Williamson	\$34,005,250	\$14,543,954	\$19,461,296
Karl-Anthony Towns	\$36,016,200	\$17,724,464	\$18,291,736
Jordan Poole	\$27,955,357	\$10,348,315	\$17,607,042
Gordon Hayward	\$31,500,000	\$15,311,898	\$16,188,102

Table 5: Ten most underpaid players according to the reduced model

	Salary	Predicted salary	Difference
LeBron James	\$47,607,350	\$79,971,554	\$32,364,204
Kevin Durant	\$47,649,433	\$76,482,185	\$28,832,752
DeMar DeRozan	\$28,600,000	\$51,810,315	\$23,210,315
Kyrie Irving	\$37,037,037	\$52,868,441	\$15,831,404
Nikola Vucevic	\$18,518,519	\$33,136,916	\$14,618,397
Jalen Brunson	\$26,346,666	\$40,868,798	\$14,522,132
Russell Westbrook	\$3,835,738	\$18,318,987	\$14,483,249
Tyrese Maxey	\$4,343,920	\$18,257,013	\$13,913,093
Kelly Oubre Jr.	\$2,891,467	\$15,308,000	\$12,416,533
Brook Lopez	\$25,000,000	\$37,263,780	\$12,263,780

Here we have a comparison between real salaries and predicted ones. The tables contain, respectively, the 10 most overpaid players and the 10 most underpaid players according to the model. The aim of this comparison is to analyse the major differences between predictions and actual salaries to understand whether, despite a big difference, the model’s predictions seem reasonable.

## MOST OVERPAID PLAYERS

The most overpaid player results to be Bradley Beal. After some brilliant seasons with Washington Wizards in which he was the league top scorer, he signed in 2022 a maximum contract (251 million \$ in 5-years). In Washington he was the best player by far, his statlines in the past years justify the huge contract. In 23-24 he was traded to Phoenix (keeping the same contract) to play with Durant and Booker (two superstars) in a team that was, on the paper, a contender for the title. Beal, being no longer the first offensive option, had a quite different statline compared to the previous years. Additionally, the whole Phoenix Suns team performed worse than expected. These facts are enough to explain that Beal’s 23-24 performance is not in line with his salary.

Darius Garland signed a big contract (near to the maximum) starting from 23-24 season. After showing superstar potential in 22-23, Cleveland Cavaliers renewed his contract with an important salary increase but Garland’s performance decreased in 23-24. He is only 24, the team bet heavily on him taking a weighted risk in order to keep with them a high potential player. This bet didn’t paid in 23-24 season.

Trae Young and Zach Lavine have superstar contracts respectively in Atlanta and Chicago, but they are not carrying their teams as expected. Both players could be traded during this summer.

Regarding Deandre Ayton, he was an amazing prospect but he repeatedly failed to meet expectations at the most important moments. He signed a big contract in 2022 but his performance were not at the same level as the salary. He was traded to Washington (keeping the same contract) but also this year in a different team he did not fulfil expectations.

Zion Williamson and Michael Porter Jr. (especially the former) are young players that in their still short careers have not shown their full potential due to injuries. Their contracts, let's say, consider their potential performance at the top of their form. Jordan Poole had an exploit in the previous seasons playing with a top team, Golden State Warriors, that somehow justifies his salary. He seemed to be ready to carry a team on his own, he was traded to Washington but his first season was a failure.

## MOST UNDERPAID PLAYERS

Lebron James and Kevin Durant have been two of the best players in the league for many years now. Even though, according to our model they should earn much more than the maximum wage. For sure their careers and their performance motivate a high salary, but equally surely they are not underpaid. We think that this overestimation depends in part on the fact that the variable AGE in the model is strongly significant, Lebron James is 39 and Kevin Durant is 35. The same reasoning could apply to Kyrie Irving (32) and especially Demar Derozan (34).

For what concerns Nikola Vucevic, his stats are always more than respectable. His salary is lower than the expected probably because he seems to lack characteristics not included in the model or generally difficult to quantify such as defense, leadership and consistency at key moments of the season.

Jalen Brunson has shown this year that he is one of the best players in the NBA after being somewhat underrated in the past years. We expected the difference between his predicted and actual salary. Very similar the situation of Tyrese Maxey, in the last year of his rookie contract. He has shown by his performances that he is worth much more than his salary says.

Russell Westbrook is in the waning phase of his career. On the expiry of his last superstar contract, no team in the league offered him a comparable salary (he earned 47 millions in 2022). Consequently, he accepted a 3.8 millions salary (veteran minimum contract) to play with Los Angeles Clippers. For sure he is no longer a player worth 47 millions, but he is not worth 3.8 millions either. Our model interprets pretty well the situation, stating that Westbrook should earn a 18.3 millions salary: not a superstar one, but not a minimum wage either.

Given the presence of correlations between the independent variables, the presence of multicollinearity is likely. For this reason, we decided to implement models that perform well when the variables are collinear such as Ridge regression and Lasso regression. In the next paragraphs we want to see if the performances of these models are better than that of the models seen so far.

## Ridge regression

The subset selection method uses least squares to fit a linear model with a subset of the predictors. On the other hand, ridge regression does not select a subset of the coefficients  $\beta_j$  of the model, but it fits a model with all  $p$  predictors adding a term  $\lambda \sum_{j=1}^p \beta_j^2$ . This term is called shrinkage penalty, since it has the effect to push the coefficient estimates towards zero. Lambda ( $\lambda$ ) is a tuning parameter that controls the impact of the penalty on the estimates. In order to determine a good value for  $\lambda$ , we used cross-validation.

```
## [1] "The best lambda is = 755675"
## [1] "The estimated test MSE with the best lambda is = 5.446854e+13"
```

In the plot of the Figure 8 the red dotted line represents the cross-validation curve along with upper and lower standard deviation curves along the  $\lambda$  sequence (error bars). We chose the value of  $\lambda$  (755675.4) that gives minimum mean cross-validated error. The mean squared error on the test set is 5.45e+13.

The final model was fitted with the best  $\lambda$  on all data. The trace plot, in the Figure 9, shows how the coefficients change if  $\lambda$  increases.

```
## [1] "R-squared = 0.693844944021397"
```

