



**SCHOOL OF INFORMATION SCIENCES &  
TECHNOLOGY**

**DEPARTMENT OF STATISTICS**

**MSc IN STATISTICS**

**Model-Based Performance Indicators for  
in-play Basketball Data**

**Argyro Damoulaki**

Supervisor: Ioannis Ntzoufras (AUEB)

Co-Supervisor: Konstantinos Pelechrinis (University of Pittsburg)

A Thesis submitted to the Department of Statistics of the Athens University of Economics and Business in partial fulfillment of the requirements for the degree of Master of Science in Statistics.

Athens, June 2023



**ΣΧΟΛΗ ΕΠΙΣΤΗΜΩΝ & ΤΕΧΝΟΛΟΓΙΑΣ ΤΗΣ  
ΠΛΗΡΟΦΟΡΙΑΣ**

**ΤΜΗΜΑ ΣΤΑΤΙΣΤΙΚΗΣ**

**ΜΕΤΑΠΤΥΧΙΑΚΟ ΠΡΟΓΡΑΜΜΑ ΣΤΗ  
ΣΤΑΤΙΣΤΙΚΗ**

**Δείκτες Αξιολόγησης παικτών βασισμένοι σε μοντέλα για  
δεδομένα αγώνων μπάσκετ**

**Αργυρώ Δαμουλάκη**

Επιβλέπων: Ιωάννης Ντζούφρας (ΟΠΑ)

Συν-Επιβλέπων: Κωνσταντίνος Πελεchrίνης (Πανεπιστήμιο του Pittsburg)

Εργασία που υποβλήθηκε στο Τμήμα Στατιστικής του Οικονομικού Πανεπιστημίου Αθηνών ως  
μέρος των απαιτήσεων για την απόκτηση Διπλώματος Μεταπτυχιακών Σπουδών στην  
Στατιστική.

Αθήνα, Ιούνιος 2023



## **Acknowledgments**

I would like to thank my supervisor, Professor Ioannis Ntzoufras, and co-supervisor, Professor Konstantinos Pelechrinis, for their full guidance, patience, and understanding during the development of this thesis. Also, I am grateful to Dionisis Mylonas for his valuable contribution through our long discussions and my family for their strong support.



## **Vita**

After completing my studies in Applied Mathematics (NTUA) I decided to continue in the Statistics field. At the same time, it was a great chance to indulge in Sports Analytics since I had decided from a young age that I could not be a professional basketball player. Hence, I found a way to combine my love for Mathematics with my love for Basketball. I am looking forward to getting started the next step of my studies.



# **Abstract**

Argyro Damoulaki

## **Model-Based Performance Indicators for in-play Basketball Data**

June 2023

The aim of this work is to estimate the offensive and defensive contributions of NBA players to their teams. As Dean Oliver, one of the most important men of Basketball Statistics, said “Teamwork is the element of basketball most difficult to capture in any quantitative sense” (“Basketball on Paper”, 2004). Rosenbaum, D. T., (2004) being inspired by this quote- and later Ilardi, S., and Barzilai, A., (2008)- proposed a Ridge Regression model for estimating the players’ contribution by using only information on their teammates and opponents in the court. The index was called Regularized Adjusted Plus/Minus (RAPM).

Firstly, we will see that low-time players should not be considered in the models since their Ridge RAPMs were found to be very high. Moreover, Lasso regularization method was used instead of Ridge for the Normal model. Results were improved since Lasso discriminates the higher and lower quality players in a more realistic and sensible way. However, the Normal distribution is not appropriate to model a discrete response variable (points per possession with values from zero to three).

As an extension of the standard Ridge Regression model, logistic models were implemented. Starting from a simplification of the initial problem, regularized Binomial models were developed with a response variable of scoring or not. An alternative interpretation for the Ridge-RAPMs is given using this method.

The final model is a Multinomial implementation of points scored per possession. The model was fitted indirectly via three separate binomial models for the different types of points scored (one, two, and three or more). A new evaluation metric is proposed for expected points per possession (EPTS-RAPM) is provided. We will see that the Multinomial model is more appropriate for the discrete variable of points and its performance is better according to some selected external validation criteria. In addition, this method seems to solve partially common problem with low time players in Plus/Minus ratings.





## Περίληψη

Αργυρώ Δαμουλάκη

**Δείκτες Αξιολόγησης παικτών βασισμένοι σε μοντέλα για δεδομένα αγώνων μπάσκετ**

Ιούνιος 2023

Στόχος αυτής της διπλωματικής εργασίας είναι να εκτιμηθεί η επιθετική και αμυντική συνεισφορά των παικτών του NBA στις ομάδες τους. Όπως είπε ο Dean Oliver, ένας από τους σημαντικότερους ανθρώπους στη στατιστική ανάλυση στην καλαθοσφαίριση, "Η ομαδική εργασία είναι το στοιχείο του μπάσκετ που είναι πιο δύσκολο να συλληφθεί με οποιαδήποτε ποσοτική έννοια" ("Basketball on Paper", 2004). Ο Rosenbaum, D. T., (2004) εμπνεύστηκε από αυτό το απόσπασμα - και αργότερα Ilardi, S., και Barzilai, A., (2008)- πρότειναν ένα Ridge μοντέλο γραμμικής παλινδρόμησης για την εκτίμηση της συνεισφοράς των παικτών χρησιμοποιώντας μόνο την πληροφορία για τους συμπαίκτες και τους αντιπάλους τους στο γήπεδο. Ο δείκτης ονομάστηκε Regularized Adjusted Plus / Minus (RAPM).

Αρχικά, θα δούμε ότι οι παίκτες με λίγο χρόνο συμμετοχής δεν πρέπει να λαμβάνονται υπόψη στα μοντέλα, καθώς τα Ridge RAPMs βρέθηκαν να είναι πολύ υψηλά. Επιπλέον, χρησιμοποιήθηκε μέθοδος κανονικοποίησης Lasso αντί για Ridge για το μοντέλο Κανονικής κατανομής. Τα αποτελέσματα βελτιώθηκαν, καθώς το Lasso διακρίνει τους παίκτες υψηλότερης και χαμηλότερης ποιότητας με πιο ρεαλιστικό και λογικό τρόπο. Ωστόσο, η Κανονική κατανομή δεν είναι κατάλληλη για τη μοντελοποίηση μιας διακριτής μεταβλητής απόκρισης (πόντοι ανά κατοχή με τιμές από μηδέν έως τρεις).

Ως επέκταση του προτεινόμενου από τη βιβλιογραφία μοντέλου παλινδρόμησης, μελετήθηκαν λογιστικά μοντέλα. Ξεκινώντας από την απλοποίηση του αρχικού προβλήματος, αναπτύχθηκαν κανονικοποιημένα (Ridge and Lasso) Διωνυμικά μοντέλα με μεταβλητή απόκρισης το αν ευστόχησαν ή όχι. Μια εναλλακτική ερμηνεία για τα Ridge-RAPMs δίνεται χρησιμοποιώντας αυτή τη μέθοδο.

Το τελικό μοντέλο είναι μια Πολυωνυμική κατανομή των πόντων ανά κατοχή. Το μοντέλο προσαρμόστηκε έμμεσα μέσω τριών ξεχωριστών διωνυμικών μοντέλων για τους διαφορετικούς τύπους πόντων. Προτείνεται ένας δείκτης αξιολόγησης για τους αναμενόμενους πόντους ανά

κατοχή (EPTS-RAPM). Θα δούμε ότι το πολυωνυμικό μοντέλο είναι πιο κατάλληλο για τη διακριτή μεταβλητή των πόντων και η αποδοτικότητά του είναι καλύτερη σύμφωνα με ορισμένα εξωτερικά κριτήρια. Επιπλέον, αυτή η μέθοδος φαίνεται να επιλύει εν μέρει το κοινό πρόβλημα των Plus/Minus δεικτών για τους παίκτες που δεν παίζουν πολύ.

## Table of Contents

Acknowledgments .....	I
Vita .....	III
Abstract.....	V
Περίληψη .....	VII
List of Tables .....	XIII
List of Figures.....	XV
Chapter 1: Review for Basketball Analytics.....	1
1.1 Introduction to Sports Statistics and Basketball Analytics.....	1
1.2 Model-Based Analytics for College Championships.....	3
1.3 The “home team” factor .....	4
1.4 Inspiration from Dean Oliver: Model based works for Prediction and Performance.....	7
1.5 Moving away from the usuals .....	11
1.5.1 Some innovative works .....	11
1.5.2 “Basketball-oriented” studies .....	13
1.5.3 Bayesian, Unsupervised Learning and Big data methods .....	14
1.5.4 Data Envelopment Analysis (DEA) .....	19
1.5.5 Plus/Minus: At first glance.....	21
1.6 Conclusion .....	23
Chapter 2: Performance Analytics via Plus/Minus Models .....	25
2.1 Early years.....	25
2.2 Model based Plus/Minus ratings.....	27
2.2.1 Adjusted Plus/Minus.....	27
2.2.2 Advanced Plus/Minus .....	30
2.2.3 Box score-Based Plus/Minus .....	32
2.2.4 Adjusted Plus/Minus Discussion and Extensions .....	35
2.3 Mixtures of Plus/Minus Ratings .....	39
2.3.1 PT-PM.....	40
2.3.2 CARMELO.....	40
2.3.3 RAPTOR.....	41

2.3.4 DARKO.....	42
2.3.5 LEBRON.....	43_Toc138710555
2.4 What about women’s basketball?.....	45
2.5 Other Sports .....	46
2.6 Conclusion .....	48
Chapter 3: A case study about Plus/Minus Ratings.....	50
3.1 Introduction.....	50
3.2 Case Study Dataset .....	50
3.2.1 Initial data.....	50
3.2.2 Data processing & final data.....	51
3.3 Explanatory Data Analysis.....	52
3.3.1 Descriptive Analysis.....	52
3.3.2 Pairwise Analysis .....	54
3.4 Calculation of RAPM: Full dataset Analysis.....	57
3.4.1 Ridge regression with all players under consideration.....	58
3.4.2 Lasso regression with all players under consideration.....	60
3.4.4 Conclusion .....	64
3.5 Calculation of RAPM: Filtered dataset Analysis .....	65
3.5.1 Filtered dataset.....	65
3.5.2 Ridge regression .....	66
3.5.3 Lasso regression .....	69
3.5.5 Assessing the effect of LTPs and the team ability .....	73
3.5.6 Validation of RAPM values: Comparison with objective criteria and ratings.....	74
3.5.7 Conclusion .....	76
3.6 Player Evaluation using RAPM.....	76
3.6.1 Comparing player’s contribution .....	77
3.6.2 Comparing traded players’ contribution.....	80
3.7 Conclusion .....	83
Chapter 4: Building Logistic Models .....	85
4.1 Introduction.....	85
4.2 Binomial models .....	85
4.3 Separate binomial models .....	90

<b>4.4 Multinomial model</b> .....	98
<b>4.4.1 Building the multinomial model via binomial models</b> .....	98
<b>4.4.2 Results of EPTS-RAPM</b> .....	100
<b>4.5 Conclusion of Chapter 4</b> .....	105
<b>Chapter 5: Comparison of the different RAPM Models</b> .....	107
<b>5.1 Introduction</b> .....	107
<b>5.2 Multinomial vs. Normal</b> .....	107
<b>5.3 Comparison of all method</b> .....	113
<b>5.3.1 Agreement between results and selected criteria</b> .....	113
<b>5.3.2 Five-lineup analysis</b> .....	118
<b>5.4 Conclusion</b> .....	124
<b>Chapter 6: Discussion</b> .....	125
<b>6.1 Summary of the methodology</b> .....	125
<b>6.2 Some highlights</b> .....	126
<b>6.3 Future work</b> .....	128
<b>APPENDIX</b> .....	131
<b>References</b> .....	133
<b>Bibliography</b> .....	133
<b>Electronic References</b> .....	139



## List of Tables

<b>Chapter 2: Performance Analytics via Plus/Minus Models .....</b>	<b>25</b>
<b>2.3 Mixtures of Plus/Minus Ratings .....</b>	<b>39</b>
Table 2.1: The discussed basketball players performance ratings of the literature.....	44
<b>Chapter 3: A case study about Plus/Minus Ratings.....</b>	<b>50</b>
<b>3.3 Explanatory Data Analysis.....</b>	<b>52</b>
<b>3.3.1 Descriptive Analysis.....</b>	<b>52</b>
Table 3.1: Frequency table of points per possessions.....	53
<b>3.3.2 Pairwise Analysis .....</b>	<b>54</b>
Table 3.2: Summary table for Teams Possessions Ratio (TPR). .....	55
<b>3.5 Calculation of RAPM: Filtered dataset Analysis .....</b>	<b>65</b>
<b>3.5.2 Ridge regression .....</b>	<b>66</b>
Table 3.3: Players with the best offensive and defensive RAPM of each NBA with more than 200 minutes played.....	67
<b>3.5.3 Lasso regression .....</b>	<b>69</b>
Table 3.4: Players with the highest offensive and defensive Lasso RAPM, with at least 200 minutes played, per team. ....	70
Table 3.5: Top 5 players with the highest offensive Lasso-RAPM per position with at least 200 minutes played.....	71
Table 3.6: Top 5 players with the highest defensive Lasso-RAPM per position with at least 200 minutes played.....	71
<b>3.5.5 Assessing the effect of LTPs and the team ability .....</b>	<b>73</b>
Table 3.7: Pearson’s linear correlation ( $r$ ) between the model with the fictional player effect ( $M_1$ ), with not the fictional player effect ( $M_2$ ) and the teams effects ( $M_3$ ).....	74
<b>3.5.6 Validation of RAPM values: Comparison with objective criteria and ratings.....</b>	<b>74</b>
Table 3.8: Percentage of agreement between RAPM and All-NBA team 2021-2022 per position (filtered dataset). ....	75
<b>Chapter 4: Building Logistic Models .....</b>	<b>85</b>
<b>4.2 Binomial models .....</b>	<b>85</b>
Table 4.1: Summary results of the ridge and binomial models.....	87
Table 4.2: Agreement between the Binomial-based high rated players and the All-NBA teams. ....	88
Table 4.3: Linear correlation between the Linear and Binomial models results.....	89



<b>4.3 Separate binomial models .....</b>	<b>90</b>
Table 4.4: Common statistically important defenders to the three separate binomial models.....	91
Table 4.5: Common statistically important offenders to the three separate binomial models. ....	92
<b>4.4 Multinomial model.....</b>	<b>98</b>
<b>4.4.2 Results of EPTS-RAPM .....</b>	<b>100</b>
Table 4.6a: Top-5 offenders per position via multinomial model's results for expected points per possession. ....	101
Table 4.7a: Top-5 defenders per position via multinomial model's results for expected points per possession. ....	101
Table 4.6b: Top-5 offenders per position via multinomial model's results for efficiency (expected points per possession based on those possessions each player was on the court). ....	101
Table 4.7b: Top-5 defenders per position via multinomial model's results for efficiency (expected points per possession based on those possessions each player was on the court). ....	101
Table 4.8: Best players per team via EPTS-RAPMs.....	103
<b>Chapter 5: Comparison of the different RAPM Models .....</b>	<b>107</b>
<b>5.2 Multinomial vs. Normal.....</b>	<b>107</b>
Table 5.1: Correlation for players RAPMs ratings of the multinomial and normal model*.....	107
Table 5.3: Starters included in the 100 highest and lowest rated RAPM (50 offensive and 50 defensive) players per model.....	110
<b>5.3 Comparison of all method .....</b>	<b>113</b>
<b>5.3.1 Agreement between results and selected criteria .....</b>	<b>113</b>
Table 5.4: Summary statistics measures for expected points per possession.....	114
Table 5.5: Agreement between model based RAPM highest rated players and the better ones in selected boxscores offensive statistics (per game).....	115
Table 5.6: Frequency of players' defensive contribution from top 50 per model that are captured in the top 50 of some in-game-statistics (per game). ....	116
Table 5.7: Percentage of agreement between the highest-rated offensive and defensive per position for each model and the All-NBA teams. ....	117
<b>5.3.2 Five-lineup analysis.....</b>	<b>118</b>
Table 5.8: Average percentage of agreement between the top players per model according to the best offensive 5-lineups per team. ....	118
Table 5.9a: Percentage of agreement between the top players per model according to the most-played 5-lineups per team. ....	118
Table 5.9b: Percentage of agreement between the top players per model according to the most-played 5-lineups per team. ....	119

## List of Figures

<b>Chapter 3: A case study about Plus/Minus Ratings.....</b>	<b>50</b>
<b>3.3 Explanatory Data Analysis.....</b>	<b>52</b>
<b>3.3.1 Descriptive Analysis.....</b>	<b>52</b>
Figure 3.1: Frequencies of points per possession.....	52
Figure 3.2: Boston Celtics possessions in which Jayson Tatum played offense.....	53
Figure 3.3: Boston Celtics defensive situations in which Jayson Tatum played. ....	53
<b>3.3.2 Pairwise Analysis .....</b>	<b>54</b>
Figure 3.4: Players' TPR based on teams' possessions. ....	55
Figure 3.5: Jason Tatum's offensive (right) and defensive (left) contribution to Boston Celtics depending on the total possessions of the season and his team's possessions. ....	56
Figure 3.6: Juancho Hernangomez's offensive (right) and defensive (left) contribution to Boston Celtics depending on the total possessions of the season and his team's possessions. ....	57
<b>3.4 Calculation of RAPM: Full dataset Analysis.....</b>	<b>57</b>
<b>3.4.2 Lasso regression with all players under consideration .....</b>	<b>60</b>
Figure 3.7: Scatterplots of Ridge vs. Lasso coefficients for the full dataset, offensive (left) and defensive (right). ....	63
Figure 3.8: Scatterplots of Ridge vs. non-zero Lasso coefficients for the full dataset, offensive (left) and defensive (right). ....	64
<b>3.5 Calculation of RAPM: Filtered dataset Analysis .....</b>	<b>65</b>
<b>3.5.2 Ridge regression .....</b>	<b>66</b>
Figure 3.9: Histograms of actual points per possession and fitted values of the Ridge regression model for the "filtered" dataset (via 1se lambda the light blue and lambda min the pink). ....	68
Figure 3.10: Minutes played of top 20 offensive and defensive RAPM players via ridge models.....	69
<b>3.5.3 Lasso regression .....</b>	<b>69</b>
Figure 3.11: Scatterplots of Ridge vs. Lasso coefficients for the filtered dataset, offensive (left) and defensive (right). ....	72
Figure 3.12: Scatterplots of Ridge vs. non-zero Lasso coefficients for the filtered dataset, offensive (left) and defensive (right). ....	72
Figure 3.13: Lasso RAPM (fitted) values vs. actual points per possession. ....	73
<b>3.5.6 Validation of RAPM values: Comparison with objective criteria and ratings.....</b>	<b>74</b>
Figure 3.14: Minutes played of top 50 offensive and defensive RAPM players via regularized models (first row: full dataset, second row: filtered dataset).....	76

<b>3.6 Player Evaluation using RAPM.....</b>	<b>76</b>
<b>3.6.1 Comparing player's contribution .....</b>	<b>77</b>
Figure 3.15: Scatterplot of Ridge Offensive and Defensive RAPM for a selection of players* .....	78
Figure 3.16: Scatterplot of Lasso Offensive and Defensive RAPM for a selection of players* .....	79
<b>3.6.2 Comparing traded players' contribution.....</b>	<b>80</b>
Figure 3.17: Radar plots for Andre Drummond and DeAndre Jordan's RAPM for their teams. ....	81
Figure 3.18: Radar plot for James Harden's RAPM for his teams. ....	81
Figure 3.19: Radar plot for Derrick White's RAPM for his teams. ....	82
Figure 3.20: Radar plot for Rajon Rondo RAPM for his teams.....	83
<b>Chapter 4: Building Logistic Models .....</b>	<b>85</b>
<b>4.2 Binomial models .....</b>	<b>85</b>
Figure 4.1: Scatterplot of coefficients from Lasso Binomial and Binomial (after the implementation of Lasso).....	88
Figure 4.2: Scatterplot of coefficients from Ridge Binomial and Ridge Normal model. ....	89
<b>4.3 Separate binomial models .....</b>	<b>90</b>
Figure 4.3: Contribution of selected players for the three separate Lasso Binomial models.....	93
Figure 4.4: Minutes played for reference group of the Multinomial model per team winning record level* . ....	96
Figure 4.5: Points scored for reference group of the Multinomial model per team winning record level* .	96
Figure 4.6: Minutes played for reference group of the Multinomial and Normal models. ....	97
<b>4.4 Multinomial model.....</b>	<b>98</b>
<b>4.4.2 Results of EPTS-RAPM .....</b>	<b>100</b>
Figure 4.7: Minutes played for reference group of the Multinomial and Normal models. ....	102
Table 4.8: Best players per team via EPTS-RAPMs.....	103
Figure 4.8: Frequency of actual vs. simulated points per possession. ....	104
Figure 4.9: Histogram of actual (data) vs. fitted points (Multinomial) per game. ....	105
<b>Chapter 5: Comparison of the different RAPM Models .....</b>	<b>107</b>
<b>5.2 Multinomial vs. Normal.....</b>	<b>107</b>
Figure 5.1: Actual (data) vs. fitted points (Multinomial & Normal) per game.....	109
Figure 5.2: Minutes played of the 100 high (blue) and low (red) rated players (50 offenders and 50 defenders) according to Normal and Multinomial models. ....	110
Figure 5.3: Minutes played per RAPM ranking position of Multinomial, Ridge Normal (first row), Lasso and Screening Lasso Normal (second row) models.....	111

Figure 5.4: Winning record of the top 100 (50 offenders and 50 defenders) and bottom 100 (blue and red respectively) rated players' teams according to Normal and Multinomial models.....	112
Figure 5.5: Points per game of the top-50 (blue) and bottom-50 (red) rated offenders according to Normal and Multinomial models. ....	113
<b>5.3 Comparison of all method .....</b>	<b>113</b>
<b>5.3.1 Agreement between results and selected criteria .....</b>	<b>113</b>
Figure 5.6: Distribution of expected points-Contribution via Multinomial and Ridge Normal model.....	114
<b>5.3.2 Five-lineup analysis.....</b>	<b>118</b>
Figure 5.7: The highest playing time 5-lineup (at the center) of Atlanta vs. the most impactful players (defense in the red and offense in the blue half-court).....	120
Figure 5.8: Lineup with the highest net rating of Atlanta (at the center) vs. the most impactful players (defense in the red and offense in the blue half-court).....	120
Figure 5.9: The highest playing time 5-lineup (at the center) of Warriors vs. the most impactful players (defense in the red and offense in the blue half-court).....	122
Figure 5.10: Lineup with the highest net rating of Warriors (at the center) vs. the most impactful players (defense in the red and offense in the blue half-court).....	122
Figure 5.11: The highest playing time 5-lineup (at the center) of Detroit vs. the most impactful players (defense in the red and offense in the blue half-court).....	123
Figure 5.12: Lineup with the highest net rating of DET (at the center) vs. the most impactful players (defense in the red and offense in the blue half-court).....	124



# Chapter 1: Review for Basketball Analytics

## 1.1 Introduction to Sports Statistics and Basketball Analytics

From time to time, studies have been carried out regarding sports statistics and analysis of data we receive from either individual or team sports. Perhaps, the most famous among them, due to its transfer to the big screen, is the book "Moneyball" by Michael Lewis (2003) that analyzes the success of the Oakland Athletics in MLB (Major League Baseball): a team that with the use of interpretive, game-related data and statistics became competitive despite its small budget. Baseball

Benjamin Alamar founded the first journal dedicated to sports statistics, the *Journal of Quantitative Analysis in Sports*. In his book, "Sports Analytics: A Guide for Coaches, Managers and Other Decision Makers", Alamar (2013) provides a survey of the practice and a detailed understanding of analytics' vast possibilities. He presented examples from professional sports teams and case studies demonstrating the use and value of analytics in the field of sports statistics. His work was a roadmap for managers and general managers.

After years, the field of sports statistics is on edge and several researchers have delved into a variety of topics. Singh (2020) offered a systematic review of research in the field of sport analytics. In various studies databases, which are classified based on business context and analytical methodology, are analyzed. There are works and papers in different application areas such as: Social media marketing, Data analytics in sport management, Technology and data analytics in sport management, Performance of teams, Payroll of players, Predictive accuracy, Performance assessment, Performance evaluation, Data analytics to predict match outcome, Betting odds, Visual analytics. Somewhat earlier, Rakshit and Babbar (2019) presented an overview of major analytical tools and technologies creating value in the sports industry with the major focus on soccer, basketball and cricket and especially on techniques to measure and improve players' performance. They, also, discussed the way big data and business intelligence could help in developing a theoretical model for tactical decision making in team sports.

Focusing on basketball, the man who established the importance and use of Statistics is Dean Oliver in 2004, who expanded on the philosophy of "Four Factors" from his book "Basketball on Paper" (2002) in an attempt to identify how four important strategies relate to success in basketball.

In essence, there are eight factors, four offensive and the corresponding four defensive. Dean Oliver determined a specific set of weights in each of the four categories to ascertain the value of each factor in relation to win. These factors, in order of decreasing importance for the result of the match, concern the shot, turnovers, rebounds, free throws and are as follows: effective field goal percentage (eFG%), turnovers per possession (TPP), percentage of offensive rebounding percentage (ORP) and number of attacks that end in free throw (free throw rate-FTR). Several times, Dean Oliver is referred as the "father" of Basketball Statistics. This, in fact, is confirmed by the research that followed the publication of "Four Factors" (2004) which gave rise to numerous studies.

Another important person for the development of basketball statistics is John Hollinger (Pro Basketball 2002, 2003, 2004, 2005) who, in his work, devoted himself to an extensive analysis and presentation of the statistics recorded in basketball games but with a more substantial look, aiming at the formation of an overall players' and teams' performance. Based on the number of possessions, Hollinger achieves to a significant extent the deepening of the defensive and offensive ability and performance of the players, not only at the individual level, but also of the team. It makes predictions based on statistics too. For example, how a team should base its attack in order to be efficient, always with the goal of winning. Hollinger is known for creating the Player Efficiency Rating (PER) statistical index. In his attempt to include all the contributions of a player in one number, he proposed the specific detailed formula that evaluates the statistical performance of each player, which is now widely used in all leagues.

In our days, one of the individuals who really pushed for Data Analytics driven basketball is Daryl Morey, General Manager of the Houston Rockets since 2007. Because of the fact that Morey's background is in statistics rather than basketball, he tried to build the play style of his team being data driven. In 2019-2020 Houston Rockets, based on statistics, decided to play without a center, known as "small ball". The fact that many teams build their offense in 3-point shoot in modern basketball is remarkable. Despite all the above, the basis of basketball analytics was laid earlier.

In one of the first studies, Stefani (1977a, b) introduced the regression-based approach with the Least Squares method in which, after initially proposing a way to evaluate college basketball and football teams using the average of the opponent's ratings with a margin of win, she then proceeded to predict the results of the games.

## 1.2 Model-Based Analytics for College Championships

NCAA has been the subject of study for several researchers, such as Schwertman et al. (1991), Carlin (1991) and Harville (2003), who, by applying linear and logistic regression, estimated the teams' exact scores, as well as their winning probability. Schwertman et al. (1991) calculated probabilities based on the assumption that the factor of team potentiality is normally distributed, while Carlin (1991) improves probability models proposed, such as Schwertman et al. (1991), taking advantage of external information for the potentiality of the teams. Harville (2003) proposes applying Least Squares to a statistical model in which the expected difference in score in each game is modeled as a difference in team effects, taking in consider whether the team plays at home or away. In fact, with some modifications he applied, the prediction model he proposed can compete with betting companies' prediction models.

Playoffs of NCAA, known as March Madness tournament, have great popularity. Jacobson and King (2009) searched about probit models with satisfactory predictive ability in the early rounds of the tournament. The authors found that the ranking of each team that enters the tournament (known as “seeds”) is a significant predictor for the earlier rounds but has less value in predicting outcomes of the final three rounds of the March MADness. As such, some models attempt to predict where seeding will go wrong (Schwertman, McCready and Howard 1991). This approach generates an upset probability for each game and then, using some threshold  $\lambda$  for all games with an upset probability  $p > \lambda$ , the higher-seeded team is predicted to lose. While this method has seen some success in recent years (Bryan, Steinke and Wilkins 2006), there is also evidence that it leads to systematic over-prediction of upsets, ultimately reducing the predictive accuracy of these models to below the seeding baseline itself (McCrea and Hirt 2009).

Based on the previous studies, Wright (2012) tried to build a model that uses more data for each team and can predict outcomes of the later round games. He analyzed data from 1986-2010 which include 1575 matchups (every matchup is one observation). He ran a probit regression and an ordinary least squares (OLS) linear regression. The dependent variable was a dummy for a win in the probit model and the margin of win in the OLS linear regression. The author developed 4 models: The 2 included variables that were only presented in the dataset from 1997 onwards, such as offensive efficiency and defensive efficiency and were developed by running a probit regression and an ordinary least squares (OLS) linear regression. The other 2 models included variables that



were presented in the entire dataset (from 1986 onwards) by using the same methods as in the previous. There were only slight differences in the results that came from each model. Crucial is that three out of four models successfully predicted the winners of the matchups in 2012 March Madness. The OLS Regression for the smaller dataset didn't perform well and it seemed to pick an unreasonable amount of upsets. A possible reason for the failure of this model might be that it heavily weighted the number of points that opponents gave up during the regular season.

Yuan et al. (2015) produced over 30 models to predict win probabilities for all possible matchups in the 2014 NCAA Men's Division I Basketball tournament that minimize loss between the predicted win probabilities and realized outcomes. Better results came by using logistic regression with lasso and ridge regularization. They, also, applied stochastic gradient boosting (SGB) and performed 10-fold cross validation on the training set of previous tournament seasons, 2003–2004 through 2012–2013, to optimize tuning parameters. Finally, they presented a neural network with a single 5-node hidden layer and a single-node output layer. Every node used a logistic activation function and model performance was evaluated by loss.

In a more recent work Brown (2019) built a logistic regression model to predict the winning team of a college basketball game, giving the historic performance metrics of the two teams. He was inspired by previous researchers, like Shanahan (1984) who built a logistic regression model to predict the probability of a win for a college basketball game. Shanahan used data from the University of Iowa men's and women's basketball teams from 1981-1983 and built a model for each team. She interestingly found that the significant variables included in the women's model were more offensive-based, while the variables in the men's model were more defensive-based. Another work came from Magel and Unruh (2013) who used logistic regression and least squares regression models with several explanatory variables that hadn't been used a lot till then, (such as home court advantage, difference in offensive rebounds, difference in defensive rebounds, difference in assists and difference in blocks), so as to determine different outcomes pertaining to a college basketball game.

### **1.3 The “home team” factor**

In some surveys, the “home” factor was also taken into account as an advantage for the home teams over the guests, like the last work above. Having got an idea from works about NCAA and

college basketball teams, below the review, continues with a reference to the NBA and less to other professional leagues.

Starting with the Pythagorean wins index (Bill James, 1980), this is an alternative method of predicting the percentage of a team's wins. It is an invention of Bill James in baseball who relied on the fact that the winning percentages of the team are generally closely related to the points scored and the points received. The Pythagorean winning percentage (Pyth) is calculated as follows:

$$PYTH_t = \frac{PTS_t^x}{PTS_t^x + PTS_o^x}$$

where the  $t$  index still indicates group, the  $o$  index indicates the opponent and the  $x$  index is an exponent determined empirically. In baseball, James (1985), empirically concluded that the exponent  $x = 2$  leads to the use of the index.

In another search, (Miller, 2006), it was found that a slightly lower exponent of about 1.8 works best. In the NBA the value of  $x$  has been empirically estimated by different studies to be between about 13 and 17 (e.g. Oliver, 1996, Coolstanding, 2006). The value varies depending on the season. In particular, the margins of win have changed slightly in recent years. These results, in smaller exponents, are necessary to correlate the points with the percentage of wins. This effect, however, if we consider the very good and very bad teams, which are only a few each season, becomes much better with a larger exponent. Typical exponent error estimates with the Least Squares method tend to yield a smaller exponent because most teams' wins accumulate between about 30 and 50 in a plurality (corresponding to percentages between about 38% and 62%). However, increasing the exponent does not significantly compromise the prediction of wins of very good or very bad teams and therefore Oliver (ESPN 2007) established the use of exponent 16.5.

In other leagues, including the NCAA and WNBA, smaller exponents have been considered to work better, as fewer possessions are in a typical game in those leagues (see Pomeroy (2006), for instance. Also, it should be noted that due to the equality of possessions for a team and its opponents,  $PTS_t$  and  $PTS_o$  can be replaced in the above equation by Offensive Rating (Ortg) and Defensive Rating (Drtg).

Compared to "Pythagorean victories", The Bell Curve method (introduced by Dean Oliver, 2004) is a more theoretical approximation of the points scored by a team and those received by the opponent, in terms of their winning percentage. Unlike the Pythagorean winning method, this method assumes that the points a team achieves and those it receives are normally distributed and, by subtracting them, they can form another normally distributed random variable, expressing "net points".

The formula for predicting the winning percentage for team  $t$  ( $Win\%$ ) of this method:

$$Win\% = NORMDIST \left[ \frac{PPG_t - PPG_o}{StDev(PPG_t - PPG_o)} \right]$$

where  $PPG_t$  is the points per game for team  $t$ ,  $PPG_o$  is the points per game for opponents  $o$ ,  $StDev(PPG_t - PPG_o)$  is the standard deviation of net points ( $PPG_t - PPG_o$ ) throughout a team's playing season and  $NORMSDIST$  is the normal Cumulative Distribution Function and represents the area under the standard normal "bell curve" to the left of the value in brackets. The advantages of this method, in relation to "Pythagorean victories", are that it does not need empirical modification for application in different leagues or different time and competition periods while, at the same time, it incorporates information about how much more or less teams score against their opponents.

Examples of such studies, including the home-team factor, are those of Harville and Smith (1994), Jones (2007, 2008) and Entine and Small (2008). In those studies, simple linear models are proposed. The results of the third study indicated that lack of rest for the road team, while not a dominant factor, is an important contributor to the home court advantage in the NBA. Being inspired by statistical works on baseball, Hu and Zidek (2004) used a weighted likelihood in order to predict the outcome of the NBA playoffs for the season 1996-1997. It is important to mention that they tried to reflect the home game advantage by applying logistic regression.

Berri et al. (2006), examined sport using statistical tools and methods and tried to explore the important factors for winning games from a statistical point of view, while, at the same time, working on the creation of statistical indicators. Berri et al. (2006), made it abundantly clear that much of the decision-making of coaches and general managers is not based on the analysis of numbers. The authors assumed that players are not racially influenced by their teammates, a fact that allows them to develop a linear formulation of weights by determining the relationship

between wins and team statistics first and then assuming a similar relationship between wins and individual player statistics.

## **1.4 Inspiration from Dean Oliver: Model based works for Prediction and Performance**

There are works based on Dean Oliver's "Four Factors" (2004), such as Küpfer's (2005), in which various weights were placed on each of the four factors mentioned above, in order to examine the importance of these in relation to the wins achieved by a team. This study, like other counterparts, simply confirmed Dean Oliver's results in the significance of the four factors.

Kubatko et al. (2007) collected the key variables of basketball analysis deriving from the key statistics and defined a framework on which subsequent investigations can be founded. Of course, Dean Oliver's "Four Factors" (2004) for which Kubatko et al. (2007) observed that their average value varied over time and even observed a general overall decrease in expected offensive indicators and apparently turnover, are very important. Also noteworthy is the introduction of advanced statistics calculated through box-score data and which have now become established in the "statistical world" of basketball. At this point, we should emphasize that Kubatko et al. (2007), was based on the statistical data of possessions following his predecessors Dean Oliver (2004) and Hollinger (2002, 2003, 2004, 2005).

As multivariate analysis makes it possible to improve the study of the effect of the four factors and not just a team's win, Doolittle and Pelton (2010) investigated the team's statistics for these factors and tried to correlate them with the win, or even to develop a method of predicting the outcome based on individual game performance predictions for each player.

Teramoto and Cross (2010), also, by using multivariate analysis, estimated the impact of four offensive and defensive factors on an NBA basketball team's winning percentage in the regular season and playoffs. In general, they confirmed the results of Küpfer (2005) under the influence of the four factors. It was found, therefore, that offensive and defensive Effective Field Goals (eFG) were the most statistically significant factors for winning both the regular season and the playoffs, with the defensive eFG the most important factor for the playoffs. On the other hand, the offensive and defensive Free Throw Rate (FTR), despite their importance in predicting the outcome (win or lose), was relatively the least important of the four factors.

Kline (2005), escaping both the analysis of Least Squares and the game-related factors considered by the majority of surveys, introduced into the study of criteria that lead a team to win is the factor of salary. By applying models of structural equation modeling (SEM), he tried to estimate the "real", substantial level of the "structure of a construction", such as in the case of the offensive or defensive quality of a team, based on various indicators. Thus, as in a normal regression model, various factors are estimated, but by integrating path analysis between the variables and those estimated factors. Thus, it is possible to correlate the factors of the study with the salary and then with the wins of a team.

A few years later, Tarek Al Baghal (2012), based on previous studies and influenced by Kline (2005), first tried to introduce and establish structural equation models (SEM) as an Effective Statistical Tool in the field of sport analytics and statistics. He worked on regression analysis for winning based on all four factors, using a larger number of match periods than Teramoto and Cross (2010). He also used SEM to determine whether the four factors could be modeled as indicators of more general individual variables, concerning the offensive and defensive quality of teams. He, then, correlated these two general characteristics of a team with its wins. Additionally, team's salaries were incorporated into the SEM models in order to assess the cost ratio with offensive and defensive quality as well as winning percentage. Baghal found that each team's total salaries have a high effect on its offensive quality and influence on a moderate degree the team's winning percentage. He concluded, therefore, that total salaries have an indirect effect on team wins, as high salaries imply a high level of offensive quality and in turn imply a high percentage of wins.

Interesting results are presented by Puranmalka (2013), who tried different models in predicting NBA outcomes such as Bayes Net, Naive Bayes, Linear and Logistic Regression, SVMs, k-nn. He proposed various new features derived from play-by-play data and many of these are important in making predictions of NBA games' outcomes. In more detail, the team-level features he proposed, namely a new measure of clutch Performance and various team-to-team synergies are important in predicting outcomes of games. Teams' rest<sup>1</sup> and distance travelled don't seem to have predictive value in the NBA. As for the player-level features, context added efficiency<sup>2</sup> and the measurement

---

<sup>1</sup> How much rest the team has gotten recently. There are different ways of measuring this, such as the number of games played in last x days or the number of days since the last game.

<sup>2</sup> A metric which incorporates the time left in the shot clock as well as general efficiency of the offensive and defensive team.

of how a player affects the shot selection of his team-mates, seems to be the most important set of player-level features.

Two studies that innovated in terms of the data they used, like the previous work mentioned, are noteworthy. First, the survey of Yuanhao (Stanley) Yang (2015), in which the individual statistics of the weightiest players were also included. In particular, Yang (2015), not departing from the trivial methodologies, applied multiple linear least squares regression to correlate, initially, the individual statistics of the players and specifically the per indicator that measures the efficiency of the player, with the result of the game for his team (win or lose). He then developed a regular season results predictive model. One year later, Eric Scot Jones (2016) proceeded to study using in-game statistics, in contrast to all the previous surveys. He was involved in estimating a model that explains the game's final point spread through in-game statistics and a predictive model of the game's outcome by analyzing the in-game factors that lead to win.

Building on the linear modeling that had been developed by 2015, Hans Manner (2015) introduced the factor of home team and back-to-back games too. The model he proposes estimates the team's capacity/potentiality in relation to the variability of time over the period of eight seasons (years). This is a dynamic Gaussian self-induced process. Manner (2015) concluded that the capacity of a team is maintained throughout all seasons and not in each individual season. He, also, empirically concluded that the capacity/potentiality of each team is significantly affected by injuries, trades and other factors that affect the composition and "chemistry" of the team.

Song et al. (2018) studied the problem of modeling and forecasting the outcomes of NBA basketball games based on performance statistics. A bivariate normal mean regression model was developed to model the scores and performance statistics.

Malarranha et al. (2013) studied the team performance in a different way. The aim of their study was to identify the intra-game variation from the four indicators (effective field goal percentage, offensive and defensive ratings, offensive rebounds percentage) that determine the outcome of basketball games. Their aim was to study and calculate the performance indicators in eight 5-minute periods (not the four quarters of the game), by analyzing all games of the Basketball World Championship (Turkey 2010). A repeated measures ANOVA was performed to identify differences in time and game outcome for each performance indicator. Also, one covariate of the fitted models was the quality of opponent teams. One special finding was that offensive rebounds

percentage had greater influence in the second half, in contrast with the fact that effective field goal percentage, offensive and defensive ratings influenced the game outcome throughout the game.

So far, most of the studies have referred to the prediction of the outcome of the game and teams' performance. Interesting, of course, is the performance evaluation of the players and the impact they have at their teams. Let's take a look at it without being transported to far back in time.

Esteller-More and Eres-Garcia (2002) presented the use of the Atkinson function to evaluate the performance of a player in any category of the game considering the inconsistency of scoring. They applied it to the ranking of scoring leaders in the NBA regular season and playoffs of 2000-2001 and they proposed an indication of how to deal with the valuation of consistency in the empirical estimation of players' compensation.