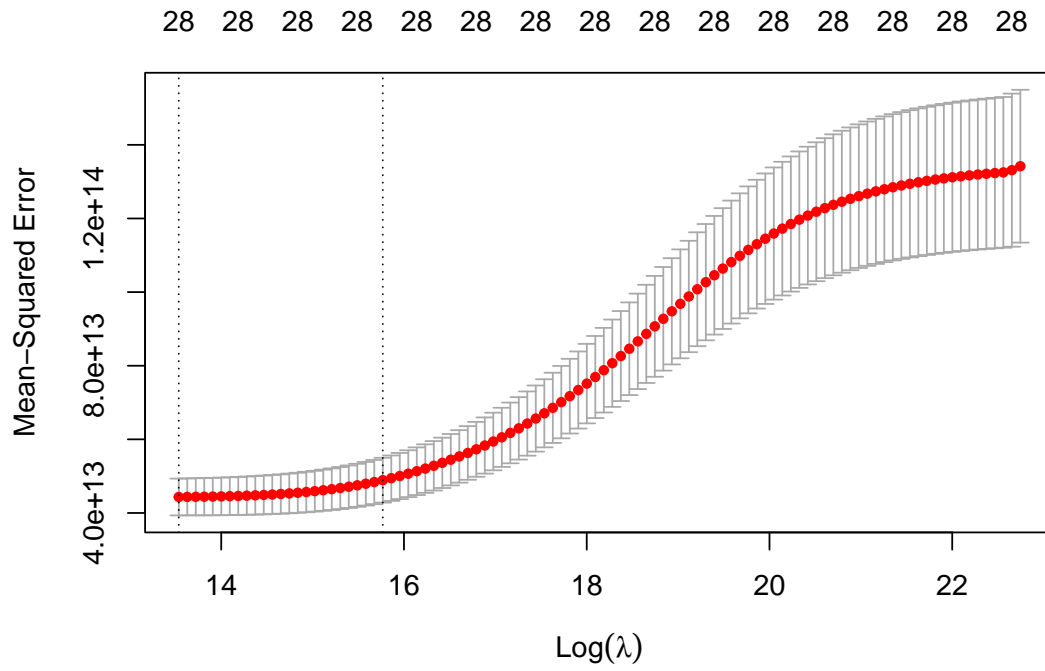


Figure 8: Plot of the MSE with respect to the value of  $\lambda$  in the Ridge regression model

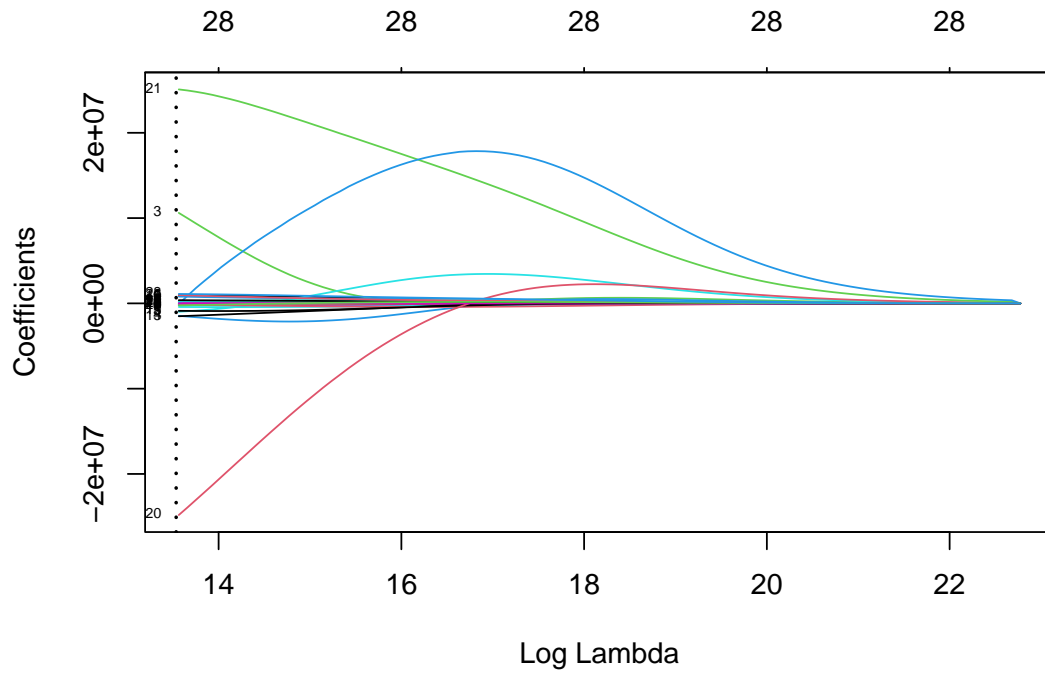


Figure 9: Plot of the values of the parameters with respect to  $\lambda$  in the Ridge regression model



```
## [1] "MSE = 4.344009e+13"
```

Once fitted the model on all data with the best lambda, we evaluated the performance. The 0.69 R-squared highlights a better fit compared to the previous model (0.64) obtained after a subset selection. Also the MSE improves: we get 4.34e+13 instead of the previous 4.77e+13.

The final step is the comparison between real salaries and predicted ones.

Table 6: Ten most overpaid players according to the Ridge regression model

	Salary	Predicted salary	Difference
Bradley Beal	\$46,741,590	\$23,835,819	\$22,905,771
Jordan Poole	\$27,955,357	\$10,348,317	\$17,607,040
Klay Thompson	\$43,219,440	\$25,696,260	\$17,523,180
Zach LaVine	\$40,064,220	\$22,726,915	\$17,337,305
Deandre Ayton	\$32,459,438	\$15,961,704	\$16,497,734
Darius Garland	\$34,005,250	\$17,803,526	\$16,201,724
Michael Porter Jr.	\$33,386,850	\$17,289,756	\$16,097,094
Tobias Harris	\$39,270,150	\$23,308,324	\$15,961,826
Fred VanVleet	\$40,806,300	\$25,497,363	\$15,308,937
Rudy Gobert	\$41,000,000	\$27,551,539	\$13,448,461

Table 7: Ten most underpaid players according to the Ridge regression model

	Salary	Predicted salary	Difference
Tyrese Maxey	\$4,343,920	\$24,414,677	\$20,070,757
Russell Westbrook	\$3,835,738	\$20,700,675	\$16,864,937
Eric Gordon	\$3,196,448	\$19,918,750	\$16,722,302
Desmond Bane	\$3,845,083	\$20,436,019	\$16,590,936
Tyrese Haliburton	\$5,808,435	\$22,145,870	\$16,337,435
Alperen Sengun	\$3,536,280	\$19,253,555	\$15,717,275
Cam Thomas	\$2,240,160	\$17,089,707	\$14,849,547
Jalen Williams	\$4,558,680	\$18,820,694	\$14,262,014
Kelly Oubre Jr.	\$2,891,467	\$16,543,402	\$13,651,935
Anthony Edwards	\$13,534,817	\$26,285,734	\$12,750,917

## MOST OVERPAID PLAYERS

Here there are a different similarities between the previous model: Beal, Poole, LaVine, Ayton, Garland and Porter Jr. still result in the most overpaid players.

Klay Thompson, after being a key piece in the Golden State Warriors dynasty, suffered a serious injury few years ago. After that, he was no longer the same player and the salary was, let's say, no longer adequate to his performance. His contract with Golden State ended after the 23-24 season and he recently signed with Dallas Mavericks for 50 millions in 3 years, thus he will earn a salary closer (even lower) to the predicted one.

Regarding Tobias Harris, this was the last contract year with Philadelphia 76ers. He signed this contract in 2019, team's situation was really different, Harris seemed to be the missing piece to build a contender for the title. After 5 years and a lot of changes, his situation is similar to Thompson's: salary not in line with

performance. In fact, he also signed recently with another team (Detroit Pistons) for 52 millions in 2 years, really close to the prediction.

Fred Vanvleet signed a big contract with Houston Rockets last year. The team has a young core, they are in a rebuilding phase so for the moment they don't have ambitions for the title. Without being a contender, teams are less attractive for the superstars. For this reason, they signed a really good player paying him like a superstar: the fact that he results as really overpaid was quite predictable.

Discussions about Rudy Gobert's value are always controversial. He doesn't shine for his technique, he is a so called hustle player: a great defender (three times best defender of the year) who gives a great contribution in terms of intangibles aspects that are really difficult to grasp with stats. It is really difficult to assess his value, especially with this type of model.

## **MOST UNDERPAID PLAYERS**

At first glance we can see that Ridge regression does not classify players with very high salaries (such as LeBron James and Kevin Durant) as the most underpaid. We believe that in this respect the prediction is better than the previous model one.

Again, we find Westbrook, Maxey and Oubre Jr. in this tier. Oubre's last contract was 30 millions in 2 years with Phoenix Suns, so in line with the predicted one. Last year he signed a small 2.89 million one-year contract with Philadelphia 76ers for several reasons: injury history, lack of performance consistency, market dynamics. Probably he will sign a new contract soon.

Eric Gordon is a veteran, he signed for a very small salary with Phoenix Suns in order to play with a contender. This move is not uncommon for good players in the final part of their career, especially if they never won a NBA title like Gordon. In the previous contract with Houston Rockets Gordon earned 75.6 millions in 4 years, perfectly in line with the prediction.

Bane, Haliburton, Sengun, Thomas and Williams have a Maxey-like situation: they are young players which are still in their rookie contracts but they clearly overperformed considering how much they earn. Maybe the model over evaluates a bit Cam Thomas, because he produces really good offensive numbers (the stats and the models capture the offensive contribution really well, much less the defensive one) when called on but his performance decrease when it comes to defense. Additionally, he could improve in leadership and understanding of the game.

Anthony Edwards had an amazing season, he carried his team to playoffs conference finals. 23-24 season was the last one in his rookie contract (he perceives a higher salary than the previous mentioned players in their rookie contracts because he was a better prospect when drafted), for this reason he perceived a lower salary compared to the model's expectation. It is really likely that he will receive a big offer in the near future.

All in all, Ridge regression has shown better results compared to the previous model: better R-squared, lower mean squared error and in some cases more meaningful predictions.

## **Lasso regression**

A disadvantage of ridge regression is that, unlike subset selection, it includes all  $p$  predictors in the final model. Also lasso regression shrinks the coefficients estimates towards zero but it has an absolute value shrinkage penalty instead of a quadratic one:  $\lambda \sum_{j=1}^p |\beta_j|$ . When  $\lambda$  is sufficiently large, some coefficient estimates become exactly equal to zero. Hence, like best subset selection, lasso performs a variable selection.

```
## [1] "The best lambda is = 86884"
## [1] "The estimated test MSE with the best lambda is = 5.419323e+13"
```

We again used the cross-validation method and we chose the  $\lambda$  value that guarantees the lower mean cross-validated error, as it can be seen in the Figure 10. Once chosen the best  $\lambda$ , the final model was fitted with that  $\lambda$  on all data. The trace plot, in the Figure 11, shows how the coefficient estimates change with increasing  $\lambda$ . Observing the coefficients, 6 are shrunk to zero. Among them, the coefficients of PIE and NET\_RATING that were significant in the best subset selection model.

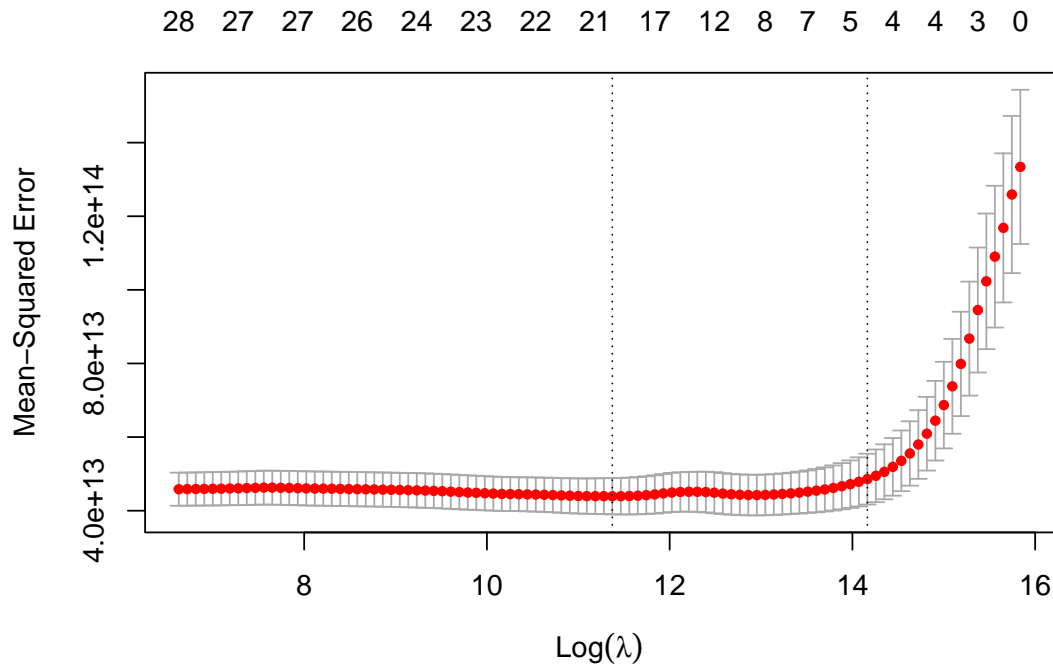


Figure 10: Plot of the MSE with respect to the value of  $\lambda$  in the Lasso regression model

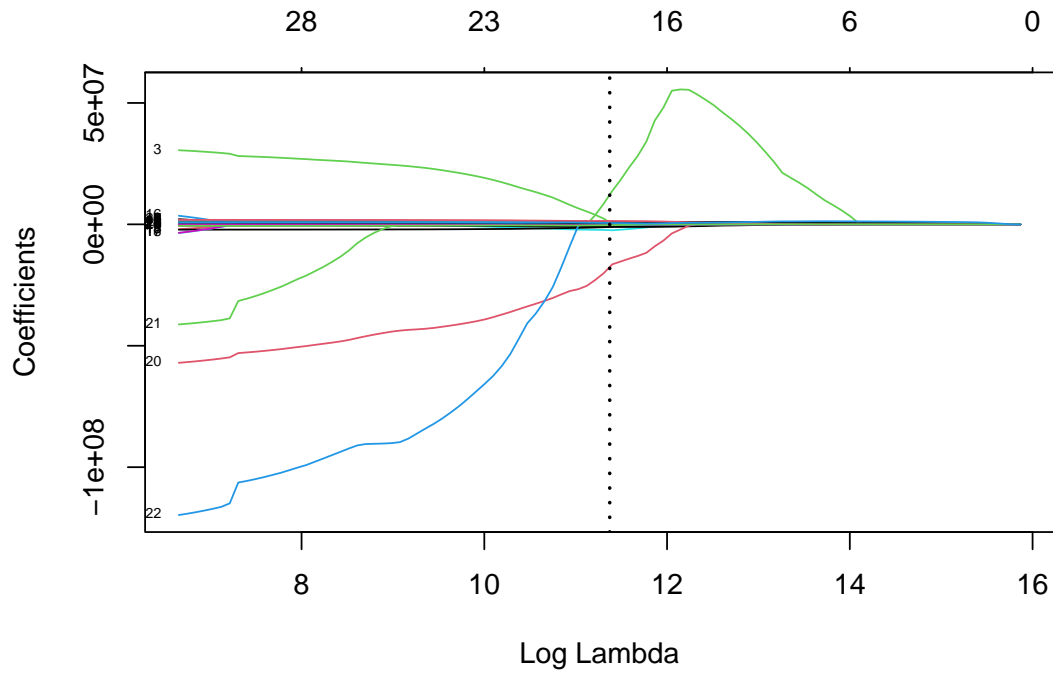


Figure 11: Plot of the values of the parameters with respect to  $\lambda$  in the Ridge regression model

```
## [1] "R-Squared = 0.69421191533625"
```

```
## [1] "MSE = 4.338802e+13"
```

Once fitted the model on all data with the best  $\lambda$ , we evaluated the performances. The R-squared is marginally better than that of the ridge, 0.6942 against 0.6938. The Mean squared error is slightly lower compared to the ridge one, 4.338802e+13 against 4.344009e+13. The performance of these two models is very similar and both outperform the model obtained with the best subset selection.