Casals and Martinez (2013) took advantage of the fact that there wasn't any study in the literature, till then, about modelling players performance through quantifiable play by play data. The authors identified variables which may affect the player's performance from one game to another and put them together in a statistical model to control their covariance. They, also, studied the significance of these variables through mixed models, showing the conditional effects on performance. It is important to mention that this research is the first attempt in the basketball analytics' literature of studying variation in player's performance in a game-by-game context.

Martinez et al. (2017) tried to study the regularity of points scored<sup>3</sup> by using data for 27 NBA players in regular season 2007-2008. In their first approach they used a full season evaluation, modelling regularity by using a robust version of the median absolute deviation for variables (in order to explain variation) and the Cochran variance outlier test in order to identify the players with the greatest variance in their performance. According to Martinez et al. (2017) some players may be more prone to influence by the variables identified by Casals and Martinez (2013), while some others may even be subject to a momentum effect. Therefore, in their second approach they analyzed the ordinal patterns of players' performance using short-term evaluations (three games per week).

---

<sup>3</sup> A measure of consistency in player performance.

An interesting study with practical significance comes from Lorenzo et al. (2019), who analyzed the changes of game-related statistics in expert players across their whole sports careers. Their study included professional basketball players from Spanish first division basketball league (ACB) in the 2017–2018 season. The following game-related statistics were studied: average points, assist, rebounds (all normalized by minute played), 3-point field goals percentage, 2-point field goals percentage and free throws percentage per season. According to the results, an increase in assist and free throw performance was shown in the players (of the sample studied) across their playing career. This information is essential for basketball coaches suggesting the use of most experienced players in the final moments of the game.

## **1.5 Moving away from the usuals**

So far, many studies based on regression (main simple linear, logistic, OLS) have been reported and only a few involved some kind of innovation. Of course, there are researchers that studied different topics based on alternative methods and theories or different types of data.

### **1.5.1 Some innovative works**

Starting with Shea and Baker (2013), in their book, “Basketball Analytics” presented original, objective and efficient strategies for understanding how teams win. New player performance statistics are proposed. Based on results of previous studies that referred to the NBA draft pick value, the authors discussed topics including the biggest draft blunders and steals, the draft success of each NBA team and the quality of each draft class dating back to 1977. This valuable understanding of the NBA Draft creates a foundation for discussing various approaches to team development and construction. Additionally, the authors discussed redefining the positions on the court, unpredictability in the game, data visualization and applications of spatial tracking technology.

One year later, Shea (2014) using new spatial tracking data collected by SportVU and others, in “Basketball Analytics” investigated game strategy, player evaluation, player types and prospect potential. He introduced new measures of a player's scoring and playmaking efficiency, quantified the spacing in an offense and the stretch of a defense and demonstrated several ways in which the NBA game has changed over the years. He presented a modern viewpoint on basketball analytics'



most fundamental principles and demonstrated the power of the industry's latest statistical breakthroughs.

Hernandez et al. (2013) had already presented their study which was based on spatial statistics that had been rarely used in the field of sports analytics till then, especially in basketball. They proposed spatial clustering techniques<sup>4</sup>, such as the Kulldroff test, for analyzing basketball data. This test detects low and high incidence clusters of shots and therefore it better characterizes the game of teams and individual players. The authors, also, used the V-test<sup>5</sup> to compare shooting maps that are very popular in basketball analytics. The application they presented was referred to the transformation of a medium-level NBA franchise into a champion team.

Completing the above, Zuccolotto et al. (2021) proposed a spatial statistical method based on classification trees to define a partition of the court in rectangles with maximally different shooting performances. Each analyzed team or player belongs to some rectangles of the court that depict their shooting performance from this spot (or “neighbour”). In this way, comparison between teams or players are available. Anyone could say that authors’ proposal is accessible to the unscientific public, especially the visualization of this.

Going back in time, Hoffman and Joseph (2003) used the Principal Component Analysis to identify the most important factors that will lead a team to qualify for the NBA playoffs. Additionally, they predicted the probability of a team finishing in one of the positions leading up to the post season. The factors they included in their research are both game-related and non-game-related: Average points scored per game, Average points allowed per game, Team percentage of field goals made, Number of years since team's establishment, Rank of team's total yearly payroll, Head Coach's NBA record, Defensive turnovers less offensive turnovers, Last season's percentage of games won, Average attendance per game, Number of rebounds per game, Ratio of new players, Median age of team.

Continuing, here are two theses using neural networks to study the factors associated with wins. Initially, Renato Amorim Torres (2013) tried to predict the results of the upcoming game by using linear regression and, in particular, the maximum likelihood method as well as through a multi-layered neural network, the results of which he compared. In this context, Jaak Uudmae (2017)

---

<sup>4</sup> which are widely employed in epidemiology

<sup>5</sup> a test based on entropy

tried to predict the results of the upcoming game (win or lose). He compared the results of different methods, including support vector machine (SVM), linear regression and neural network regression (NNR).

In the previous three studies of Lori Hoffman and Maria Joseph (2003), Renato Amorim Torres (2013) and Jaak Uudmae (2017), where an attempt at methodological innovation was made, it is observed that they are more result-oriented and not much attention is paid to the statistical structure of the data. In the application of the linear model, it was considered that the data are independent, which is not true since the teams participate in different games (observations). Thus, Xiao Zhang (2019), based on the hypothesis of correlation between observations, introduces a statistical model with the method of Generalized Estimating Equation (GEE) to estimate the results of the games and explain the most important factors.

In a recent work by Huang and Lin (2020) a different type of regression is proposed. They analyzed game data of the Golden State Warriors and their opponents in the 2017–2018 season of the NBA. They developed a regression tree model for score prediction of each player on two teams for every game and they compared the predicted total scores to obtain the predicted results (lose or win) of the team of interest. The predictive accuracy of the model was satisfactory in high degree.

### **1.5.2 “Basketball-oriented” studies**

Before we continue with methodological alternative studies, it is worth mentioning two works in which women’s basketball data are analyzed. The first is a statistical analysis of momentum in basketball that analyzes the distribution of time between scoring events for the BGSU Women’s Basketball team from 2011 to 2017. According to this study, the scoring events within a game follow the Poisson process, while the time between scoring events for each game can typically be modeled by an Exponential distribution with a mean equal to the reciprocal of the average time between successive BGSU scoring events within that game. However, there is a significant number of games in which the exponential model does not fit the data well. This may mean that these games have larger or smaller gaps between scoring events.

The second study Noivo et al. (2022) identified the predictors of success in ball screens considering time, space, player, task, and contextual related variables in elite Women’s Portuguese Basketball League. Classification tree analysis with Chi Squared Automatic Interaction Detection method

was used to identify the set of variables that best predicted the success in ball screens. Results showed that quarter half, time possession remaining, finishing action, defensive strategy and offensive system were the best predictors of ball screen effectiveness. Several conclusions had been reached with practical significance. For example, playing ball screens in transition at the beginning of each quarter it is very important. In addition, a drive from the ball handler, as well as a pop out from the screener were the finishing actions that ensured greater success after the ball screen. This type of information would be helpful for basketball coaches to make safer decisions when planning tactical strategies for their teams, although, there are not many studies with such a practical significance.

One could argue that the above work has more basketball content than all the previous publications we have cited above. Taking advantage of this, we should mention that NBA draft picks are of scientific interest and maybe more for operational research and less for statistics. However, in his research paper, Parker (2018), aimed to understand the statistical analysis behind the draft pick success in the NBA. The three discussed methodologies are Player Efficiency Rating (PER) and standard statistical categories, Roland Beech's Rating System and Win Shares. Due to the paper, the most efficient method is to combine the three above.

Moreover, the study of Martinez and Caro (2011) who examined the opinions of basketball stakeholders regarding several questions of special interests to value players, is remarkable. After analyzing opinions of players, coaches, agents, journalists, editors, bloggers, researchers, analysts, fans and chairs, by using the content analysis methodology, they found that player's evaluation systems "ignore" intangibles.

### **1.5.3 Bayesian, Unsupervised Learning and Big data methods**

Coming back, emphasizing the methodological part, Bayesian's theory doesn't absent from the literature. Indicatively, in the work of Deshpande and Jensen (2016) a Bayesian linear regression model is proposed, so as to estimate an individual player's impact. According to the authors there are highly paid players with low impact relative to their teammates, as well as players whose high impact is not captured by existing metrics. Maybe intangibles is an important factor that not everyone is able to see.

Shi and Song (2018) proposed a discrete-time and a finite-state Markov chain model was developed to fit the NBA basketball data. It can be used to produce in-play prediction for basketball matches. The authors proposed a model to calculate probabilities of the final score difference, which performs well.

Zuccolotto et al. (2019) proposed a different method. In their work the variability “holds a central position”. Shooting performance variability was modeled with a Markov Switching dynamic, assuming the existence of two alternating performance regimes. Then, the relationships between each player’s variability and the layup composition were modeled as an ARIMA process with covariates and was described with network analysis tools, to extrapolate positive and negative interactions between teammates, helping the coach to decide the best substitution during the game.

Brown (2017) applied the idea of Google PageRank to rate and rank the basketball, soccer and hockey players of a game. The foundations of PageRank lie in Markov chain theory. There were interesting results from the model, where some players, who had impressive stat-lines, had lower ranks and others, who had less impressive stat-lines, had higher ranks. The model’s ranking and ratings reflect the flow of the game more compared to traditional sports statistics.

Unsupervised Learning methods is one of the topics and methods that could not be missing from the relevant literature. There is a large increase of studies, over the last five years, that have been applied clustering methods. Only some of them will be mentioned here for brevity.

Patel (2017), in his thesis, attempted to reclassify NBA players into new groups based on personal performance in the 2016-2017 regular season. By using k-means clustering, he revealed four groups of players with similar playing styles in contrast to the five typical positions. The four clusters are presented below:

- The Paint Protectors (moderate scorers with great 2-point shots, rebounds, and blocks but poor 3-point shooting ability)
- The Supporters (relatively low scorers with good assists, stealing ability and decent 3-point shooting ability)
- The Shooters (moderate scorers with great free throw, 2-point, and 3-point shooting ability)
- The Insiders (high scorers who excel at free throws, 2-point shots, rebounds and blocks)

In the same path, Armanious (2019) in his thesis categorized the players of the regular seasons from 2015 to 2019 of the men's U Sports teams for the Ontario University Athletics (OUA) Division into 8 clusters by using k-means. The types of players he proposed are:

- Efficient Playmakers & Scorers: This cluster of players have the most assists and the second most points per game. They have a big defensive impact through the number of steals they get and can control the tempo well and score.
- All-Around Players: These players can get rebounds, pass and score well.
- Dominant Big Men: These players are the most dominant big men in the league with the most rebounds (defensive and offensive), blocks and points.
- Smart Catch & Shoot Players: These players make the best decisions and turnover the ball the fewest. They do not dribble the ball much and are the most efficient shooters.
- Aggressive Defenders: These players are aggressive and foul the most out of all the other clusters. They have a bigger impact on defense since they do not shoot well.
- Role Players: These players contribute to many plays and work both offensively and defensively.
- Second Tier Playmakers: These players are less dominant playmakers that can still score efficiently.
- Second Tier Small Players: These players play small but do not shoot as efficiently as the other players or create as many plays.

Hussian (2019) worked to find a better approach to define players roles, based on the value they bring to their team. By applied k-means for statistics of every NBA player from 2011 to 2018, he ended up to the next clusters:

- Perimeter Wing/Scorers
- Three & D
- Do it All, Elite Wings
- Backup Bigs (Inside)

- Elite Bigs (Inside)
- Star Bigs (Inside)
- All Stars
- Superstars.

Jyad (2020) used hierarchical clustering in order to redefine NBA players classifications. The 9 clusters he proposed by studying season 2018-2019 are:

- Elite Modern Big Men
- Traditional Big Men
- Elite 3-point shooters
- Role players
- 3 and D-players
- Level Scorers
- Decent Ball Handlers
- Elite All Stars
- 2-way Perimeter players.

A slightly different procedure comes from Duman et al. (2021) who grouped NBA basketball players into similar clusters, according to their playing styles, for each of the traditionally defined five positions. This is helpful for teams' recruiting and building because coaches are able to focus on players' characteristics in the game. In this work, 17 game-related statistics from 15 seasons of the NBA were analyzed using a hierarchical clustering method. Based on this analysis, four clusters were identified for PG, SG and SF positions, while five clusters for PF position and six clusters for C position were established. In addition to the definition of the created clusters, their individual achievements were examined based on three performance indicators: adjusted plus-minus (APM), average points differential and the percentage of clusters on winning teams. This study contributes to the evaluation of team compatibility, which is a significant part of winning, as it allows one to determine the playing styles for each position, while examining the success of position pair combinations.

Khobden et al. (2021) tried to cluster players based on their individual abilities by using Artificial neural network (ANN). NBA players' statistics from 2011-2018 were considered as features (such

as points, rebounds, goal passes, ball stealing, defenses, etc.). For this purpose, Self-Organizing Map (SOM) Neural Networks were used. The result of the study showed that the performance of the SOM in clustering basketball players was higher than the k-means algorithm.

Furthermore, Hong (2021) processed big data in order to group players. He used Spark framework (based on memory computing to enhance the effect of basketball data analysis) and k-means algorithm. The shooting training effect in the active area has the most measurable impact and a good influence on the training effect.

In a very recent study, Hu et al. (2022) tried to conduct a descriptive analysis of the anthropometric features of the line-ups of strong teams (top 16) in the 2019 FIBA Basketball World Cup, to group the line-ups mentioned above into different clusters, based on their average height, weight and body mass index (BMI) and to explore the performance variables that discriminate between various line-up clusters. They processed big data, like the previous referenced researcher. Also, they analyzed play-by-play statistics using two-step cluster and discriminant analysis. Line-ups were classified into four groups: low average height and weight with middle BMI, high average height and low average weight with low BMI, low average height and high average weight with high BMI, high average height and weight with middle BMI. The results of the discriminant analysis demonstrated that LowH–LowW–MiddleBMI line-ups had the least time played and the lowest offensive rating, but the best offensive rebounds, turnovers and fastest game pace performance; HighH–LowW–LowBMI line-ups demonstrated the best defensive rating, but performed poorly with a low value of assists and a high value of turnovers; the LowH–HighW–HighBMI group achieved the best time played statistics, but had the lowest number of free throws made; the HighH–HighW–MiddleBMI group had a higher number of assists and a higher offensive rating and 2-point field goal performance, while also achieving the lowest number of offensive rebounds and ball possessions. The main results had been mentioned by Zhang (2019).

From devising game plans, to improving player and team performance, big data is changing the way basketball is being played. Basketball teams and businesses have similarities. Nowadays, they use big data in order to attract and select the right people for their teams and grow a competitive advantage in the championship/on the market. Branga (2021) compared teams and businesses and found some differences between them. The results, also, showed that big data analytics can help basketball teams' and businesses' recruiting.

Lee and Page (2021) studied advanced data analytics of men's professional basketball statistics of the last 16 seasons in more than 25 professional leagues and 71 FIBA tournaments. The complete database consists of more than 37,000 games and upwards of 20,000 players. Thus, we can imagine a huge dataset. By analyzing those data, they determined the players' performance curve, an optimal age in professional men's basketball, a rating correction factor for different basketball leagues, which accounts for intra-league and cross-league variability as well as for player characteristics. They, also, studied the factors for predicting the professional career of players. Finally, they provided statistical models to evaluate the players' performance based on position, age, skills, league and other characteristics and their influence in the game.

A comprehensive study has been published by Zuccolotto and Manisera (2020), related to their last year's work (2019). Using data from one season of NBA games, "Basketball Data Science" is a book-guide for basketball analytics, with applications in R. The book provided statistical and data mining methods for the growing field of analytics in basketball. Tools for modelling graphs and figures to visualize the data are presented. In fact, the authors provided their own programming package for analyzing basketball data and they presented many real-data examples.

#### **1.5.4 Data Envelopment Analysis (DEA)**

A completely different method compared to those already mentioned is Data Envelopment Analysis (DEA), which is an econometric method. There are not a lot of studies about this method. Cooper et al. (2009), who measured the effectiveness from 172 players participating in the Spanish Basketball League during the 2003-2004 season, argued that the DEA can be used interchangeably instead of the indicator of the evaluation of the players' performance. In addition, they focused on multiplier values for player effectiveness and on how these values can be used to find the advantages and disadvantages of each individual player.

Lee and Worthington (2013) examined the output effectiveness of 62 key regionals competing in the NBA's "1" and "2" positions for the 2011-2012 season. It should be noted that this period coincides with the outstanding performance of Jeremy Lin, known as "*Linsanity*". The results of their study showed that between 29% - 42% of the NBA's key guards were fully effective, including Jeremy Lin. However, Jeremy Lin rarely served as a benchmark for ineffective DMUs, which shows his unique style of play, including altruistic play and leadership ability.



Radovanovic et al. (2013) measured the performance of 26 NBA players during the 2011-2012 season using DEA and distance-based analysis (DBA). According to the results All NBA players were in the performance range from 70% - 116% and 7 players had a performance rating of more than 100% showing actual performance. The DEA and DBA results had a positive correlation, indicating that the two methods are in harmony with each other.

In contrast to the previous studies, Bartholomew and Collier (2011) presented a team-level study of the Collegiate League of America. In fact, he examined the defensive efficiency of the teams using DEA with DMU the 20 halves of matches, outflows the defensive rebounds (DR), the contested shots and inputs total opposing points (TOP). In fact, they studied a second model with the same inputs and outputs, forced turnovers (FTO), defensive rebounds (DR), total fouls (TF), contested shots (CS) and defensive steals (DS). DEA. Also, a third model they studied has as inputs the total opposing points (TOP) and field goal percentage (FG%) and outputs the FTO, DR, TF and CS.

Aizenberg et al. (2014) applied the DEA to measure the efficiency of NBA basketball teams in the years 2006-2010. In this context, each team constitutes a DMU, while the total payroll and the average attendance are selected as inputs, while the wins and average points per game are selected as output.

Hai Yang et al. (2014) proposed a 2-stage DEA implementation. In particular, they decomposed the overall performance of NBA teams into two pieces: one concerns salaries (1<sup>st</sup> DEA stage) and the other concerns the playing (on court) efficiency (2<sup>nd</sup> DEA stage). According to their results, NBA teams perform better in terms of salaries compared to the on court-part, as the effectiveness on the field is influenced by many uncontrollable factors. In addition, NBA teams, on average, tend to place more weight on the first stage, suggesting that management's decisions in player selection are quite critical to the team's development.

An alternative approach was presented by Moreno and Lozano (2014). With DEA network application (Network DEA) their goal is to evaluate the efficiency of NBA teams for the 2009-2010 regular season and compare the results with those of the simple DEA method. Both methods use a Slack-Based Measure (SBM) to assess the potential reduction in consumed inputs, teams' budget and output, the number of wins. The study examines the distribution of the budget between the players of the first-team and the payroll of the rest of the players. The proposed approach

consisted of five stages, which evaluate the performance of the first-team's players and the players of the "bench", the offensive and defensive systems and the ability to "transform" the points achieved and allowed into wins.

Villa and Lozano (2016), considering that the number of points scored in a basketball game greatly influences the appeal of a game, proposed a new approach that focused on measuring two teams' efficiency in scoring at one game. To do this, the performance of each team in each quarter must be considered. This results in a dynamic network DEA model with two subprocesses (corresponding to home and away teams), performed in each quarter. The data processed was related to the 2014-2015 NBA season.

Finally, in a recent study Assani et al. (2021) proposed models to measure players' overall, offensive and defensive efficiencies based on a non-homogeneous parallel data envelopment analysis (DEA) network. They, also, introduced input-output oriented network models to estimate the marginal returns from salary on the outcomes of both offensive and defensive activities.

### **1.5.5 Plus/Minus: At first glance**

To conclude up to this point of the review, the field of basketball analytics is on edge. In our days, many teams make use of statistics and analytics for their recruiting, training their players, building their tactical plan, scouting their opponents. In the above perhaps there is a lack of studies referred to the impact of players in their teams with more practical significance. The aim of this thesis is to study the players impact through a plus/minus model and use it to manage players' rest for the upcoming games. Thus, previous works for plus/minus models performing basketball data will be useful. There is no large number of such studies, but the necessary basis has been laid.

The starting point is Rosenbaum's (2004) work, in which he introduced Adjusted Plus-Minus to estimate the average number of points a player scores per 100 possessions after controlling for his opponents and teammates. He computed plus/minus ratings that measure how point differentials change when a particular player is in the game, versus when he is not. The logic of this approach is straightforward; teams should perform better when their good players are playing, versus when they are not. Because, according to Rosenbaum, these plus/minus ratings measure the value of the player relative to the players that substitute in for him and there are differences in the quality of players that play with and against which aren't included, he proposed the "adjust" plus/minus

ratings. These would account for home court advantage and for clutch time/garbage time play. Contributions for individual players are isolated statistically. He also proposed a hybrid approach by using both, pure and adjusted plus/minus ratings and presented an application by analyzing data from two seasons 2002-2004.

Ilardi and Barzilai (2008) presented the most accurate (low-noise) adjusted plus/minus rating till then, by using data for five seasons (from 2003-2004 to 2006-2007). In addition, they modeled separately each player's impact on offense and defense, treating these as completely independent variables.

Sill (2010) presented a framework for evaluating adjusted plus/minus models in terms of their ability to predict the outcome of future games. He used Bayesian regularization (specifically ridge regression) to improve accuracy of the model.

Some years later, Sisneros and Moer (2013) applied the concept of measuring impact through plus/minus ratings to all box score statistics and expanded the player performance to team performance by using data per game and considering traditional plus/minus numbers at the team level as a measure of the quality of a win/loss. They presented the plus-minus plot and examples found in 2012-2013 NBA box score data.

Deshpande and Jensen (2016) used plus/minus to produce a retrospective measure of individual player contributions and not to measure a player's latent ability or talent. Their proposal is context-dependent, thus they supported that their procedure provides a more appropriate accounting of what actually happened, than existing player evaluation metrics like PER and Real Plus/Minus (RPM) from ESPN. Their estimates of player effect could evaluate, if the playing time is divided effectively, and help to understand which player's individual performance is translated to wins for his team.

In some blogs anyone could find discussions about plus/minus models or ratings and stats. As for the practical significance of these studies, of course, they could help coaches and their stuff to decide about the players' playing time, rest or which combination of players is more efficient in different parts of a game.

## **1.6 Conclusion**

From all the above anyone could say that sports and, especially, basketball analytics have been developed over time. This field generates many interesting questions and case studies. Due to them basketball and scientific community should collaborate and exchange opinions and knowledge. Some believe that now statistics have gone deep in basketball analytics but without practical significance. This is something that those involved in the field may be concerned about.

In this study, we will try to present an analysis that will be useful for teams in a practical way. Play by play data from season 2021-2022 are analyzed. The goal of the research is to support teams' load management, by applying plus/minus methodology and optimizing the probability of wins or the number of wins in the upcoming games. In the second chapter more works for plus/minus ratings and models are presented.



## Chapter 2: Performance Analytics via Plus/Minus Models

### 2.1 Early years

Advanced statistics have grown up in the field of basketball analytics and of course in the field of sport statistics in general. Especially in team sports, one of the most important and useful metrics is the (positive or negative) effect a player has on his own team. As for basketball, impact player ratings were developed several years ago. Starting with the most popular index, Player Efficiency Rating (PER) and Win Shares have been established in basketball analytics and they are box score statistics.

The introduction of PER by Hollinger (2007)<sup>6</sup> through ESPN was the first big step for welcome advanced statistics in basketball and it also belongs to the Worldwide Leader of Sports. In 2007, PER was a unique pace-adjusted, per-minute rating of player productivity. The calculation includes 12 different statistics which are weighted differently (he gave more weight to the categories that he “feels” are more important).

Win Share was developed by Kubatko (2007) who, not only was inspired by Bill James who had created a similar metric for baseball (MLB) earlier, but also used ideas from Dean Oliver’s (2004) “Basketball on Paper”. The measure assigns a score that estimates exactly how many wins a player contributed to his team’s record. The Win Share is based on the seasonal numbers (i.e., for example, scoring a buzzer-beater or blocking the opponent’s buzzer-beater does not lead to higher Win Share values). In other words, the importance of a contribution within a game is not weighed differently. According to Kubatko: “A win share is worth one-third of a team win. If a team wins 60 games, there are 180 ‘Win Shares’ to distribute among the players.”. However, according to Ben Taylor (2019) for Nylon Calculus at the “FANSIDED” blog, Win Shares yielded particularly volatile results.

This work will be based on another type of metric, the Plus/Minus rating, which has several dimensions. Plus/Minus ratings have great popularity in team sports and they are widely used around the world. The most common approach is to compute plus/minus ratings that measure how

---

<sup>6</sup> Hollinger had used PER earlier, but he presented officially in 2007.

point differentials change when a particular player is in the game, compared to the period when he is not in the game.

Hockey experts pioneered the use of such a plus/minus system for many years in the past. Historically, the first team that used the plus-minus approach to track the performance its players sometime in the 1950s was NHL's Montreal Canadiens. Other teams followed in the early 1960s and the NHL started officially compiling the statistics for the 1967–1968 season. While Emile Francis, an ice hockey player, coach and general manager of NHL, was often credited with devising the system, he only popularized and adapted the system in use by the Canadiens. In those early years, plus-minus was used to evaluate a player's offensive and defensive contributions to his team at even strength. According to the NHL, it is calculated as follows: *“A player is awarded a ‘plus’ each time he is on the ice when his club scores an even-strength or shorthanded goal. He receives a ‘minus’ if he is on the ice for an even-strength or shorthanded goal scored by the opposing club.”*. Despite the dubiousness and discussions around the plus/minus, it is the most popularized statistic in hockey.

The metric plus/minus was used primarily with corresponding logic in other sports, until 2004, when Rosenbaum (2004) in “Measuring How NBA Players Help Their Teams Win”, 82games<sup>7</sup>, introduced the application and use of adjusted plus/minus (APM) ratings for the NBA. His work has inspired many researchers (not necessarily from the scientific community) to study and develop plus/minus models for basketball data till nowadays.

It is worth mentioning that the adjusted plus/minus technique was first developed by Wayne Winston and Jeff Sagarin in the form of their WINVAL<sup>8</sup> software system in 2002. For each NBA player, it starts with the team's average point differential for each possession when they are on the court. This gives a number showing how effective the player's team was when they were in the game. Its main problem is bias in favor of players that usually are in the court alongside good teammates or/and against weak opponents.

---

<sup>7</sup> <http://www.82games.com/comm30.htm>

<sup>8</sup> One of the first of its kind analytics software program, measures team chemistry and analyzes a player's impact on his team's ability to produce points.

## 2.2 Model based Plus/Minus ratings

### 2.2.1 Adjusted Plus/Minus

Rosenbaum (2004) said that basketball is not like baseball, a game structured around repeated one-on-one contests between pitchers and batters, where the contributions to winning of any given player can be measured well by individual game statistics. “Basketball is much more of a team game”, as Rosenbaum emphasized. He also mentioned that there are differences in the quality of players that players play with and against. According to him, “weak starter on a team with exceptionally good starters (relative to bench players) will generally get a very good unadjusted plus/minus rating, regardless of their actual contribution to the team”. In this context, Rosenbaum proposed adjusted plus/minus ratings to account, firstly, for the quality of players that a given player plays with and against and, secondly, the home court advantage and clutch time/garbage time play. On the other hand, the “unadjusted” plus/minus ratings measure the value of the player relative to the players that substitute him.

Let us now present the methodological framework of Rosenbaum (2004), “Measuring How NBA Players Help Their Teams Win”, 82games<sup>9</sup>. In his study, he considered all players with playing time higher than 250 minutes in seasons 2002-2004. Every observation was a segment of time (or period) within a game with no substitutions (i.e. with the same 5-player lineup). Observations were weighted by the number of possessions (with observations in 2003-04 weighted twice as high as those in 2002-03) and higher weights during critical “crunch” instants and lower (or zero) weights during “garbage” time. In the first attempt, the effects of the other players on the floor are accounted for as covariates. The coefficients measure the difference of the point margin (measured per 100 possessions) of the player under study relative to a set of reference players in the same lineup<sup>10</sup>. In other words, coefficients are plus/minus statistics adjusted for the other players on the floor.

However, those ratings as mathematical estimates, each plus-minus rating contains measurement noise. To solve this problem, Rosenbaum (2004), “Measuring How NBA Players Help Their

---

<sup>9</sup> <http://www.82games.com/comm30.htm>

<sup>10</sup> i.e. holding constant all of the players that shared the floor with that player and with the reference players.



Teams Win”, 82games<sup>11</sup>, regressed the pure plus/minus scores (i.e. the coefficients from the regression described in the previous paragraph) against the box scores of the players. Then, by using the effect of each box score (i.e., the coefficients of the second model) in point margin, Rosenbaum calculated the statistical plus/minus (SPM) for each player of the dataset. Finally, he combined the two plus/minus ratings to obtain an overall rating called the adjusted plus/minus (APM). The complementary weights of each plus/minus score in APM were selected in such a way that they minimize the standard error of the overall rating.

Likewise, Ilardi (2007) presented the first in-season computation of adjusted plus/minus ratings for the NBA 2006-2007 season, included all the playoffs’ games. To help improve the accuracy of the 2006-2007 adjusted plus-minus estimates, Ilardi (2007) added data to the model from the preceding season (2005-2006) and weighted those data much less heavily than last season’s data. This had a net effect of yielding much better (less noisy) player estimates. He aimed to provide a gist-level sense of what the APM model does and a means of understanding the accompanying rank-ordered listings of each player in the league who logged at least 400 minutes during the 2006-2007 season. Also, he ran the model for players that played more than 1640 minutes, while he was trying to minimize the estimations’ noise.

A year later, Ilardi and Barzilai (2008) presented their work, which was aimed at reducing the adjusted plus/minus estimations’ noise (error). They used data from five seasons, including players with at least 300 minutes in the 2007-2008 season and more than 2000 minutes in total in the regular season. They also weighted the ratings of each season with the highest weight set in the last season (2007-2008) to distinguish the individual effects of teammates who frequently appear on the court at the same time. They calculated the players’ total impact on their team and they modelled separately their impact on offence and defence, too. As Ilardi and Barzilai (2008) said, they presented “the most accurate (low-noise) adjusted plus-minus ratings ever to appear in the public domain” till then. In this context, they mentioned that Eli Witus had used a two-step estimation technique to estimate offensive and defensive adjusted plus-minus ratings for the season 2007-2008. However, his estimations had larger standard errors than Ilardi’s and Barzilai’s, who used five seasons’ data, whereas Eli used just one and they incorporated player offensive and

---

<sup>11</sup> <http://www.82games.com/comm30.htm>

defensive effects directly into their model, rather than utilizing a two-step (indirect) procedure to derive them from a net (offence + defense) estimate.

Adjusted plus/minus (APM) ratings have been a subject of study for many researchers and were the starting point for various works. Sill (2010) discussed about the issues of multicollinearity<sup>12</sup>, overfitting, the selected threshold for the playing time of the dataset's players, the number and weights of seasons were taken into account in previous studies. The author, firstly, proposed cross validation (CV) for the evaluation of APM model's out-of-sample prediction, by using the root-mean-squared-error (RMSE) of predicted vs. actual margin efficiency (the difference of points per 100 possessions). Smaller the RMSE better the accuracy of the model. CV was used, as well as, to determine the optimal "reference" player's minutes cutoff and weightings of past years<sup>13</sup> for the standard APM linear regression technique, which were chosen to succeed in the lowest CV RMSE.

Sill (2010) used data from three seasons (2006-2009). The data from March and April 2009 were the test set of the analysis. They used a weighting scheme in which data from  $k$  years ago is weighted with weighting param  $D^k$ , for various values of  $D$ . He also ran the above only for the last season 2008-2009 and he compared the results with those that came by including all three seasons. To solve the overfitting problem and improve the accuracy, Sill (2010) applied ridge regression, which is presented through a Bayesian interpretation. The second term can arise out of a Bayesian prior distribution over the vector of player ratings ( $w$ ), which are to be estimated, by independent Gaussian distributions. The penalty ( $\lambda$ ) corresponded to the ratio of the variance of the inherent, unpredictable noise to the variance of this gaussian prior and it was chosen by CV. The above technique is known, in general, as Regularized Adjusted Plus/Minus (RAPM). RAPM is one of the most dominant plus/minus ratings used in academic literature (see for example in Grassetti *et al.*, 2020, Janeczko *et al.*, 2022) or professional analytics websites and companies (see for example Squared Statistics<sup>14</sup>, The Spax<sup>15</sup>, NBAstuffer<sup>16</sup>, Evolving Hockey<sup>17</sup>, Hockey Graphs<sup>18</sup>)

---

<sup>12</sup> Mention that multicollinearity in terms of players refers to those players that play together most frequently.

<sup>13</sup> Sill (2010) mentioned these as meta-parameters.

<sup>14</sup> <https://squared2020.com/2017/09/18/deep-dive-on-regularized-adjusted-plus-minus-i-introductory-example/>

<sup>15</sup> <https://www.thespax.com/nba/calculating-regularized-adjusted-plus-minus-for-25-years-of-nba-basketball/>

<sup>16</sup> <https://www.nbastuffer.com/analytics101/regularized-adjusted-plus-minus-rapm/>

<sup>17</sup> <https://evolving-hockey.com/glossary/regularized-adjusted-plus-minus/>

<sup>18</sup> <https://hockey-graphs.com/2019/01/14/reviving-regularized-adjusted-plus-minus-for-hockey/>

Some drawbacks of APM rating were discussed by Berri (2011) at “*Wages of Win*” Journal. He talked about the relatively large standard error and, due to this, coefficients couldn’t be differentiated from zero. By adding more data (i.e. more seasons) the standard errors are reduced but there are still statistically insignificant coefficients. Berri (2011) talked about inconsistent measurements across time, too. He claimed that APM model can’t give predictions in case players change teams, because APM of one year is not statistically related to the previous season’s performance. In the same network, Galletti (2011) made an extensive discussion about the APM model based on the insignificant and inconsistent estimations that Berri (2011) mentioned. According to his analysis, APM model calculated two variables with a low correlation to wins ( $R^2 < 5\%$ ) and adds them up to minimize the error and guarantee at least 90%  $R^2$  for the overall model. In this context, he claimed that a simple model using the percentage of minutes played for a team to assign wins to substitute true plus/minus and wins produced or win shares for statistical plus/minus (SPM) would be much more consistent to team wins prior to the error correction.

### **2.2.2 Advanced Plus/Minus**

Many researchers have worked to create a box score-based player evaluation statistic. In this context, Myers (2012) proposed Advanced Statistical Plus/Minus (ASPM), which followed the concept of Statistical Plus/Minus (SPM). It would be remarkable at this point to mention Paine’s (2009) opinion about SPM. In more detail, he noted that SPM couldn’t be a box score-based method, but it is more a complement of the pure plus/minus. On the other hand, it has an effective predictive ability for teams’ performance based on the weighted projected SPM of their players<sup>19</sup>.

Before going back to Myers’ (2012), one of his previous works should be referred. Well, Myers (2011) presented a review of Adjusted Plus/Minus (APM)<sup>20</sup> and its derivatives (the drawbacks of APM were the center of many studies). He said that APM has an issue with stability, when one-season data are used. Perhaps, the biggest problem is collinearity due to coaches’ decisions about players’ rotation and especially when there are players that play at the same position (it always happens at all teams). As Myers (2011) said, collinearity is related to reliability. Specifically, evidence of collinearity greatly decreases APM’s reliability. Another characteristic, that Myers

---

<sup>19</sup> Mayers (2012) refer readers to Paine’s work about SPM.

<sup>20</sup> <http://godismyjudgeok.com/DStats/2011/nba-stats/a-review-of-adjusted-plusminus-and-stabilization/>

(2011) noted, is validity. According to him, APM is completely valid, because it directly measures the desired result. However, it is not as valid for players who change teams (something that has been already referred to previously). Thus, the main problem, which was referred to by Myers (2011), is collinearity (or reliability in the opposite direction). He presented four different ways to solve this problem.

Firstly, Myers (2011) proposed a larger dataset than the ones used in previous publications. To be more specific, he mentioned a four-season (2007-2010) analysis of Engelmann (2011) and he called this approach a “long-term APM”. However, the collinearity problem, which was present also in previous studies, hadn’t been solved. Additionally, there was an issue with players that had played only a part of the total period under study. In the second approach, Myers (2011) followed the approach of Iladi (2008) weighted the ratings of each year assigning a higher weight to the most recent seasons. Although the validity of APM was not improved with this approach, the reliability of the estimations for the season in question was better.

The third approach he noted is statistically stabilized APM. Myers (2011) talked about Rosenbaum’s SPM, which was developed when he tried to create a box score metric and used that to stabilize the regression. Rosenbaum combined that SPM with his un-stabilized APM after each were run and the result was weighted toward either SPM or APM depending on which had the lower standard error. Myers (2011) underlined that he hadn’t seen another approach like Rosenbaum’s till then. The last approach he discussed was the Regularized APM (RAPM), which had been proposed by Sill (2010)<sup>21</sup> for the first time. Among others, Myers (2011) mentioned that the penalty factor was chosen such that the maximum out-of-sample accuracy was attained (or the minimum RMSE as we saw before). In this way, collinearity (or noise) was reduced, but a small amount of validity was lost.

Based on all these, one year later, Myers (2012)<sup>22</sup> introduced the box score statistic ASPM (Advanced and Not Adjusted). He used team-by-team advanced statistics as data from eight seasons 2003-2011 and he was based on an unweighted RAPM. Players with more than 3000 possessions had been included in the final dataset. By using long-term APM and RAPM, he had, as a result, less random error and slightly more bias. Also, ASPM was more accurate than previous

---

<sup>21</sup> His work has been already referred above.

<sup>22</sup> <http://godismyjudgeok.com/DStats/aspm-and-vorp/>

works on SPM and it wasn't restricted to measuring simple points and field goals attempted (FGA). In addition, modelling some non-linear interactions helped to succeed more accurately. ASPM was built as a regression equation like this:

$$ASPM = a*MPG + b*TRB\% + c*BLK\% + d*STL\% + e*USG\%[TS\%^2*(1-TOV\%) - f*TOV\% - g + h*AST\% + i*USG\%]$$

where:

- i. *Minutes Per Game (MPG)*, *Total Rebounds Percentage (TRB%)*, *Blocks Percentage (BLK%)*, and *Steals Percentage (STL%)* are included as simple linear terms.
- ii. The last term is a scoring term (SC) given by this general form:

$$SC = Usage*(points\ per\ possession - threshold\ points\ per\ possession)$$

where *Usage Percentage (USG%)* is an estimate of the percentage of team plays used by a player while he was on the floor (basically shooting possessions plus turnover possessions),  $TS\%^2*(1-TO\%)$  is a points-per-shot term,  $-f*TOV\%$  gives a negative value to turnovers,  $-g$  is the raw threshold value,  $+h*AST\%$  gives a positive value to assists and  $+i*USG\%$  gives a positive value to shooting more—if the player uses a ton of possessions, they don't have to be quite as efficient to benefit the team overall.

In the above measure, if the *points per possession* are above a certain threshold, then SC is considered positive, and the *Usage* is not involved in the equation. Otherwise, the player is considered to be negative for his own team.

Furthermore, Myers (2012) introduced a new rating called *Value Over Replacement Player (VORP)*, which is helpful to estimate the overall value of a team's ASPM and is given by

$$VORP = (ASPM + Replacement\ Player\ Level) * \%min.$$

### 2.2.3 Box score-Based Plus/Minus

One year later, Sisneros and Van More (2013) published their work on expanding plus/minus for visual analysis of NBA box score data. They used data per game, as they believe that this type of data and data per possession help compare teams and players. They measured the players' impact on their teams through differentials to all box score statistics and then presented visualization tools. In more detail, Sisneros and Van More (2013) proposed plus/minus plot (pluMP) for all players'

and teams' statistics as a valuable tool for simple comparisons. The home or away team information was considered.

They also proposed a new metric for evaluating the win-contribution (WC) of a player based on statistics most relevant at the team level. In short, they calculated the sum of a team's per-game box score differential over all games played in the season and the seasonal differential value for a specific box score. This represents the total value for that box score, statistic or rating of the team. To divide this among the players of the team, they simply assigned each player the percentage of value for a state exactly corresponding to his contributing percentage of the team-wide, season-long accumulation of that state. A player's total number of blocks for the season is divided by the team's total number. The win contribution metric is the sum of all a player's percentages of status values.

According to Sisneros and Van More (2013), "*the WC metric is an interesting first step in evaluating players whose presences are believed to be underrepresented in box score data*". Both PluMP and WC are not model-based indices since they are based on simple differential values of various box score statistics.

In 2014, the *Sports Reference* blog<sup>23, 24</sup> introduced a new statistical rating: the Box Plus/Minus (BPM) metric, which was developed by Myers<sup>18</sup> (2014). BPM is an advanced metric which records the total contribution of each player calculated using box score data that is available from season 1973-1974 like Usage Percentage (USG%), True Shooting Percentage (TS%), Steals Percentage (STL%) and other advanced box score statistics<sup>25</sup>, as well as the statistical interactions between usage, rebounds and assists. The study about BPM is similar to those that were developed for Statistical Plus/Minus (SPM) and ASPM, and also long-term data was used. BPM measures the contribution of each player with respect of the excess of points in comparison to the average points per 100 possessions played. Note that we obtain a separate measure for the offensive and defensive components of a player's BPM: OBPM (Offensive Box Plus/Minus) and DBPM (Defensive Box Plus/Minus). Further, BPM is scaled so that zero represents a decent starter or solid 6th player and -2 represents a theoretical "*replacement level*"<sup>26</sup>. Thus, this concept is easily extended to permit

---

<sup>23</sup> The blog has been renamed to *Basketball Reference*.