Table 8: Ten most overpaid players according to the Lasso regression model

	Salary	Predicted salary	Difference
Bradley Beal	\$46,741,590	\$24,401,038	\$22,340,552
Jordan Poole	\$27,955,357	\$10,171,469	\$17,783,888
Klay Thompson	\$43,219,440	\$25,865,270	\$17,354,170
Zach LaVine	\$40,064,220	\$23,367,732	\$16,696,488
Darius Garland	\$34,005,250	\$17,735,816	\$16,269,434
Michael Porter Jr.	\$33,386,850	\$17,520,241	\$15,866,609
Deandre Ayton	\$32,459,438	\$16,626,254	\$15,833,184
Tobias Harris	\$39,270,150	\$23,634,443	\$15,635,707
Fred VanVleet	\$40,806,300	\$26,335,599	\$14,470,701
Trae Young	\$40,064,220	\$26,643,866	\$13,420,354

Table 9: Ten most underpaid players according to the Lasso regression model

	Salary	Predicted salary	Difference
Tyrese Maxey	\$4,343,920	\$24,351,859	\$20,007,939
Desmond Bane	\$3,845,083	\$21,627,917	\$17,782,834
Eric Gordon	\$3,196,448	\$20,056,953	\$16,860,505
Russell Westbrook	\$3,835,738	\$20,610,467	\$16,774,729
Tyrese Haliburton	\$5,808,435	\$22,042,488	\$16,234,053
Alperen Sengun	\$3,536,280	\$18,956,778	\$15,420,498
Jalen Williams	\$4,558,680	\$19,044,255	\$14,485,575
Kelly Oubre Jr.	\$2,891,467	\$17,305,075	\$14,413,608
Cam Thomas	\$2,240,160	\$16,357,848	\$14,117,688
Kevin Love	\$3,835,738	\$16,608,865	\$12,773,127

## MOST OVERPAID PLAYERS

It is interesting to note that 9 out of 10 players in this table are the same as in the corresponding table for ridge. Also the difference between real salaries and predicted ones is very similar to that of the previous model. The only change is the presence of Trae Young here (he was one the most overpaid players in the best subset selection model) instead of Rudy Gobert in the ridge.

## MOST UNDERPAID PLAYERS

Also in this case, 9 out of 10 players are the same as in ridge and the differences are really small. The only change is the presence of Kevin Love. As Eric Gordon, he is a veteran and he signed for a small salary with Miami Heat.

After analyzing these three models, it emerges that ridge and lasso regression outperform the model obtained through the best subset selection. Additionally, the predictions for the most underpaid players seem to be more reasonable in ridge and lasso regression. In conclusion, lasso regression results (even only slightly compared to the ridge regression) the model that best fits our data.

## Salaries analysis by position

In the last part of the study, we wanted to analyse salaries by grouping players with respect to their playing positions. As positions, we considered the classic split of centers, forwards and guards. First of all, we wanted to check whether players earn, on average, the same salary regardless of their role. To do so, we used the ANOVA to compare the means of the different groups. Secondly, we implemented 3 lasso regression models to explore the relationship between salaries and performance for each position. The objectives are:

- observe the differences between the coefficients of the role-specific models among themselves and with respect to the lasso general model (the one which considers all the position)
- compare position-specific models' performances between each other and with the lasso general model
- compare the predictions on the most overpaid and most underpaid players between position-specific and general models

We used lasso regression because it turned out to be the best model in the previous phase.

In our dataset the division was somewhat different, with 5 positions (PG, SG, PF, SF, C). We considered point guards (PG) and shooting guards (SG) as guards (G), power forwards (PF) and shooting forwards (SF) as forwards (F). The centers (C) are the same.

## ANOVA

```
##
## Bartlett test of homogeneity of variances
##
## data: Salary by Pos
## Bartlett's K-squared = 0.054132, df = 2, p-value = 0.9733

##           Df      Sum Sq   Mean Sq F value Pr(>F)
## Pos         2 2.599e+12 1.299e+12  0.009  0.991
## Residuals  357 5.108e+16 1.431e+14
```

The ANOVA is a hypothesis test of equal means in different groups in which it is assumed that the variance is the same for every group; we used it to verify if players with different roles have, on average, different salaries. The null hypothesis states that the average salary is the same for every position; on the other hand, the alternative hypothesis states that at least one average salary is different. Firstly, we tested the hypothesis of variances homogeneity with Bartlett's test in order to see if it was possible to proceed with the ANOVA. Looking at the output, the Bartlett's K-squared was 0.054132. This small value indicates that the difference between the observed variances between the groups is small, as would be expected under the null hypothesis of equality of variances. The p-value is 0.9733: considering a significance level of 0.05, we do not have sufficient evidence to reject the null hypothesis. Therefore, Bartlett's test shows no evidence of unequal variances, so we can confidently proceed to the ANOVA considering satisfied the hypothesis of homogeneity of variances. Proceeding with ANOVA, the test statistic F results equal to 0.009. This quantity indicates that the variability of salaries between positions is rather small compared to the variability within positions. The p-value is much greater than 0.05 (chosen again as level of significance), so we don't have sufficient evidence to reject the null hypothesis. In conclusion, there is no significant evidence to suggest that average salaries differ significantly between different positions.

## Models

**Centers** We now considered only the centers, 67 players.

Starting from the complete model, we found the best  $\lambda$  through a cross-validation.