<sup>24</sup> <https://www.sports-reference.com/blog/2014/10/introducing-box-plusminus-bpm-2/>

<sup>25</sup> <https://www.basketball-reference.com/about/glossary.html>

<sup>26</sup> i.e., a bench player

calculations of one player's value over that theoretical threshold. This rescaled BPM is named VORP<sup>27</sup>, which corresponds to BPM per 100 team possessions.

More recently, Myers (2020) presented at the *Sports Reference* blog<sup>28</sup> a new version of BPM which is calculated by using player box score statistics per 100 possessions, the team adjusted efficiency per 100 possessions and player positions (position and offensive role are estimated from box score data, unless the player has very few minutes) and the regression's results. The new BPM is fully obtained by a regression model using a larger number of observations and large number of features/covariates. Nevertheless, it totally ignores the playtime per game (minutes-per-game).

Moreover, BPM and VORP were discussed in the *Hack a Stat* blog<sup>29</sup> the same year, in 2020. As a general conclusion from the article, BPM is a fascinating advanced statistic, but one must know its limits in order to use it wisely. The VORP, on the other hand, is a useful statistic to get a first indication of which are the best players in the league.

Engelmann and Ilardi (2014) introduced a new advanced metric in ESPN, inspired by Taj Gibson, who was awarded as the “*NBA Sixth Man of the Year*” helping his team win by doing “*all the little things*” that never show up in statistics. It is about the Real Plus/Minus (RPM), which estimates how many points each player adds or subtracts, on average, from his team's net scoring margin for every 100 possessions played. Also, offensive RPM (ORPM) and defensive RPM (DRPM) could be estimated/calculated separately. It follows the development of adjusted plus-minus (APM) by several analysts and regularized adjusted plus-minus (RAPM) by Sill (2010). Via the use of prior distributions with the Bayesian framework, and by considering as model covariates the aging curve of each player and the game score, Engelmann and Ilardi (2014) obtained RPM for which they obtained empirical evidence that it has better out-of-sample predictive accuracy. Hence, RPM is considered as an important improvement of Engelmann's (2011) RAPM. This metric isolates the unique plus-minus impact of each NBA player by adjusting for the effects of each teammate, opposing player and coach, in contrast with the simple plus/minus metric which does not take into account any of these features<sup>30</sup>.

---

<sup>27</sup> the formula is  $[BPM - (-2.0)] * (\% \text{ of minutes played})$

<sup>28</sup> <https://www.basketball-reference.com/about/bpm2.html>

<sup>29</sup> <https://hackastat.eu/en/learn-a-stat-box-plus-minus-and-vorp/>

<sup>30</sup> RPM statistics are provided by Jeremias Engelmann in consultation with Steve Ilardi. RPM is based on Engelmann's xRAPM (Regularized Adjusted Plus-Minus).

As any newcomer, RPM was heavily criticized and originally was not widely accepted (*The Data Jokes* blog, 2021). For instance, Erler (2014) in the *Pounding the Rock* blog<sup>31</sup>, an online San Antonio Spurs' community, expressed his doubts about RPM, which are related to the fact that players who come from the bench have specific limitations and role with limited playing time. He also related RPM with win-shares in baseball. Nevertheless, Erler (2014) did not present scientific evidence about his doubts, but only his opinion which was more based on specific game facts and situations.

## 2.2.4 Adjusted Plus/Minus Discussion and Extensions

The same year, Clemens (2014) presented his review about plus/minus ratings in the “*Fansided*” blog<sup>32</sup>. He discussed APM, RAPM, and SPM, he underlined the advantages and disadvantages of APM and explained his thoughts about the comparison of Hollinger’s PER and APM. In more detail, he strongly criticized PER which is not a model-based metric (i.e. it is simply calculated by an arbitrary but intuitive formula). He underlined that we have no idea about the statistical significance of the PER difference among players or among the PER-value for different numbers of games. On the other side, he supported the use of APM because it is obtained through a scientific model using empirical evidence and data.

Deshpande and Jensen (2016), based on Rosenbaum’s (2004) work, proposed an alternative method for adjusted plus/minus model. Once the shifts (i.e., periods of play between substitutions) are determined and both the point differential and total number of possessions are measured in each shift then, the point differential per 100 possessions is regressed against indicators corresponding to the ten players on the court. Deshpande and Jensen (2016) proposed to regress the change in the home team’s win probability (instead of regressing the point differential) during a shift against indicators corresponding to the five home team players and five away team players, in order to estimate each player’s partial effect on his team’s chances of winning.

For the practical part of the study of Deshpande and Jensen (2016), they estimated the home team win probability as a function of its lead and the time elapsed with data from season 2006-2007 to 2013-2014. To solve the collinearity problem, they used the Bayesian approach with a Laplace

---

<sup>31</sup> <https://www.poundingtherock.com/2014/4/8/5594238/problem-with-real-plus-minus>

<sup>32</sup> <https://fansided.com/2014/09/25/glossary-plus-minus-adjusted-plus-minus/>



prior distribution with mode at zero on each partial effect. The Laplace prior has been chosen in favor of the common choice of a normal prior since it shrinks more abruptly smaller partial effects towards zero than the normal prior. These two approaches (Laplace vs. normal prior) correspond to implementing lasso over ridge regression in the classical statistics framework. Finally, they have used play-by-play data from the 2013-2014 season to update their priors and finally obtain the posterior distribution of the partial effects.

Deshpande and Jensen (2016) compared their impact metric with other performance ratings and reached some interesting conclusions. They underlined that the measure they proposed is not suitable for predictions, but it is more accurate than PER or RPM. They also found a small positive correlation ( $r=0.22$ ) between their metric and PER and a larger one with RPM ( $r=0.65$ ). However, they mentioned that it would be more informative to identify cases where these metrics are in agreement or disagreement.

One year later, Jacobs (2017) presented in his blog<sup>33</sup> the APM, and the RAPM. He discussed the ability of RAPM to control the impact of multicollinearity in an attempt to get stable and comparable numbers between players and implemented it in 2016-2017 season data.

The same analyst, Jacobs (2018) walked through a “*vanilla-flavored*” methodology for building a RAPM model for NBA data. He focused on the data necessary, the required data manipulation process, and the methodology for determining required hyper-parameters. The two posts of Jacobs (2017, 2018) about RAPM were focused on offensive and defensive versions of the original model developed by Sill (2010).

In another post, Jacobs (2018) focused on unweighted stints with offensive or defensive ratings. He aimed to rehash some of the key points in an effort to understand the pitfalls of RAPM. In his article, turned the crank and showed explicitly why we need to be very careful when using this analysis. He tried to solve some problems by pushing in weights, but confidence bounds weren’t improved. Also, the Gaussian assumption did not hold.

The same year, Goldstein (2018), with his article in the *Fansided* blog<sup>34</sup>, introduced the Player Impact Plus-Minus (PIMP) rating. His aim was to create a metric not just to measure a player’s

---

<sup>33</sup> *Squared Statistic: Understanding Basketball Analytics* blog <https://squared2020.com/2017/>

<sup>34</sup> <https://fansided.com/2018/01/11/nylon-calculus-introducing-player-impact-plus-minus/>

performance but also to predict it. PIMP combines luck-adjusted plus-minus data with the value of the box score and a handful of interaction terms to estimate a player's value over the course of a season. As Goldstein (2018) said the three components: box score component, luck-adjusted on-off data and luck-adjusted net rating can provide a descriptive and predictive function for players' performance.

For the box score component, the author used pace adjusted per 36-minute statistics to calculate an initial estimate of offensive and defensive contribution. The offensive and defensive components estimated through a weighted regression against 15-season RAPM with independent variables games started percentage (GS%) in a squared term, points (PTS), total rebounds (TRB), assists (AST), steals (STL), blocks (BLK), turnovers (TOV), personal fouls (PF), free throw attempted (FTA), 2-point attempted (2PA) and 3-point attempted (3PA).

As for the other two components, on-off data are the most basic data available for tracking how a team performs when a player is on-court versus off-court. Luck-adjusted data, designed by Nathan Walker earlier, are used to adjust for factors that are out of an individual team or player's control and ratings allow for a clearer view of team performance with and without a player. In short, the luck-adjusted methodology uses more statistically predictive factors of a team's offensive and defensive rating (points scored or allowed per 100 possessions) to estimate what the team's Offensive Rating (ORTG) or Defensive Rating (DRTG) should be without the variance that some statistics have. The clearest example of this variance comes from opponent 3-points percentage (3P%), which a team has relatively limited control over but can drastically shift how a team looks in raw ORTG and DRTG. Thus, Defensive-PIMP and Offensive-PIMP are calculated by a formula that included the player's luck-adjusted on-off net ratings and the player's luck-adjusted on-court ORTG and DRTG relative to the league average.

PIPM is one of the most accurate publicly available impact metrics in terms of predicting future results (Goldstein, 2018). For seasons prior to 1997-1998, when plus-minus data weren't available league-wide, an estimate of on-off data is calculated using simple interaction terms. The final portion of Player Impact Plus-Minus is converting the per 100 possession estimates of value into cumulative Wins Added based on playing time, by estimating Pythagorean wins<sup>35</sup> (as an exponent

---

<sup>35</sup> Pythagorean wins were mentioned at the previous chapter (1).

is used 13.98). The replacement level for the NBA is set at a PIPM of approximately -2.3, below average but in basketball there are often below-average players on the court who still help their teams more than they negatively affect it.

Grasseti et al. (2020)<sup>36</sup>, working around RAPM, applied a Bayesian approach with regularization by using play-by-play data from season 2018-2019 of Euroleague<sup>37</sup>. Their approach had 2-directions:

- i. The classical response variable, points scored, was replaced by a comprehensive score combining a set of box score statistics.
- ii. This approach was extended to the implementation of the entire lineup rather than individual players.

Another study about APM was developed by Ghimire et al. (2020) and published in the “PLOS ONE” journal. They studied the relationship between one player’s RPM and his teammate’s quality. They used two approaches in which potential endogeneity in the relationship of interest was taken into account: i) linear fixed effect regression to explain the variation if RPM across players-seasons and ii) 2-stage least square (2SLS) for robustness check.

In the first approach, the RPM regressed against players and seasons, which capture the player specific and time effects and, as well as other players’ RPM and age (linear and squared term). In the second approach, to solve endogeneity, the authors applied 2SLS instrumental variable estimation using the lagged values of the minutes weighted average APM of lineup teammates as an instrument. They used data about the previous season’s performance of present season line up teammates as an instrument for the present quality of lineup teammates and they used the estimated values in the regression they ran the dependent variable the RPM of a player and independent variables: the other players’ RPM (i.e. the estimated values from the 1<sup>st</sup> stage regression) and the age (linear and squared term). In this way, they addressed such a reverse causality, as they were based on the possibility that a player may influence his team’s current period performance. From their results, they found evidence that RPM is related to the on-court teammates’ quality.

---

<sup>36</sup> <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0237920#sec008>

<sup>37</sup> That’s something we haven’t seen till now. Data not from NBA!

Snarr (2020) introduced another metric, the Estimated Plus/Minus (EPM), which is a new kind of Adjusted Plus/Minus using Bayesian priors. It is a hybrid model that uses player-tracking data as part of its evaluation. He described EPM as the “most accurate” all-in-one NBA player metric and the lowest RMSE.

Cheema (2021) discussed how RAPM could be improved by using Bayesian priors with box scores and tracking statistics to give a better prior estimate of player value for the model. Instead of regressing to zero, the regression could be run to this prior value. Thus, he presented two versions of RAPM:

- i. non-prior informed RAPM (NPI RAPM) which uses nothing but the lineup metrics. It shows weaknesses in the most of its uses. For small sample sizes, especially for one-season data, there isn't suitable prior and there is also a lot of variances because of factors like the three-point shot.
- ii. prior informed RAPM (PI RAPM) which improves the accuracy of the model in small samples by incorporating box score numbers.

The author noted that a five-year RAPM has often more reasonable results than a one season. His study is based on 25-year data from regular and post-season 1997-2021 (we talk about 5,972,736 possessions). He doubled playoff possessions in order to give larger weight to the post-season games. He had planned to utilize a simple prior based on age. The problem is that all players do not age on the same curve. Inspired by previous works for RAPM, Cheema (2021) decided to apply a Bayesian approach with a prior based on playtime and team strength and use two external statistics in the regression: player's minutes per game and team's net rating. He, finally, presented his final player rankings and concluded that this approach indeed identifies top-quality players.

## **2.3 Mixtures of Plus/Minus Ratings**

The last years some fresh metrics have been proposed especially for predicting the future performance of a team or the final outcomes of the games. Most of them are a combination of one plus/minus rating (of the referred) and a different type of data (e.g., box scores, player tracking statistics).

### 2.3.1 PT-PM

Johnson (2014)<sup>38</sup> presented an alternative version of statistical plus/minus, the Player Tracking Plus/Minus (PT-PM). It is based on some box score metrics, basic and advanced, NBA's SportVU player tracking statistics and RAPM from Jerry Engelmann's Stat's for the NBA (2011). RAPM is split into offense and defense in terms of estimates of a player's impact on the court, so the two scores were used to perform two separate regressions (Weighted Least Squares) weighted by the number of possessions on the court for each player.

He calculated the Offensive PT-PM (OPT-PM) through a formula which includes Points, Field Goals Attempted (FGA), Minutes Per Game (MPG), Free Throw Attempted (FTA), Rebounds multiplied by the Three-point Rate and the new measures: i) Passing Efficiency<sup>39</sup>, ii) Turnovers (TOV) per 100 touches<sup>40</sup> and Contested Rebounds Percentage<sup>41</sup>.

The Defensive PT-PM (DPT-PM) formula includes as independent variables Steals and Personal Fouls (PF) per 100 possessions and the new measures: i) Opponent Field Goal percentage (FG%) Rim (the percentage a players opponent shot at the basket when they were in a position to contest the shot, the higher the opponent's shot percentage the less effective the player is on defense), ii) Opponent Field Goals Attempted (FGA) Rim (the number of times per 100 possessions the player was in a position to contest a shot at the rim).

### 2.3.2 CARMELO

Another plus/minus metric is Career-Arc Regression Model Estimator with Local Optimization (CARMELO), which obviously refers to the known player. In fact, the backronym was developed after creating the metric. CARMELO was developed by Silver et al. (2015)<sup>42</sup> in the *FiveThirtyEight* blog<sup>43</sup>. Silver (2015) mentioned that this metric originated out of his work about New York Knick's Carmelo Antony the previous year (2014). The real inspiration for CARMELO

---

<sup>38</sup> *Counting the Baskets* blog: <https://counting-the-baskets.typepad.com/my-blog/2014/09/introducing-player-tracking-plus-minus.html>

<sup>39</sup> SportVU: points created per pass attempts.

<sup>40</sup> SportVU: the number of turnovers committed per 100 touches of the ball.

<sup>41</sup> SportVU: the percentage of rebounds a player gets, that are up for grabs with the opposing team.

<sup>42</sup> Silver (2015) in the *FiveThirtyEight*'s article referred that Nail Paine and Allison McCann (one of our visual journalists) helped developing CARMELO.

<sup>43</sup> <https://fivethirtyeight.com/features/introducing-raptor-our-new-metric-for-the-modern-nba/>

is Player Empirical Comparison and Optimization Test Algorithm (PECOTA), a system that Silver (2003) built for Baseball Prospectus to forecast the careers of baseball players. According to Silver (2015), CARMELO is considerably simpler than PECOTA and it has “*fewer bells and whistles*”. It predicts each player’s overall value (contribution) on offence and defence, but not his component statistics.

CARMELO is able to answer how good can any player be in the near future. In other words, it is a predictive metric for the players’ career performance. It compares a player to similar players of the past and generates future predictions in terms of wins above replacement (WAR). It is calculated through a 50-50 percent combination of Real Plus/Minus (RPM) and Box Plus/Minus (BPM), where the advanced plus-minus metrics give a measure of the number of points per 100 possessions that a player contributed to his team, relative to an average NBA player.

### 2.3.3 RAPTOR

Silver et al. (2019)<sup>44</sup> introduced a new plus-minus metric, called Robust Algorithm using Player Tracking and On/Off Ratings (RAPTOR), which replaced CARMELO. This metric measures the number of points a player contributes to team offense and team defense per 100 possessions, relative to a league-average player, by using modern NBA data, specifically player tracking and play-by-play data, in contrast to other plus/minus ratings like BPM or RAPM.

RAPTOR does not account for coaching tactic and plans or synergies between teammates. It consists of two major components:

- i) a “*box*” (as in “box score”) component
- ii) an “*on-off*” component

The box component uses individual statistics (player tracking and play-by-play data) and it weighs more highly than the second component (*on-off*) which evaluates a team’s performance when the player and various combinations of his teammates are on or off the court. The “*on-off*” element of RAPTOR evaluates how a player’s team performed while he was on the court, how the player’s court-mates performed while they were on the court without the player and, finally, how those

---

<sup>44</sup> *FiveThirtyEight* blog: <https://fivethirtyeight.com/features/introducing-raptor-our-new-metric-for-the-modern-nba/>

court-mates' other court-mates performed when they were on the court without the player's court-mates, all adjusted for the strength of competition they were facing.

Silver et al. (2019) noted that they found the on-court/off-court (On/Off) statistics that come from the out-of-sample testing of RAPTOR noisy, especially in comparison with the individual measures of the “box” component. They also referred that descriptive RAPTOR does not use priors based on a player's characteristics like height, weight, age, position, or any other factor. However, RAPTOR could be used for predictions. For this, height, age, draft position, and whether a player recently appeared on an All-NBA team (and other data like these) are used extra. Predictions also weigh variables slightly differently than descriptive RAPTOR does, as certain statistics are more subject to luck than others. They called this predictive version of RAPTOR as PREDATOR (PREdictive rApTOR).

### 2.3.4 DARKO

Medvedovsky (2021)<sup>45</sup> introduced a composite predictive metric that uses box score and plus-minus statistics, the Daily Adjusted and Regressed Kalman Optimized (DARKO) or Daily Plus/Minus (DPM). It is similar in concept to baseball projection systems such as PECOTA, i.e., not to only measure player performance, but also to predict it. The main difference between DPM and the other player impact metrics, in accordance with the literature, is that “*DPM solely looks forward by giving the results in a Bayesian model which predicts all elements of the box score*”.

According to the creator, DARKO is built to provide player box score and impact predictions which are updated daily. NBA box scores, play-by-play data, tracking data, and other game-level data from season 1998-1999 are used. DARKO, by considering all games of players' career and recency (the more recent games are weighted more heavily), explains how much each stat day-to-day performance is due to random noise or talent and it estimates the change of players' talent across the time. Based on root mean square error (RMSE), DPM is better than any all-in-one metric in terms of predictive analytics, while EPM and LEBRON (it is presented below) follow it in the ranking.

---

<sup>45</sup> The Athletics blog, <https://theathletic.com/2613015/2021/05/26/introducing-darko-an-nba-playoffs-game-projection-and-betting-gui>

Moreover, DARKO can be used for not only players' performance. Regarding game-level predictions, DARKO is generating player impact predictions for every game, based on a combination of their historic box score statistics, their on/off data, and some team-level statistics. An estimation for a player's total plus/minus rating per game is calculated after pairing the previous DPM with theoretical (researcher's projection) minutes played at a specific game. Adding all that up, the in-home court gives a predictive point spread. Furthermore, Medvedovsky (2021) worked on the estimated change of players' impact during the playoffs via DARKO. This is based on each player's own career history in the playoffs and on how much their teams have over or underperformed in the playoffs with them in the court.

### **2.3.5 LEBRON**

One more fresh metric created by the *BBall-Index.com* team, Krishna Narsu and Tim Cranjic (2021). It is the Luck-adjusted player Estimate using a Box prior Regularized On-Off (LEBRON). It's a kind of prior informed RAPM. This rating uses box score results and on-off calculations (specifically, luck-adjusted RAPM) for an impact score measured per 100 possessions.

The box score component of LEBRON uses weights from PIPM (Player Impact Plus-Minus). Players with less playing time are problem for the plus/minus indices and analysis. Thus, in order to deal with those players and identify if a high performance in a small sample is noise or not, they stabilized data by combining the padding technique of Medvedovsky (2020)<sup>46</sup> and Narsu, K., and Cranjic, T., (2016)<sup>47</sup> offensive archetypes about players' offensive role. Through this approach, the volume at which each box score statistic stabilizes is determined. Additionally, that box score prior is used for the RAPM calculations to derive LEBRON values. The on-off data are referred to the difference between points or percentages of box score statistics when a player is on the court or not.

Consequently, LEBRON is a measure of impact, not talent, and it is expected that players get better (more rapidly when younger) over time and then drop over time (more rapidly the older they get) later in their careers.

---

<sup>46</sup> <https://kmedved.com/2020/08/06/nba-stabilization-rates-and-the-padding-approach/>

<sup>47</sup> <https://www.bball-index.com/offensive-archetypes/>



In addition, Narsu, K., and Cranjis, T., (2021) added calculations to estimate the total wins a player has added for a season (Wins Added). They used the weights for the minutes and total impact of Goldstein's (2018) formula instead of PIMP. At the end of their work, estimated a per 100 possession impact that is role-adjusted, stabilized, and utilizes luck-adjusted values along with RAPM calculations: zero value is for average players while the value of -2.7 corresponds to players that usually play limited time (substitutes) and more often for G-league replacement added to the roster.

Table 2.1: The discussed basketball players performance ratings of the literature.

METRIC	Full Name	Formula/Description	Reference
APM	Adjusted Plus/Minus	<i>Pts per 100 possessions ~ Offenders+ Defenders</i> (Regression model)	Winston and Sagarin (2002), Rosenbaum (2004)
RAPM	Regularized Adjusted Plus/Minus	<i>Pts per 100 possessions ~ Offenders+ Defenders</i> (Ridge Regression model)	Sill (2010), Ilardi (2007), Iradi and Barzilai (2008)
SPM	Statistical Plus/Minus	$APM \sim Pts + EFGA + EFGA^2 + 3PA + FTA + AST + OREB + DREB + TOV + STL + BLK + PF + (PTS \times AST \times REB)^{1/3} + MP$ Game statistics per 40 minutes (Regression model)	Rosenbaum (2004)
ASPM	Advanced Statistical Plus/Minus	$ASPM = a * MPG + b * TRB\% + c * BLK\% + d * STL\% + e * USG\% * [TS\% * 2 * (1 - TOV\%) - f * TOV\% - g + h * AST\% + i * USG\%]$	Myers (2012)
VORP	Value Over Replacement Player	$VORP = (ASPM + Replacement Player Level) * \%min$	Myers (2012)
pluMP	Plus/Minus Plot	Differentials to all box score statistics & visualization tools	Sisneros and Van More (2013)
WC	Win Contribution	The sum of a team's per-game box score differential for all games vs. the seasonal differential value for a specific box score	Sisneros and Van More (2013)
BPM	Box Plus/Minus	The excess of points vs. the average points per 100 possessions played.	Myers (2014)
RPM	Real Plus/Minus	<i>Pts per 100 possessions ~ Offenders</i> (Regression model)	Engelmann and Ilardi (2014)
PIMP	Player Impact Plus-Minus	PIMP~luck-adjusted plus-minus + box score + interactions	Narsu, K., and Cranjis, T., (2021)
EPM	Estimated Plus/Minus	A kind of Adjusted Plus/Minus. Player-tracking data are used.	Snarr (2020)
PT-PM	Player Tracking Plus/Minus	<i>Offensive PT-PM ~ Pts + FGA + Passing Efficiency + TOV per 100 touches + Contested Rebs % + Mi per game + FTA + Rebs * 3Prate, Defensive PT-PM ~ OppFG%Rim + Steal100 + OppFGARim + PF100</i>	Johnson (2014)
CARMELO	Career-Arc Regression Model Estimator with Local Optimization	$CARMELO = 0.50 * RPM + 0.50 * BPM$	FiveThirtyEight blog, Silver et al. (2015)
RAPTOR	Robust Algorithm using Player Tracking and On/Off Ratings	~ box (player tracking & play-by-play data) + on-off data	FiveThirtyEight blog, Silver et al. (2019)
DARKO/DPM	Daily Adjusted and Regressed Kalman Optimized/Daily Plus/Minus	~ box scores + play-by-play + tracking + other game-level data	Medvedovsky (2021)
LEBRON	Luck-adjusted player Estimate using a Box prior Regularized On-Off	~ box score (weights of PIMP) + on-off (luck-adjusted RAPM), measured per 100 possessions.	Krishna Narsu and Tim Cranjis (2021)

## 2.4 What about women's basketball?

Up to this point, we have seen various models and ratings for the players' impact on their teams, which are based on plus/minus. All of them are referred to men's basketball data and some of them have been inspired by NBA players. We have not found any relevant in women's basketball. So, it is important to mention here one recent study which deals with the development of a rating system for WNBA.

The *Athlytics* blog<sup>48</sup> (2022) presented a WNBA model to evaluate players' performance, which is built by two components that are combined to a regularized linear regression. The first component is a box score player rating that serves as a prior for the final ratings. To create the first component, the offensive and defensive box plus/minus (OBPM and DBPM respectively) are estimated through a linear regression like those applied for NBA (because there wasn't something to follow from the literature). For both OBPM and DBPM, a variety of explanatory variables were used: (a) a binary indicator about the player position (power forward/center (PF/C) versus the rest), the playing time in minutes per game, and (c) some offensive and defensive box score statistics. The second part of the total performance rating is an on/off component based on adjusted plus/minus regression. It is important to note that ridge regression was applied in order to improve the accuracy of the coefficients.

To estimate the overall player rankings, the box score prior and the on/off regression components are combined via ridge regression. Instead of shrinking every coefficient to zero, as the ridge regression does by default, it shrinks every coefficient to the box score-based rating. In this way, any collinearity that might exist between players, is alleviated. Also, it is worth mentioning that the penalty ( $\lambda$ ) is chosen by cross-validation (CV).

Possession data for 2021-2022 season in WNBA were used. In fact, there is not any innovation in the applied methodology. However, more studies in women basketball are needed in order to provide better measures for women performance and possibly identify differences in the game between the two genders.

---

<sup>48</sup> <https://412sportsanalytics.wordpress.com/2022/08/18/a-detailed-guide-for-developing-player-ratings-for-wnba-and-other-leagues/>

## 2.5 Other Sports

From the previous section, we can assume that plus/minus ratings in basketball are quite popular. This may be the result of the “nature” of the game, which has a lot of sport relate events at every possession. Of course, plus/minus ratings have been used to other sports and it would be useful to refer to some indicative relative studies.

A starting point to our quest is hockey, in which the Plus/Minus is also very popular. Although plus/minus ratings were used, for hockey before the implementation in basketball, new works and studies by hockey analysts were introduced after the work of Rosenbaum (2004) in basketball. Macdonald (2011), inspired by basketball’s APM models and developed two weighted least squares regression models to estimate an NHL (National Hockey League) player's effect on his team's success in scoring and preventing goals, independent of that player's teammates and opponents. He modeled power play and shorthanded situations to estimate a player's offensive and defensive contributions. Also, for those shifts that begin with a faceoff, he had accounted for the zone on the ice in which a shift begins. Next year, Macdonald (2012), made another step by using Regularized APM. Specifically, as he said, ridge regression helped solve problems of large errors of the estimations and collinearity, which may exist due to the frequency some players play together (which is a problem that appears in all team sports), exactly as RAPM used in basketball analytics’ studies. He also presented a review of adjusted plus/minus models, with references to the classic basketball literature such as Rosenbaum (2004) and Ilardi (2007), accompanied by results on the *Arctic Ice Hockey* blog<sup>49</sup>.

The Plus/Minus approach was also used in football (or soccer). It identifies a player’s implied effect on his team’s goal difference while he is on the field of play. Hamilton (2010) wrote about some starting documents on the plus/minus problem in the *Soccer Metrics Research* blog<sup>50</sup>. He also talked about APM and RAPM four years later (2014) in the same blog. Matano et al. (2018), mentioned that APM hasn't had the same impact in football, like basketball and hockey, since soccer games are low scoring with a low number of substitutions. In soccer, perhaps the most comprehensive player value statistics come from video games and in particular FIFA ratings. Thus, they developed the Augmented APM metric, by combining FIFA ratings and APM, with better

---

<sup>49</sup> <https://www.arcticicehockey.com/2011/9/22/2441898/nhl-adjusted-plus-minus-part-01>

<sup>50</sup> <https://www.soccermetrics.net/>

predictive ability than standard APM or a simple regression model using only FIFA ratings. This metric also considers the problem of collinearity among players. The idea is to recast APM into a Bayesian framework and incorporate FIFA ratings into the prior distribution. The same year, Schultze and Wellbrock (2018), presented a simpler work about a weighted plus/minus metric for individual soccer player performance of three Bundesliga teams. Kharrat et al. (2019), compared European football leagues and their players, by applying the usual goal-differential plus/minus and proposing two new variations. The first is referred to the evaluation of an expected goals plus/minus rating. In the second, in-play probabilities of match outcome are used for evaluating expected goals plus/minus rating. To note that win probability, used earlier for plus/minus models', works in basketball analytics. Moreover, Hvattum (2020) proposed the separate calculation of Offensive and Defensive Plus–Minus Player Ratings for soccer.

Talking about Hvattum, he is quite popular at the field of academic sports analytics. Some of his latest works refers to the: i) relationship between plus/minus ratings and event-level performance statistics on football, Gelade and Hvattum (2020), ii) Comparison of the bottom-up and top-down ratings for individual football players, Hvattum and Gelade (2021). Also, in his YouTube channel “Football Player Ratings”, videos describing and analyzing mathematical models for evaluating individual football players are presented.

For other sports, there are fewer studies regarding to plus/minus. Hass and Craig (2018) implemented a plus/minus approach to obtain volleyball players' evaluation metrics. They introduced a methodology to recover court performance information from standard play-by-play data for one NCAA team. The recovery is in the form of a posterior distribution of player presence, which can then be used to not only calculate the plus/minus metric but also quantify the uncertainty of the metric due to the incomplete information.

Sabin (2021) studied plus/minus models for American football. He mentioned that models like APM, RAPM and BPM from basketball analytics and Augmented APM from football analytics are useful in coming up with a results-oriented estimation of each player's value. Since in American football many positions (such as offensive lineman) have no recorded statistics, which hinders the ability to estimate a player's performance, Sabin (2021) provided a fully hierarchical Bayesian plus/minus (HBPM) model framework that extends RAPM to include position-specific

penalization. According to cross-validated results the out-of-sample predictive ability of BPM is better than RAPM or APM models, which do not fit American football.

Last but not least, Hvattum (2019) presented a comprehensive review of the: i) plus/minus ratings for players' performance in basketball, ice hockey, football, and other team sports, ii) mathematical models, that have been developed, their difficulties and problems which occurred, as well as the proposed solutions. An important conclusion is that *“the literature on plus/minus ratings is quite fragmented, but that awareness of past contributions to the field, should allow researchers to focus on some of the many open research questions related to the evaluation of individual players in team sports.”*

## 2.6 Conclusion

There is a great number of studies for plus/minus models and ratings in team sports and especially in basketball and hockey, which are characterized by faster pace and more events. We have reviewed many works based on basketball in which different metrics have been developed. We may keep in mind the fact that basketball statistics leading the plus/minus models and ratings. Since a lot of people are spending time in the field of basketball analytics (with scientific view or not), it may be important from now on to try using the findings as practical as possible in the game.

The plus/minus ratings were introduced in hockey through a very simple form (goal-difference when a player is on or off the court) but their usage is quite popular in basketball. After the work of Rosenbaum (2004) the plus/minus ratings became a popular topic of study.

These ratings are performance metrics more often for players and less for teams. Some of the plus/minus ratings are model-based while some others are intuitive-based. Also, some of them have been introduced in very popular basketball analytics blogs while others are popular in the academic community.

All works have a common aim: to reduce the plus/minus flaws in an effort to get more accurate results regarding a player's impact. In addition, the effect of teammates and opponents should be considered. Some of the most popular plus/minus ratings for players' performance in basketball is the Adjusted Plus/Minus, by Winston and Sagarin (2002), Statistical and Adjusted Plus Minus by Rosenbaum (2004), Regularized Adjusted Plus/Minus by Sill (2010), Real Plus/Minus by ESPN,

Ilardi and Engelmann (2014), and Box Plus/Minus by Myers (2014). These ratings reflect players' contributions to their teams on offense or defense.

Despite the studies on the plus/minus ratings via statistical models, players with less playing time create noisy estimations for the players' impact, as they usually seem to be better than they really are. Although, researchers try to deal with this problem in different ways. For example, Ilardi and Barzilai (2008) developed separate models for players with only a few minutes played and those that played a lot.

Except for basketball, plus/minus ratings are used in other team sports as well. For example, Macdonald (2011, 2012) publish a very useful work on RAPM for ice hockey, while more researchers studied plus/minus ratings for football such as Hamilton (2010), Schultze and Wellbrock (2018) or Hvattum (2020), who also discuss in a simpler way footballer's performance models in online videos.

As for our work, having reviewed a variety of literature studies, regularized models will be developed for evaluating the NBA players' contributions to their teams based on the teams' points scored/allowed. Regularized Adjusted Plus/Minus ratings will be the starting point. The teammates and opponents of each player will be considered in the models. Finally, a different regularization method and a type of model, not only the simple version of the linear regression, will be applied.

## Chapter 3: A case study about Plus/Minus Ratings

### 3.1 Introduction

One of the basic goals of this thesis is to construct a plus/minus rating via a statistical model, in order to estimate the contribution of each player to his team. According to the relevant literature, we will obtain a Regularized Adjusted Plus/Minus (RAPM), separately for offensive and defensive contributions.

In many RAPM models, the home-team factor has been included. In our RAPM model, the home-team factor and the type of season (regular season or playoffs) will be considered. Also, according to previous studies, for the calculation of RAPM, a sparse matrix is necessary as we will see below.

Section 3.2 describes the data we have used in this thesis in order to construct model based RAPM measures. Section 3.3 presents the explanatory analysis of the available data (NBA season 2021-22). Sections 3.4 & 3.5 present the fitted models in the whole dataset and in a subset of data (after excluding low-time players). Finally, Section 3.6 presents the ranking using the model based RAPM we obtained using the model-based approach of Section 3.5.

### 3.2 Case Study Dataset

#### 3.2.1 Initial data

In this thesis we have used data from NBA season 2021-2022. The data involve possessions (study unit) and their characteristics (points and playing lineups). The dataset involved 322,852 observations (possessions) and 14 features. Specifically, the available variables are:

- **home\_off**: A binary variable that shows if the home team is on offense (1) or the away team is on offense (0).
- **pts**: Points per possession with values 0-6. The case of 5 or 6 points at one possession is extremely unusual. So, it makes sense that there are only 2 observations with 5 points and only 1 with 6 points. This variable will be considered as integer or categorical.

- **season\_type**: A categorical variable that refers to the type of season the game took place (regular season or playoffs).
- **time**: A numerical variable that shows the dataset's duration of each possession (observation). This variable, probably, will not be used for the plus/minus model as the RAPM will be calculated per possession based on the literature.
- **O1, O2, O3, O4, O5**: These five variables referred to the offensive players per possession (observation). The form of the character is team's name & player's name, e.g. LAC22 Ivica-Zubac.
- **D1, D2, D3, D4, D5**: These five variables referred to the defensive players per possession (observation). The form of the character is like the offensive players calculate RAPM it is necessary to transform the form of the above data.

### 3.2.2 Data processing & final data

The first step before the analysis is to transform our dataset into a format suitable to implement the required models. The basic part of the data comprises from the information of the player into the game, and it is an  $n \times m$  sparse matrix where  $n=322,852$  is the number rows/possessions in the training set and  $m=1434$  is equal to two times the number of the players (here 717) in the training set –each player is split into their offensive and defensive coefficient.

For each player we have two dummy variables: one indicating its presence in the lineup of the offensive team and one indicating its presence to the lineup of the defensive team. In the first, the indicator takes the value of one when he is present, while in the second the value of minus one.

Points per possession is a vector of length  $n$ , where each value is the points per possession in that possession (0, 1, 2, 3). Each possession's duration (time) is not taken into consideration here. Alternatively, at the final interpretation we consider points per 100 possession instead of points per possession.

The final training set is a sparse matrix  $322,852 \times 1434$  plus 3 extra columns: home\_off (binary), season\_type (factor with 2 levels) and pts (points per possession).



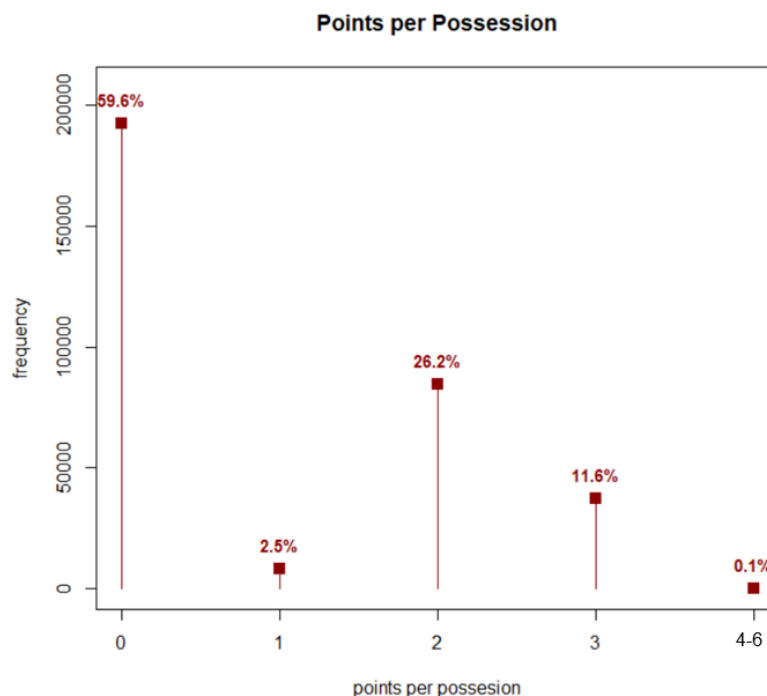
### 3.3 Explanatory Data Analysis

In this section, we present the descriptive and explanatory analysis of the NBA season 2021-22 possession dataset. We focus on the team possessions ratios of offensive and defensive contribution of each player, and we present indicatively the performance of specific cases of players.

#### 3.3.1 Descriptive Analysis

From Figure 3.1 and Table 3.1, it is obvious that the majority of possessions (59.6%) end up with no points. From the rest, 26.2%, ended with two points scored while 11.6% with three points. Finally, only 0.1% of possession resulted in more than three points while one-point possessions are limited to 2.5% of the total possessions since possessions with two successful free throws do account for two-point possessions. Although, more than three-point possession sound unnatural for basketball (since the points are 1-2-3), such possessions are referred to goal-foul cases from three-point shoot or extra free throws earned by the attacking team from technical fouls. Note that there are only two possessions with five points and only one with six points.

Figure 3.1: Frequencies of points per possession.



\*4-6 point category include the 3 observations with 5 and 6 points per possession.

Table 3.1: Frequency table of points per possessions.

Points	Frequency	Cumulative Frequency	Relative Frequency	Cumulative Relative Frequency
0	192492	192492	59.6%	59.6%
1	8145	200637	2.5%	62.1%
2	84461	285098	26.2%	88.3%
3	37521	322619	11.6%	99.9%
4-6	233	322852	0.1%	100.0%

\*2 observations for 5 points per possession and 1 for 6 points per possession.

Concerning the two opponent teams, the number possessions are almost equally split between the home and away teams (278 possessions more for the away than the home team over 322852 possessions) which was expected due to the sequential nature of the game. Regarding the type of season, the majority of possessions, 93.2%, belong to the regular season. This is due to the NBA league format since each team plays 82 games in the regular season and 4-28 in playoffs.

Concerning the player indicator/dummy variables, if we try to obtain the percentage of the participation of the players in the offensive and defensive possessions then the corresponding percentage is quite low since every player can participate only in their team. Generally, it is more sensible to look at the percentage of participation of the possessions of their team. For example, Jason Tatum, who is a starter and crucial player for Boston Celtics (BOS), played only about 3% of the total possessions of the season. However, his offensive and defensive appearances on the court for his team are about 72% of Boston possessions in 2021-2022 (Figure 3.2).

Figure 3.2: Boston Celtics possessions in which Jayson Tatum played offense.

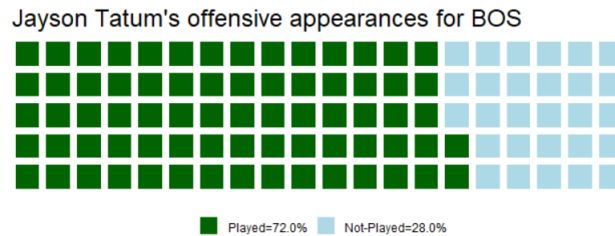
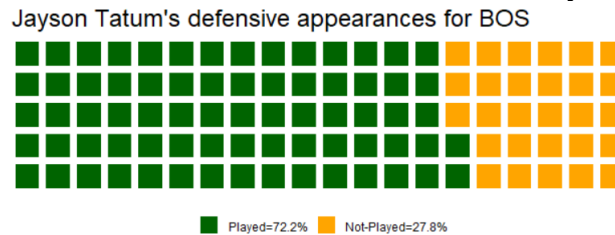


Figure 3.3: Boston Celtics defensive situations in which Jayson Tatum played.



The separate presentation of all players' descriptive plots about their appearances in the offensive and defensive lineups may not be useful. Indicatively, we present the top10 players with respect to their overall playing time (expressed as % of participation in possessions and the player with the highest participation per team).

Therefore, we focus on the team possessions ratios of offensive and defensive contribution of each player as part of pairwise analysis.

### 3.3.2 Pairwise Analysis

For the needs of this thesis, the interest lies in the association between points per possession and the players presence in offensive or defensive possessions.

A way to study the team efficiency in terms of points between i) different players in offensive and defensive lineups, ii) game type (regular season or playoffs) and iii) home and away teams.

Hence, we apply the non-parametric Kruskal test in order to test the hypothesis of equality of the median point efficiency (in terms of points per possession) at every level of categorical variables. According to our findings the median of points per possession differs between the above home and away teams. On the other hand, no difference was found between the median of points per possession in the regular season and the post-season (playoffs). Both results were expected according to the nature of the game.

As for the player appearances, by testing the null hypothesis of the Kruskal test with a significance level of 10% at most, only for 162 defenders (out of 717) and 222 offenders (out of 717) the median value of points per possession differs in case of a player appears on the defensive or offensive lineup respectively and in the opposite case (Figures 3.21 & 3.22 in Appendix I).

It may be more insightful to focus on the teams' possessions ratio (TPR) of points scored and conceded while a player is on the court. If a player has a "positive" contribution to his team's offense, it is expected the probability of scoring more than one point per possession to be higher when he plays compared to his team performance when he does not play. In this case, the Offensive TPR would be higher than one. In this context, Defensive TPR is higher than one when the probability of **not scoring** is larger if a defender is on the court instead of not playing.

$$OTPR = \frac{\#scored\ possessions\ with\ player}{\#scored\ possessions\ without\ player}$$

$$DTPR = \frac{\#not-scored\ possessions\ with\ player}{\#not-scored\ possessions\ without\ player}$$

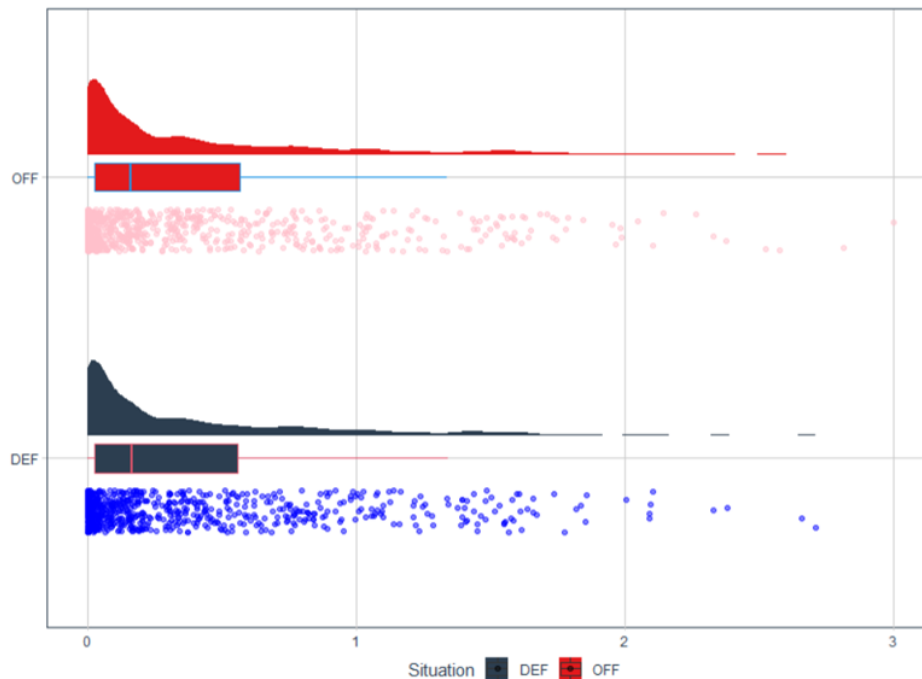
When TPR is calculated over all season possessions, the corresponding values are expected to be smaller ( $< 0.05$ ). This is reasonable since players participate in a very small percentage ( $< 5\%$ ) of the total possessions of all NBA teams in the season. On the other hand, it is more appropriate to calculate TPRs over the possessions of the player's teams. Under this approach we identified 93 offenders and 86 defenders with  $TPR > 1$  which is an evidence of “positive” or “negative” impact of these players.

Table 3.2: Summary table for Teams Possessions Ratio (TPR).

	Total Possessions		Team Possessions	
	Off TPR	Def TPR	Off TPR	Def TPR
Min.	0.000	0.000	0.000	0.000
1st Qu.	0.001	0.001	0.027	0.029
Median	0.005	0.005	0.161	0.167
Mean	0.007	0.007	0.389	0.381
3rd Qu.	0.012	0.012	0.567	0.561
Max.	0.030	0.031	2.999	2.712

Figure 3.4: Players' TPR based on teams' possessions.

Players' TPR based on teams possessions



For example, Jayson Tatum (Boston Celtics) according to TPR, is an influential player since both OTPR and DTPR are larger than 2.5 (2.82 and 2.66 respectively). This means that the probability of the Boston Celtics scoring is almost tripled if Tatum plays than when he does not play. Similarly, the probability of the Boston Celtics opponents fail to score when Tatum is in-play is 2.6 times as higher when Tatum is not playing.

The offensive and defensive contribution of Tatum to his team is presented at Figure 3.5. Figure 3.6 provides the comparison between all Boston Celtics players with respect to TPR. As we can see, Tatum is the better player of Boston with the highest offensive and defensive TPR. Also, there is no difference between the offensive and defensive ranking.

Figure 3.5: Jason Tatum's offensive (right) and defensive (left) contribution to Boston Celtics depending on the total possessions of the season and his team's possessions.

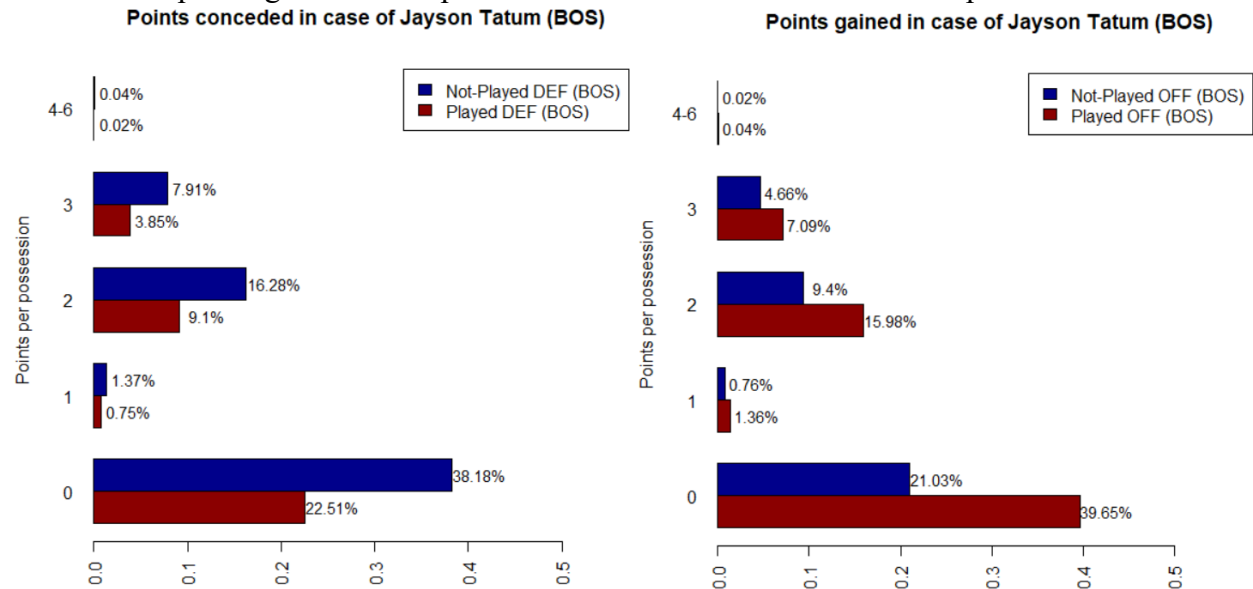
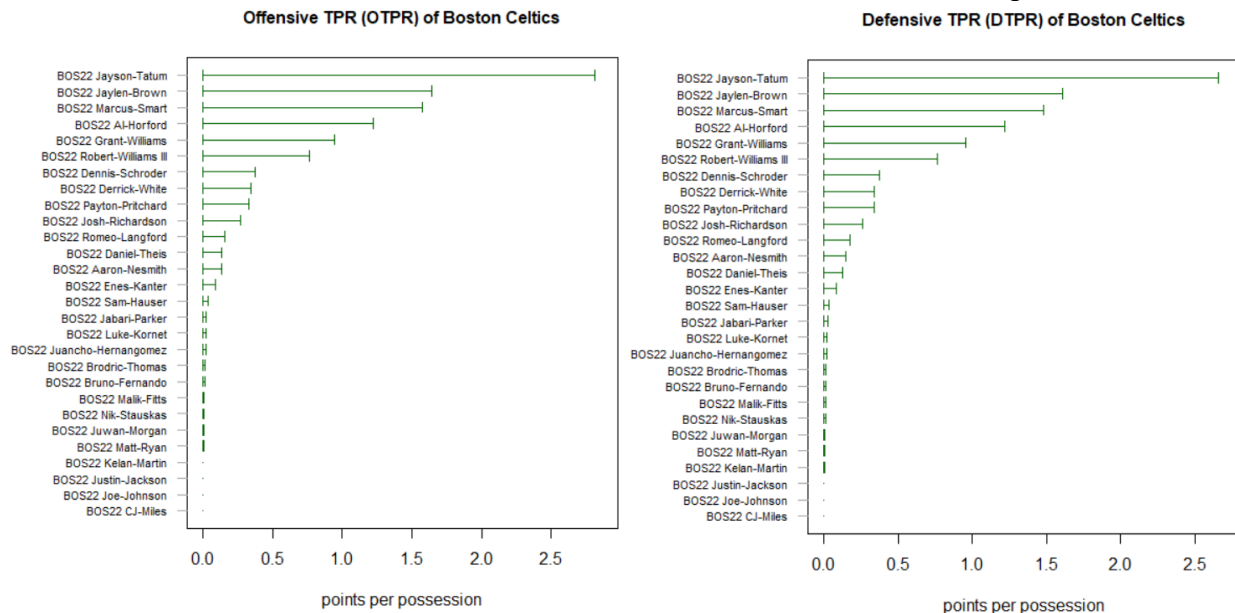


Figure 3.6: Juancho Hernangomez’s offensive (right) and defensive (left) contribution to Boston Celtics depending on the total possessions of the season and his team’s possessions.



To conclude, it is difficult to study in detail all players (as they are 717!). Although, having an idea about the dataset, the next step is calculating the plus/minus index via a statistical model.

### 3.4 Calculation of RAPM: Full dataset Analysis

This section deals with the construction plus/minus ratings via a normal linear regression model. The basis of our study is the estimation of Regularized Adjusted Plus/Minus (RAPM) according to the literature.

According to Jacobs (2017)<sup>51</sup>, “ridge regression is a Bayesian filter with a particular goal in mind: if a slight perturbation to the player interaction matrix is applied, the matrix is invertible”. While a slight bias in the results is introduced, low variance estimates for each player are obtained. Therefore, it is not necessary to eliminate some out of the evaluation model. Moreover, the shrinkage parameter controls for possible multicollinearity between the players in the court-lineup<sup>52</sup>. Finally, the interpretation of linear models is straight forward and easy to communicate to sports community.

<sup>51</sup> [Deep Dive on Regularized Adjusted Plus-Minus I: Introductory Example | Squared Statistics: Understanding Basketball Analytics \(squared2020.com\)](#)

<sup>52</sup> Multicollinearity in terms of players refers to those players that play together most frequently.

Hence, the linear model implemented here is given by:

$$EFFICIENCY_i = b_0 + \sum_{j=1}^K b_j O_{ij} + \sum_{j=1}^K b_{K+j} D_{ij} + \varepsilon_i, \varepsilon_i \sim N(0, \sigma^2)$$

where  $EFFICIENCY_i$  refers to points per  $i$  possession (home/away team),  $\forall i=1, \dots, n$ ,  $O_{ij}$ ,  $\forall j=1, \dots, K$ , takes the value of one when player  $j$  is on offense in  $i$  possession,  $\forall i=1, \dots, n$ , and zero otherwise,  $D_{ij}$ ,  $\forall j=1, \dots, K$ , takes the value of minus one when player  $j$  is on defense in  $i$  possession,  $\forall i=1, \dots, n$ , and zero otherwise and  $\varepsilon$  is the error term. In many literature works, points per possession are transformed into points per 100 possessions by multiplying the estimated coefficients with 100. Hence, the interpretation of estimations in this chapter refers to 100 possessions.

The estimated coefficients of the offensive and defensive contribution of player are the so called plus/minus ratings of each player. More specifically, the constant term  $b_0$  here is an antisense parameter with no direct meaningful interpretation,  $b_j$  measures the offensive contribution of player  $j$ , and  $b_{K+j}$  measures defensive contribution of player  $j$ , for  $j=1, \dots, K$ . Differences of  $b_j - b_m$ ,  $\forall j, m=1, \dots, K$ ,  $j \neq k$ , show the average difference in the expected number of points between the players  $j$  and  $k$ , when their teammates and opponents remain the same.

Most of the related literature proposes Ridge regression for obtain RAPM ratings. Ridge regression shrinks coefficients close to zero but not exactly to zero. This is the main reason why this regularization method is preferable on some occasions instead of Lasso, which shrinks the coefficients of some players exactly to zero in a massive group.

### 3.4.1 Ridge regression with all players under consideration

By applying ridge regression starting from the full linear model, many coefficients will be quite close to zero but not zero. The tuning parameter which corresponds to minimum error's (CV MSE) lambda is found to be  $\lambda_{\min}=241.07$  while  $\lambda_{1se}=1068.07$ , which corresponds to the one standard error (se) away from the minimum value.

The shrinkage percentage of the coefficients against the linear regression estimations is quite high (>99%) for both selected lambda ( $\lambda$ ). Considering: i) the theory, the model with the  $\lambda_{1se}$  is more parsimonious and ii) the literature, most RAPM estimations, particularly on ESPN, were produced

with a lambda value of 2000, which is confirmed also by Sill (2010)<sup>53</sup>. In addition, the results of the two procedures do not differ a lot.

More specifically, the ranking of the players differs in some cases, but the total results are almost the same. For example, the top 50 of each ranking is composed by the same players, that did not play a lot (see Figure 3.26 in Appendix I). Also, the linear correlation between the model based ridge RAPM with the minimum error and the one-standard error tuning parameter is 98.6%.

The top 50 of the three ridge RAPM models are consisted of the same players (see Table 3.12 in Appendix I), while the minutes played of those are presented in Figure 3.30 (Appendix I). The results of the ridge regression models with tuning parameter  $\lambda=2000$  and the one-standard error  $\lambda_{1se}=1068.07$  are quite close (linear correlation  $r = 99.9\%$ ). The model with the minimum error's tuning parameter raises in the ranking about the half of the players in the ranking compared to the other two models. Undoubtably, the lower-ranking position of those players, that compose the best 5-lineups of the NBA, is the common observation of the three ridge linear models, since the agreement of the RAPM high rated players and the three best All-NBA teams players is zero.

As there is no clear and strong evidence to use the one-standard error approach or the literature's proposal, the minimum error's lambda is used for the presenting results. Some first observations of the results:

- Ahmad-Caver (IND) has the best offensive RAPM, equal to 35.49, while he has the 14<sup>th</sup> better in the row defensive RAPM (5.97). Those values of offensive and defensive performance raise him in the first place of the total RAPM (41.46).
- Sekou Doumbouya (LAL) has the best defensive RAPM, equal to 15.82, while he has the 384<sup>th</sup> better in the row offensive RAPM (-0.37). Those values of offensive and defensive performance raise him in fourth place of the total RAPM (15.46).
- Sam Dekker (TOR) has the worst offensive RAPM, equal to -26.20, while he has the 712<sup>th</sup> better in the row defensive RAPM (-25.41). Those values of offensive and defensive performance rank him in the 717<sup>th</sup> place of the total RAPM (-51.62).

---

<sup>53</sup> lambda  $\approx$  2222 for one season.



- Carlik-Jones (DEN) has the worst defensive RAPM, equal to -30.33, while he has the third better in the row offensive RAPM (17.78). Those values of offensive and defensive performance rank him in the 696<sup>th</sup> place of the total RAPM (-12.55).

For the model interpretation two lineups, one offensive and one defensive, of the 2021-2022 NBA finalist teams are used. Considering the offensive lineup of Golden State Warriors (GSW) with Curry, Thompson, Green, Wiggins and Looney. and the defensive lineups of Boston Celtics (BOS) with Tatum, Smart, Horford, Brown and Williams III, the average expected points scored by GSW per 100 possessions are 87.19 while the average actual points are 91.67. If Porter Jr. (1396 minutes played in the season) substitute Looney (1732 minutes played in the season) and the rest of the player remain the same, then the average expected points are 77.27.

For many teams, the best RAPM performers are players that did not have a basic role in their team. That is probably due to the result that most of the teams use second unit players in time periods where the score differential is large or playing against weaker teams. For example, for a player that plays only two minutes at every game and his team scores about four points on average, his final RAPM will be positive. For detailed presentation of the best RAPM based players per team, see Table 3.5 in Appendix I.

The best NBA players per position are presented in tables 3.6 and 3.7 in Appendix I. As you may observe, the top players are not the one expected with the exception of Kevin Durant. From the top five players, only the following can be considered as first lineup players (but not top):

- the shooting guards DeJon Jarreau and Malcolm Hill
- the small forwards Rodney Hood and Anthony Lamb
- the point guard Isaiah Thomas
- the power forward Malik Flits
- the center Deandre Ayton

### **3.4.2 Lasso regression with all players under consideration**

The second approach for the construction of RAPM can be based on Lasso regularization method, which is not as popular in basketball as its usage in Statistics and Machine Learning. Lasso

regularization shrinks regression coefficients towards zero and additionally sets an automatic cutoff for the good or worse players.

The idea here is to apply cross-validated Lasso in order to identify which players have a non-zero coefficient in at least 20% of the times. In more detail, after running 20 times the Lasso, players with no zero estimated RAPM coefficients in more than 4 repetitions are kept in the model. A total of 184 out of the 717 offenders and 143 defenders have identified with a non-zero coefficient. The rest of the players are summed up and recorded in the intercept of the model. To mention that the minimum error penalty  $\lambda_{\min} \approx 0.27$  is used for the model implementation, as the one-standard error  $\lambda_{1se} = 1.17$  shrinks all coefficients towards zero.

Let us reinstate the example of section 3.4.1. Considering the offensive lineup of Golden State Warriors (GSW) with Curry, Thompson, Green, Wiggins and Looney and the defensive lineups of Boston Celtics (BOS) with Tatum, Smart, Horford, Brown and Williams, the average expected points scored by GSW per 100 possessions are 86.92. If Porter Jr. substitutes Looney and the rest of the players remain the same, then there is no difference in expected points as only Curry has a non-zero contribution.

Both ridge and lasso identify the same best players. Regarding the best performers of each team (Table 3.20 in Appendix I), for some teams the highest contribution was found to be zero. More, in lasso, there is an important number of “top star” players that also have high/top-rated RAPM values such as Kevin Durant, Jason Tatum, and Ja Morant. This is in contrast with ridge regression which totally fails in this task.

Under this approach a number of all-stars or very good players is picked up as top players in these terms (Tables 3.21 and 3.22 in Appendix I present the lasso top 5 offenders and defenders per position) in contrast with ridge regression, where the method failed to pick up good players. For example:

- the shooting guards Donovan Mitchell, Buddy Hield, and Derick White
- the small forwards Jayson Tatum, and Michael Porter Jr.
- the point guards Steph Curry, Jrue Holiday, and James Harden
- the power forwards Kevin Durant, Aaron Gordon, DeMar DeRozan, and Paul George

- the centers Joel Embiid, Deandre Ayton, Karl Anthony Towns, and Rudy Gobert are stars of NBA.

From the above results it is clear that the Lasso method gives better and more realistic results than Ridge since it identifies all-star as the best players for each team and position in contrast to ridge which failed. A notable difference is that the estimation of the intercept of Lasso is larger than Ridge regression results value.

When we study closer the top 50 players of Lasso RAPM we reach the following notable observations and conclusions<sup>54</sup>:

- ✓ Considering the **offensive RAPM**, Lasso regression identifies at least 20 all-stars in the top 50, like Tatum, Durant, Harden, Curry, Antetokounmpo, Jokic, Embiid, Irving, Trae Young, Jimmy Butler, Booker, LaVine, DeRozan, etc. On the other hand, Ridge regression identifies only 10 very good players and only five superstars (Durant, Harden, Booker, Embiid, Irving) in the top 50 players but with a lower ranking than in Lasso. Considering the All-star votes for the 2021-2022 NBA season, from the 28 higher-voted for All-Star players 16 were included in the top-50 and 13 in the top-28 of Lasso while only four and zero in Ridge top players respectively.
- ✓ Considering the **defensive RAPM**, in top-50 if Lasso ratings there are at least 20 very good players, such as Curry, Tatum, Paul George, Bam Adebayo, Derrick White, Ricky Rubio, Caruso, Iguodala, Chris Paul, Ntilikina, Gobert, Obi Toppin, VanVleet. Ridge regression includes in the top 50 the more mediocre Anthony Lamp, Isaiah Thomas, Malik Flits, Romeo Langford, Thomas Satoransky, Michael Porter Jr. and Iguodala. The two last NBAers who are good players belong also to the Lasso top 50 performances but on higher-ranking positions. Considering the two All defensive teams (5-lineups) of the season 2021-2022, six out of 10 players belong at the Lasso RAPM top 50 with the rest four players belong at the reference level. On the other hand, the ridge RAPM top 50 does not include anyone of these players.

Moreover, there are players with zero Lasso RAPM, but they were ranked at a high position according to the Ridge regression. For example, DeJon Jarreau (IND) which have the fourth best

---

<sup>54</sup> See Appendix I Tables 3.17 and 3.18.

offensive RAPM performance via Ridge, in Lasso his RAPM was shrunk to zero. Moses Wright (LAC) has the second highest defensive contribution according to ridge regression results, but Lasso was set his RAPM to zero. The most important here is that both played only one minute in the season, something that Lasso identified by shrinking their contribution to zero.

In this context, by focusing on the total results, it seems that the ridge coefficients are quite different from the Lasso (See Figure 3.7). The main difference is of the zero-coefficients of Lasso. However, a stronger linear relationship is observed when zero coefficients are not considered (Figure 3.8), with 79% (62.6% considering the zero values of Lasso) Pearson's linear correlation between the lasso and ridge offensive RAPM and 69.3% (54.6% considering the zero values of Lasso) for the defensive case.

Figure 3.7: Scatterplots of Ridge vs. Lasso coefficients for the full dataset, offensive (left) and defensive (right).

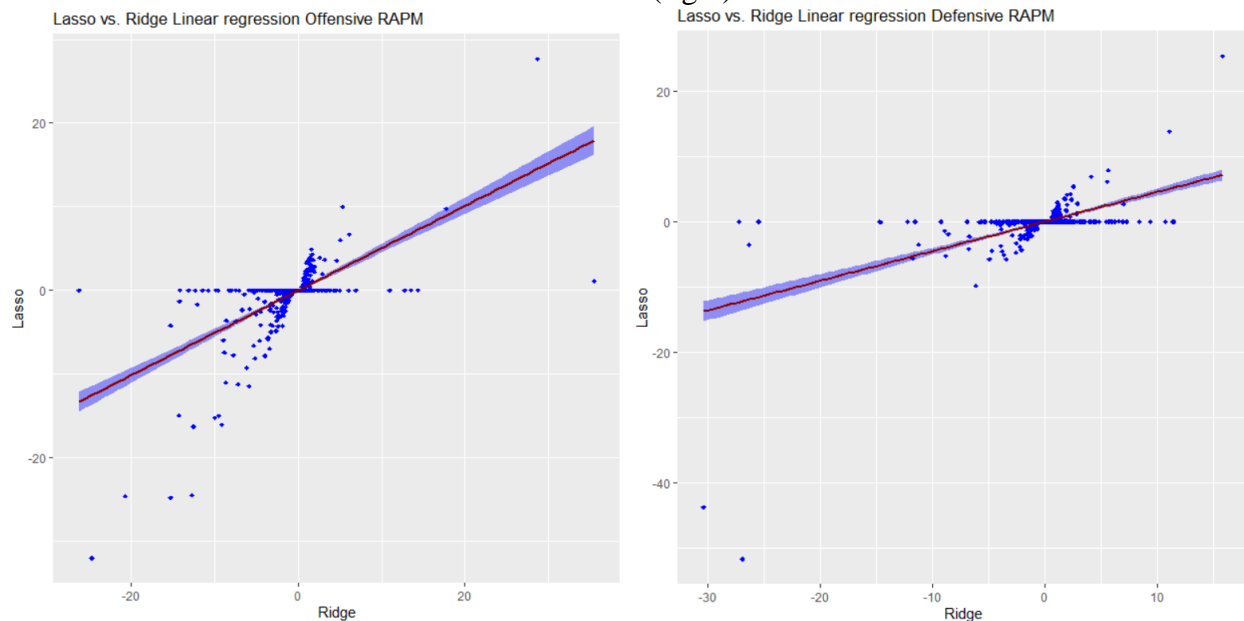
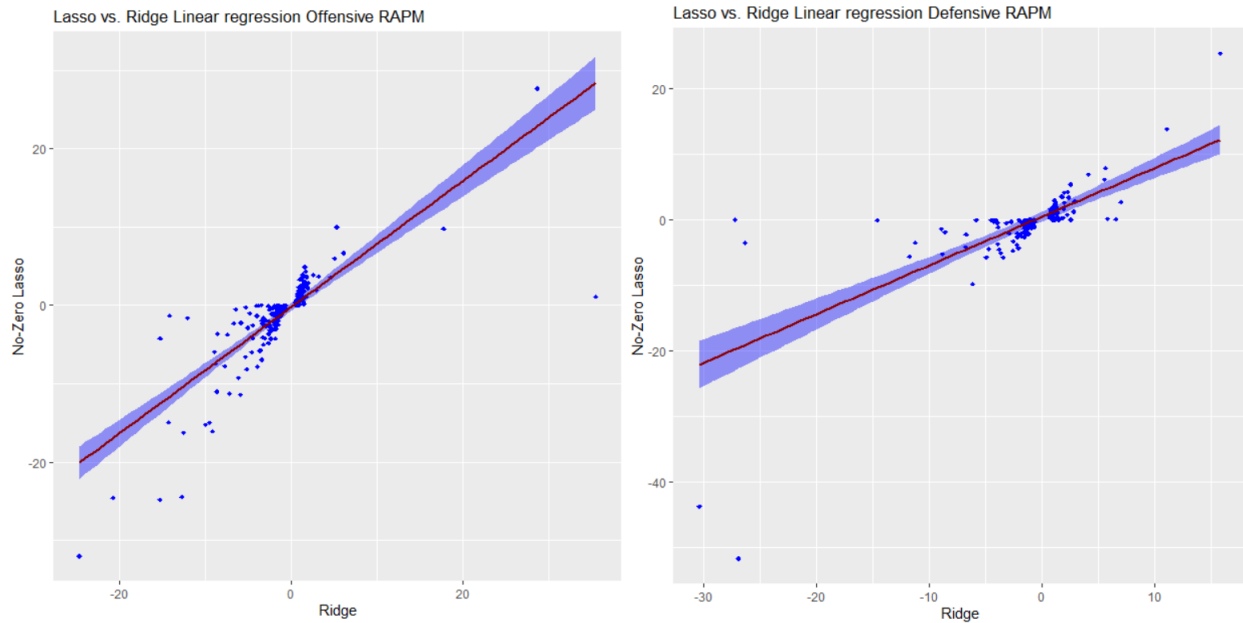


Figure 3.8: Scatterplots of Ridge vs. non-zero Lasso coefficients for the full dataset, offensive (left) and defensive (right).



### 3.4.4 Conclusion

To conclude with regularization regression methods (Ridge and Lasso) were used to calculate Adjusted Plus/Minus, which is a useful index for players rating, comparing players per team or position and their performance among the different teams they played.

Ridge regression for all players did not provide reliable ratings, since players with a low playing time appear with high contributions, higher than the ones expected and in comparison to better players with a considerable playing time. For this reason, we applied Lasso regression, which provided better results with ratings well-respected to good players and players with increased playing time. Although, Lasso is not used often in the basketball RAPM bibliography since it shrinks some coefficients exactly to zero. For those players, we cannot make further comparisons concerning their contribution or their ratings.

Although lasso provides better results than ridge, we still have the problem that some low time players (LTP) appear to be more impactful. This problem is common in the recent basketball bibliography. On the other hand, the basketball community is probably more interested in the better players and not the average, hence the players with zero RAPMs might not be of great difference.

A possible pathway is to exclude low playing time players. This practice was followed by other researchers (for example Rosenbaum (2004), and Ilardi (2007)). In the next section, we present the updated results without low playing time players.

### **3.5 Calculation of RAPM: Filtered dataset Analysis**

#### **3.5.1 Filtered dataset**

Not all players are important for the performance of a team. Some players are crucial, other players are starters, some are role-players or rookies, which play their first years in the NBA, and there is a number of low time players that constitute the “third unit” of the teams. Especially, NBA teams gives the chance to many players to try their abilities in the league (G-league, 2-way contracts, etc.).

As we have seen in section 3.4, dataset includes an impactful number of players with less playing time that influence the results of the ridge regression. Due to this, we decided to filter the dataset and not consider players with less than 200 minutes played (20 games \* 10 minutes). This value is compatible with similar choices in the literature.

For example, Rosenbaum (2004) used data for players with more than 250 minutes played in seasons 2002-2004. Ilardi (2007) set a stricter limit by including players with more than 400 minutes played in season 2006-2007, while he ran the regression separately for those players with higher playing time than 1640 minutes played in order to decrease the noise of estimations. Ilardi and Barzilai (2008) used a smaller threshold as they ran their model for all players with more than 300 minutes played in the 2007-2008 season.

The last work is closer to our implementation. In our work, the threshold is selected to be 200 minutes, not too strict but slightly lower than the choice of Rosenbaum’s (2004) and Ilardi and Barzilai’s (2008).

Table 3.25 (Appendix I) presents the excluded 232 players (with less than 200 total minutes played). It is obvious that there are players with one game and players with more than 20 games

but with few minutes on the court. For example, Petr Cornelie played for Denver Nuggets 31 games with 176 minutes played, which means less than 6 minutes per game of the 31.

All low playing time players of our analysis are used to form a “fictional” player that will be used as a reference player level in some models.

### **3.5.2 Ridge regression**

Ridge regression for the calculation of the plus/minus index is applied in the new filtered dataset and it is presented in the previous section of this paragraph. The values of coefficients are now smaller than the RAPM ridge ratings with all players. Again, the players with the best offensive and defensive performance are not popular and all-star players.

The shrinkage percentage of the linear regression estimations via the ridge method is higher when the one standard error lambda ( $\lambda_{1se} = 1422.86$ ) compared to the minimum's error lambda ( $\lambda_{min} = 166.16$ ), but with no great difference (shrinkage at 99.93% and 99.74% for the absolute terms, and 99.47% and 92.78% for the squared terms, respectively). Therefore, we present the results of the ridge regression model with the minimum error penalty.

Considering the offensive lineup of Golden State Warriors (GSW) with Curry, Thompson, Green, Wiggins and Looney and the defensive lineups of Boston Celtics (BOS) with Tatum, Smart, Horford, Brown and Williams, the average expected points scored by GSW per 100 possessions are 111.53. If Porter Jr. substitute Looney and the rest players remain the same, the average expected points are 112.00.

Considering now those results that came by applying the Ridge regression with tuning parameter the minimum error's value ( $se0$ ):

- Usman-Garuba (HOU) has the best offensive RAPM, equal to 4.13.
- Michael Porter Jr. (DEN) has the best defensive RAPM, equal to 5.29.
- Semi Ojeleye (MIL) has the worst offensive RAPM, equal to -5.19.
- Alfonzo McKinnie (CHI) has the worst defensive RAPM, equal to -4.83.

Here we introduce an extra dummy indicator if any of the low time is in the playing 5-lineup. The coefficient of this fictional player is equal to -1.62 for offense while it is slightly larger for the defense, -0.40.

In Table 3.3 players with the best offensive and defensive RAPM per team are presented. Under this approach (ridge with low playing time players form a “fictional” player), important players are observed in the list of the highest rated players. Even for players do not appear, their difference from the top-rated players in terms of RAPM is low. For example, Giannis Antetokounmpo did not have the highest value of RAPM for Milwaukee Bucks, but he belongs in the top-3 performances of his team.

Regarding the goodness of fit of Ridge regression in the filtered dataset the overall performance is not satisfactory. It seems like the model does not fit well with the data<sup>55</sup>. This is expected since here we try to model an asymmetric discrete response with a Gaussian distribution (see Figure 3.9).

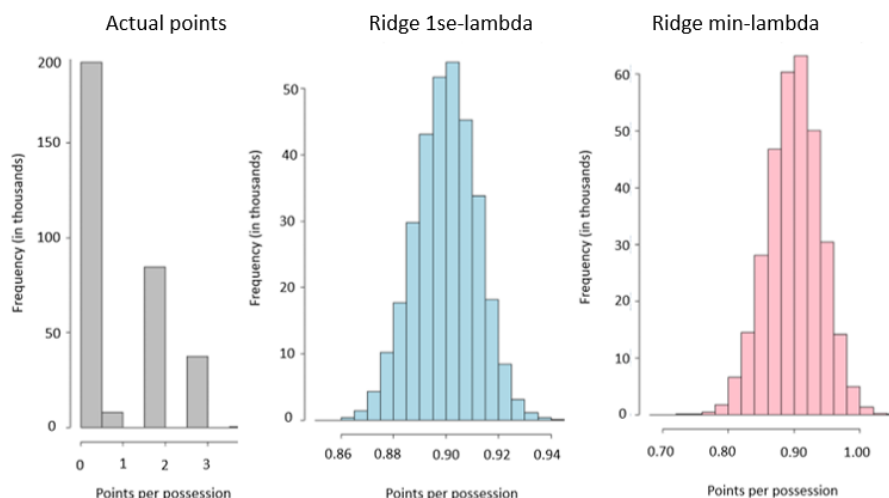
Table 3.3: Players with the best offensive and defensive RAPM of each NBA with more than 200 minutes played.

TEAMS	OFFENSE	DEFENSE	TEAMS	OFFENSE	DEFENSE
ATL	Bogdan-Bogdanovic	Onyeka-Okongwu	MIA	Markieff-Morris	Gabe-Vincent
BKN	Kyrie-Irving (s1) Kevin Durant (s0)	Paul-Millsap	MIL	Jevon Carter (s1) Jrue-Holiday (s0)	George-Hill
BOS	Daniel Theis	Robert-Williams III	MIN	Jaylen-Nowell	Nathan Knight
CHA	P.J. Washington	JT Thor	NOP	CJ-McCollum	Jose-Alvarado
CHI	Javonte Green	Tony Bradley	NYK	Miles-McBride	Derrick-Rose
CLE	Caris-LeVert	Ed Davis	OKC	Lindy Waters III	Kenrich-Williams
DAL	Spencer Dinwiddle	Frank-Ntilikina	ORL	Markelle Fultz (s1) Moritz-Wagner (s0)	Tim-Frazier
DEN	Bones-Hyland	Michael Porter Jr.	PHI	James Harden	Isaiah Joe (s1) Andre-Drummond (s0)
DET	Isaiah Livers	Rodney-McGruder	PHX	Deandre Ayton	Jalen Smith
GSW	Stephen-Curry	Andre-Iguodala	POR	Damian-Lillard (s1) Justise-Winslow (s0)	Nassir-Little (s0)
HOU	Usman Garuba	Daniel-Theis	SAC	Terence Davis	Terence Davis
IND	Buddy-Hield	Kelan-Martin	SAS	Josh-Richardson	Derrick-White
LAC	Robert-Covington	Justise-Winslow	TOR	Pascal-Siakam	Chris-Boucher
LAL	Wenyen-Gabriel	Rajon-Rondo	UTA	Juancho-Hernangomez	Joe-Ingles (s1) Danael-House Jr. (s0)
MEM	Brandon-Clarke (s0)	Dillon-Brooks	WAS	Brad-Wanamaker	Brad-Wanamaker

<sup>55</sup> Moreover, note that  $R^2$  is very small again (0.00105 for the Ridge regression with 1se lambda and 0.00338 for the case of lambda min).



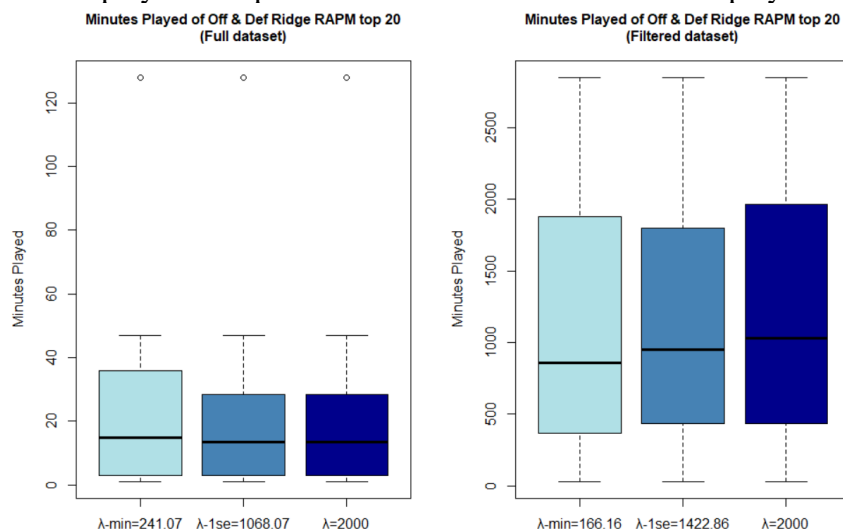
Figure 3.9: Histograms of actual points per possession and fitted values of the Ridge regression model for the “filtered” dataset (via 1se lambda the light blue and lambda min the pink).



Overall, the results of the ridge on the filtered dataset are rather improved in comparison to the results of the full dataset. The improvement is summarized in the following observations:

- In the top 20 offensive RAPM performed players of the filtered dataset, some very exceptional players appear, such as James Harden, Kevin Durant, Joel Embiid, Kyrie Irving, Steph Curry, Devin Booker, and DeAndre Ayton.
- Similarly, in the top 20 defensive RAPM performed players of the filtered dataset, some very good players appear, such as Paul George, Andre Iguodala, George Hill, Gary Payton II, Ricky Rubio.
- There is not top star player or even a very good player in the top-20 of the offensive or defensive RAPM that was top-rated in the full dataset. This is validated by also looking at the playing time of the RAPM top-rated players. Specifically, except the mediocre Rodney Hood who played 128 minutes for Lakers (he was traded later to Cleveland Cavaliers) the rest of the players with the highest offensive contribution played less than 47 minutes in the season (also including players with one and two minutes played). Similar is the case of the defensive RAPM, from the top 20 players all of them had playing time less than 28 minutes in the season (apart from RJ Nembhard Jr. of Cleveland Cavaliers who played 63 minutes).

Figure 3.10: Minutes played of top 20 offensive and defensive RAPM players via ridge models.



- Considering the minutes played by the basic players of each team (the six high minutes played by players per team), the 37% of the top 50 offensive and defensive ridge RAPMs from the filtered dataset (independently of the  $\lambda$  choice) is a starter player (he played at least the same time as the lowest-minute starter). On the other hand, only 8% of the top 50 offensive and defensive ridge RAPMs from the whole dataset. A notable observation here is that the Lasso model estimates the 56% starters of the top 50 offensive and defensive RAPMs from the full dataset.

### 3.5.3 Lasso regression

In this section, we proceed by implementing Lasso in the filter dataset. Cross-validated Lasso is applied in order to estimate the RAPM ratings. After running 10-fold cross-validated Lasso, we keep players with non-zero estimated coefficients in more than four iterations (out of 20), a number of 158 out of the 489 offenders and 138 defenders have been identified with non-zero RAPMs. The minimum error's tuning parameter is selected again,  $\lambda_{\min}=0.25$ , as the one-standard error tuning parameter shrink all coefficients towards zero.

Returning to the illustrated case of the offensive lineup of Golden State Warriors (GSW) with Curry, Thompson, Green, Wiggins and Looney and the defensive lineups of Boston Celtics (BOS) with Tatum, Smart, Horford, Brown and Williams, the average expected points scored by GSW

per 100 possessions are 86.90. If Porter Jr. substitute Looney for GSW, the expected points do not change, as both appear with zero offensive RAPM.