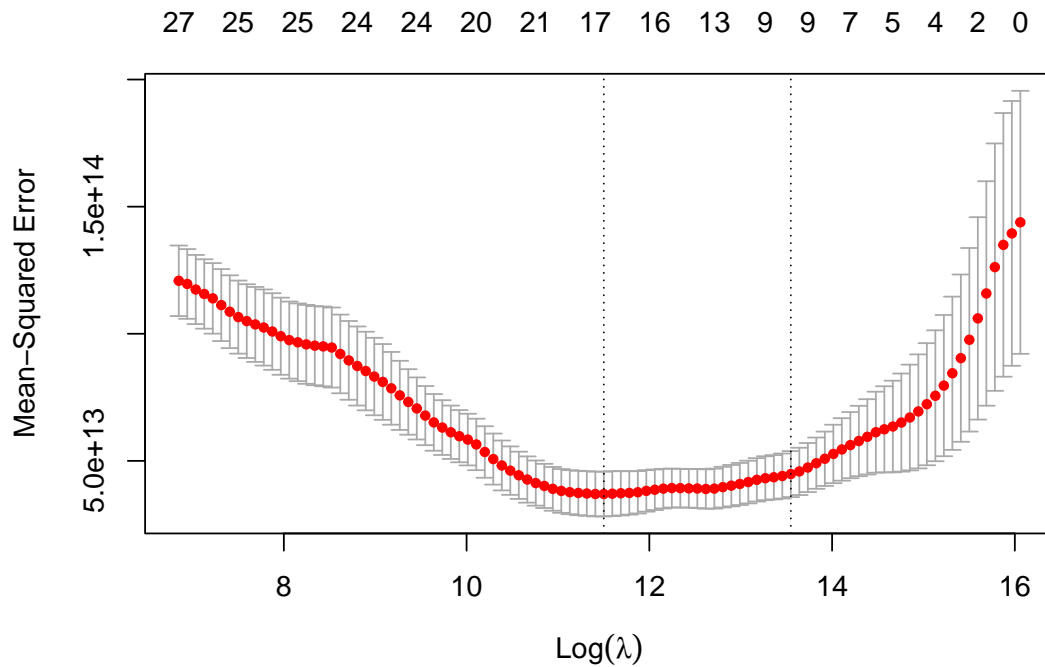


Figure 12: Plot of the MSE with respect to  $\lambda$  in the Lasso regression model for centers

```
## [1] "The best lambda is = 98893"
## [1] "The estimated test MSE with the best lambda is = 7.335062e+13"
```

We chose the  $\lambda$  value that guarantees the lower mean cross-validated error, as it can be seen in the Figure 12.

```
## 29 x 1 sparse Matrix of class "dgCMatrix"
##              s1
## (Intercept) 45780670.697
## AGE         826686.157
## GP          -134932.924
## FG_PCT      5366103.429
## FG3_PCT     -2050472.111
## FT_PCT      -106867.223
## OREB         .
## DREB        -182595.066
## REB         -47579.185
## AST          .
## TOV         -203905.609
## STL          .
```

```

## BLK          .
## BLKA        -2172425.478
## PF          .
## PTS         507593.849
## OFF_RATING  -130630.748
## DEF_RATING  -384576.335
## NET_RATING  .
## AST_TO     -1513719.318
## TS_PCT     -24792521.628
## USG_PCT    .
## PIE        .
## PFD        283721.386
## MIN        .
## MIN_G      513743.638
## WS         1584267.130
## BPM        3195.267
## VORP       188998.839

```

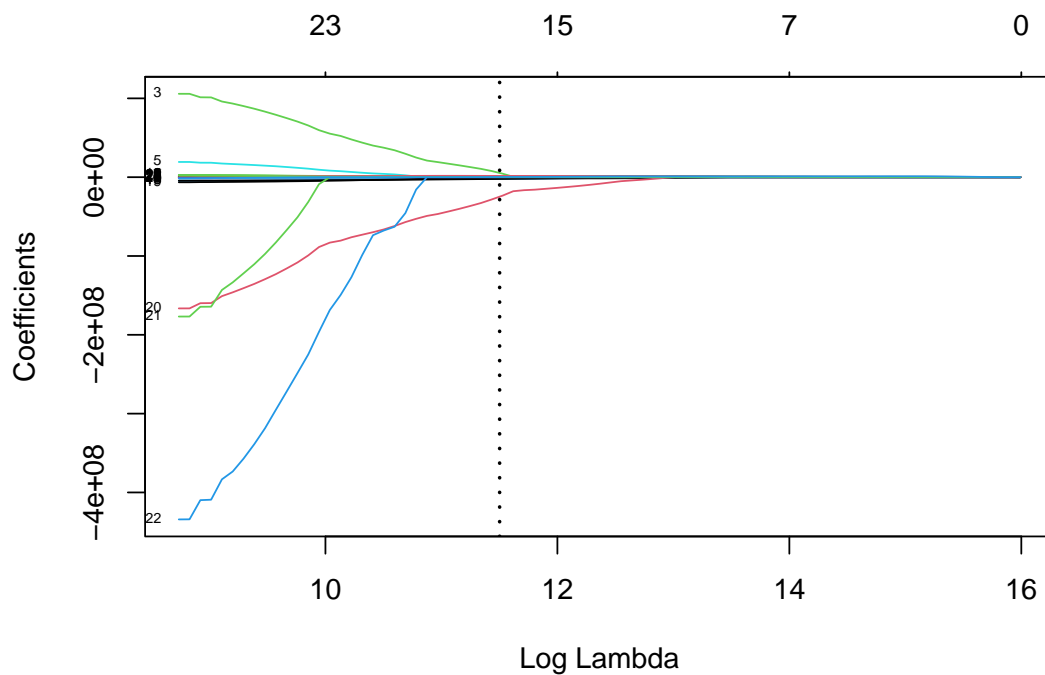


Figure 13: Plot of the values of the parameters with respect to  $\lambda$  in the Lasso regression model for centers

Once chosen the best  $\lambda$ , we fitted the model on all data. The Figure 13 shows how the coefficient estimates change increasing  $\lambda$ . It is interesting that, in the particular case of centers, 9 variables are shrunk to zero: the model considers less variables compared to the general lasso model (where only 6 variables were shrunk to zero). Unexpectedly, despite the fact that blocks are typically a center play, the variable BLK is excluded from the model.

```
## [1] "R-Squared = 0.79769165156743"
```

```
## [1] "MSE = 2.749367e+13"
```

The models' performance is really good: approximately 0.8 R-squared and 2.749367e+13 MSE. By far, the best performance till now. It is possible to infer that the relationship between salaries and performance is very well represented from the model; however, it is important to consider that the sample here is quite small, only 67 observations. For this reason, the results must be considered carefully.

Table 10: Three most overpaid centers according to the Lasso regression model by position

	Salary	Predicted salary	Difference
Deandre Ayton	\$32,459,438	\$17,949,767	\$14,509,671
James Wiseman	\$12,119,440	\$1,404,043	\$10,715,397
Jaren Jackson Jr.	\$27,102,202	\$18,324,760	\$8,777,442

Table 11: Three most underpaid centers according to the Lasso regression model by position

	Salary	Predicted salary	Difference
Alperen Sengun	\$3,536,280	\$20,207,681	\$16,671,401
Al Horford	\$10,000,000	\$19,101,403	\$9,101,403
Andre Drummond	\$3,360,000	\$12,006,476	\$8,646,476

Looking at the most underpaid and most overpaid centers it emerges that the differences between actual and predicted salaries are smaller than in the models seen above.

## MOST OVERPAID CENTERS

Among the overpaid centers we found again Deandre Ayton.

Regarding James Wiseman, he was a great prospect but he was really unlucky with injuries. Although he was traded from Golden State Warriors to Detroit Pistons (to a weaker team), he has not yet managed to fulfil its potential. He is still in his rookie contract: like Edwards, his salary is higher compared to average rookie contracts because he was the 2nd pick in 2020 draft. But unlike Edwards, according to the model his performance did not match his salary.