Joel Embiid has found with the highest offensive contribution, 4.79 points, and Michael-Porter Jr. with the highest defensive contribution, 7.81 points. Also, the coefficient of the fictional player's offensive plus/minus rating was found to be equal to -3.15, while the lasso coefficient for the defensive plus/minus rating is zero.

Below the best performers of each team are presented (Table 3.4). Best players with their coefficient shrunk to zero are fewer than previously (for the full dataset). Based on these results, some very good players of each team have the best RAPM, in contrast with the Ridge regression's results for players with more than 200 minutes played and the results from Lasso for the full dataset.

Table 3.4: Players with the highest offensive and defensive Lasso RAPM, with at least 200 minutes played, per team.

TEAMS	LASSO		TEAMS	LASSO	
	OFFENSE	DEFENSE		OFFENSE	DEFENSE
ATL	Danilo-Gallinari	Onyeka-Okongwu	MIA	Jimmy-Butler	Bam-Adebayo
BKN	Kevin-Durant	Jevon-Carter	MIL	Jrue-Holiday	George-Hill
BOS	Jayson-Tatum	Jayson-Tatum	MIN	Jaylen-Nowell	Josh-Okogie
CHA	P.J. Washington	JT-Thor	NOP	Brandon-Ingram	Jose-Alvarado
CHI	DeMar-DeRozan	Tony-Bradley	NYK	Miles-McBride	Immanuel-Quickley
CLE	Darius-Garland	Ricky-Rubio	OKC	Shrunk to zero	Kenrich-Williams
DAL	Dorian-Finney-Smith	Frank-Ntilikina	ORL	Shrunk to zero	Jalen-Suggs
DEN	Bones-Hyland	Michael Porter Jr.	PHI	Joel-Embiid	Matisse-Thybulle
DET	Isaiah-Livers	Shrunk to zero	PHX	Deandre Ayton	Cameron-Payne
GSW	Stephen-Curry	Andre-Iguodala	POR	Damian-Lillard	Shrunk to zero
HOU	Usman Garuba	Shrunk to zero	SAC	Terence Davis	Terence Davis
IND	Buddy-Hield	Chris-Duarte	SAS	Jakob-Poeltl	Derrick-White
LAC	Isaiah-Hartenstein	Justise-Winslow	TOR	Pascal-Siakam	Chris-Boucher
LAL	Wenyen-Gabriel	Rajon-Rondo	UTA	Donovan-Mitchell	Rudy-Gobert
MEM	Ja-Morant	Dillon-Brooks	WAS	Corey-Kispert	Shrunk to zero

In addition, in tables 3.5 and 3.6 the best offenders and defenders per position via Lasso regression-RAPM are presented. The existence of more all-stars or very good players is distinct in contrast with Ridge regression results again. In comparison with the Lasso top-5 for the full dataset, there

are not many differences. We only observe that this time most of the players are in higher ranking positions while some fewer good players have been replaced with others just as good or better.

Table 3.5: Top 5 players with the highest offensive Lasso-RAPM per position with at least 200 minutes played.

Rank	Shooting Guard (SG)	Small Forward (SF)	Point Guard (PG)	Power Forward (PF)	Center (C)
1	DEN Bones-Hyland	BOS Jayson-Tatum	GSW Stephen-Curry	BKN Kevin-Durant	PHI Joel-Embiid
2	UTA Donovan-Mitchell	PHX Mikal-BRidges	MIL Jrue-Holiday	LAL Wenyen-Gabriel	PHX Deandre-Ayton
3	MIN Jaylen-Nowell	WAS Corey-Kispert	PHI James-Harden	HOU Usman-Garuba	CHA P.J.-Washington
4	IND Buddy-Hield	NOP Brandon-Ingram	CLE Darius-Garland	DEN Aaron-Gordon	MIN Karl-Anthony-Towns
5	SAC Terence-Davis	MIA Jimmy-Butler	BKN Kyrie-Irving	CHI DeMar-DeRozan	UTA Hassan-Whiteside

Table 3.6: Top 5 players with the highest defensive Lasso-RAPM per position with at least 200 minutes played.

Rank	Shooting Guard (SG)	Small Forward (SF)	Point Guard (PG)	Power Forward (PF)	Center (C)
1	GSW Gary-Payton II	DEN Michael-Porter Jr.	MIL George-Hill	LAC Justise-Winslow	MIA Bam-Adebayo
2	SAS Derrick-White	OKC Kenrich-Williams	CLE Ricky-Rubio	LAC Paul-George	CHI Tony-Bradley
3	MIN Josh-Okogie	GSW Andre-Iguodala	PHX Cameron-Payne	GSW Draymond-Green	UTA Rudy-Gobert
4	OKC Ty-Jerome	BOS Jayson-Tatum	NYK Immanuel-Quickley	NYK Obi-Toppin	CLE Ed-Davis
5	DAL Frank-Ntilikina	BOS Jaylen-Brown	MIA Gabe-Vincent	MEM Jaren-Jackson Jr.	BOS Robert-Williams III

Estimations from both ridge and lasso in the filtered dataset are closer compared to the ones in the full dataset, with 78.8% linear correlation between the ridge and lasso offensive RAPM and 73.9% for the defensive case (see Figure 3.11). Also, again, not considering the zero RAPMs of Lasso, the values of the two regularized modeling procedures are pretty close with 91.5% linear correlation and 90.2% for the offensive and defensive RAPMs respectively (see Figure 3.12).

Figure 3.11: Scatterplots of Ridge vs. Lasso coefficients for the filtered dataset, offensive (left) and defensive (right).

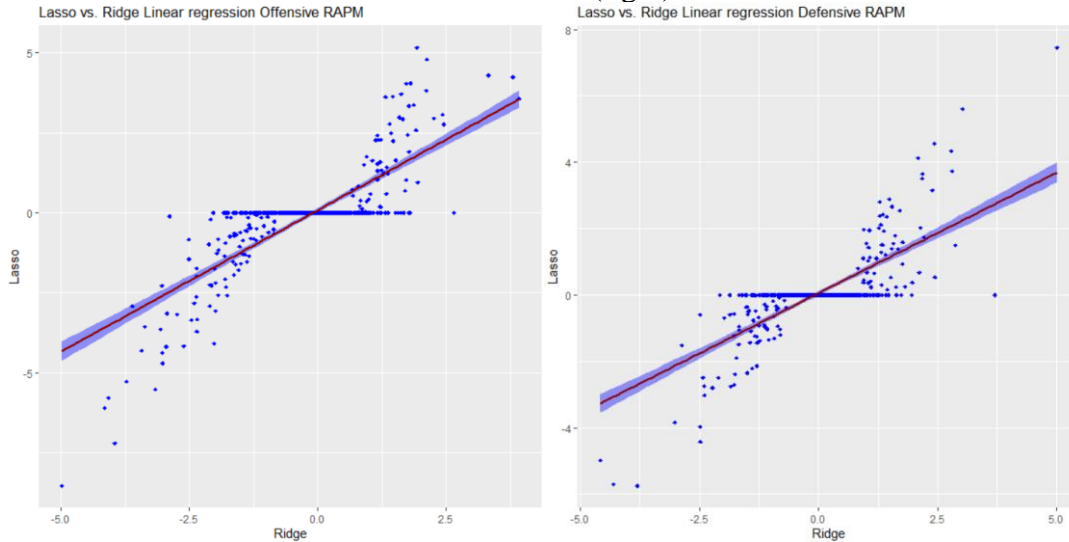
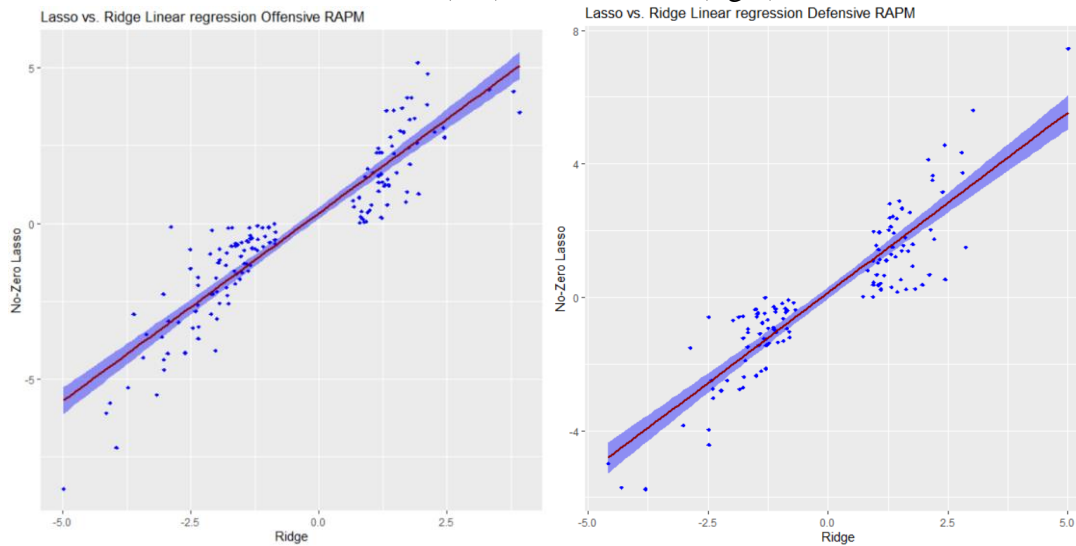
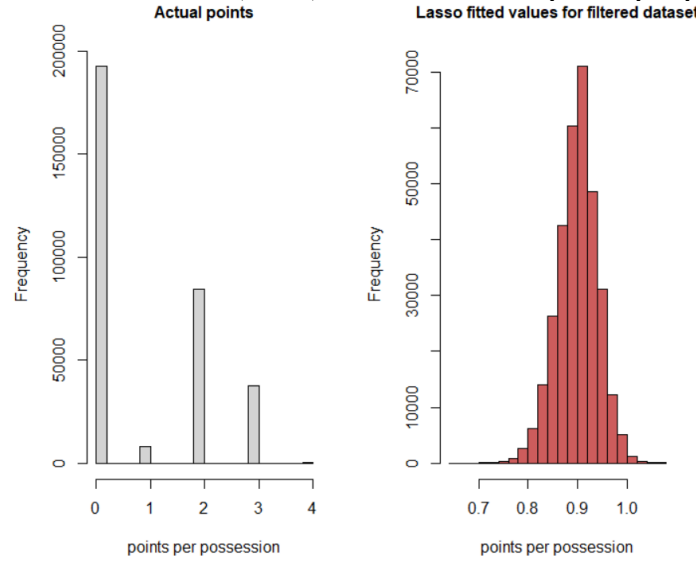


Figure 3.12: Scatterplots of Ridge vs. non-zero Lasso coefficients for the filtered dataset, offensive (left) and defensive (right).



Regarding the goodness of fit of Lasso regression in the filtered dataset the overall performance is not satisfactory. It seems like the model does not fit well with the data just like the Ridge model (see Figure 3.13).

Figure 3.13: Lasso RAPM (fitted) values vs. actual points per possession.



### 3.5.5 Assessing the effect of LTPs and the team ability

Here we study the case of not eliminating the low time players (LTP) from the model by removing the “fictional” player. Also, we will study the teams effect by adding the 30 dummies of the teams in our RAPM models.

First, let us devote the model with the “fictional” player effect by  $M_1$ , while the model without the “fictional” player effect by  $M_2$  and finally the model with not the “fictional” player effect but with the team effects as  $M_3$ .

The results of the new two models ( $M_2$  and  $M_3$ ) are similar to the ones of  $M_1$ . No notable differences at the top five RAPMs for each position or the best rated players per team are observed. More differences are observed for the defensive RAPM<sup>56</sup>.

Moreover, the correlations of the coefficients between all three models are high (84.35% - 99.98%). Specifically see Table 3.7.

<sup>56</sup> See Appendix I (I.4).

Table 3.7: Pearson's linear correlation ( $r$ ) between the model with the fictional player effect ( $M_1$ ), with not the fictional player effect ( $M_2$ ) and the teams effects ( $M_3$ ).

$r$	Ridge Off RAPM			Ridge Def RAPM		
	M1	M2	M3	M1	M2	M3
M1	1.0000	-	-	1.0000	-	-
M2	0.9998	1.0000	-	0.9997	1.0000	-
M3	0.9519	0.9470	1.0000	0.8435	0.8526	1.0000
$r$	Lasso Off RAPM			Lasso Def RAPM		
	M1	M2	M3	M1	M2	M3
M1	1.0000	-	-	1.0000	-	-
M2	0.9908	1.0000	-	0.9953	1.0000	-
M3	0.9893	0.9941	1.0000	0.9949	0.9919	1.0000

It seems that the three models do not give different information and we could not say that one of them is better or more suitable compared to the others.

### 3.5.6 Validation of RAPM values: Comparison with objective criteria and ratings

From the comparison of model based RAPMs from the “filtered” dataset, we observe that Lasso RAPMs are better than the rest of the approaches. In this section we are going to validate the results by comparing model based RAPM to actual facts and “objective” external criteria.

The All-NBA Team is an annual NBA honor bestowed on the best players in the league following every NBA season. All-NBA Team is composed of three five-lineups - a first, second, and third team. The voters select two guards (G), two forwards (F) and one center (C) for each team. Thus, there are six positions for guards, six for forwards and three for a center player to be selected as one of the best.

Hence, we validate our models by accounting the number of All-NBA teams players that was highly rated by model based RAPM. Results per position (Shooting Guard, Point Guard, Small Forward, Power Forward and Center) and per “general” position (Shooting Guards and Point Guards are considered as Guards, Small Forwards and Power Forwards are considered as Forwards, Center) are provided in the following. From Table 3.8 we can conclude that Lasso models have higher percentage of agreement than Ridge.

Table 3.8: Percentage of agreement between RAPM and All-NBA team 2021-2022 per position (filtered dataset).

	Ridge ( $M_1$ )	Ridge ( $M_2-M_3$ )	Lasso ( $M_1$ )	Lasso ( $M_2-M_3$ )
Per Position	40.00%	33.33%	53.33%	46.67%
Per “General” Position	26.67%	20.00%	33.33%	33.33%

A second validation is made through the percentage of agreement between the best model based RAPM player for each team and the All-NBA teams players. Considering the top three Ridge and Lasso RAPM players the percentage of agreement is 73.33%, while this percentage is the maximum (100%) considering the top five per team.

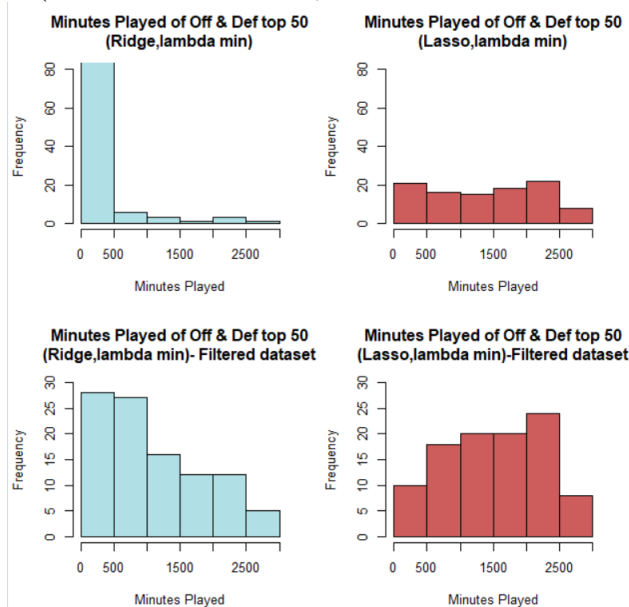
Similarly, we compare how many from the All-NBA defensive team which is composed of two five-man lineups, i.e., ten players were identified by our RAPM models. The percentage of predicting the best defenders of the 2021-2022 NBA season per position is only 10% for Ridge models and 60% for Lasso. Note that the performance of both approaches was much worse for the full dataset (only 10% for Lasso and zero for Ridge).

A last notable observation refers to the minutes played of the top offensive and defensive RAPM players. From Figure 3.14 we can see that:

- Regarding the full dataset, the top-ridge RAPM players had obviously less playing time than the ones of Lasso (about 80% with less than 500 minutes played and 80% with more than 500 minutes played respectively).
- Regarding the filtered dataset: the high rated Lasso RAPM players appear with more playing time and in fact the distribution of the playing time of those players is left skewed.



Figure 3.14: Minutes played of top 50 offensive and defensive RAPM players via regularized models (first row: full dataset, second row: filtered dataset).



### 3.5.7 Conclusion

From the results of section 3.5, we have seen that when we consider low time players (with less than 200 minutes played) as a single player, the RAPM rating improve in the sense that models track better players according to external criteria. Moreover, the RAPM ratings are less affected by low time NBAers. In addition, RAPMs of both Ridge and Lasso are closer to each other than the corresponding of the full dataset.

Results did not improve when the low time players were totally excluded from the model (i.e., their effect was embodied in the model formula constant) or when team effects were included.

However, the calculated RAPM ratings from the “filtered” dataset captures some actual facts/events. In the following paragraph, examples of the RAPM usage are presented.

## 3.6 Player Evaluation using RAPM

As the RAPM is a player evaluation index, it is useful for comparing players in a league, within each team or position. In the following, we focus on the RAPM values for specific players.

### 3.6.1 Comparing player's contribution

The contribution of a player depends on his team and its opponents. In this study, both the teammates and opponents have been considered via the model formula.

Scatterplots depict the 10 top players according to offensive RAPM, the 10 top players according to defensive RAPM and some of the top NBAers (such as Giannis Antetokounmpo, Steph Curry, LeBron James, Kevin Durant, Joel Embiid and Jason Tatum) in Figures 3.15 and 3.16. Each player's position appears in a different color, while the offensive RAPM is presented next to the players' name.

Starting from Figure 3.15 (Ridge RAPM), players like LeBron James, Giannis Antetokounmpo, Klay Thompson, Russell Westbrook, Damian Lillard, Rudy Gobert and Luka Doncic have not the top offensive performances as we would expect. Nevertheless, their RAPM value have no great difference from the maximum values, some comparisons could be done. For example, Steph Curry and Jason Tatum have a higher positive contribution to their team offense and defense, in contrast with other NBA stars that have high offensive contributions and a low or negative defensive impact to their teams, such as Damian Lillard or LeBron James.

RAPMs are also a useful tool for comparisons of players of the same position. From the results we can see that Durant's and Antetokounmpo's plus/minus ratings of offense and defense have not a large difference. Jason Tatum, after a great season, has a higher offensive and defensive RAPM than Jimmy Butler. Also, Karl Anthony Towns has a higher offensive and lower defensive RAPM than Rudy Gobert. This may suggest that a team that wishes to build a stronger defense may take an eye on the later player.

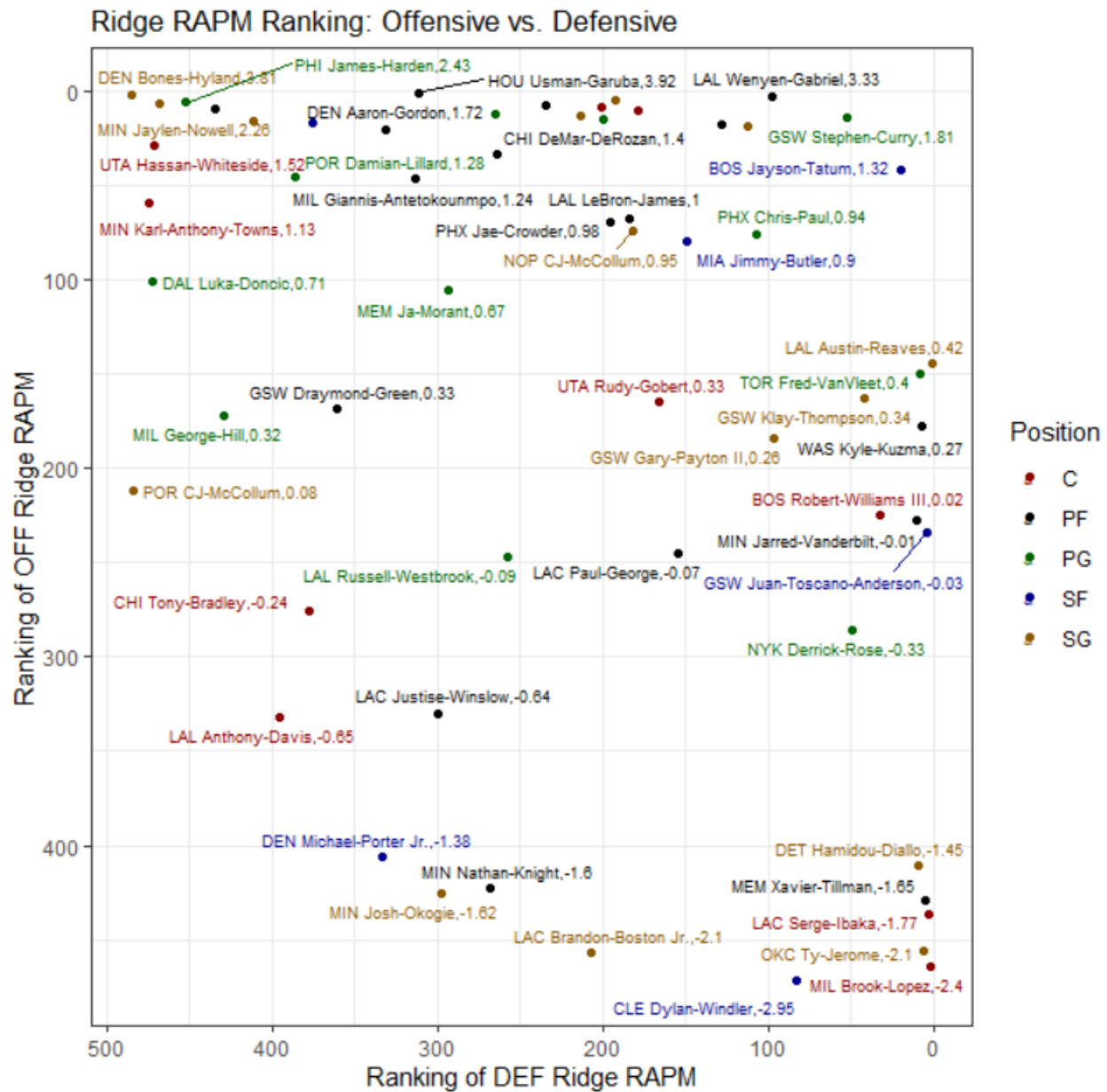
Generally, RAPM based comparisons like the ones above can help managers make decisions concerning the roster or management of a team.

Lasso RAPM rank some very good players in the top of performances. On top of the offensive RAPM ratings we find players like Joel Embiid, DeAndre Ayton, James Harden, Donovan Mitchell and Jrue Holiday. Note that in the case of Ridge, the rank of these players was lower.

Most of the good offenders presented in Figure 3.16 have smaller defensive than offensive contributions to their team. Steph Curry and Jason Tatum are exceptions as since they help their

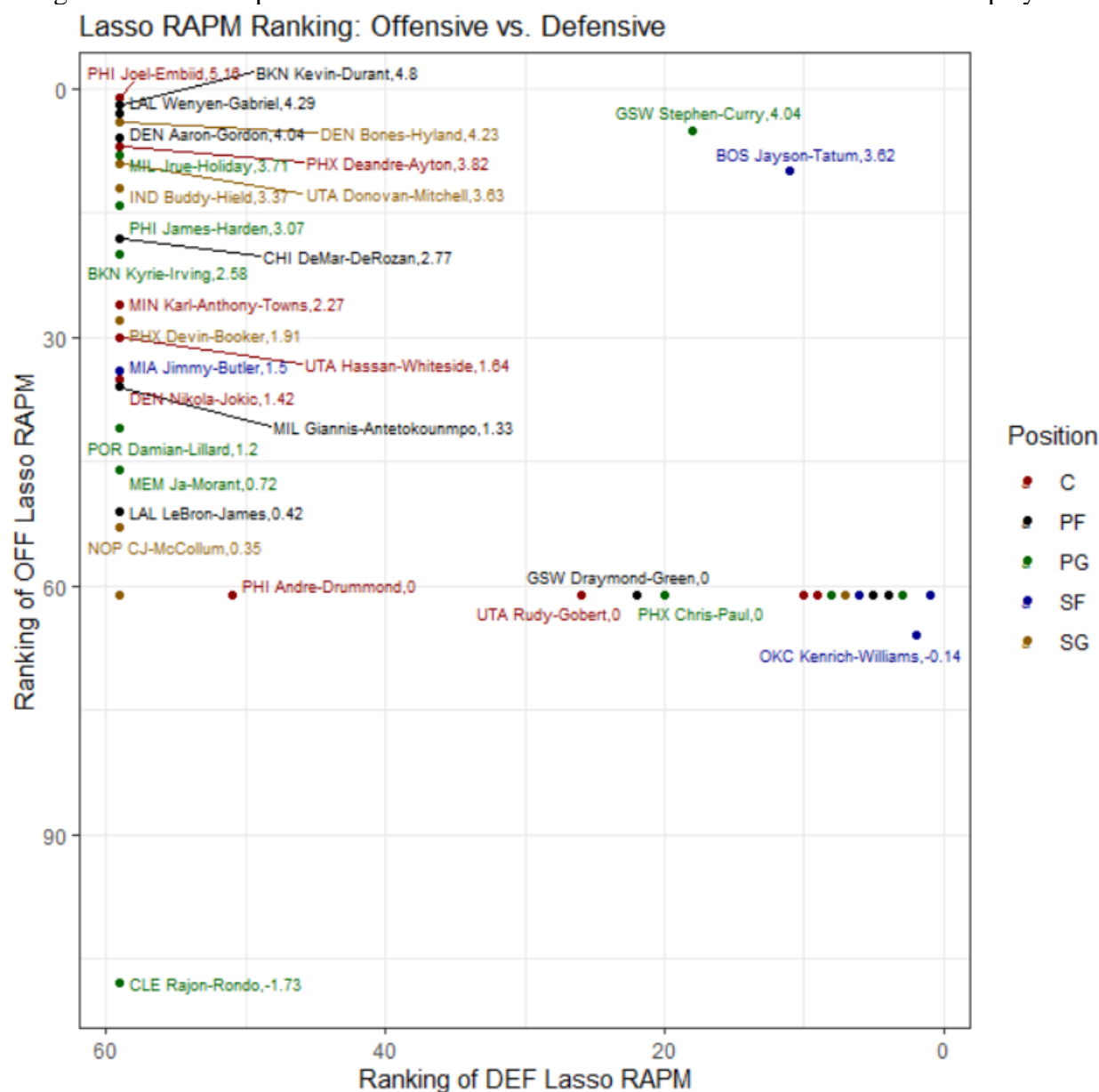
team in both defense and offense. However, Lasso does not allow us to compare a group of players who shrunk to zero. These can be considered as average players that do not increase or decrease the points per possession significance from the overall mean.

Figure 3.15: Scatterplot of Ridge Offensive and Defensive RAPM for a selection of players\*.



\*Consist of the top 10 Ridge RAPM offenders, top 10 Ridge RAPM defenders along with selected top NBAers. Each position is in a different color.

Figure 3.16: Scatterplot of Lasso Offensive and Defensive RAPM for a selection of players\*.



\*Consist of the top 10 Lasso RAPM offenders, top 10 Lasso RAPM defenders along with selected top NBAers. Each position is in a different color.

The previous analysis illustrated the usefulness of RAPM for comparing players of the same team or position. A remarkable finding is that James Harden was not selected for the All-NBA teams (three 5-lineups players) although his high offensive Lasso and Ridge RAPM.

Moreover, comparison of the performance of a player who played two different teams could be done. For example, CJ McCollum played for New Orleans Pelicans (NOP) and Portland Blazers

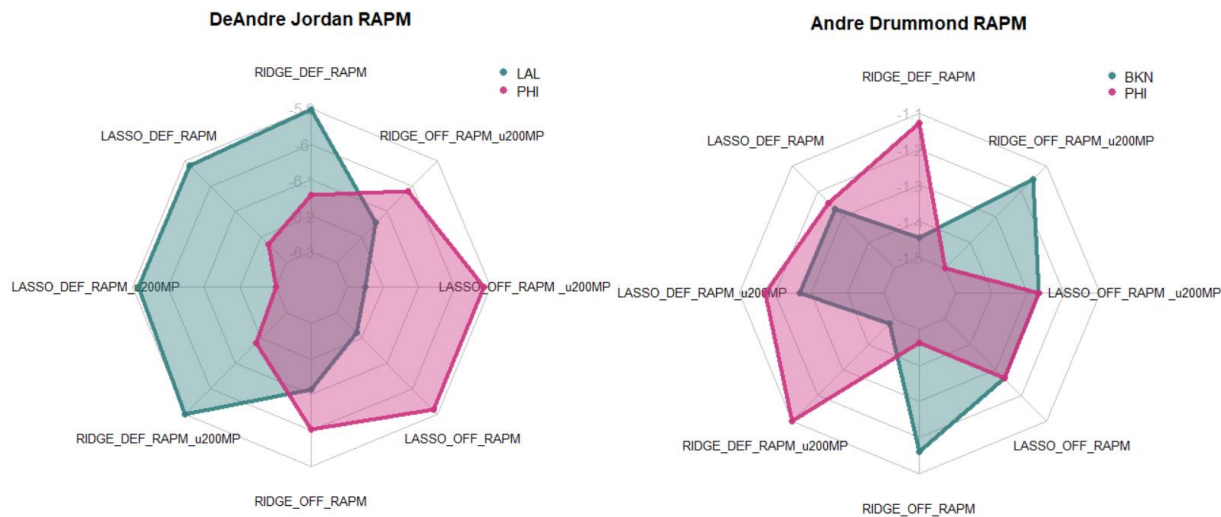
(POR) during the season 2021-2022. From Figure 3.15, we observe that McCollum had a higher offensive contribution when he was playing in NOP (offensive RAPM=0.954) than in POR (offensive RAPM=0.078). This is confirmed by Lasso RAPM as well since his offensive RAPM when he was playing in POR shrunk to zero while it was positive in NOP (0.351). This is probably because in Portland, McCollum had less participation in offense, with Damian Lillard having the largest participation, than NOP. During his participation in NOP, McCollum helped his new team with a more important role in offense. More details about some selected traded players are provided in the Section 3.6.2 which follows.

### **3.6.2 Comparing traded players' contribution**

NBA teams usually trade players within season based on specific regulations, in order to get stronger in the mid-season. Hence, here we focus on selected traded players evaluations for different teams they played during the season 2021-2022. During this season, 42 players were traded in total. Naturally, for some of them, the trade affected positively their RAPM performance while for others negatively.

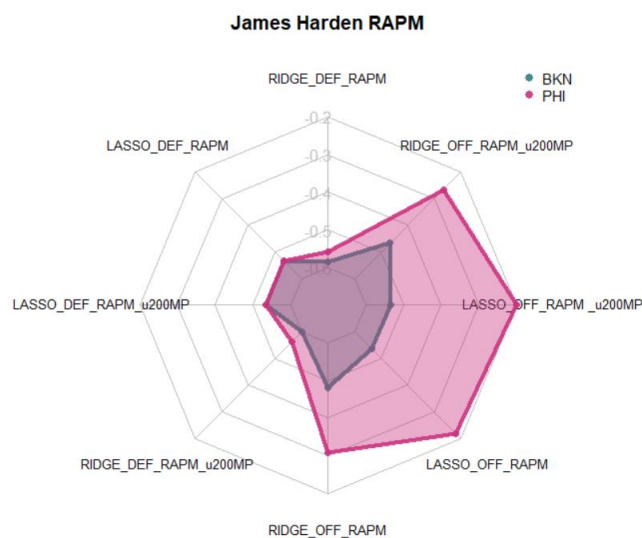
For example, in Figure 3.17, Andre Drummond appears with higher offensive RAPM when he was playing in Brooklyn Nets (BKN) than in Philadelphia 76ers (PHI). As for his defensive RAPM it is observed exactly the opposite. Since Drummond is not considered as a top player, he may have less impact on a higher quality NBA team. Another example is DeAndre Jordan (see Figure 3.17) who had different role in Philadelphia 76ers (PHI) and Los Angeles Lakers (LAL). In the first, his offensive RAPM is higher while in the second his defensive RAPM is improved. This is reasonable, since Philadelphia overall performed better than Lakers and therefore scored more, while the later team was more important to have better defensive performance.

Figure 3.17: Radar plots for Andre Drummond and DeAndre Jordan's RAPM for their teams.



Let us now assess the case of a top NBAer, James Harden, who started the season to Brooklyn Nets (BKN) along with Kevin Durant and Kyrie Irving. Although there were high expectations for this team, its performance was poor. Therefore, Harden soon moved to Philadelphia 76ers (PHI) where he met another superstar teammate, Joel Embiid. From Figure 3.18, we can see that Harden had higher offensive RAPM when he was playing in Philadelphia than Brooklyn, where his offensive role was more limited since the Brooklyn offensive plan was based on (and therefore divided between) three top players (Durant, Harden, and Irving). On the other hand, his defensive RAPM is similar for both teams.

Figure 3.18: Radar plot for James Harden's RAPM for his teams.



As a second demonstrative example, let us consider Derrick White (Figure 3.19). He appears to have higher offensive RAPM and lower defensive RAPM in Boston Celtics (BOS) than San Antonio Spurs (SAS). This may be the result of the necessities and strategies of each team. Since in the 2021-2022 season Spurs had lower performance than Boston, their players had to pay more attention to defense. Thus, Derrick White, being one of the starter players appears with increased defensive impact on his first team (Spurs) compared to the later one.

Finally, there are players whose model based RAPM performance was not affected in obvious way their transfer, see for example Rajon Rondo whose performance is depicted in Figure 3.20. The radar plots for traded players performance are presented in Appendix I.5.

Figure 3.19: Radar plot for Derrick White's RAPM for his teams.

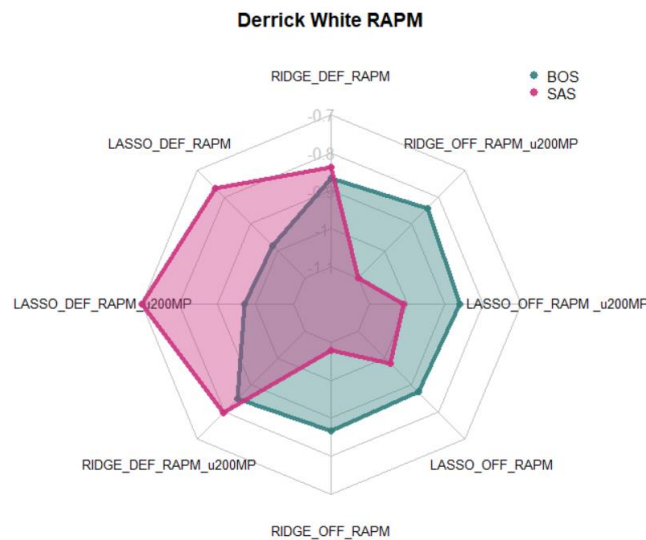
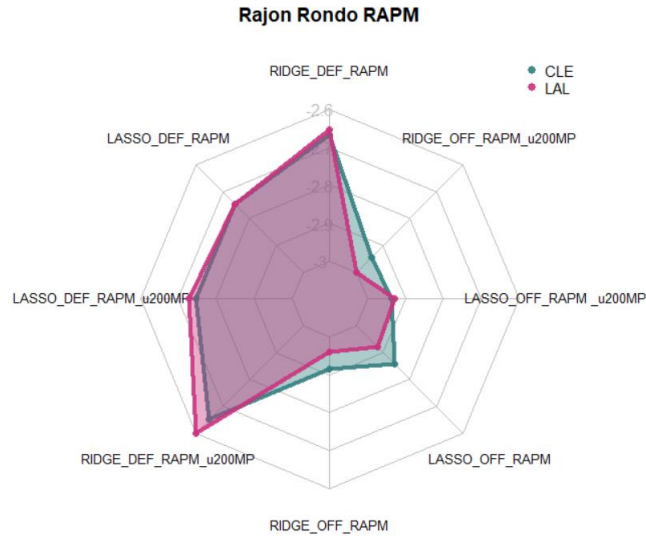


Figure 3.20: Radar plot for Rajon Rondo RAPM for his teams.



### 3.7 Conclusion

This chapter deals with the implementation of regularization linear regression models for obtaining RAPM performance ratings. Following the literature, a linear ridge model was used to construct initial RAPM ratings. Concerning the tuning parameter of ridge, the most common choice in the basketball RAPM literature is to set the penalty ( $\lambda$ ) equal to 2000. This choice was higher than the one-standard error  $\lambda$  in terms of the minimum RMSE found in our data and is the suggested choice in the several regularization literatures. Nevertheless, the two choices gave similar results in terms of the ranking.

Finally, we decided to use the choice of  $\lambda = \lambda_{\min}$ , which corresponds to the minimum RMSE. This was considerably smaller than the previous two values for obtaining RAPM performance ratings. The reason for this choice was that were similar to the rankings of RAPMs uses  $\lambda_{1se}$  (correlation close to one) and the scale was higher (due to less shrinkage) leading to better discrimination of players with.

However, the results of the ridge model are not sufficient, as low time players appear in the top RAPM rated players. Hence, a Lasso regression was implemented by using  $\lambda_{\min}$ . Lasso RAPM



ratings rank the players in a more sufficient way and capture in a higher degree best players of each team and position (in contrast to ridge RAPM). Although, low time players still appear more impactful than we expected.

To avoid this, players with less than 200 minutes played were initially considered as a single “fictional” player. This time, Ridge and Lasso model based RAPMs were considerably improved. The agreement of the RAPM rated players and some external criteria have been improved (40% and 53% agreement of All-NBA teams and higher rated players per position for Ridge and Lasso respectively, while the agreement reach 73% when considering the high rated players per team).

Although its popularity in the estimation of RAPM from the statistical point of view, the Normal distribution is not suitable to model the response variable in this problem since it is clearly discrete taking mainly values from zero to three. A model, which accounts for the discrete nature of the response, such as a multinomial, is clearly more suitable for this problem. For this reason, next chapter deals with the calculation of RAPMs from a Multinomial model.

## Chapter 4: Building Logistic Models

### 4.1 Introduction

In Chapter 3, the Regularized Adjusted Plus/Minus (RAPM) ratings were constructed via regularized linear models. Although the popularity of the RAPM based on normal model in basketball, the Normal distribution is not suitable for modelling the variable of points per possession, which is discrete taking values from zero to three, and rarely until six.

Hence the multinomial distribution for such a response. For this reason, seems more suitable, this chapter deals with the implementation of multinomial logistic models after removing the low time players (with less than 200 minutes played) from the dataset. In the first stage we implement a binomial logistic model for the modelling the probability of scoring or not, which is presented shortly in Section 4.2. More details and results about the binomial model are provided in the Appendix II.

In order to overcome computational problems, the multinomial model is implemented indirectly through three separate binomial models for the three main scoring situations:

- i. Scoring 1 point per possession or not scoring at all.
- ii. Scoring 2 points per possession or not scoring at all.
- iii. Scoring 3 or more points per possession or not scoring at all.

These three separate binomial models are presented in Section 4.3 while the final multinomial model, which is composed from these three models, is presented and discussed in Section 4.4

### 4.2 Binomial models

For the implementation of the binomial model the discrete variable of points per possession is transformed into a binary with two levels: scoring with value equal to one and not scoring with value equal to zero. The equation of the binomial model is the following:

$$Y_i \sim \text{Bernoulli}(P_{i1})$$
$$\log \frac{P_{i1}}{1 - P_{i1}} = b_0 + \sum_{j=1}^K b_j O_{ij} + \sum_{j=1}^K b_{K+j} D_{ij}$$

where  $Y_i$  is the response random variable which takes value one when the attacking team scores at possession  $i$  and zero otherwise,  $P_{i1}$  is the probability of scoring in  $i$  possession,  $P_{i0} = 1 - P_{i1}$  is the probability of not scoring in  $i$  possession,  $O_j$ , for  $j=1, \dots, K$ , takes the value one, when player  $j$  plays offense in  $i$  possession, and 0, otherwise,  $D_j$ , for  $j=1, \dots, K$ , takes the value minus one, when player  $j$  plays defense in  $i$  possession, and zero, otherwise.

The estimated coefficients of the offensive and defensive players' variables are measures of their contribution to their team's scoring and for the opponent team not to score. More specifically, the constant term  $b_0$  refers to the log-odds probability to score when only low-time players (with less than 200 minutes played) are on the court (which is an unrealistic scenario),  $b_j$  measures the difference in offensive log-odds for a team to score between player  $j$  and reference low-time players, for  $j=1, \dots, K$ , and  $b_{K+j}$  measures the difference in defensive log-odds for the opponent team not to score between player  $j$  and reference low-time players, for  $j=1, \dots, K$ . Differences of  $b_j - b_m$ , for  $j, m=1, \dots, K, j \neq m$ , show the average difference in log-odds between the players  $j$  and  $m$ , when their teammates and opponents remain the same. Hence,  $e^{b_j - b_m}$  will show the relative change in the odds of scoring, when player  $j$  substitute player  $m$  (and the rest of the players remain the same).

Following the methodology of the linear models, regularization methods are also applied for the binomial model. For both ridge and lasso methods, the results between the models with tuning parameter ( $\lambda$ ) in terms of the minimum error ( $\lambda_{\min}$ ) and the one-standard error ( $\lambda_{1se}$ ) do not differ with respect to the rank of the players while the correlation between the estimated RAPMs with the two values of lambda ( $\lambda$ ) is high (95% for ridge and 88% for lasso). The Lasso method with  $\lambda_{\min}$  penalty shrinks to zero 200 players less than the  $\lambda_{1se}$  penalty (617 players in overall, 299 offenders and 318 defenders). Therefore, we selected  $\lambda = \lambda_{\min}$  for the regularized binomial models, in order to take advantage of the bigger number of players with non-zero RAPMs, and the larger scale of in the values of RAPM.

Some results of Lasso and Ridge binomial models are presented in Table 4.1. The probability of a team to score in one possession is about 39% and 34% according to the ridge and lasso binomial models respectively for scenario where the lineups of the two opponents are composed only by the low-time players.

Table 4.1: Summary results of the ridge and binomial models.

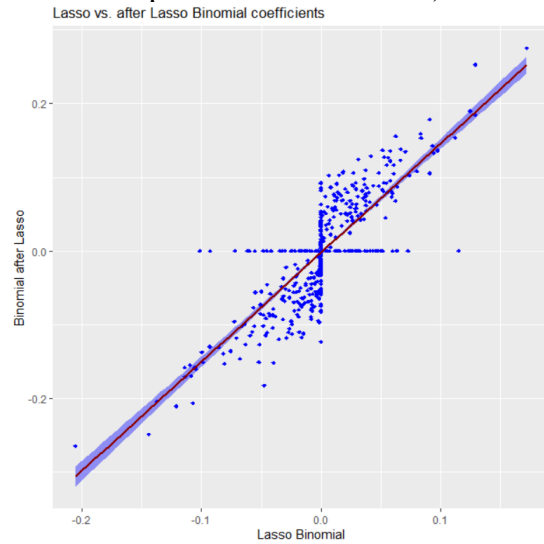
	<b>Ridge Binomial</b>	<b>Lasso Binomial</b>
$\lambda_{\min}$	0.301	0.001
<b>Intercept</b>	-0.495	-0.415
<b>Max OFF</b>	0.093 (HOU, Usman-Garuba)	0.129 (PHI, Joel Embiid)
<b>Min OFF</b>	-0.117 (MIL, Semi-Ojeleye)	-0.205 (MIL, Semi-Ojeleye)
<b>Max DEF</b>	0.101 (DEN, Michael-Porter Jr.)	0.172 (DEN, Michael-Porter Jr.)
<b>Min DEF</b>	-0.106 (CHI, Alfonzo-McKinnie)	-0.144 (CHI, Alfonzo-McKinnie)

According to the highest offensive contribution, the odds probability of a team to score is about 9.75% and 13.77% higher than the odds of the low-time players for the ridge and lasso model respectively. These odds refer to the case of the highest performance rated players being on offense, who are Usman Garuba for ridge and Joel Embiid for lasso. To note that the first player only played in 24 games (about 10 minutes per game). Hence, he was not expected to such a high RAPM rating. Similarly, according to the lowest defensive contribution, the odds of a team to score is about 11.0% and 19.5% smaller than the odds of low time players, when Semi Ojeleye plays offense, for the ridge and lasso model respectively.

Regarding the defensive contribution of the players, according to the highest one the odds probability of a team not to score is about 10.6% and 18.7% larger than the odds of the low time players when Michael Porter Jr. is on defense for ridge and lasso respectively. Moreover, due to the lowest defensive contribution, the odds probability of a team not to score is about 10.1% and 13.4% smaller than the odds low time players, when Michael Porter Jr. is on defense for ridge and lasso respectively.

As the results of the regularized model with the penalty are usually biased, Lasso method should be used as a screening method. Therefore, a binomial model was developed only with the players whose Lasso-coefficients were found to be non-zero. The unbiased OLS coefficients, and Lasso RAPMs are linearly correlated, 81.3%, with the players ranking were found to be 78.3% correlated (see Figure 4.1).

Figure 4.1: Scatterplot of coefficients from Lasso Binomial and Binomial (after the implementation of Lasso).



Regarding the fit of the implemented binomial models the accuracy under the ROC curve is not too high, 56%, while the Hosmer-Lemeshow goodness-of-fit test the equality of observed and expected proportions is rejected ( $p\text{-value} < 0.05$ ). Moreover, the misclassification error of the points levels (zero or one, i.e., not score or score) prediction is about 39%.

Adding the team effect in the logistic model resulted in the inclusion of 60 extra dummy variables (30 offensive and 30 defensive teams). The results are quite similar to the results of the models with no team effects as the NBA teams effect was not found to be high. The correlation of the RAPM coefficients under the two approaches is very high (97.6%).

Closing the discussion about the binomial logistic models, let us consider the three best All-NBA teams and the two defensive All-NBA teams of the 2021-2022 season. Table 4.2 presents the percentage of agreement between the high rated players from each binomial model and those players that compose the All-NBA teams. The Lasso model clearly captures the high rated players in a greater degree.

Table 4.2: Agreement between the Binomial-based high rated players and the All-NBA teams.

	Ridge	Lasso	Binomial (after Lasso)
All-NBA teams	26.67%	<b>46.67%</b>	26.67%
defensive All-NBA teams	0.00%	<b>50.00%</b>	30.00%

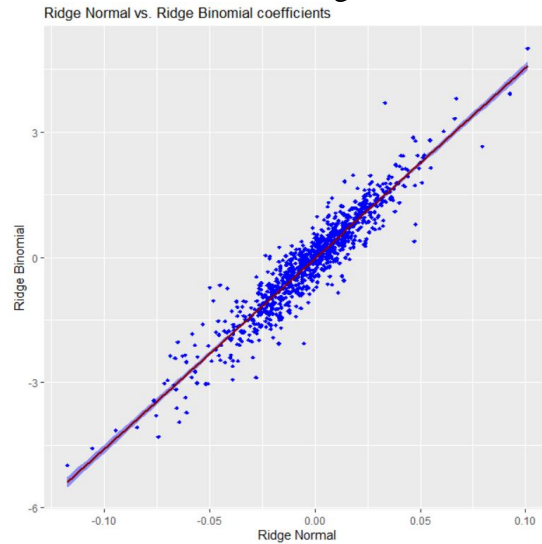
The correlations between the Ridge RAPMs of the Normal RAPM ratings developed in Chapter 3 and the Binomial RAPM ratings provided in this section are shown in Table 4.3.

Table 4.3: Linear correlation between the Linear and Binomial models results.

	Ridge Normal	Lasso Normal	Lasso Binomial	Ridge Binomial
Lasso Normal	0.767	-	-	-
Lasso Binomial	0.761	<b>0.882</b>	-	-
Ridge Binomial	<b>0.940</b>	0.733	0.816	-
Binomial (after Lasso)	0.716	0.709	0.813	0.775

From all the comparisons, it is obvious that the highest correlations are between the Ridge Binomial and Ridge Normal. The later is the one which is used frequently in the basketball literature (see Figure 4.2). We can exploit this correlation to offer an extra interpretation in the Normal RAPM ratings.

Figure 4.2: Scatterplot of coefficients from Ridge Binomial and Ridge Normal model.



In fact, the RAPM ratings of the Ridge Binomial and Ridge Normal are related via the below simple linear model:

$$Binomial_j = 0.02 \times Normal_j + \varepsilon, \varepsilon \sim N(0, 0.008^2)$$

for all  $j=1, \dots, K$  players (offender or defender), with  $R^2 \approx 0.89$ . This equation is extremely useful, since we can calculate the Binomial RAPMs from the Normal ones. By this way we can offer an

extra interpretation to the Normal RAPMs and connect them to RAPMs based on a correctly fitted modes with respect to the distribution used for the response.

### 4.3 Separate binomial models

Due to the size of the data, we had computational problems when trying to fit the multinomial. To overcome these problems, the Multinomial model was implemented in an indirect way through the three separate binomial models for the three scoring situations (one, two and three or more points per possession against possessions without scoring).

For each of the three models Lasso regularization is used as a screening method with the penalty parameter set to the minimum RMSE ( $\lambda_{\min}$ ). Then, we have fitted the binomial models including only those players whose coefficient was not shrunk to zero on Lasso. Note that, using a Lasso penalty equal to  $\lambda_{\text{lse}}$ , then all coefficients were shrunk to zero for all the three models. Hence, it resulted to a useless model with respect to the RAPM estimations.

Here we present a discussion based on the overall results of the three binomial models. Details about the fit of the three separate binomial models are given at the Appendix II.4 aloud with the top players per position and per type of points.

The three binomial models are described as following:

$$Y_{i1} \sim \text{Bernoulli}(P_{i1}), Y_{i2} \sim \text{Bernoulli}(P_{i2}), Y_{i3} \sim \text{Bernoulli}(P_{i3})$$

$$\begin{aligned} \log \frac{P_{i1}}{1 - P_{i1}} &= a_0 + \sum_{j=1}^K a_j O_{ij} + \sum_{j=1}^K a_{K+j} D_{ij} \\ \log \frac{P_{i2}}{1 - P_{i2}} &= b_0 + \sum_{j=1}^K b_j O_{ij} + \sum_{j=1}^K b_{K+j} D_{ij} \\ \log \frac{P_{i3}}{1 - P_{i3}} &= c_0 + \sum_{j=1}^K c_j O_{ij} + \sum_{j=1}^K c_{K+j} D_{ij} \end{aligned}$$

where  $Y_{i1}$ ,  $Y_{i2}$ , and  $Y_{i3}$  is the response random variable of each model which take value one when the attacking team scores one, two and three points at possession  $i$  respectively and zero otherwise,  $P_{i1}$ ,  $P_{i2}$ , and  $P_{i3}$  is the probability of scoring one, two and three points in  $i$  possession respectively,  $O_j$ , for  $j=1, \dots, K$ , takes the value one, when player  $j$  plays offense in  $i$  possession, and 0, otherwise,

$D_j$ , for  $j=1, \dots, K$ , takes the value minus one, when player  $j$  plays defense in  $i$  possession, and zero, otherwise.

Let us consider the offensive and defensive 5-lineups of Golden State Warriors and Boston Celtics, respectively, as an example. When Steph Curry, Klay Thompson, Draymond Green, Andrew Wiggins and Kevon Looney are playing offense and Jason Tatum, Marcus Smart, Jaylen Brown, Robert Williams III and Al Horford are playing defense, the probability of Warriors to score one, two or three points is 4.61%, 28.53% and 22.22% respectively. If Kevon Looney substitute by Otto Porter Jr. the probability of Warriors scoring one point is increased by 0.26% (4.87%) and for three points is decreased by 4.89% (23.64%).

Only 37 players out of the 970 in total (28 offenders and 9 defenders, see Tables 4.4 and 4.5) were found with non-zero Lasso RAPMs for all three binomial fitted models. Specifically, these players can be thought as players with contribution (positive or negative) in all types of scoring.

Table 4.4: Common statistically important defenders to the three separate binomial models.

PLAYERS	Position	log-odds of 1 pt per possession probability	log-odds of 2 pts per possession probability	log-odds of 3 or more pts per possession probability	G	MP
CHI Tony-Bradley	C	0.085	-0.147	<b>0.206</b>	55	549
CLE Jarrett-Allen	C	-0.317	-0.327	<b>0.192</b>	56	1809
DAL Davis-Bertans	PF	0.116	0.083	-0.074	22	306
DAL Kristaps-Porzingis	PF	<b>0.358</b>	<b>0.131</b>	0.135	34	1002
DEN Austin-Rivers	SG	0.059	-0.125	0.110	67	1480
DET Cory-Joseph	PG	-0.146	<b>0.149</b>	-0.060	65	1600
LAC Ivica-Zubac	C	-0.885	-0.151	-0.214	76	1852
MIA Bam-Adebayo	C	<b>0.466</b>	0.115	-0.099	56	1825
POR Keljin-Blevins	SF	0.161	-0.062	-0.252	31	349

Among these discriminated the following ones for their top-class offensive quality. The superstars Anthony Davis, Giannis Antetokounmpo, Ja Morant, Joel Embiid and, Tobias Harris (a very good power forward) appear in the relevant list. Among defenders with non-zero RAPMs three high quality Centers are included: Bam Adebayo, Ivica Zubac and Jarrett Allen, see Figure 4.3 for the corresponding RAPMs.

Giannis Antetokounmpo has the highest contribution in the two offensive situations, with the odds of scoring one point and two points per possession (against not scoring) is about 54% and 23%,



respectively, higher when he plays than the case he does not play. Also, his appearance in the court increases the odds of scoring three or more points per possession by 25% (vs. not scoring).

Table 4.5: Common statistically important offenders to the three separate binomial models.

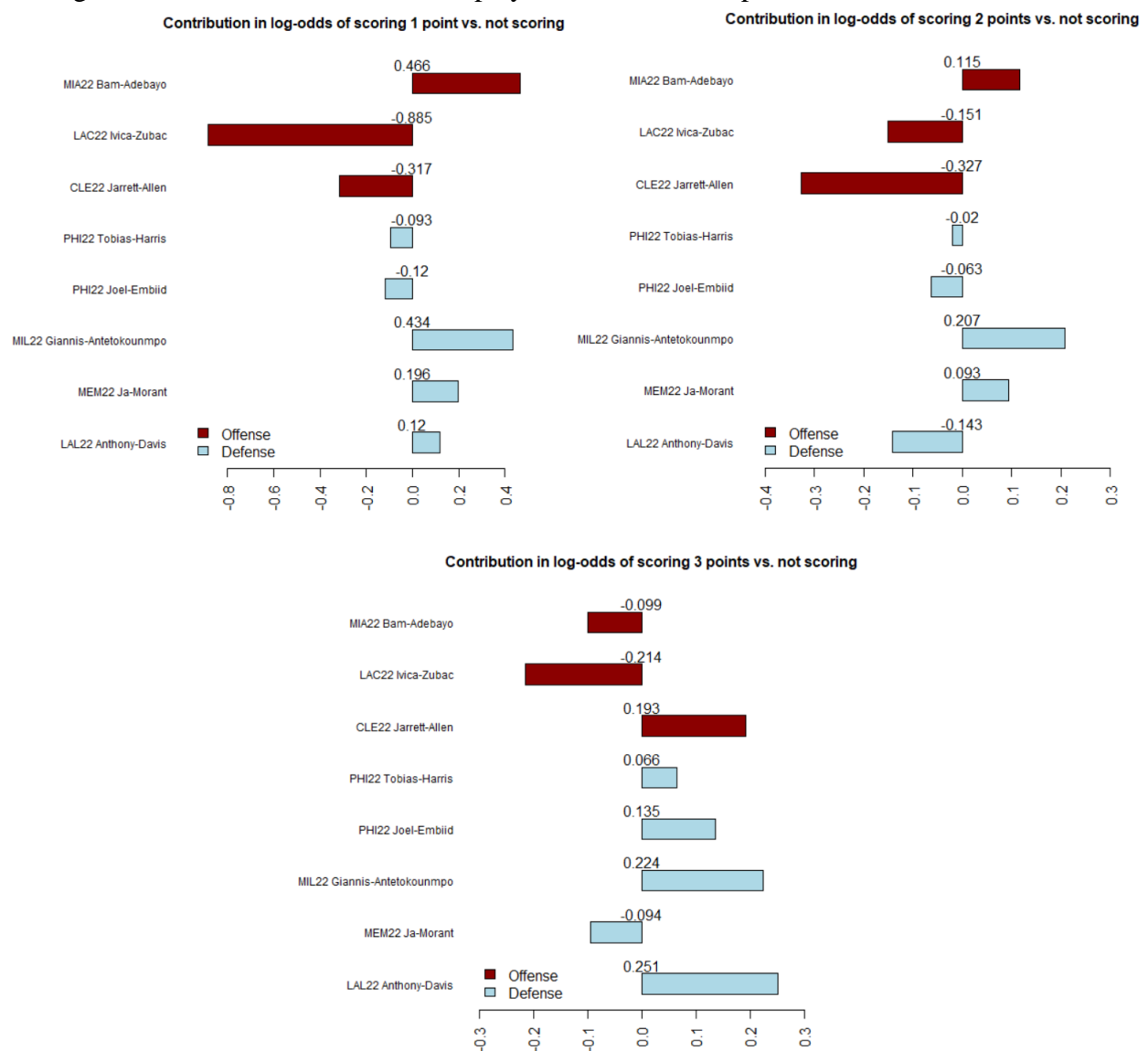
PLAYERS	Position	log-odds of 1 pt per possession probability	log-odds of 2 pts per possession probability	log-odds of 3 or more pts per possession probability	G	MP
PHI Tobias-Harris	PF	-0.092	-0.020	0.065	73	2543
NOP Herbert-Jones	PF	-0.267	-0.060	-0.038	78	2335
PHI Joel-Embiid	C	-0.120	-0.063	0.135	68	2297
MIL Giannis-Antetokounmpo	PF	<b>0.434</b>	<b>0.207</b>	<b>0.224</b>	67	2204
MEM Ja-Morant	PG	0.196	0.093	-0.094	57	1889
NYK Immanuel-Quickley	PG	0.226	0.036	0.162	78	1802
HOU Garrison-Mathews	SG	0.228	<b>0.155</b>	-0.228	65	1712
OKC Josh-Giddey	SF	0.398	-0.081	<b>0.177</b>	54	1700
WAS Daniel-Gafford	C	-0.173	-0.152	-0.217	72	1444
LAL Anthony-Davis	C	0.120	-0.143	<b>0.251</b>	40	1404
LAC Isaiah-Hartenstein	C	0.335	-0.077	-0.121	68	1216
UTA Joe-Ingles	SF	0.055	0.048	0.176	45	1122
NOP Garrett-Temple	SG	-0.389	-0.192	-0.064	59	1098
BOS Payton-Pritchard	PG	-0.324	0.101	0.074	71	1001
MIN Jaylen-Nowell	SG	0.471	0.119	0.085	62	975
CHI Derrick-Jones Jr.	PF	0.221	-0.104	-0.096	51	899
LAL Wayne-Ellington	SG	-0.356	-0.050	0.055	43	810
ATL Cam-Reddish	SF	0.114	-0.071	-0.321	34	797
HOU Daniel-Theis	C	-0.351	-0.149	0.118	26	584
NYK Jericho-Sims	PF	0.073	0.083	0.053	41	555
MIN Josh-Okogie	SG	<b>0.517</b>	0.153	-0.183	49	516
OKC Lindy-Waters III	SG	0.173	0.091	-0.113	25	465
BKN Joe-Harris	SF	0.209	-0.022	0.168	14	423
CHI Tristan-Thompson	PF	0.108	<b>0.168</b>	-0.033	23	376
NOP Kira-Lewis Jr.	PG	-0.542	0.118	-0.173	24	341
LAL Rajon-Rondo	PG	<b>0.800</b>	0.083	<b>0.184</b>	18	289
HOU Usman-Garuba	PF	0.293	0.144	0.098	24	239
CHI Tyler-Cook	PF	-0.111	-0.162	-0.193	20	200

It is worth saying that Anthony Davis appears for the first time with a positive effect (one-point and three-points scoring situations). The odds probability of Lakers scoring three or more points per possession increased by 28.5% when Davis is in the offensive lineup. However, he has a negative effect on the log-odds of two points, which means that when he plays the odds of scoring two points vs. not scoring is lower 13.3% compared to the case that a “reference” player is included

in the lineup instead of the Davis (for easier interpretation we assume that the studied player does not play in this situation).

Moreover, Joel Embiid affects negatively but in a less severe way the odds of scoring one or two points per possession (lower by 11.3% and 6.1% respectively). On the other hand, he has a positive impact on the odds of three or more points increase by 14.4%. Almost the opposite is true of Ja Morant. His one-point RAPM when he is playing is about 0.2, twice as high as the corresponding RAPM for the two-points. For the three points his contribution is negative and at the same magnitude as the two-points RAPM (in absolute value).

Figure 4.3: Contribution of selected players for the three separate Lasso Binomial models.



As for the defenders (presented in the Figure 4.3), Bam Adebayo was one of the best defenders of the 2021-2022 NBA season according to the two best defensive 5-lineups. His contribution in defense according to the three separate binomial models is analyzed as follows: When he does not play on defense, the odds of the opponent to score one or two points per possession is increasing by 59.6% and 12.2%, respectively. On the other hand, when Adebayo is not on the court for his team's defense, the odds of three point scored by the opponent team to score is decreasing by 9.4%.

Moreover, Jarrett Allen's defensive effect is exactly the opposite of Bam Adebayo's. His role and contribution in defense are not so remarkable. His total plus/minus in the season of the study is 2.4 according to the NBA official site. At the same time, Ivica Zubac's contribution in defense is worse according to the models' results. His coefficient effect was not shrunk to zero through the Lasso, because this player does not belong to the mediocre/average players, but probably to the one below average, according to our fitted Lasso model.

Generally, the non-zero Lasso RAPM players can be divided into players with positive or negative contributions. Indicatively, Immanuel Quickley had one of his best seasons and Josh Giddey is the first rookie with the most triple doubles in one season in NBA. Also, Isaiah Hartenstein is a high-quality center who had a basic role in Los Angeles Clippers offensive plan. Nevertheless, according to the model, his effects on log-odds of scoring two and three points per possession was found negative.

In addition, some low time players have non-zero RAPMs for the three models. Rajon Rondo of the Lakers played about 16 minutes (average time) in 18 games. His contribution is positive for all types of scoring when his team is attacking. Let us now take the opportunity to explain in more detail what RAPM in terms of contribution express. Rajon Rondo (with 3.1 points per game, 15.7% usage of Lakers, -0.3 Win Share) was not as crucial as LeBron James for Lakers (with 56/56 games started, 37 minutes per game, 30.3 points per game, 32.3% usage of Lakers and 7.5 Win Share). Players with less playing time usually are those that play in specific situations and when the points difference between the opponents is large and therefore the winning team is almost finalized. Thus, anyone could say that in the case of Rondo, the effect on the Lakers' offense is not reflecting his actual contribution. Hence, we need to consider a metric that takes into consideration time since a high-level player will affect his team (whatever his effect is) due to his increased playing time.

Regarding the defensive contribution, we observe that seven common players of the binomial models out of nine are “front-court” players (Centers and Power Forwards) while the five have positive effect. This is something close to real facts as the best defenders are usually Centers or Power Forwards. This is validated through some external criteria. Evidently, for the 2021-2022 NBA season six out of 10 of the best two-defensive lineups are Centers or Power Forwards. Also, using boxscores statistics, seven of the 10 highest Defensive Win Share players and 13 of the 25 highest Defensive Rating players are Power Forwards or Centers.

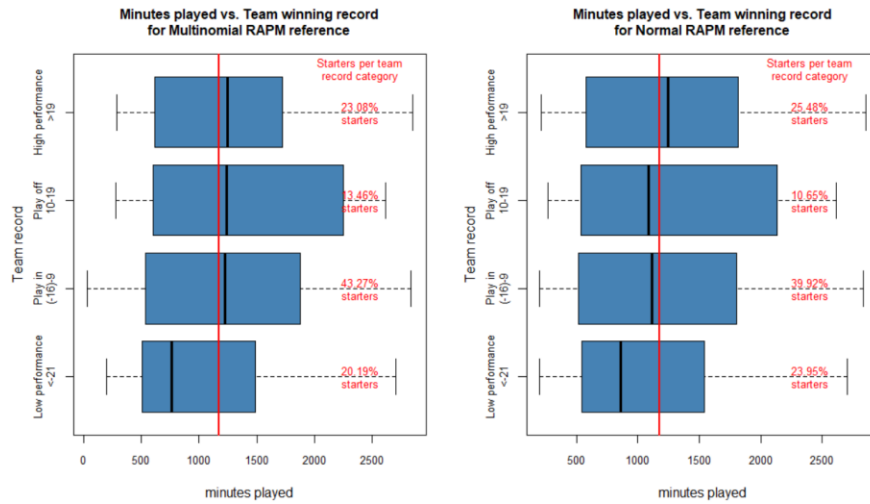
The profile of zero-RAPM Lasso (for all the three binomial models) is reasonable. Specifically, the “reference” group is composed by:

- Players of the lower performance playoffs teams (the lower winning record teams which proceeded to playoffs) with more minutes played. From the reference players group 43.3% and 39.9%, for Multinomial and Normal respectively, played for teams with the second lower winning record (Figure 4.4). Players of this winning record level teams also are 40.8% of those scorers with more than 10 points per game (Figure 4.5).
- Players of the lower winning record teams (more than 20 losses than wins) with less playing time. This group of players are the 20.19% and 23.95% of the starters included in the reference group (Figure 4.4). Players of this winning record level teams also are 36.7% of those scorers with more than 10 points per game (Figure 4.5).

In overall the zero-Lasso group includes the average players of the teams who were not responsible for the high or low performance of the teams. For the higher winning record teams, we can see players with significant playing time but low scorers. On the other hand, for the lower winning teams players with less playing time and average scorers were found.

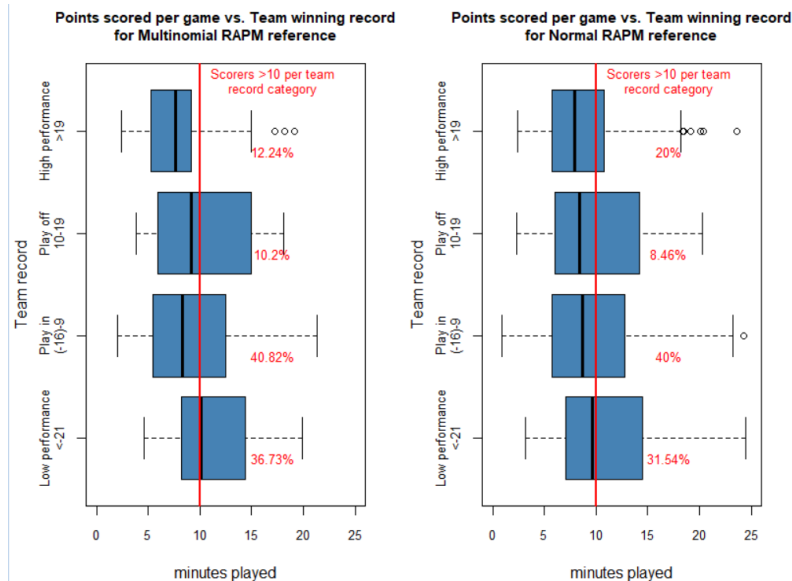
In the defenders’ reference group, some of the top NBAers are also included. For example, James Harden (in Philadelphia and Brooklyn), Nikola Jokic, Ja Morant, Damian Lillard, Russell Westbrook, Devin Booker, Kevin Durant, and more Joel Embiid. All of them are known for their offensive skills and contribution, and the Lasso RAPM confirms that their defensive contribution is on average level. That, also, explains the highest values of minutes played for the reference group of players.

Figure 4.4: Minutes played for reference group of the Multinomial model per team winning record level\*.



\* Play offs teams: winning record >10, Play in winning record >-16

Figure 4.5: Points scored for reference group of the Multinomial model per team winning record level\*.



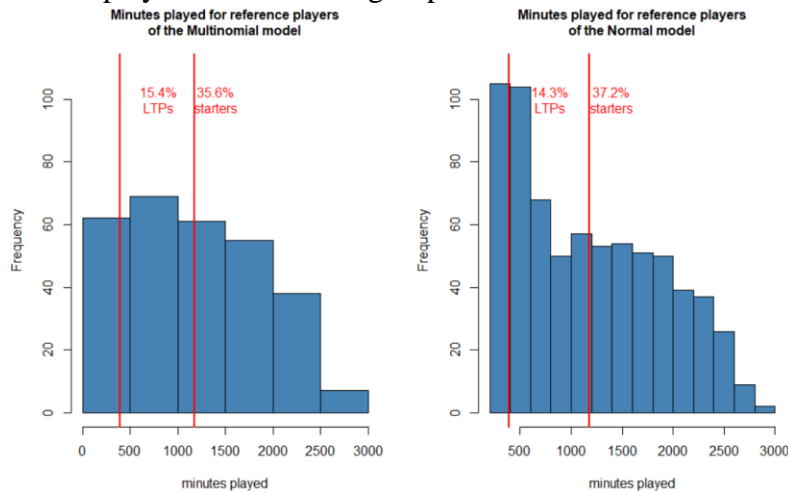
\* Play offs teams: winning record >10, Play in winning record >-16

One difference between the reference group of Multinomial and Normal zero-RAPM players is founded in the points per game of the offensive ones. In the first group, players score less than the second, where we can see some with more than 20 points. Also, from the higher-performed teams (with >20 wins than losses) higher percentage of reference players belong to scorers with more

than 10 points per game (20% for Normal against 12.24% for Multinomial). Nevertheless, for both groups, it is obvious that players from low-performed teams score more than the other teams.

Figure 4.6 summarizes all the above findings. More than half of the Multinomial reference players (64.82%) played less than 1500 minutes in the season (and 65.9% for the zero-Lasso Normal). Moreover, 22.7% of all NBA starters (according to the definition of Section 3.5.2, first six players per team in playing time) belong to the Multinomial reference group players respectively. On the other hand, the Normal reference group includes much more starters since 57.4% of the NBA 2021-2022 season starters have zero-Lasso Normal RAPM. In addition, 45 and 101 low-time players (LTPs based on the 50 lowest playing time) appear in the reference group of Multinomial and Normal respectively.

Figure 4.6: Minutes played for reference group of the Multinomial and Normal models.



Finally, regarding the estimation of the reference players impact (intercept), the probability of the attacking team to score when all players in the court are from this group: i) one point per possession is about 4%, ii) two points per possession is about 42%, and iii) three or more points per possession is about 19%.

After getting all initial idea about the results of the three separate binomial models, which compose the Multinomial, Section 4.4 discusses the final model implementation of our analysis.

## 4.4 Multinomial model

### 4.4.1 Building the multinomial model via binomial models

The goal of the thesis is to propose an improved model for the estimation of players contribution as an extension of the RAPM, which is used in the basketball literature. As we mentioned, a Multinomial logistic distribution is more appropriate to model the discrete response variable of interest which is the number points per possession (with values from zero to three). This Multinomial model is implemented through the fit of three binomial models as described in the following  $4 \times 4$  system of equation:

$$\begin{cases} \frac{P_{i1}}{P_{i0}} = \exp(a_0 + \sum_{j=1}^K a_j O_{ij} + \sum_{j=1}^K a_{K+j} D_{ij}) = L_{i1} \\ \frac{P_{i2}}{P_{i0}} = \exp(b_0 + \sum_{j=1}^K b_j O_{ij} + \sum_{j=1}^K b_{K+j} D_{ij}) = L_{i2} \\ \frac{P_{i3}}{P_{i0}} = \exp(c_0 + \sum_{j=1}^K c_j O_{ij} + \sum_{j=1}^K c_{K+j} D_{ij}) = L_{i3} \\ P_{i1} + P_{i2} + P_{i3} = 1 \end{cases}$$

where  $L_{i1}$ ,  $L_{i2}$  and  $L_{i3}$ , are the link functions of the three separate binomial models, for all  $i=1, \dots, n$ . The solution of the above system offers the scoring probabilities, and it is the following:

$$\begin{cases} P_{i0} = \frac{1}{1+L_{i1}+L_{i2}+L_{i3}} \\ P_{i1} = \frac{L_{i1}}{1+L_{i1}+L_{i2}+L_{i3}} \\ P_{i2} = \frac{L_{i2}}{1+L_{i1}+L_{i2}+L_{i3}} \\ P_{i3} = \frac{L_{i3}}{1+L_{i1}+L_{i2}+L_{i3}} \end{cases}$$

where  $P_{i1}$ ,  $P_{i2}$ , and  $P_{i3}$  is the probability of scoring one, two and three points in  $i$  possession respectively. Since we are interesting to obtain a single RAPM (instead of three), expected points per possession, for all players  $j=1, \dots, 485$  are calculated by using the scoring probabilities:

$$EPTS_{ij} = P_{ij1} + 2 \times P_{ij2} + 3.01 \times P_{ij3}$$

where  $EPTS_{ij}$  is the expected points scored from: 1) the team of a player  $j$  (for  $j=1, \dots, 485$ ) in  $i$  possession, for all  $i=1, \dots, n$ , when the  $j$  offensive player is on the court while his teammates and opponents are players from the reference group (intercept), 2) the opponent team of  $j$  player in  $i$

possession, for all  $i=1, \dots, n$ , and  $j=1, \dots, 485$ , when the  $j$  defensive player is not on the court while players from the reference group (intercept) play. The mean of three or more points scored (3.01) is used for the last term to account also the limited possessions of four to six points where scored. This value is simple the average points for all possessions with points more or equal to three. Finally, the probabilities are given by the following equations:

$$\begin{cases} P_{ij1} = \frac{\exp(a_0 + a_j O_{ij})}{1 + \exp(a_0 + a_j O_{ij}) + \exp(b_0 + b_j O_{ij}) + \exp(c_0 + c_j O_{ij})} \\ P_{ij2} = \frac{\exp(b_0 + b_j O_{ij})}{1 + \exp(a_0 + a_j O_{ij}) + \exp(b_0 + b_j O_{ij}) + \exp(c_0 + c_j O_{ij})} \\ P_{ij3} = \frac{\exp(c_0 + c_j O_{ij})}{1 + \exp(a_0 + a_j O_{ij}) + \exp(b_0 + b_j O_{ij}) + \exp(c_0 + c_j O_{ij})} \end{cases} \cdot \begin{cases} P_{ij1} = \frac{\exp(a_0 + a_j D_{ij})}{1 + \exp(a_0 + a_j D_{ij}) + \exp(b_0 + b_j D_{ij}) + \exp(c_0 + c_j D_{ij})} \\ P_{ij2} = \frac{\exp(b_0 + b_j D_{ij})}{1 + \exp(a_0 + a_j D_{ij}) + \exp(b_0 + b_j D_{ij}) + \exp(c_0 + c_j D_{ij})} \\ P_{ij3} = \frac{\exp(c_0 + c_j D_{ij})}{1 + \exp(a_0 + a_j D_{ij}) + \exp(b_0 + b_j D_{ij}) + \exp(c_0 + c_j D_{ij})} \end{cases} .$$

Finally, for all players' expected points are calculated as the average points considering all the possessions of each team and the possessions in which the player was on the court (a kind of efficiency). Hence, the final metric is:

$$EPTS_j = P_{j1} + 2 \times P_{j2} + 3.01 \times P_{j3}$$

where  $EPTS_j$  is the expected points scored from: 1) the team of a player  $j$  (for all  $j=1, \dots, 485$ ) per possession when the  $j$  offensive player is on the court while his teammates and opponents are players from the reference group (intercept), 2) the opponent team of  $j$  player per possession (for all  $j=1, \dots, 485$ ), when the  $j$  defensive player is not on the court while players from the reference group (intercept) play. The probabilities are described as following:

$$\begin{cases} P_{j1} = \frac{\exp(a_0 + a_j)}{1 + \exp(a_0 + a_j) + \exp(b_0 + b_j) + \exp(c_0 + c_j)} \\ P_{j2} = \frac{\exp(b_0 + b_j)}{1 + \exp(a_0 + a_j) + \exp(b_0 + b_j) + \exp(c_0 + c_j)} \\ P_{j3} = \frac{\exp(c_0 + c_j)}{1 + \exp(a_0 + a_j) + \exp(b_0 + b_j) + \exp(c_0 + c_j)} \end{cases} .$$

For the above calculations of scoring probabilities and expected points, the results of the three separate binomial models no team effects were included, since there were not found to be significant. In the final Multinomial model, 362 offensive and 316 defensive players were actively contributed. These were players with non-zero RAPM at least one of the binomial components.



#### 4.4.2 Results of EPTS-RAPM

The interpretation of the Multinomial model results appears obviously with more difficulties than the normal, but it is not such complicated:

- The more the expected points per possession when a player is on the offensive lineup, the higher his contribution.
- For the defenders, the lower the probability for a team to concede points, the higher the contribution. But if we would like to use the same measure for offensive and defensive contribution, the higher expected points per possession the opposite teams score when a specific player does not play defense, the higher his contribution. The ranking of players based on their contribution is about the same, thus anyone can use with safety one of the two procedures.

The player's contribution is estimated based on his team possessions and not on the total possessions of the season (see Table 4.25 in Appendix II). In addition, expected points and scoring probabilities per player are calculated considering the player's possessions on the court. This is a kind of efficiency contribution (from now on we call it "efficiency"), which expresses the expected points per possession of a team according to each player's participation.

Firstly, offenders and defenders with the higher contribution based on expected points per possession according to their teams' total possessions and players' attendance in offense and defense as well, are presented in tables 4.19 and 4.20 for each position.

According to the top five per position for expected points, the results are very close to real facts. There are some of the best players in the best offenders (noted in bold, table 4.5). For the defenders, there is not anything special to mention, except for some players like Bam Adebayo or Jason Tatum who are impactful for their teams' defense.

A notable observation, here, is that considering the efficiency (i.e. expected points per possession based on the participation of the players and not his team possessions), not high-quality players seem better than others. This metric is more useful for comparing mediocre or worse players with only a few minutes playing in the season.

This phenomenon (lower quality players raise the ranking due to their efficiency-contribution) is not observed in the same degree for the best players per team (Table 4.8).

Table 4.6a: Top-5 offenders per position via multinomial model's results for expected points per possession.

Rank	Shooting Guard (SG)	Small Forward (SF)	Point Guard (PG)	Power Forward (PF)	Center (C)
1	DAL Jalen-Brunson	PHX Mikal-Bridges	CLE Darius-Garland	BKN Kevin-Durant	DEN Nikola-Jokic
2	UTA Donovan-Mitchell	MIA Jimmy-Butler	ATL Trae-Young	CHI DeMar-DeRozan	MEM Steven-Adams
3	ATL Bogdan-Bogdanovic	ORL Franz-Wagner	MIL Jrue-Holiday	DEN Aaron-Gordon	WAS Daniel-Gafford
4	CHI Zach-LaVine	BOS Jayson-Tatum	MEM Ja-Morant	CHA Miles-Bridges	LAC Ivica-Zubac
5	PHX Devin-Booker	NOP Brandon-Ingram	OKC Shai-Gilgeous-Alexander	PHI Tobias-Harris	PHX Deandre-Ayton

Table 4.7a: Top-5 defenders per position via multinomial model's results for expected points per possession.

Rank	Shooting Guard (SG)	Small Forward (SF)	Point Guard (PG)	Power Forward (PF)	Center (C)
1	BKN Patty-Mills	NOP Brandon-Ingram	MIL George-Hill	DAL Dorian-Finney-Smith	NOP Jonas-Valanciunas
2	LAL Malik-Monk	BOS Jayson-Tatum	OKC Shai-Gilgeous-Alexander	POR Robert-Covington	CLE Kevin-Love
3	CHI Ayo-Dosunmu	GSW Andre-Iguodala	PHI Tyrese-Maxey	MIN Jarred-Vanderbilt	LAC Ivica-Zubac
4	ORL Gary-Harris	CHI Troy-Brown Jr.	OKC Tre-Mann	LAC Nicolas-Batum	MIA Bam-Adebayo
5	MEM Desmond-Bane	ATL Cam-Reddish	NYK Kemba-Walker	MIA P.J.-Tucker	ORL Mo-Bamba

Table 4.6b: Top-5 offenders per position via multinomial model's results for efficiency (expected points per possession based on those possessions each player was on the court).

Rank	Shooting Guard (SG)	Small Forward (SF)	Point Guard (PG)	Power Forward (PF)	Center (C)
1	DAL Jalen-Brunson	PHX Mikal-Bridges	WAS Brad-Wanamaker	BKN Kevin-Durant	WAS Montrezl-Harrell
2	DEN Bones-Hyland	MIA Jimmy-Butler	CLE Darius-Garland	HOU Usman-Garuba	CHA Montrezl-Harrell
3	MIN Jaylen-Nowell	ORL Franz-Wagner	POR Damian-Lillard	DEN Aaron-Gordon	MIL DeMarcus-Cousins
4	SAC Terence-Davis	MEM Dillon-Brooks	POR Kris-Dunn	CHI DeMar-DeRozan	WAS Daniel-Gafford
5	IND Buddy-Hield	BOS Jayson-Tatum	MEM Ja-Morant	LAL Wenyen-Gabriel	NOP Willy-Hernangomez

Table 4.7b: Top-5 defenders per position via multinomial model's results for efficiency (expected points per possession based on those possessions each player was on the court).

Rank	Shooting Guard (SG)	Small Forward (SF)	Point Guard (PG)	Power Forward (PF)	Center (C)
1	BKN Patty-Mills	GSW Andre-Iguodala	MIL George-Hill	POR Robert-Covington	CLE Ed-Davis
2	ATL Skylar-Mays	DEN Michael-Porter Jr.	GSW Chris-Chiozza	DAL Dorian-Finney-Smith	NOP Jonas-Valanciunas
3	CLE Denzel-Valentine	ATL Cam-Reddish	NYK Kemba-Walker	LAC Nicolas-Batum	CLE Kevin-Love
4	LAL Malik-Monk	NOP Brandon-Ingram	OKC Tre-Mann	LAL Stanley-Johnson	MIL Brook-Lopez
5	ORL Gary-Harris	CHI Troy-Brown Jr.	OKC Shai-Gilgeous-Alexander	LAC Paul-George	LAC Ivica-Zubac

In the highest offensive rated players based on EPTS-RAPMs for each team, a number of top NBAers of the season 2021-2022 appear (such as Trae Young for Atlanta, Kevin Durant for Brooklyn, Jason Tatum for Boston, DeMar DeRozan for Chicago, Steph Curry for Golden State, Russell-Westbrook for Lakers, Ja Morant for Memphis, Jimmy Butler for Miami, Jrue Holiday for Milwaukee, Damian Lillard for Portland, Pascal Siakam for Toronto and Donovan Mitchell for Utah). In fact, 90% (27/30) of the highest EPTS-RAPM-rated offensive players for each team are starters and 40% (12/30) of these are the highest-time players in their teams.