It is quite surprising to find Jaren Jackson Jr. here. After good seasons with Memphis Grizzlies, probably this year he suffered the drop in performance of the entire team. He is a really good defender and an amazing blocker: as said before, the defensive aspect of basketball is difficult to grasp with stats (and consequently with models). Moreover, we saw that the variable BLK is shrunk to 0 in this model, this may partly explain why he is classified as overpaid.

## MOST UNDERPAID CENTERS

The most underpaid center results to be Alperen Sengun and it makes sense for what we said above.

Al Horford is a veteran, he contributed excellently to the Boston Celtics' title victory, performing above the expectations.

Andre Drummond is a great rebounder, capable of consistent performances in terms of statistics. However, he is lacking in aspects of the game that are not considered by the model, such as basketball IQ: maybe this is why he results underpaid according to the model.



**Forwards** For what concerns forwards, we have 145 observations in our dataset.

```
##
## Call:
## lm(formula = Salary ~ +., data = fd_forward_nopos)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -17924383 -3474249  -153153   3096674 16232473
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -54258789  49966494  -1.086   0.2798
## AGE          1088269    153648    7.083 1.17e-10 ***
## GP          -117364     202185   -0.580   0.5627
## FG_PCT       34571222   35260552    0.980   0.3289
## FG3_PCT      28460591   14605848    1.949   0.0538 .
## FT_PCT       1063850    9814807    0.108   0.9139
## OREB        -8059593   12574105   -0.641   0.5228
## DREB        -9601894   12414671   -0.773   0.4408
## REB          8902973   12522833    0.711   0.4785
## AST          565889    1371158    0.413   0.6806
## TOV          1273614    3030615    0.420   0.6751
## STL          123653     1745694    0.071   0.9437
## BLK          2065023    1390173    1.485   0.1401
## BLKA        -2310592    2082275   -1.110   0.2694
## PF          -326377     1068712   -0.305   0.7606
## PTS          717063     1111652    0.645   0.5202
## OFF_RATING   12824426   11208353    1.144   0.2549
## DEF_RATING  -12370276   11216024   -1.103   0.2723
## NET_RATING  -12549273   11223242   -1.118   0.2658
## AST_TO      -2469170     2804182   -0.881   0.3804
## TS_PCT      -86632178   57325355   -1.511   0.1334
## USG_PCT       3613122   140382585    0.026   0.9795
## PIE        -71974525   186272354   -0.386   0.6999
## PFD          -54292      640701   -0.085   0.9326
## MIN           2572        8082    0.318   0.7509
## MIN_G         61398     522267    0.118   0.9066
## WS           1241370    1249206    0.994   0.3224
## BPM          -242631    1293090   -0.188   0.8515
## VORP          2096275    2511873    0.835   0.4057
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6206000 on 116 degrees of freedom
## Multiple R-squared:  0.7842, Adjusted R-squared:  0.7321
## F-statistic: 15.05 on 28 and 116 DF,  p-value: < 2.2e-16

## [1] "The best lambda is = 310190"
## [1] "The estimated test MSE with the best lambda is = 4.484197e+13"
```

As for centers, we started from the complete model and we chose the best  $\lambda$ , as it can be seen in the Figure 14.

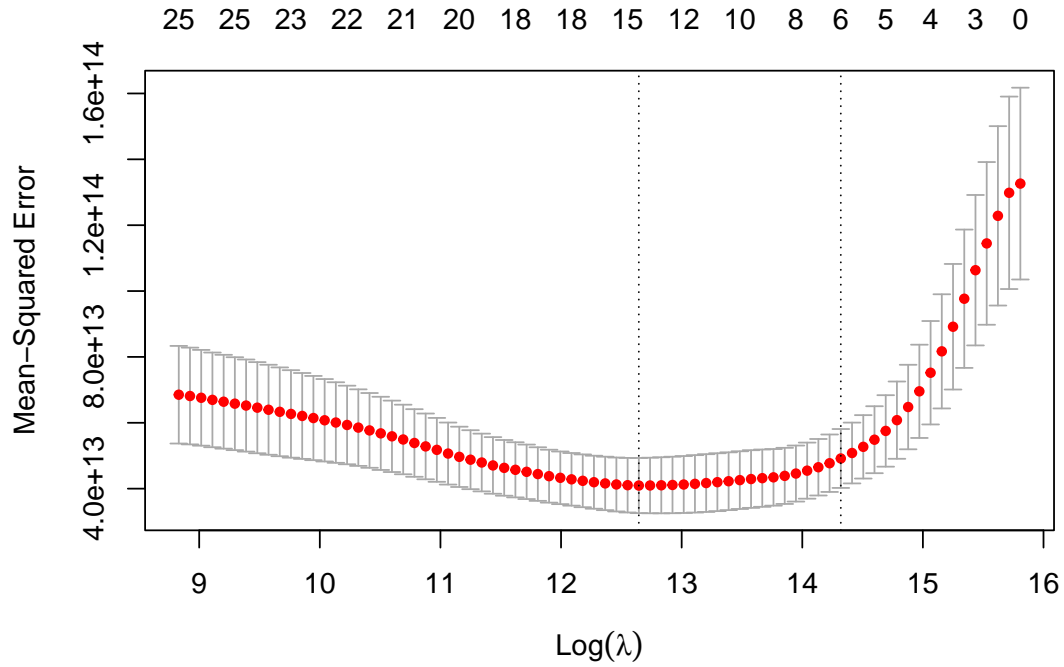


Figure 14: Plot of the MSE with respect to  $\lambda$  in the Lasso regression model for forwards

```
## 29 x 1 sparse Matrix of class "dgCMatrix"
##              s1
## (Intercept) -58430282.29
## AGE         1017388.93
## GP          .
## FG_PCT      .
## FG3_PCT     3745346.20
## FT_PCT      .
## OREB         .
## DREB        -358763.59
## REB          .
## AST          .
## TOV          486089.61
## STL          .
## BLK          88052.79
## BLKA         .
## PF           .
## PTS          .
## OFF_RATING   175806.62
## DEF_RATING   84909.00
## NET_RATING   .
## AST_TO      -1593815.38
## TS_PCT       -5887473.76
## USG_PCT      60152767.65
## PIE          .
## PFD          .
```

```
## MIN .
## MIN_G 298566.17
## WS 163860.04
## BPM .
## VORP 2681879.79
```