As for the defenders, most of the top-rated players are average ones or role players, whose participation in defense appears to have a positive impact on their team, with their opponents scoring less points per possession. That is explained by the fact that 70% (21/30) of the highest-rated EPTS-RAPMs are starters and the same time only 20% (6/30) of these players are the highest time played for each team.

Regarding the above observations, in Figure 4.7 the percentage of agreement between the highest EPTS-RAPMs for each team and: i) starters, ii) highest-time players per team (High-MP), iii) highest points scored (High-PTS) and defensive rebounds (High-DRB) for the offensive and defensive players respectively.

Figure 4.7: Minutes played for reference group of the Multinomial and Normal models.

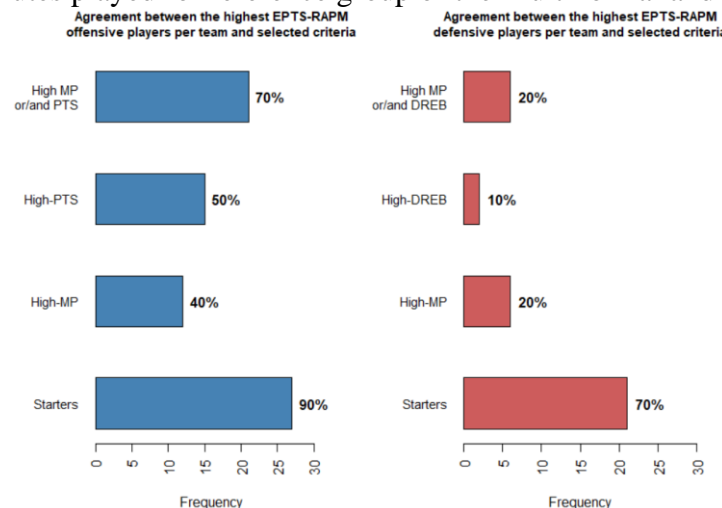


Table 4.8: Best players per team via EPTS-RAPMs.

TEAM	OFFENSE	DEFENSE	TEAM	OFFENSE	DEFENSE
ATL	Trae-Young	Cam-Reddish (pts) Skylar-Mays (efficiency)	MIA	Jimmy Butler	P.J.-Tucker
BKN	Kevin Durant	Patty Mills	MIL	Jrue Holiday (pts) DeMarcus Cousins (efficiency)	George-Hill
BOS	Jason Tatum	Jason Tatum	MIN	D'Angelo-Russell (pts) Jaylen-Nowell (efficiency)	Jarred-Vanderbilt
CHA	Miles-Bridges (pts) Montrezl Harrel (efficiency)	LaMelo Ball	NOP	Herbert-Jones (pts) Willy Hernangomez (efficiency)	Jonas-Valanciunas
CHI	DeMar DeRozan	Ayo-Dosunmu (pts) Troy-Brown Jr. (efficiency)	NYK	Immanuel-Quickley	Kemba-Walker
CLE	Darius-Garland	Kevin Love (pts) Ed Davis (efficiency)	OKC	Shai-Gilgeous-Alexander	Shai-Gilgeous-Alexander
DAL	Jalen-Brunson	Dorian-Finney-Smith	ORL	Franz-Wagner	Gary-Harris
DEN	Aaron-Gordon	Monte Morris (pts) Michael-Porter Jr. (efficiency)	PHI	Tobias-Harris	Tyrese-Maxey
DET	Isaiah-Livers	Josh-Jackson	PHX	Mikal-Bridges	Chris Paul (pts) Elfrid-Payton (efficiency)
GSW	Stephen-Curry	Andre-Iguodala	POR	Damian-Lillard	Robert-Covington
HOU	Usman Garuba	Josh-Christopher	SAC	De'Aaron-Fox (pts) Terence Davis (efficiency)	Buddie Hield (pts) Chimezie-Metu (efficiency)
IND	Caris LeVert (pts) Buddie Hield (efficiency)	Chris-Duarte	SAS	Dejounte-Murray (pts) Josh-Richardson (efficiency)	Doug-McDermott
LAC	Ivica Zubac (pts) Isaiah-Hartenstein (efficiency)	Ivica-Zubac	TOR	Pascal-Siakam (pts) Thaddeus-Young (efficiency)	OG-Anunoby
LAL	Russell-Westbrook (pts) Wanyen Gabriel (efficiency)	Malik-Monk	UTA	Donovan-Mitchell	Jared-Butler
MEM	Ja Morant	Desmond-Bane	WAS	Daniel-Gafford (pts) Brad-Wanamaker (efficiency)	Ish-Smith

For a limited number of teams (4/30), the top-rated player is not their star player or the best-performed player according to boxscore statistics. For those anyone expected. For example, the player with the highest expected points per possession of Milwaukee Bucks (MIL) is Jrue Holiday and not Giannis Antetokounmpo. Although Holiday is an important player for Milwaukee, Antetokounmpo is objectively considered as the most valuable. Such players usually have more responsibilities than their teammates, and many times their performance affects both wins and losses since they are presented in all crucial times of the games. In addition, Antetokounmpo makes the most decisions for the offense, as he is the basic player for Milwaukee. Most of the teams' stars

have a similar profile with Antetokounmpo. Similar is the case of Denver, Philadelphia, and Dallas, since their top players appear without having the highest contribution.

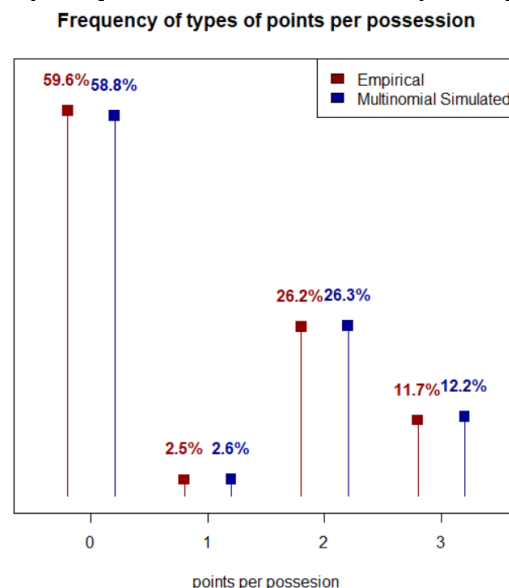
Regarding the three best 5-lineups for the NBA 2021-2022 season, according to the expected points per possession, the multinomial model identifies 8/15 players in the EPTS-RAPMs ratings top 6 per position. Again, studying the top three EPTS-RAPMs players per team identified 13/15 players of the top 5-lineups.

It is interesting and very useful to compare the multinomial EPTS-RAPMs with the corresponding RAPMs of the previous fitted models. More details are presented in the Chapter 5. However, this model has a statistical advantage over the Normal, which is the standard choice in the basketball literature.

To test the fit for the multinomial, the probabilities of scoring (from zero to three or more points) per possession were calculated as described in the corresponding equations in Section 4.4. Then, 1000 samples were simulated from a multinomial distribution with these probabilities per possession.

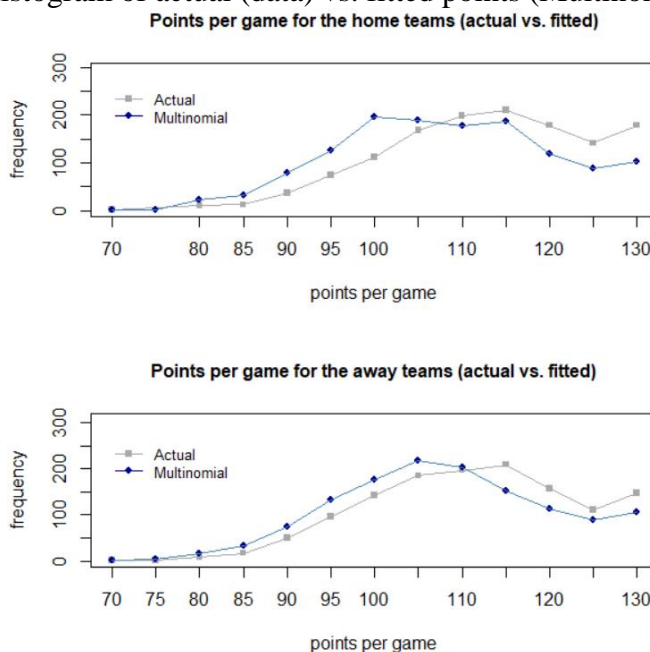
The relative frequencies of points per possession are very close to those of the actual data (Figure 4.8). The observed data do not differ significantly from the fitted ones (Bootstrap Chi-square test,  $p\text{-value} = 0.55 > 0.05 \Rightarrow$  the null hypothesis is not rejected), which means a good marginal fit of the model.

Figure 4.8: Frequency of actual vs. simulated points per possession.



The model fit was also visually tested in a more detailed level via the fitted points per game<sup>57</sup>. The fitted points were calculated by the probabilities to score one, two, and three or more points per possession for the Multinomial case. Points were then categorized into classes of five points range: <70, 70-75, 75-80, ..., 125-130, >130. As we observe in Figure 4.9, the fitted values for both home and away teams are close to observed ones.

Figure 4.9: Histogram of actual (data) vs. fitted points (Multinomial) per game.



## 4.5 Conclusion of Chapter 4

This chapter deals with the implementation of logistic regression models. The first step was a binomial model for the probability of scoring/conceding points against not scoring/conceding, when a player is in the offensive or in the defensive lineup. Regularization methods were applied and one of our most important findings is the high linear relationship between the Ridge Binomial and Ridge Normal RAPM ratings. This allows us to obtain Binomial RAPM by standard Ridge regression RAPMs, which are used in the relative bibliography. Moreover, it offers an alternative interpretation of the Ridge regression RAPMs.

<sup>57</sup> The playoffs weren't include because we hadn't the Game at the dataset and we fill it by hand.

Finally, we proposed a full Multinomial model which is composed by three separate binomial models. Due to this indirect approach, the scoring probabilities and the expected points were calculated based on the teams' possessions and the players' participation as well. Hence, we have obtained a new measure which we call it EPTS-RAPM.

Although the Multinomial model is more complicated than the simple regression RAPMs, the extra benefits we earn are the following:

- We have separate RAPMs and contribution per type of points.
- Simulated and fitted data are close to the actual, in contrast with the fitted values of the normal models which indicated a worse fit. Hence the model is more appropriate resulting in a better fitted approach, which explains the game in a more realistic way.

Therefore, the multinomial model offers a considerably improved and realistic evaluation approach. Moreover, there is statistical evidence that fits better the data than the Normal distribution model. To support more the superiority of this approach, we proceed in Chapter 5 with a comparison between the different RAPM methodology.

## Chapter 5: Comparison of the different RAPM Models

### 5.1 Introduction

In this chapter we compare the fit and predictive ability of the RAPM models presented in Chapters 3 and 4, i.e., the Ridge and Lasso regularization methods applied in the normal regression, binary logistic and multinomial models.

Moreover, we validate the RAPM ratings with some external objective criteria. In the Section 5.2 we focus on the comparison between the Multinomial and Normal models developed, while in the Section 5.3 the comparison includes all the implemented models and is based on external, game-oriented criteria.

### 5.2 Multinomial vs. Normal

Starting with a comparison between the Normal and Multinomial models, from Table 5.1 we observe that the RAPM ratings, are linearly correlated in a higher degree ( $r=0.82$ ) when each player's participation (efficiency) was considered (and not each teams possessions). The Ridge and Normal (after screening Lasso) RAPMs are higher correlated with the Multinomial ratings, while the Lasso RAPMs has obviously lower correlation with Multinomial RAPM (20.8% and 28.6%).

Table 5.1: Correlation for players RAPMs ratings of the multinomial and normal model\*.

	Multinomial	Multinomial (efficiency)
Lasso Normal	0.208	0.286
Normal (screening Lasso)	0.600	0.765
Ridge Normal	0.630	0.823

*\*All players were considered.*

Moreover, it is useful to study the fitted values for assessing the goodness of fit of each model. The sample mean points per 100 possessions is 89.9. At the same time the mean of expected points (fitted) of the:

- Multinomial model is 92.0.
- Ridge Normal model is 90.0.



- Lasso Normal model is 90.1.
- Normal model after screening Lasso is 89.9.

By developing the normal model for a train set (with 70% random observations of the full dataset) the in-sample and out-of-sample prediction for a test dataset (the rest 30% of the observations)<sup>58</sup> were studied. The RMSE of the train set was slightly smaller than the test set (1.145 and 1.147 respectively). Also, the RMSE between the fitted models of the constant normal model and the actual data points is quite close (equal to 1.147). In addition, the RMSE of the ridge and lasso regression estimations are almost equal to the normal after screening Lasso. For the Multinomial model, considering the fitted values to be the expected points per possession as a function of the probabilities per possession (like it was determined in Chapter 4), the RMSE is slightly higher (1.149).

For the Multinomial estimations, we studied the simulated samples (as mentioned in Chapter 4) using the estimated probabilities of scoring per possession in comparison with the constant model (using as probabilities the frequencies of points scored). Note that simulation based on the Normal model is not appropriate since negative values are generated but it is included in the comparisons.

The bootstrap Chi-squared test does not reject the null hypothesis of independence of the simulated and empirical sample (for Multinomial:  $p\text{-value}=0.55 > 0.05$ , for Normal:  $p\text{-value}=0.49 > 0.05$ ). Especially for the Normal simulated samples, the result of the Chi-squared test implies that the difference between data points and simulated is not large. In fact, this is confirmed since the majority of data points is zero and the mean of the simulated samples takes values in the range of  $(-0.4, 0.4)$ . Also, negative values do not affect the Chi-squared result due to the squared term. On the other hand, the same test rejects the null hypothesis for the ridge, lasso, and normal after screening lasso samples of fitted values against the data points sample ( $p\text{-value}<0.05$ ).

The RMSE of the fitted Multinomial model is slightly smaller than the constant ( $1.62 < 1.63$ ) while the RMSE of the Normal simulated samples is obviously higher than the Multinomial ( **$1.86 > 1.62$** ).

In addition, the percentage of agreement between the Multinomial simulated and data samples per points category is not high. The fitted model appears with higher percentage of agreement compared to the constant model for the three out of four types of points scored (see Figure 5.1 and

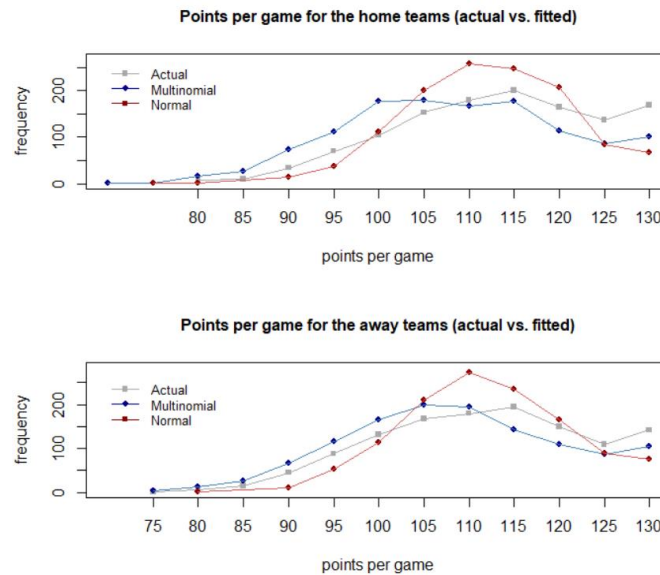
---

<sup>58</sup> As test set the payoffs observations were used and the result was the same.

Table 5.10 in Appendix III). These findings imply that the fitted multinomial has slightly better performance than the constant.

We, also, examine the fit of both interesting models (Multinomial and Normal) through the points scored per game, as it was described at Chapter 4. From Figure 5.1<sup>59</sup>, it is obvious that the Multinomial model is slightly closer to the actual data than the Normal (after screening Lasso). Additionally, if we do not classify points per game and compare the samples of Normal and Multinomial models with the empirical one, it is observed that the independence hypothesis is not rejected (Chi-squared test with p-value equal to 0.45 for Multinomial and 0.42 for Normal).

Figure 5.1: Actual (data) vs. fitted points (Multinomial & Normal) per game.



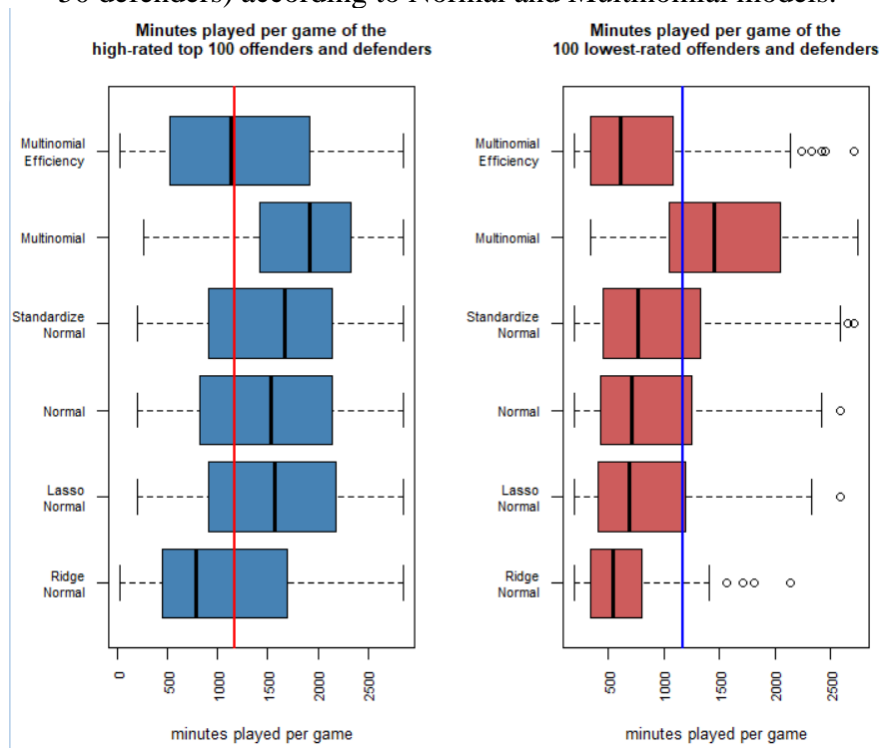
From Table 5.3 and Figure 5.2, we can confirm that the Multinomial model rank high time players in the top RAPM ratings (with 83% being starters) in a clearly better way than the Normal models (identified  $\leq 64\%$  starters). In addition, the 52% of the low-rated EPTS-RAPM players can be considered as starters according to their playing time. This is in contrast to 100 lowest rated Ridge RAPM players were only 7 starters are included (Table 5.3). Although, initially, this sounds as a surprising result, it seems that this model separates successfully high-time players according to their performance (Figure 5.3).

<sup>59</sup> The playoffs weren't included because we hadn't the Game at the dataset and we fill it by hand. Points are categorized into levels differ by 5 points: <70, 70-75, 75-80, ..., 125-130, >130

Table 5.3: Starters included in the 100 highest and lowest rated RAPM (50 offensive and 50 defensive) players per model.

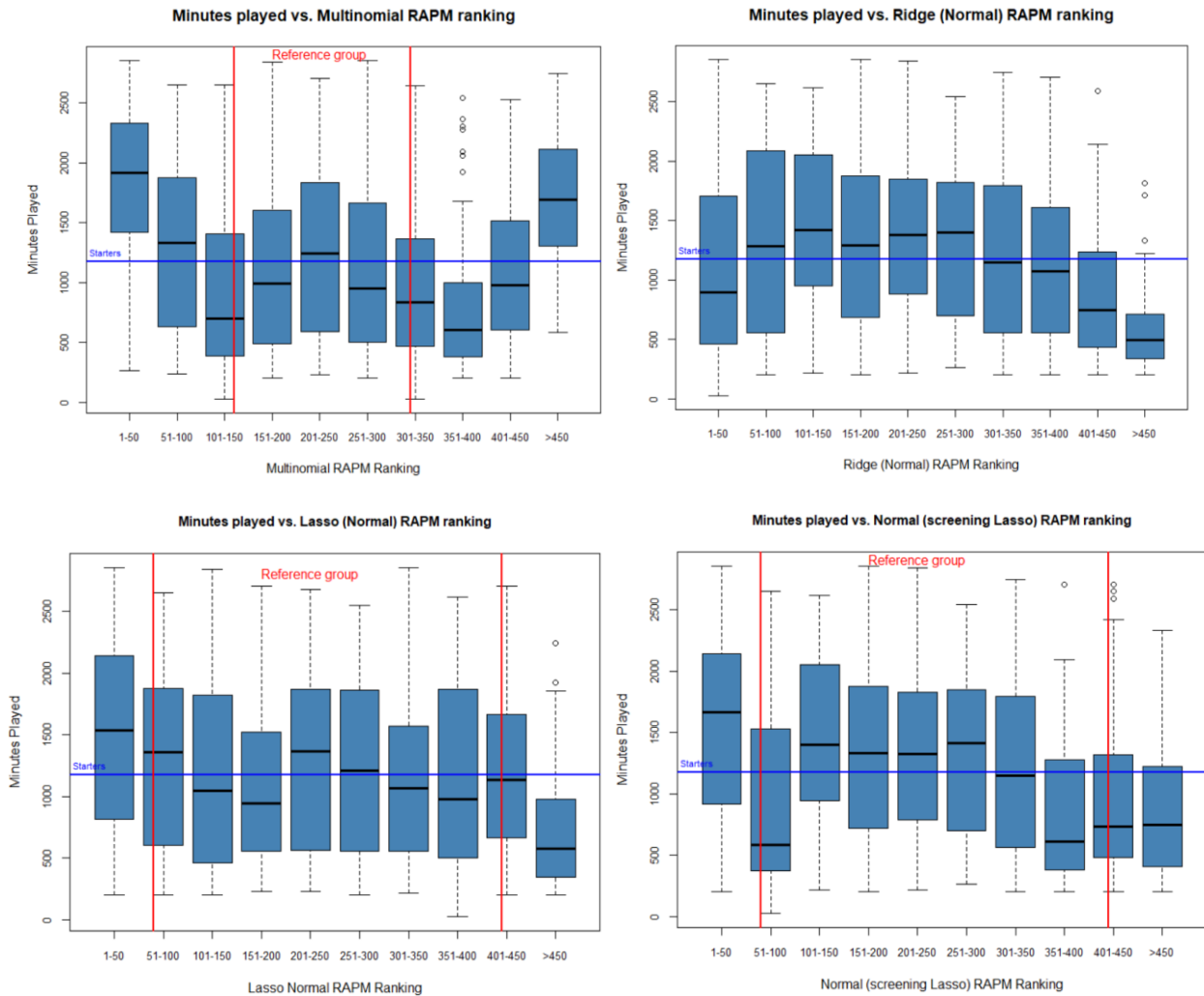
Model-based RAPM	100 high-rated	100 low-rated
Ridge Normal	27%	7%
Lasso Normal	50%	19%
Normal	51%	23%
Standardized Normal	64%	32%
Multinomial	83%	52%
Multinomial (Efficiency)	49%	17%

Figure 5.2: Minutes played of the 100 high (blue) and low (red) rated players (50 offenders and 50 defenders) according to Normal and Multinomial models.



On the other hand, players are separated in a different way with Ridge RAPM (Figure 5.3), where low time players are rated in the higher and lower ranking positions. Moreover, in Ridge-normal RAPMs, starters and higher-time players are ranked in the middle of the list which is not desirable. In the multinomial model, on the other hand, separate players low time players appear either in the reference level (with their RAPM shrunk to zero) or in the mid-ranking positions.

Figure 5.3: Minutes played per RAPM ranking position of Multinomial, Ridge Normal (first row), Lasso and Screening Lasso Normal (second row) models.

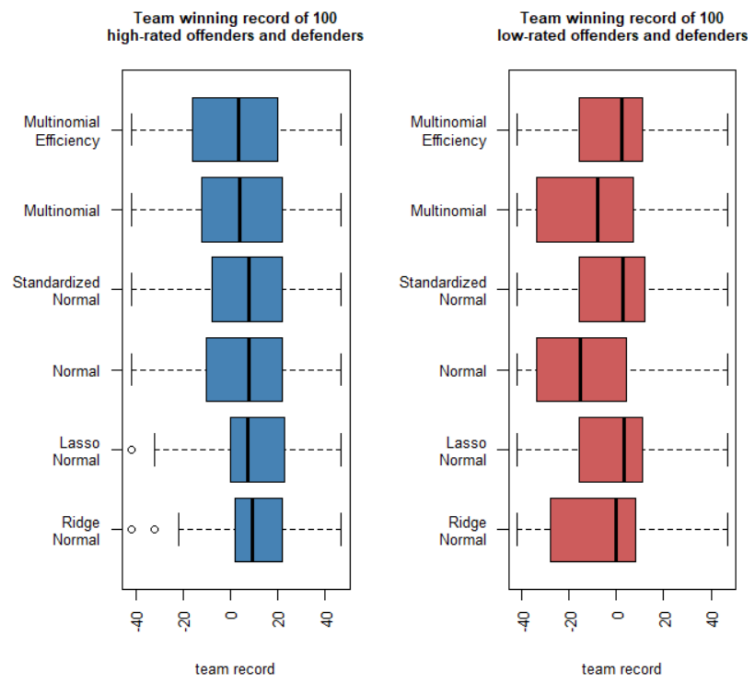


As an extension, in Figure 5.4 the team winning record top-100 (in blue) and bottom-100 (in red) rated players according to the implemented models are depicted. For the top players the mean and median values of the winning record is positive for all models, but the multinomial models appear with obviously larger range. This implies that the Multinomial recognizes players with positive or negative contribution, independently of their team winning record.

Furthermore, the points scored per game for the top-50 and bottom-50 players according to the model based RAPMs are presented in Figure 5.5. The 90% of the top Multinomial RAPM offensive players are scorers with more than 10 points. The corresponding performance is lower for the rest of the models RAPM ratings (82% for Normal after Lasso, 78% for Lasso Normal, 60% for Ridge Normal, and 72% for the Multinomial-Efficiency).

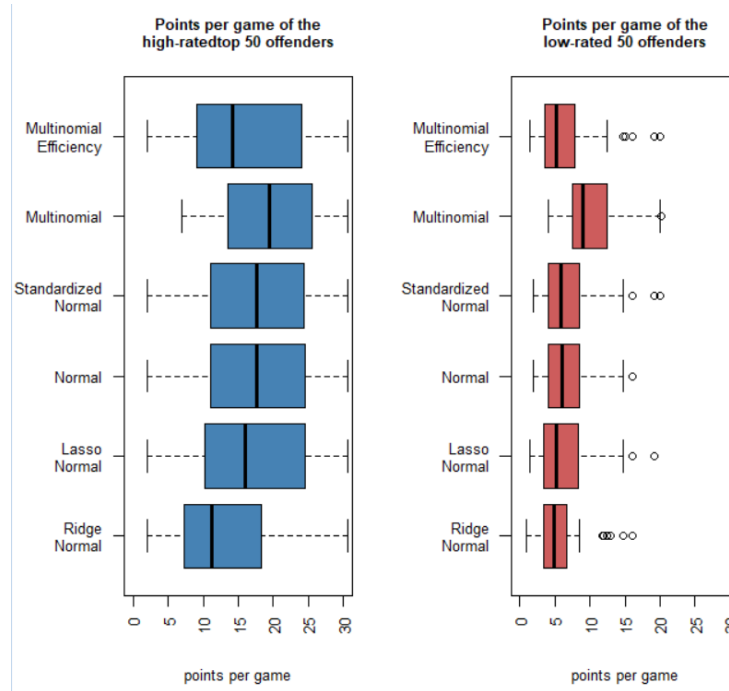
Regarding the lower RAPM rated players, their scoring is obviously lower than top rated players as well as the corresponding range of points (mostly from zero to 10 with the exception of the multinomial where the range is larger). Moreover, the 78% of the Multinomial low-rated players scored less than 15 points and the 44% more than 10 points per game, while the range of points for the low EPTS-RAPM ratings is larger than the RAPMs obtained by other models. This strengthens the hypothesis that is not sufficient to quantify the contribution of a player through points scored and playing time measurements.

Figure 5.4: Winning record of the top 100 (50 offenders and 50 defenders) and bottom 100 (blue and red respectively) rated players' teams according to Normal and Multinomial models.



Finally, the profile of zero-Lasso RAPM players for the Multinomial and Normal models was discussed in Chapter 4. According to our findings, the reference group of Normal includes higher number of starters than the Multinomial one (57.4% and 22.7% of all NBA starters respectively). This is clearly a disadvantage of the Normal-Lasso since many starters end up without essential rating squeezed in the reference group.

Figure 5.5: Points per game of the top-50 (blue) and bottom-50 (red) rated offenders according to Normal and Multinomial models.



### 5.3 Comparison of all method

In this section, a number of external validation criteria are used in order to study the efficiency of the models according to information that is not included in the data. Also, we will examine the expected points that are estimated through all implemented methods.

#### 5.3.1 Agreement between results and selected criteria

First, let us present the summary statistics of expected points per possession for each method. From Table 5.4 it is obvious that the discrimination of the players is not appropriate through the Lasso Normal and Normal after Lasso where the first quantile, the third quantile and the median are almost equal. Also, the median value of Multinomial, Lasso Normal and Normal after Lasso estimation is equal to the reference group (intercept) expected points.

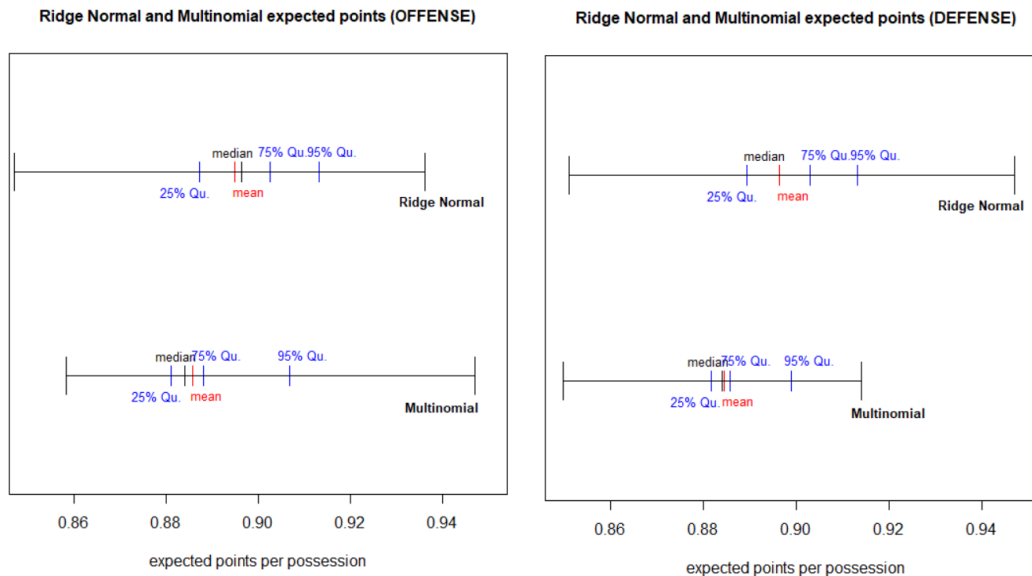
Moreover, while the range of the Multinomial estimations is higher than the Ridge, in the second case the standard deviation is higher than the first and the distribution is quite close to Normal (as it is expected). The defensive-expected points based on Multinomial model take values in a smaller range than the Ridge Normal, in contrast to the offensive estimations.

Also, the range from the third quantile to the maximum value is larger for the offensive Multinomial estimations than the scale of the 75% of the players contributions. This fact implies a more difficult discrimination for an important percentage of players. On the other hand, for the defensive contributions the distribution appears more symmetric. In Figure 5.6 the above observations are visualized.

Table 5.4: Summary statistics measures for expected points per possession.

Offensive							
Model	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	95% Qu.
Multinomial	0.858	0.881	0.884	0.886	0.888	0.947	0.907
Multinomial-Efficiency	0.778	0.868	0.884	0.884	0.899	1.001	0.939
Ridge Normal	0.847	0.887	0.896	0.895	0.902	0.936	0.913
Lasso Normal	0.820	<b>0.902</b>	<b>0.902</b>	<b>0.901</b>	<b>0.902</b>	0.949	0.921
Normal after Lasso	0.753	<b>0.888</b>	<b>0.888</b>	<b>0.886</b>	<b>0.888</b>	1.003	0.935
Defensive							
Model	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	95% Qu.
Multinomial	0.850	0.882	0.884	0.884	0.886	0.914	0.899
Multinomial-Efficiency	0.781	0.873	0.884	0.885	0.891	1.061	0.934
Ridge Normal	0.851	0.889	0.896	0.896	0.903	0.947	0.913
Lasso Normal	0.842	<b>0.902</b>	<b>0.902</b>	<b>0.902</b>	<b>0.902</b>	0.977	0.919
Normal after Lasso	0.771	<b>0.888</b>	<b>0.888</b>	<b>0.888</b>	<b>0.888</b>	1.049	0.927

Figure 5.6: Distribution of expected points-Contribution via Multinomial and Ridge Normal model.



Now we focus on the comparison between the top 50 model based RAPM rated players with the top 50 of selected boxscore statistics per game, such as:

- Points (PTS), assists (AST), and offensive rebounds (OREB) for offensive contribution.
- Defensive rebounds (DREB), steals (STL) and blocks (BLK) for defensive contribution.
- A bottom 50 in terms of minutes played (MP) is used for both offensive and defensive contribution.

Tables 5.5 and 5.6 present the agreement percentage of the players in the top 50 contributions per model and the ones for each boxscores statistics. We can see that Multinomial RAPM ratings appear with the highest percentage agreement while the more popular Ridge RAPM ratings (frequently in the literature) with the lowest percentage in terms of both offensive and defensive boxscore statistics.

Table 5.5: Agreement between model based RAPM highest rated players and the better ones in selected boxscores offensive statistics (per game)<sup>60</sup>.

OFFENSIVE RAPM	MP	PTS	AST	OREB
Multinomial	<b>0%</b>	<b>58%</b>	<b>68%</b>	<b>52%</b>
Binomial after Lasso	4%	46%	58%	42%
Lasso Binomial	4%	42%	54%	34%
Normal after Lasso	6%	50%	58%	30%
Standardized Normal	6%	50%	58%	32%
Lasso Normal	6%	46%	54%	32%
Ridge Binomial	12%	28%	36%	26%
Multinomial (Efficiency)	16%	40%	52%	38%
Ridge Normal	16%	26%	32%	20%

---

<sup>60</sup> MP refer to those 50 players with the less minutes played in the season. PTS refer to points scored per game, AST to assists per game and OREB to offensive rebounds per game. Methods are listed in terms of the low-time players higher efficiency (lower percentages for MP are desirable).



Table 5.6: Frequency of players' defensive contribution from top 50 per model that are captured in the top 50 of some in-game-statistics (per game)<sup>61</sup>.

DEFENSIVE RAPM	MP	DREB	BLK	STL
Multinomial	<b>2%</b>	<b>60%</b>	<b>20%</b>	<b>46%</b>
Binomial after Lasso	<b>2%</b>	52%	14%	28%
Lasso Binomial	<b>2%</b>	38%	14%	24%
Standardized Normal	8%	38%	12%	36%
Multinomial (Efficiency)	12%	38%	14%	28%
Normal after Lasso	8%	36%	12%	34%
Lasso Normal	8%	36%	14%	34%
Ridge Normal	20%	24%	8%	22%
Ridge Binomial	16%	22%	8%	20%

Especially for the low time players none of them appears in the Multinomial RAPM top rated offensive players, in contrast to the rest of the implemented models. Similar is the picture for the defensive players where only one low-time player appears in the Multinomial RAPM higher ratings. However, the defensive contribution of each player cannot be evaluated in a straightforward manner since many defensive actions and moves are not captured by boxscore or other game statistics. Also, it seems to be easier for low-time defensive players to have higher offensive ratings, as they usually play in the garbage time (when the outcome of the game has been almost finalized). In addition, Ridge-RAPM top 50 of offensive and defensive players includes the highest number of players with the lowest-time players in the season.

We have already provided the All-NBA teams (see Section 3.5.6) as one of the external validation criteria for model comparison. Each one of the three 5-lineups are composed of one Center, two Forwards and two Guards. Therefore, the percentage of agreement between the highest rated players per position in terms of model-based RAPM is examined, since each player position has six “chances” to be predicated among the best 5-lineups.

---

<sup>61</sup> MP refer to those 50 players with the less minutes played in the season. PTS refer to points scored per game, AST to assists per game and OREB to offensive rebounds per game. Methods are listed in terms of the high-DREB (and secondly the low-time) players higher efficiency.

Table 5.7: Percentage of agreement between the highest-rated offensive and defensive per position for each model and the All-NBA teams.

<b>RAPM</b>	<b>Offensive agreement</b>	<b>Defensive agreement</b>
Multinomial	<b>60.0%</b>	<b>50.0%</b>
Lasso Normal	<b>53.3%</b>	<b>60.0%</b>
Lasso Binomial	<b>53.3%</b>	<b>50.0%</b>
Standardized Normal	<b>46.7%</b>	<b>40.0%</b>
Binomial After Lasso	46.7%	20.0%
Ridge Normal	40.0%	20.0%
Normal after Lasso	40.0%	20.0%
Multinomial (efficiency)	40.0%	20.0%
Ridge Binomial	33.3%	10.0%

From Table 5.7, it is obvious that the highest percentage of agreement between the RAPM and All-NBA teams (offensive agreement) is observed for Multinomial model (after Lasso) with 60% and in Lasso-Normal and Lasso-Binomial with 53%. The Multinomial appears with the lowest percentage agreement in case of considering each player's participation (efficiency) and not the team's total possessions. Finally, based on the top three players per team, all models have estimated more than 11 players who belong to the best three 5-lineup: 73% for Ridge Normal, 87% for the Multinomial and 80% for the rest of the methods. In addition, all All-NBA teams' players are included in the top 5 of each team with exception of the Ridge Normal where the agreement is 93%.

Another criterion is the agreement between the Team Possessions Rating (TPR), which was introduced in Chapter 3, and the model-based players' performance. This can be used as an extra measure of marginal goodness of fit. In fact, players with a positive impact on their teams' offense and defense are considered (i.e.,  $TPR > 1$ ). The 80.7% (75/93), 29.3% (27/93) and 35.5% (33/93) of those offenders remain at the Lasso approach for Multinomial, Binomial, and Normal models respectively. However, their contribution was positive for 64.4% (59/93) for the Multinomial, 18.3% (17/93) for the Binomial, and 32.3% (30/93) for the Normal model.

For the positive D-TPR (defensive) players, 60.9%, 21.8%, and 28.7% of them "survived" after Lasso application at the Multinomial, Binomial, and Normal models respectively. However, the models are lower efficient this time since 28.7% (25/87) of the Multinomial and 14.9% (13/87) of Normal positive-rated and 10.3% (9/87) of the Binomial appear in the results.

### 5.3.2 Five-lineup analysis

It is also useful to consider higher-level criteria based on the 5-lineups (and not stay at the lower level of players) in order to study the efficiency of the methods. For this reason, Tables 5.10a and 5.10b present the agreement between the highest playing time 5-lineups for each team and the 5-lineups composed by the model-based RAPM highest-rated offenders and defenders per position (Multinomial, Normal and Binomial after screening Lasso).

It seems that the Multinomial model appears with the highest agreement compared to the rest of the models (see summary Table 5.8). The Binomial model again appears with lower percentage of agreement, which is probably due to the accumulation of the whole scoring information in one case.

Table 5.8: Average percentage of agreement between the top players per model according to the best offensive 5-lineups per team.

Average Percentage	Multinomial	Normal	Binomial
Most played	67%	48%	27%
High Net rating	50%	35%	18%

Table 5.9a: Percentage of agreement between the top players per model according to the most-played 5-lineups per team.

TEAM	Model	Percentage	TEAM	Model	Percentage	TEAM	Model	Percentage
ATL	Multinomial	40%	HOU	Multinomial	40%	OKC	Multinomial	40%
	Normal	60%		Normal	20%		Normal	20%
	Binomial	40%		Binomial	0%		Binomial	40%
BKN	Multinomial	60%	IND	Multinomial	60%	ORL	Multinomial	60%
	Normal	20%		Normal	40%		Normal	60%
	Binomial	20%		Binomial	0%		Binomial	0%
BOS	Multinomial	100%	LAC	Multinomial	80%	PHI	Multinomial	100%
	Normal	80%		Normal	20%		Normal	60%
	Binomial	80%		Binomial	0%		Binomial	40%
CHA	Multinomial	40%	LAL	Multinomial	60%	PHX	Multinomial	100%
	Normal	20%		Normal	80%		Normal	80%
	Binomial	20%		Binomial	0%		Binomial	80%
CHI	Multinomial	80%	MEM	Multinomial	100%	POR	Multinomial	40%
	Normal	60%		Normal	80%		Normal	20%
	Binomial	40%		Binomial	40%		Binomial	0%

Table 5.9b: Percentage of agreement between the top players per model according to the most-played 5-lineups per team.