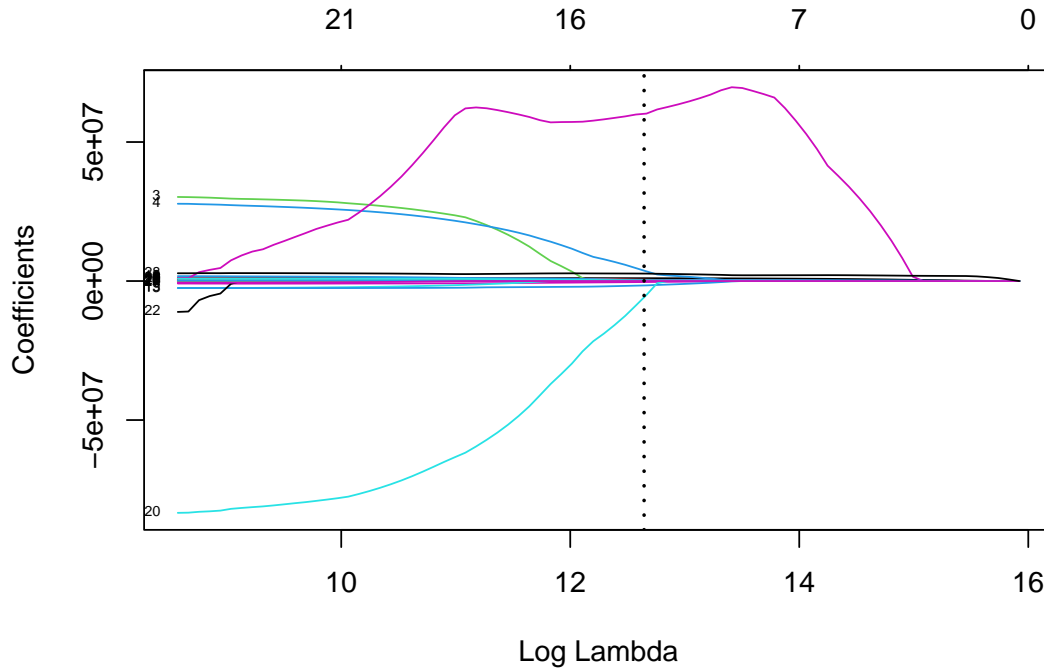


Figure 15: Plot of the values of the parameters with respect to  $\lambda$  in the Lasso regression model for forwards

Once chosen the best  $\lambda$ , we fitted the model on all data. In this specific lasso regression 15 variables are reduced to zero (15), far more than in the lasso regression for centers (9) and the general lasso regression (6).

```
## [1] "R-Squared = 0.753408923705351"
```

```
## [1] "MSE = 3.520197e+13"
```

The model has a 0.75 R-squared and a  $3.520197e+13$  MSE. The performance is very good, also considering that this model has only 13 predictors as we mentioned earlier. Also for forwards we can say that the relationship between salaries and player performance is well captured by the model.

Table 12: Three most overpaid forwards according to the Lasso regression model by position

	Salary	Predicted salary	Difference
Michael Porter Jr.	\$33,386,850	\$17,418,864	\$15,967,986
Tobias Harris	\$39,270,150	\$23,551,577	\$15,718,573

	Salary	Predicted salary	Difference
Klay Thompson	\$43,219,440	\$27,607,531	\$15,611,909

Table 13: Three most underpaid forwards according to the Lasso regression model by position

	Salary	Predicted salary	Difference
Jalen Williams	\$4,558,680	\$19,113,798	\$14,555,118
Kelly Oubre Jr.	\$2,891,467	\$17,431,407	\$14,539,940
Kevin Love	\$3,835,738	\$18,212,706	\$14,376,968

The 3 most overpaid forwards according to this model are Michael Porter Jr., Tobias Harris and Klay Thompson, who were also among most overpaid players in the overall Lasso regression. As we said before, it is reasonable to find these players in this tier

The same applies to Williams, Oubre Jr. and Love in terms of the 3 most underpaid forwards.

**Guards** Regarding guards, we consider 148 observations.

```
##
## Call:
## lm(formula = Salary ~ +., data = fd_guard_nopos)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -17491704 -4133266  151945  4251578 20559797
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -917653    51713161  -0.018  0.9859
## AGE             1098768     176097    6.240 6.97e-09 ***
## GP              24348      214655    0.113  0.9099
## FG_PCT         40005355    42523403    0.941  0.3487
## FG3_PCT       -18746248    19317481   -0.970  0.3338
## FT_PCT        -2662700    13386120   -0.199  0.8427
## OREB           18079410    11922943    1.516  0.1321
## DREB           19684785    11923256    1.651  0.1014
## REB          -17554017    11814813   -1.486  0.1400
## AST             1113406     1284659    0.867  0.3879
## TOV            2024924     2938466    0.689  0.4921
## STL             327282      1895837    0.173  0.8632
## BLK            -239655      2588200   -0.093  0.9264
## BLKA          -1130455      2334111   -0.484  0.6290
## PF            -2076987      1359115   -1.528  0.1291
## PTS             1861414      1141662    1.630  0.1057
## OFF_RATING     19372187    12487097    1.551  0.1235
## DEF_RATING    -19721268    12490871   -1.579  0.1170
## NET_RATING    -19537267    12485173   -1.565  0.1203
## AST_TO           940733      1470565    0.640  0.5236
## TS_PCT        -12856606     54842278   -0.234  0.8151
```

```
## USG_PCT      -37203927  137922483  -0.270  0.7878
## PIE          -528843189  274444067  -1.927  0.0564 .
## PFD          -105777    887236    -0.119  0.9053
## MIN          -8945      8316     -1.076  0.2843
## MIN_G        914696    519847    1.760  0.0811 .
## WS           3120012    1643765    1.898  0.0601 .
## BPM          -81320    1578527   -0.052  0.9590
## VORP         -594952    2976839   -0.200  0.8419
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7652000 on 119 degrees of freedom
## Multiple R-squared:  0.6724, Adjusted R-squared:  0.5954
## F-statistic: 8.724 on 28 and 119 DF,  p-value: < 2.2e-16
```

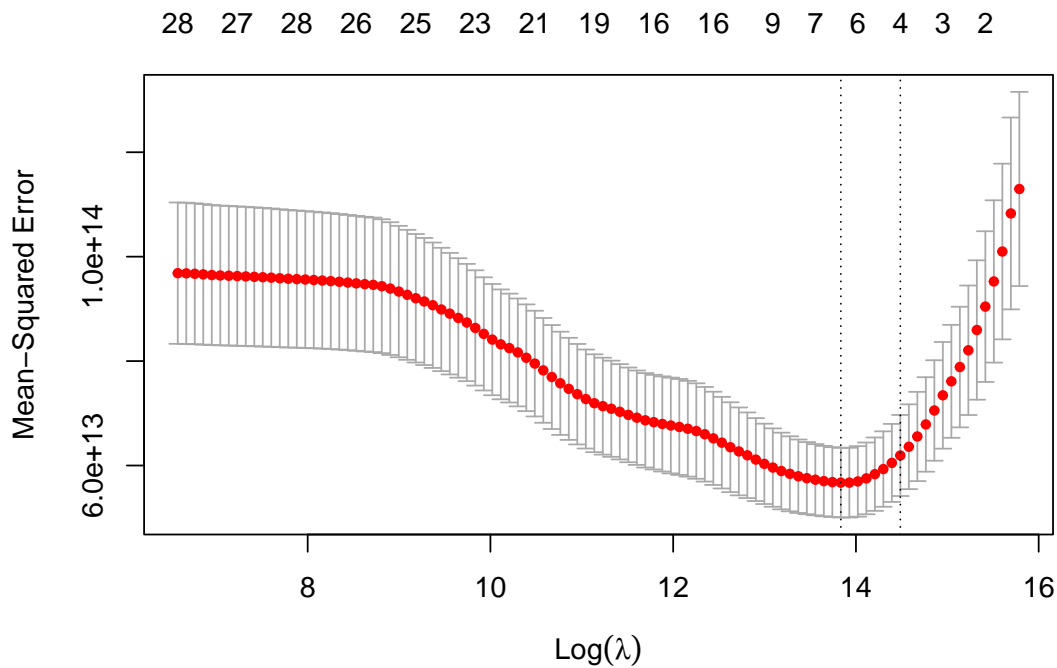


Figure 16: Plot of the MSE with respect to  $\lambda$  in the Lasso regression model for guards

```
## [1] "The best lambda is = 1019523"
## [1] "The estimated test MSE with the best lambda is = 8.220954e+13"

## 29 x 1 sparse Matrix of class "dgCMatrix"
##              s1
## (Intercept) -29556114.8
## AGE         760281.0
## GP          .
## FG_PCT      .
## FG3_PCT     .
```

## FT_PCT	.
## OREB	.
## DREB	.
## REB	.
## AST	.
## TOV	.
## STL	.
## BLK	.
## BLKA	.
## PF	.
## PTS	247361.5
## OFF_RATING	.
## DEF_RATING	.
## NET_RATING	.
## AST_TO	.
## TS_PCT	.
## USG_PCT	7217391.1
## PIE	.
## PFD	.
## MIN	.
## MIN_G	572478.0
## WS	.
## BPM	.
## VORP	683547.4

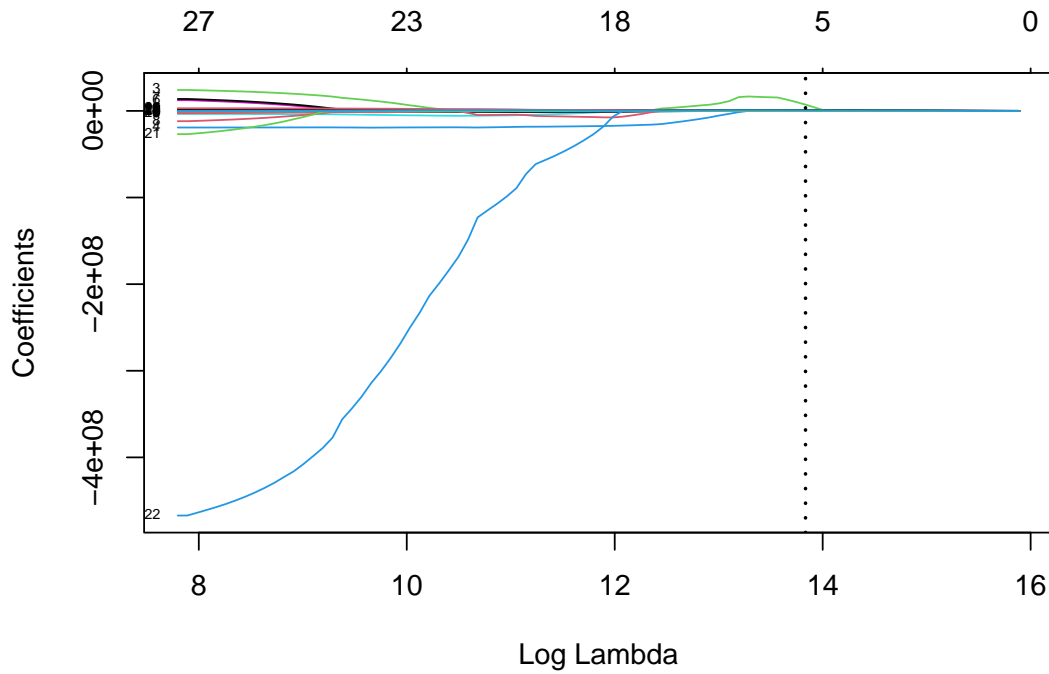


Figure 17: Plot of the values of the parameters with respect to  $\lambda$  in the Lasso regression model for guards

After selecting the best value of  $\lambda$  also in this case (16), we obtained a lasso regression model for guards. The

result in terms of coefficients is very interesting (17): 23 variables are shrunk to zero, the model considers only 5 predictors. Given that the guard role tends to be the most difficult to interpret as the creative component plays a major role, we expected a more complex model than those obtained for centers and forwards. Instead, we are faced with a very simple model.

```
## [1] "R-Squared = 0.58850980282926"
```

```
## [1] "MSE = 5.914565e+13"
```

As expected, given the low number of predictors, here we have a lower R-squared and a fairly higher MSE compared to models of the other positions. Probably the relationship between salaries and performance for guards is harder to grasp and it is best explained by variables not considered in our model.

Table 14: Three most overpaid guards according to the Lasso regression model by position

	Salary	Predicted salary	Difference
Bradley Beal	\$46,741,590	\$21,185,105	\$25,556,485
Stephen Curry	\$51,915,615	\$31,332,367	\$20,583,248
Zach LaVine	\$40,064,220	\$21,145,701	\$18,918,519

Table 15: Three most underpaid guards according to the Lasso regression model by position

	Salary	Predicted salary	Difference
Tyrese Maxey	\$4,343,920	\$21,911,682	\$17,567,762
Desmond Bane	\$3,845,083	\$20,462,332	\$16,617,249
Eric Gordon	\$3,196,448	\$18,924,824	\$15,728,376

We again find Beal and LaVine in the most 3 overpaid. Maxey, Bane and Gordon also here are classified among the most underpaid players. The predictions are very similar to those of the overall lasso regression model.

The only but very interesting difference is that Stephen Curry results very overpaid in this model. Curry is one of the biggest talents in the NBA: he won 4 titles as the absolute star and 2 MVPs thanks to his unique style of play. After several years of dominance, the performance of his Golden State Warriors has been declining in recent seasons. Despite this, Curry overall performance in 23-24 was in line or even better compared to the two previous seasons. For this reason, we think that in this case the prediction of the model does not make much sense.

## Conclusion

### Models performances summary

Summarizing, we started from a linear regression model and we applied a logarithmic transformation to the dependent variable salary. Then, in order to exclude unnecessary or redundant variables and to obtain a simpler model, we performed a variable selection. We got a simpler model with a more than acceptable decrease in performance. However, given the probable presence of multicollinearity between the predictors, we implemented a Ridge regression model and a Lasso regression model. Indeed, Ridge regression showed an improvement in both R-squared and MSE that was further enhanced by the Lasso regression. Since the

latter was the best model in terms of performances, we used it to analyse separately players with different roles. The models for centers and forwards have shown a really good fit with the data. Moreover, the variables considered changed both in number and type compared to the model fitted on all players. For what concerns guards, a very simple model emerged but with significantly worse performance.

### **Models' limitations**

It is necessary to remember that a lot of factors concur to determine how much a player should earn. Regarding performances, the statistics we used can only give a partial idea of a player's defensive contribution; additionally, it is very difficult to understand how players perform mentally, i.e. to grasp characteristics such as attitude, leadership and basketball IQ. Factors outside the basketball court also greatly influence the determination of salaries: player's potential, player's career, player's fit into the team, market dynamics, and so on. We considered only the 2023/2024 Regular Season: to have more complete models one would have to consider more seasons and, as mentioned at the beginning, also consider the playoffs.

### **Final comments**

In the end, although their obvious limitations, some of our models performed pretty good and provide a good basis for studying the relationship between NBA players' salaries and performances. In particular, Lasso regressions (except in the case of guards) fitted the data quite well with an acceptable mean squared error. As seen progressively, most of the biggest differences between actual salaries and predicted salary were quite reasonable.