TEAM	Model	Percentage	TEAM	Model	Percentage	TEAM	Model	Percentage
<b>CLE</b>	Multinomial	80%	<b>MIA</b>	Multinomial	80%	<b>SAC</b>	Multinomial	20%
	Normal	60%		Normal	40%		Normal	40%
	Binomial	40%		Binomial	20%		Binomial	0%
<b>DAL</b>	Multinomial	80%	<b>MIL</b>	Multinomial	80%	<b>SAS</b>	Multinomial	80%
	Normal	40%		Normal	40%		Normal	80%
	Binomial	40%		Binomial	60%		Binomial	20%
<b>DEN</b>	Multinomial	80%	<b>MIN</b>	Multinomial	100%	<b>TOR</b>	Multinomial	80%
	Normal	60%		Normal	60%		Normal	40%
	Binomial	40%		Binomial	0%		Binomial	0%
<b>DET</b>	Multinomial	80%	<b>NOP</b>	Multinomial	60%	<b>UTA</b>	Multinomial	40%
	Normal	60%		Normal	20%		Normal	40%
	Binomial	0%		Binomial	40%		Binomial	40%
<b>GSW</b>	Multinomial	60%	<b>NYK</b>	Multinomial	20%	<b>WAS</b>	Multinomial	60%
	Normal	40%		Normal	0%		Normal	60%
	Binomial	40%		Binomial	20%		Binomial	0%

We continue our analysis with a comparison between the lineup with the highest playing time and the highest net rating<sup>62</sup> for each team are compared with the ones composed of the highest-rated offenders and defenders per position according to the Multinomial and Normal results. Three selected teams are discussed as case study.

Starting with Atlanta Hawks (ATL), differences between the highest playing time and top RAPM 5-lineups are given in Figure 5.7. The two most valuable guards of Atlanta, Young and Bogdanovic, are also included to the RAPM-based lineups. The Multinomial captures the third guard, Huerter, as well, in contrast to the Normal lineup. From the power forwards, Gallinari appears with highest offensive RAPM ratings for both Normal and Multinomial model. Collins has lower offensive contribution than Gallinari. However, his participation in the highest playing-time lineup is not unexpected, therefore he tolerates a limited playing time for which he can be impactful for his team. Also, Collins, who belongs at the highest playing time lineup of Atlanta, appears in the highest defensive lineup for the Multinomial one.

<sup>62</sup> An advanced boxscore statistic that measures a team's point differential per 100 possessions. On player level this statistic is the team's point differential per 100 possessions while they are on court. Formula: OFFRTG - DEFRTG

Figure 5.7: The highest playing time 5-lineup (at the center) of Atlanta vs. the most impactful players (defense in the red and offense in the blue half-court).

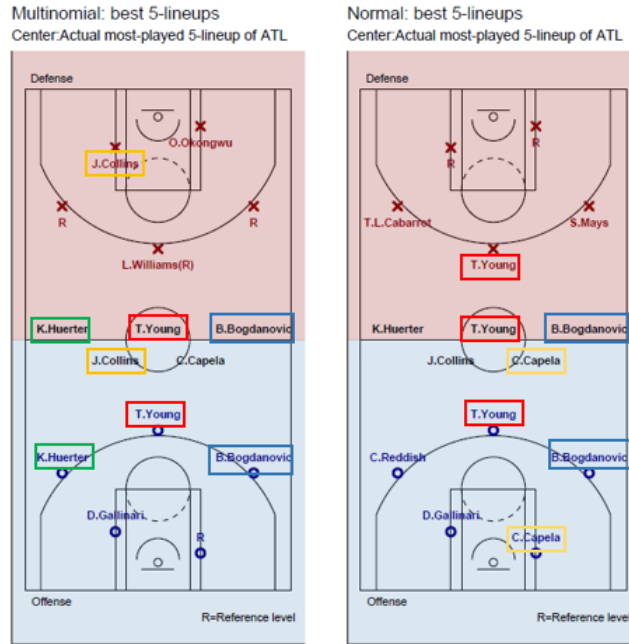
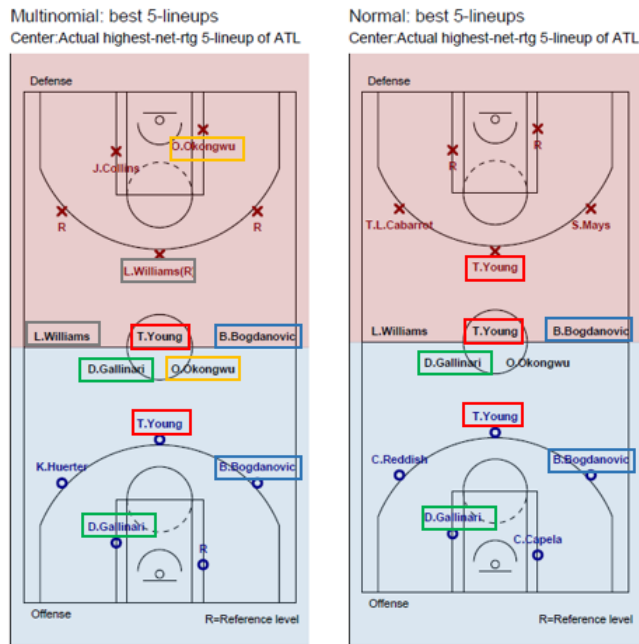


Figure 5.8: Lineup with the highest net rating of Atlanta (at the center) vs. the most impactful players (defense in the red and offense in the blue half-court).



With respect to the centers of Atlanta, Capela and Okongwu appear with the highest offensive EPTS-RAPM, who belong at the reference group of the Multinomial model, while Capela has the highest offensive Normal RAPM. With respect to the external criteria, the Multinomial model agrees with both the net rating and playing time while the Normal agrees only with the playing

time. In fact, Capela is one of the starters who appear with negative EPTS-RAPM contribution even his overall participation is less than Okongwu (2042 minutes played against 992).

This observation agrees with our findings for the discrimination of players in terms of their playing time according to the models, since Okongwu playing time per game is seven minutes less, the boxscore statistics PER (Player Evaluation Rating) is quite close to Capela (19.8 and 21.4) and the normalized boxscore statistics are slightly higher than Capela's (for example Win Share per game which is 0.202 for the first player and 0.195 for the second).

The second example is Golden State Warriors (GSW), the NBA champions of 2021-2022 season. From Figure 5.9 we can see that only the point guard Curry has positive offensive Normal RAPM (Normal) while the highest contribution for the rest of the positions belongs is equal to zero (i.e., reference group).

However, the Multinomial is more informative in this case (see Figure 5.9). In fact, the offensive lineup and the most-played lineup agree in four positions. The two lineups do not agree in the small forward (left guard in the figure). The highest EPTS-RAPM player of the Multinomial model is Kuminga for offense. The highest offensive Normal-RAPM appears in four power forwards of the reference group. With respect to the defensive highest RAPM, Iguodala belongs to both models-lineup. This might be the result of Klay Thompson's injury, one of the best players of Warriors who was a starter.

The 5-lineup with the highest net rating of Warriors, is composed by Payton, Toscano-Anderson, Thompson, Wiggins, and Looney. The last three players are good players and starters while Payton and Toscano cannot be considered as starters, but they played mostly in the garbage time.

Figure 5.9: The highest playing time 5-lineup (at the center) of Warriors vs. the most impactful players (defense in the red and offense in the blue half-court).

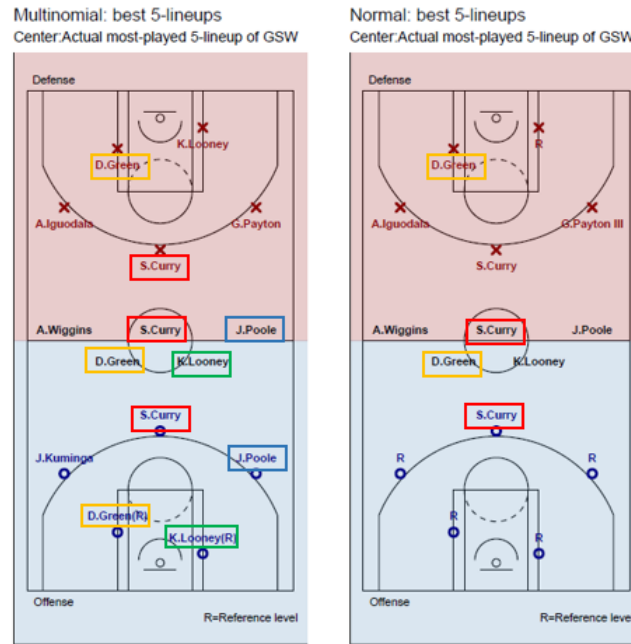
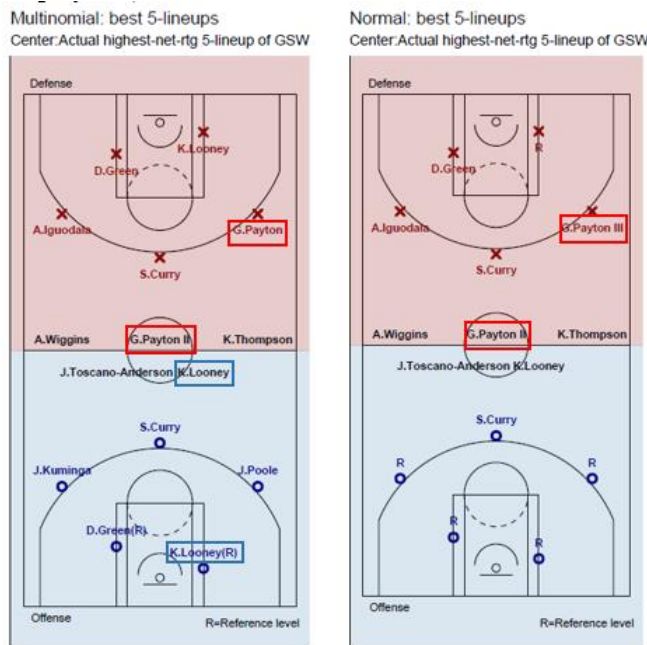


Figure 5.10: Lineup with the highest net rating of Warriors (at the center) vs. the most impactful players (defense in the red and offense in the blue half-court).



Detroit Pistons (DET) is the final team we are going to examine as case study; This team was one of the worst performed team in the 2021-2022 NBA season (finished 14<sup>th</sup> in the regular season among the 15 teams in the East conference). Generally, the two lineups do have a number of differences (Figure 5.11). Let us first consider the case of power forwards, Jerami Grant and Trey

Lyles. Regarding the offensive EPTS-RAPM, Lyles is more impactful than Grant, in contrast with the defensive EPTS-RAPM. However, Grant belongs to the lineup of Detroit with the highest playing time and his box-score statistics are better than the ones of Lyles. These ratings can work as an alarm that would make the manager think if he needs to exploit more the talents of this player and optimize his impact.

A last observation here is that the lineup with the highest net rating and the lineup with the highest playing time is about the same for a low performed team like Detroit. However, these two lineups differ for high-quality teams (like Warriors). In the latter case, low-time players may belong to the highest net rating lineup.

Figure 5.11: The highest playing time 5-lineup (at the center) of Detroit vs. the most impactful players (defense in the red and offense in the blue half-court).

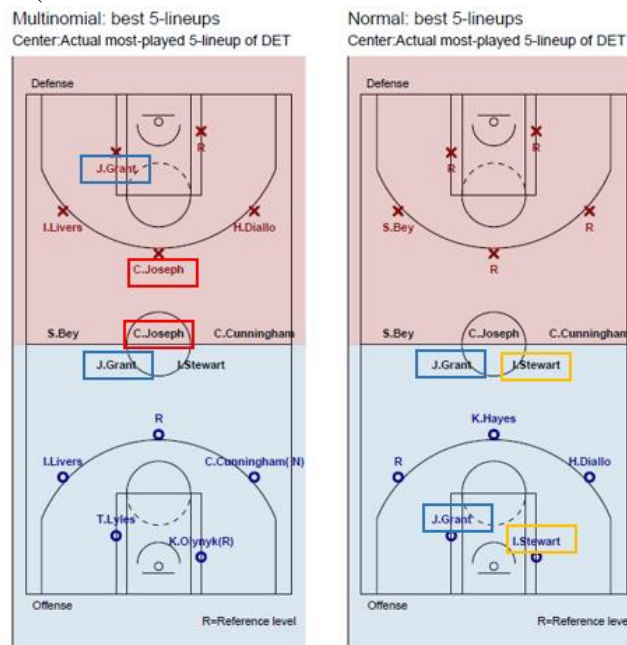
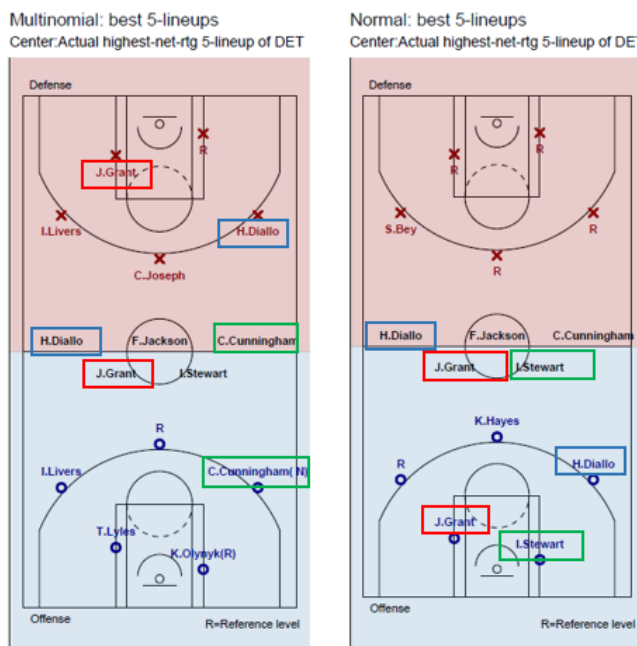




Figure 5.12: Lineup with the highest net rating of DET (at the center) vs. the most impactful players (defense in the red and offense in the blue half-court).



## 5.4 Conclusion

From all aforesaid in Section 5.2, we conclude that the fit of the proposed Multinomial model is mainly marginal but much better than the common normal-based model. Since the aim of RAPMs is not the prediction from the results we may keep the bad fit of the Normal models and the satisfactory simulated tests for specific characteristics of the Multinomial distribution.

The Multinomial model performed better in terms of the external validation criteria provided in Section 5.3. These validation measures were slightly worse for the defensive RAPM ratings. This might be due to the fact that defense is based on teamwork and very difficult to be measured.

Moreover, the indirect implementation approach we have used in this thesis of the Multinomial model, rates more fairly the players who may affect positively their team in one of the three types of points scored/conceded (one, two, or three points). Additionally, we have observed that the Multinomial model rate in the lower ranking positions players of low-winning-record teams and the lower-quality role-players (of better teams with important participation in the games). In this way, the low-time players that had a positive contribution in the court are rated higher than the role players who did not perform well. Also, high-time players with high performance are correctly awarded with the highest EPTS-RAPM ratings.

## Chapter 6: Discussion

### 6.1 Summary of the methodology

“Teamwork is the element of basketball most difficult to capture in any quantitative sense” (Dean Oliver, *Basketball on Paper*, 2004). This quote inspired Rosenbaum (2004) to study players’ performance based on their teammates and opponents on the court. Hence, the starting point of our thesis was the Regularized Adjusted Plus/Minus index estimated via a Ridge regression model as suggested by the related literature.

Ridge regression was proposed for the following reasons: i) low variance of the estimated coefficients for each player, ii) all players can be considered, iii) regularization works as a dimensionality reduction technique, iv) shrinkage parameter controls multicollinearity between the players, and v) simple interpretation.

Hence, the first method implemented in this thesis was the Ridge Normal model for the whole possession dataset of NBA season 2021-2022. Low-time players considerably affected the results since they appear with higher than expected contributions. To minimize this problem, we decided to: i) apply Lasso instead of Ridge and ii) implement the Ridge Regression without considering low-time players.

Lasso resulted in RAPMs with better discrimination among the evaluated players since better players have higher RAPM contributions. When LTP players (i.e., NBAers who played less than 200 minutes in the season) are removed the results were considerably improved. The main drawback (but also an advantage at the same time) of the Lasso method is that does not discriminate zero-RAPM players. This problem is more severe for the normal model where 674 out of 970 were set to zero by Lasso.

Note that, a more appropriate model than the Normal one should be fitted since the response variable (points per possession) is discrete taking values in the range of 0-3 (mainly). Thus, logistic regression models were the next natural step in our analysis. Firstly, we consider as a response the binary of scoring against not scoring. Regularized methods were applied resulting to accuracy of about 55% (which is not very high). This is a simplification of the original problem since the number of points is not considered in this treatment.

The final method implemented in this thesis was a multinomial model on the number of points scored. In order to avoid computational problems and be more flexible concerning which coefficients were set equal to zero by Lasso, we used three separate binomial implementations. Finally, expected points (EPTS-RAPM) were calculated, in order to evaluate each player.

## **6.2 Some highlights**

Since we have made various observations up to this point in the main Chapters of the thesis, it is useful to discuss some special conclusions that stand out.

The agreement between the Ridge RAPM high-rated players and the best-performed players in selected boxscore statistics (minutes played, points, assists, offensive rebounds, defensive rebounds, blocks, steals) or the All-NBA teams' players is lower than Lasso-RAPM.

Moreover, an interesting finding is the interpretation that is induced by the relationship between the Ridge Binomial and Ridge Normal RAPMs. Specifically, ratings of Ridge-Normal RAPMs were found to be linearly related to the binomial ones. This gives us the opportunity to calculate Binomial RAPMs by the Normal RAPMs via a more valid interpretation for the latter ones which suffer by the fact that the underlying normal model is not appropriate for the response at hand (i.e., points per possession).

The final Multinomial model is more appropriate for such a response. The EPTS-RAPM obtained by this model behaves much more better than RAPMs obtained by the Normal or the simple binomial model.

An advantage of the Multinomial model is that offers more detailed information concerning the player contribution in different types of scoring (one, two, or three points). For example, if a team needs a three-point shoot at a crucial time of a game, higher-rated players in terms of three points scored could be selected to play. It is not about players with high shooting percentages but those that help their team to increase their odds to score the specific type of points. Such player can be a role player who drives to the paint easily and pass the ball to a suitable player to shoot or a center that can pass the ball.

Furthermore, Expected Points (EPTS-RAPM) were provided as one player evaluation metric through the Multinomial model. This is an easily interpreted index, especially for offenders, in contrast to the log-odds interpretation of the standard Multinomial model.

In the Normal-Lasso regression, the non-zero RAPM players were to be 296 offensive and defensive players (30.5%). In the reference group (with zero-RAPMs), a number of top performed players of the 2021-2022 NBA season appear. On the other hand, 678 (69.9%) players have non-zero EPTS-RAPM, which is about double (more than double) compared with the corresponding zero-RAPM players in Normal-Lasso. Players with fewer responsibilities are found in the reference group, which are actually: i) average performed offensive players and some top with no significant defensive contribution of the higher-quality teams and ii) average performed with less playing time players of the lower-quality teams.

For the non-zero RAPM players, we observe three different situations:

- Higher and lower-rated Ridge RAPM players are usually low-time players. Hence, Ridge-Normal fails in ranking top players at the top of its list.
- The higher the Lasso-RAPMs the higher the minutes played. This partially solves the problem of the high evaluation of LTPs which is the main problem of plus/minus ratings.
- Higher and lower-rated EPTS-RAPM players are players with increased playing time. Hence, the Multinomial treats the contribution of players in more sensible way.

Team effects were not considered in the analysis, since Lasso eliminated them from the model (for both normal and binomial models).

Regarding the players' contribution, it is observed that the high-quality players of high performed teams affected their teammates' performance positively and made them seem more important than they appear according to their personal statistics. On the other hand, for lower-quality teams, role players who played less than others appear with a positive impact.

In all the top 50 offensive ratings, most of the players are Shooting Guards, Power Forwards, or Point Guards frequently appear. This can be considered as a confirmation of the style of the game NBA teams play over the last seasons (fast attacking with 3-pointers). Also, from the group of Small Forwards, only a few superstars appear in the top-50 ranking (with best examples to be Jason Tatum and Jimmy Butler).

For the defensive contribution, there is no great difference in the different types of players positions in the top-50 list. Centers and Small Forwards (with the less offensive contributions) are important for the defensive part by playing close or in the paint, where it is easier for the opponents to score.

The RAPM ratings within each team can provide useful information and valuable comparisons. Two important points are the following:

- For a number of teams Ridge RAPMs fail to identify the most valuable players in contrast with the Lasso RAPMs. Representative examples are Bogdan Bogdanovic, Goran Dragic, Daniel Theis, Tyler Cook, Jevon Carter, and Juancho Hernangomez, who appear to be the best offenders based on the Ridge models, although their teammates are top NBAers (Trae Young, Kevin Durant, Jason Tatum, DeMar DeRozan, Giannis Antetokounmpo and Donovan Mitchell, respectively).
- For an important number of low performed teams Lasso models as the highest-rated offensive and especially defensive players had RAPM equal to zero. Specifically, there is no player with a positive contribution to the defense of Houston Rockets, Detroit Pistons, Orlando Magic, Oklahoma City Thunder, and the offense of Portland Trail Blazers, and Washington Wizards.

To conclude, a variety of different topics can be studied using the player's contribution we have implemented in this thesis. Some of them have been discussed in this thesis. For example, there is no doubt that an expert in this field (such as a basketball coach) can use the estimated RAPM ratings for specific problems. Our proposed methodology offers a more detailed rating system to study players' impact on offense and defense for different scoring situations. We also proposed an overall through the alternative proposed evaluation index (EPTS-RAPM).

## **6.3 Future work**

This work focuses on calculating the players' contribution by taking into account their teammates and opponents with different types of models. Extensions of this methodology can be employed in a variety of different ways.

A first methodological choice is to use the shrinkage of the standard Ridge regression model (which is most commonly used in the literature) in an alternative way than the one used until now in practice, by setting the penalty term to be a function of the playing time in order to reduce the effect of the low-time players in the estimated RAPMs.

Another idea is to implement a two-stage regression model, with the first stage of modeling to be on the estimating of RAPM ratings at the team level and the second stage on the lineup of the team. Hence, in the second level, the team coefficient will be split into the contribution of each player. By this way, some boxscore statistics might be useful in the second stage.

Nevertheless, the study of the data for a specific team is of interest. Load management is of primer importance for team sports. Hence, we might be able to build RAPM models for the players of a specific team which will help managers to optimize the team load and use appropriate players in the right games time. This can be a potentially useful tool for the teams by estimating the impact of the “second unit” (or bench) players and the contribution of the substitutes at specific time points of the game compared to the first lineup players.

Another interesting point is the study of specific lineups. For example, lineups for playing “small ball”. Such an analysis can provide valuable extra choices for a team. For example, when one or more players are injured.



## **APPENDIX**

You can find the full part of Appendix in the electronic supplement which can be found here:

<https://cloud.aueb.gr/index.php/s/WwkKxZFLKePTp99>

R-Shiny Application for players contribution according to Multinomial fitted model results:

[https://bballstats.shinyapps.io/nba21-22\\_Players\\_Contribution/](https://bballstats.shinyapps.io/nba21-22_Players_Contribution/)





## References

### Bibliography

- Aizemberg, L., Soares de Mello, J.C., and Alves, A., (2014). “Measuring the NBA Teams’ Cross-Efficiency by DEA Game”, *American Journal of Operations Research*, 4, 101-112.
- Alamar, B.C., (2013). *Sports Analytics: A Guide for Coaches, Managers, and Other Decision Makers*. USA, Columbia University Press.
- Amorin, A.M., and Guimarães, E., (2022). “Ball screen effectiveness in elite women's basketball”, *Journal of Physical Education and Sport*, DOI: 10.7752/jpes.2022.03095.
- Armanious, M., (2019). *Men’s U-Sports Basketball Analysis*, Ottawa, Carleton University (Master of Science Thesis).
- Assani, S., Mansoor, M., Li, Y., and Yang, F., (2021). “Efficiency, RTS, and marginal returns from salary on the performance of the NBA players: A parallel DEA network with shared inputs”, *Journal of Industrial and Management Optimization*, DOI: 10.3934/jimo.2021053.
- Baghal, T.A. (2012). “Are the “Four Factors” Indicators of One Factor? An Application of Structural Equation Modeling Methodology to NBA Data in Prediction of Winning Percentage”. *Journal of Quantitative Analysis in Sports*, Vol. 8(1), Manuscript 1355.
- Bartholomew, J.T., and Collier, D.A., (2011). “A defensive basketball efficiency score using data envelopment analysis”, *Journal of Sport Management Research*, <https://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.648.9054&rep=rep1&type=pdf>.
- Berri, D., (2011). “Deconstructing the Adjusted Plus-Minus Model”, *Wages of Win Journal*, <https://dberri.wordpress.com/2011/03/05/deconstructing-the-adjusted-plus-minus-model/>.
- Berri, D.J., Schmidt, M.B. and Brook, S.L. (2006). *The Wages of Wins: Taking Measure of the Many Myths in Modern Sport*. Stanford, CA: Stanford University Press.
- Bhat, Z. U. H., Sultana, D., & Dar, Q. F. (2019). “A comprehensive review of data envelopment analysis (DEA) in sports”, *Journal of Sports Economics & Management*, 9(2), p. 82-109.
- Bhatnagar, R., and Babbar, M., (2019). “A systematic review of sports analytics”, *Research Gate, Conference paper, International Conference of Business and Management, Delhi School of Management*.
- Braga, V.A., (2021). “Big Data Analytics in Basketball Versus Business”, *Sciend*, Volume 16 (3), p. 24 – 31, DOI: <https://doi.org/10.2478/sbe-2021-0042>.
- Brown, B., (2019). *Predictive Analytics for College Basketball: Using Logistic Regression for Determining the Outcome of a Game*. USA, Peter. T Paul College of Business and Economics, University of New Hampshire (Undergraduate Honor Thesis).
- Brown, S., (2017). “A PageRank Model for Player Performance Assessment in Basketball, Soccer and Hockey”, *Research Gate*, McGill University.
- Carlin, B. P. (1996) “Improved NCAA basketball tournament modeling via point spread and team strength information”, *The American Statistician*, Vol. 50, pp. 39–43.

- Casals, M., and Martinez, J.A., (2013). “Modelling player performance in basketball through mixed models”, *International Journal of Performance Analysis in Sport* 13(1):64-82, DOI:10.1080/24748668.2013.11868632.
- Deshpande, S., and Jensen, S., (2016). “Estimating an NBA player’s impact on his team’s chances of winning”, *Journal of Quantitative Analysis in Sports*, DOI:10.1515/jqas-2015-0027.
- Doolittle, B. and Pelton, K. (2010) *Pro Basketball Prospectus*. United States (California): CreateSpace Independent Publishing Platform.
- Duman, E., Sennaroglu, B., and Tuzkaya, G., (2021). “A cluster analysis of basketball players for each of the five traditionally defined positions”, *Proceedings of the Institution of Mechanical Engineers, Part P: Journal of Sports Engineering and Technology*, DOI: 10.1177/17543371211062064.
- Entine, O. A., and Small, D.S., (2008). “The role of rest in the NBA home-court advantage”, *Journal of Quantitative Analysis in Sports*, Vol. 4(2), Article: 6.
- García, J., Ibáñez, S.J., Martinez De Santos, R., Leite, N., and Sampaio, J., (2013). “Identifying Basketball Performance Indicators in Regular Season and Playoff Games”, *Journal of Human Kinetics*, volume 36, p. 163-170.
- Gelade, G., and Hvattum, L.M., (2020). “On the relationship between +/- ratings and event-level performance statistics.”, *Journal of Sports Analytics*. 6. 1-13. 10.3233/JSA-200432.
- Ghimire, S., Ehrlich, J., and Sanders, S., (2020). “Measuring individual worker output in a complementary team setting: Does regularized adjusted plus minus isolate individual NBA player contributions?”, *PLOS ONE*, <https://doi.org/10.1371/journal.pone.0237920>.
- Goldstein, J., (2018). “Nylon Calculus: Defining and calculating luck-adjusted ratings for the NBA”, *FANSIDED blog*, <https://fansided.com/2018/01/08/nylon-calculus-calculating-luck-adjusted-ratings/>.
- Harville, D. A. (2003). “The selection of seeding of college basketball or football teams for postseason competition”, *Journal of the American Statistical Association*, Vol. 98, pp. 17–27.
- Harville, D. A. and M. H. Smith (1994). “The home-court advantage: How large is it and does it vary from team to team”, *The American Statistician*, Vol. 48, pp. 22–29.
- Hass, Z., and Craig, B.A., (2018). “Exploring the potential of the plus/minus in NCAA women’s volleyball via the recovery of court presence information”, *Journal of Sports Analytics*, DOI: 10.3233/JSA-180217.
- Hernández, F.A.L., Martinez, J.A., and Marin, M.A., (2013). “Spatial Pattern Analysis of Shot Attempts in Basketball”, *Revista Internacional de Medicina y Ciencias de la Actividad Fisica y del Deporte* 13(51).
- Hoffman, L. and Joseph, M. (2003). *A multivariate statistical analysis of the NBA*. University of Wisconsin River Falls (River Falls, WI) and Kentucky State University Frankfort, KY
- Hollinger, J. (2002). *Pro Basketball: Prospectus*. United States: Brassey’s
- Hollinger, J. (2003). *Pro Basketball: Prospectus*. United States: University of Nebraska Press
- Hollinger, J. (2004). *Pro Basketball: Forecast*. United States: Potomac Books
- Hollinger, J. (2005). *Pro Basketball: Forecast*. United States: Potomac Books Inc.

- Hong, X., (2021). "Basketball Data Analysis Using Spark Framework and K-Means Algorithm", *Hindawi Journal of Healthcare Engineering*, Volume 2021, Article ID: 6393560, <https://doi.org/10.1155/2021/6393560>.
- Hu, F., and Zidek, J.V., (2004). "Forecasting NBA basketball playoff outcomes using the weighted likelihood", *Research Gate*, DOI: 10.1214/lnms/1196285406.
- Huang, M.L., and Lin, Y.J., (2020). "Regression Tree Model for Predicting Game Scores for the Golden State Warriors in the National Basketball Association", *Symmetry*, 12(5), 835, <https://doi.org/10.3390/sym12050835>.
- Hussain, H., (2019). "Using K-Means Clustering Algorithm to Redefine NBA Positions and Explore Roster Construction", *Towards Data Science*, <https://towardsdatascience.com/using-k-means-clustering-algorithm-to-redefine-nba-positions-and-explore-roster-construction-8cd0f9a96dbb>.
- Hvattum, L.M., (2019). "A comprehensive review of plus-minus ratings for evaluating individual players in team sports", *International Journal of Computer Science in Sport*, Volume 18, Issue 1, DOI: 10.2478/ijcss-2019-0001.
- Hvattum, LM., (2020) "Offensive and Defensive Plus–Minus Player Ratings for Soccer", *Applied Sciences*, 10(20):7345. <https://doi.org/10.3390/app10207345>.
- Hvattum, LM., and Gelade, G., (2021). "Comparing bottom-up and top-down ratings for individual soccer players.", *International Journal of Computer Science in Sport*. 20. 23-42. 10.2478/ijcss-2021-0002.
- Jin, Y., (2021). "Analysis of NBA Business Strategy", *Education and Humanities Research*, volume 543, *Proceedings of the 2021 6th International Conference on Social Sciences and Economic Development (ICSSED 2021)*.
- Jones, M.B. (2007). "Home advantage in the NBA as a game-long process", *Journal of Quantitative Analysis in Sports*, Vol. 3 (4), Article: 2.
- Jones, M.B. (2008). "A note on team-specific home advantage in the NBA", *Journal of Quantitative Analysis in Sports*, Vol. 4 (3), Article: 5.
- Jones, S.E., (2016). *PREDICTING OUTCOMES OF NBA BASKETBALL GAMES*, USA, North Dakota State University (Master of Science Thesis).
- Jyad, A., (2020). "Redefining NBA Player Classifications using Clustering: Using Hierarchical Clustering to define NBA Players", *Towards Data Science*, <https://towardsdatascience.com/redefining-nba-player-classifications-using-clustering-36a348fa54a8>.
- Kharrat, T., Pena, J.L., and McHale, I.G., (2017). "Plus-Minus Player Ratings for Soccer", *European Journal of Operational Research* 283(2), DOI:10.1016/j.ejor.2019.11.026.
- Knorr-Held, L. (2000). "Dynamic ratings of sports teams", *The Statistician*, Vol. 49, pp. 261–276.
- Kubatko, J., Oliver, D., Pelton, K. and Rosenbaum, D.T. (2007). "A Starting Point for Analyzing Basketball Statistics", *Journal of Quantitative Analysis in Sports*, Vol. 3, Iss. 3, Article 1.
- Küpfer, E. (2005). Team Similarity. *APBRMetrics Forum*.

- Lee, D.J., and Page, G.L., (2021). “Big Data in Sports: Predictive Models for Basketball Player’s Performance”, *Mathematics in Industry Reports (MIIR)*.
- Lewis, M. (2003). *Moneyball: The Art of Winning an Unfair Game*. United States: W.W. Norton & Company.
- Lorenzo, J., Lorenzo, A., Conte, D., and Giménez, M., (2019). “Long-Term Analysis of Elite Basketball Players’ Game-Related Statistics Throughout Their Careers”, *Frontiers in Psychology*, 10:421. DOI: 10.3389/fpsyg.2019.00421.
- MacDonald, B., (2011). “NHL Adjusted Plus-Minus Part 1”, *Arctic Ice Hockey*, <https://www.arcticicehockey.com/2011/9/22/2441898/nhl-adjusted-plus-minus-part-01> .
- MacDonald, B., (2011). *An Improved Adjusted Plus-Minus Statistic for NHL Players*, Boston, MA, USA, Sloan Sports Analytics Conference 2011, March 4-5, 2011.
- MacDonald, B., (2012). “Adjusted Plus-Minus for NHL Players using Ridge Regression with Goals, Shots, Fenwick, and Corsi”, *Journal of Quantitative Analysis in Sports* 8(3), DOI:10.1515/1559-0410.1447.
- MacDonald, B., (2012). “NHL Adjusted Plus-Minus Part 2: Even Strength Offense”, *Arctic Ice Hockey*, <https://www.arcticicehockey.com/2012/1/16/2711979/nhl-adjusted-plus-minus-part-2-even-strength-offense> .
- MacDonald, B., (2012). “NHL Adjusted Plus-Minus Part 3: Even Strength Offense”, *Arctic Ice Hockey*, <https://www.arcticicehockey.com/2012/1/29/2756500/nhl-adjusted-plus-minus-part-3-even-strength-defense> .
- MacDonald, B., (2012). “NHL Adjusted Plus-Minus Part 4: Even Strength Offense”, *Arctic Ice Hockey*, <https://www.arcticicehockey.com/2012/5/20/3032244/nhl-adjusted-plus-minus-part-4-even-strength-overall> .
- Malarranha, J., Figueira, E.B., Leite, N., and Smpaio, J., (2013). “Dynamic Modeling of Performance in Basketball”, *International Journal of Performance Analysis in Sport* 13(2):377-387, DOI:10.1080/24748668.2013.11868655.
- Manner, H. (2015). *Modeling and forecasting the outcomes of NBA Basketball games*. Institute of Econometrics and Statistics, University of Cologne.
- Martinez, J.A., and Caro, L.M., (2011). “A stakeholder assessment of basketball player evaluation metrics”, *Journal of Human Sport and Exercise*, 6(1), DOI: 10.4100/jhse.2011.61.17.
- Martinez, J.A., Marin, M.A., Casals, M., and Hernández, F.A.L., (2017). “Regular point scoring by professional basketball players”, *Electronic Journal of Applied Statistical Analysis* 10(3):759-772, DOI:10.1285/i20705948v10n3p759.
- Matano, F., Richardson, L.F., Pospisil, T., Eubanks, C., and Qin, J., (2018). “Augmenting Adjusted Plus-Minus in Soccer with FIFA Ratings”, *Research Gate*, [https://www.researchgate.net/publication/328380476\\_Augmenting\\_Adjusted\\_Plus-Minus\\_in\\_Soccer\\_with\\_FIFA\\_Ratings](https://www.researchgate.net/publication/328380476_Augmenting_Adjusted_Plus-Minus_in_Soccer_with_FIFA_Ratings).
- Moreno, P., and Lozano, S., (2014). “A network DEA assessment of team efficiency in the NBA”, *Annals of Operations Research*, 214 (1), p. 99-124, DOI: 10.1007/s10479-012-1074-9.

- Noivo, A., Amorim, A., Guimarães, E., and Janeira, M., (2022). "Ball screen effectiveness in elite women's basketball", *Journal of Physical Education and Sport* 22(3):757-766, DOI:10.7752/jpes.2022.03095.
- Oliver, D. (2002). *Basketball on Paper: Rules and Tools for Performance Analysis*. United States (Washington D.C.): Brassey's, Inc.
- Parker, C., (2018). "NBA Draft Pick Valuation", *Research Gate*, University of Pennsylvania.
- Parunmalka, K., (2012). *Modelling the NBA to Make Better Predictions*. USA, Massachusetts Institute of Technology (Master of Science thesis).
- Patel, R., (2017). *Clustering Professional Basketball Players by Performance*, Los Angeles, University of California (Master of Science Thesis).
- Rosenbaum, D.T., (2004). "Measuring How NBA Players Help Their Teams Win", *82games*, <http://www.82games.com/comm30.htm>.
- Sabin, R. P., (2021). "Estimating player value in American football using plus-minus models," *Journal of Quantitative Analysis in Sports*, De Gruyter, vol. 17(4), p. 313-364, DOI: 10.1515/jqas-2020-0033.
- Schmidt, M.B. and Berri, D.J., (2007). "Does One Simply Need to Score to Score?", *International Journal of Sport Finance*, 2007, 2, p. 142-148.
- Schultze, S.R. and Wellbrock, C.M. (2018). "A weighted plus/minus metric for individual soccer player performance", *Journal of Sports Analytics*, 4, p. 121–131, DOI 10.3233/JSA-170225.
- Schwertman, N. C., McCready, T.A. and Howard, L. (1991). "Probability models for the NCAA regional basketball tournaments", *The American Statistician*, Vol. 45, pp. 35–38.
- Shea, S.M., (2014). *Basketball Analytics: Spatial Tracking*. Createspace Independent Publishing Platform.
- Shea, S.M., and Baker, C.E., (2013). *Basketball Analytics: Objective and Efficient Strategies for Understanding How Teams Win*. CreateSpace Independent Pub. Platform.
- Sill, J., (2010). *Improved NBA Adjusted +/- Using Regularization and Out-of-Sample Testing*. MIT Sloan Sports Analytics Conference.
- Singh, N., (2020). "Sport Analytics: A Review", *The International Technology Management Review*, Vol. 9(1), p. 64–69, DOI: <https://doi.org/10.2991/itm.r.k.200831.001>; ISSN 2213-7149; eISSN 1835-5269.
- Sisneros, R., and Van More, M., (2013). "Expanding Plus-Minus for Visual and Statistical Analysis of NBA Box-Score Data", *Research Gate, Conference Paper, Conference: The 1st Workshop on Sports Data Visualization*.
- Snarr, T., (2020). "New EPM (estimated plus minus) metric", *Real GM*.
- Song, K., and Shi, J., (2019). "A discrete-time and finite-state Markov chain based in-play prediction model for NBA basketball matches", *Communications in Statistics - Simulation and Computation*, 50:11, p. 3768-3776, DOI: 10.1080/03610918.2019.1633351.
- Song, K., Zou, Q., and Shi, J., (2018). "Modelling the scores and performance statistics of NBA basketball games", *Communications in Statistics - Simulation and Computation*, 49:10, p. 2604-2616, DOI: 10.1080/03610918.2018.1520878.



- Soroush, B.K., Mohamad, R.Y., and Siavash, K.S., (2021). "Clustering of basketball players using self-organizing map neural networks", *Journal of Applied Research on Industrial Engineering*, E-ISSN: 2676-6167 / P-ISSN: 2538-5100.
- Stefani, R. T. (1977a). "Football and basketball prediction using least squares", *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-7, p. 117–121.
- Stefani, R. T. (1977b). "Improved least squares football, basketball, and soccer predictions", *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-7, p. 117–121.
- Stern, H. S. (1994). "A brownian motion model for the progress of sports scores", *Journal of the American Statistical Association*, Vol. 89, pp. 1128–1134.
- Stump, M., (2017). *Statistical Analysis of Momentum in Basketball*, USA, Bowling Green State University, Honors College.
- Teramoto, M. and Cross, C.L. (2010). "Relative Importance of Performance Factors in Winning NBA Games in Regular Season versus Playoffs", *Journal of Quantitative Analysis in Sports*, Vol. 6, Iss. 3, Article: 2.
- Torres, R. A. (2013). *Prediction of NBA games based on machine learning methods*. University of Wisconsin-Madison
- Uudmae, J. (2017). *Predicting NBA game outcomes*. Leland Stanford Junior University
- Villa, G., and Lozano, S., (2018). "Dynamic Network DEA approach to basketball games efficiency", *Journal of the Operational Research Society*, 69:11, p.1738-1750, DOI: 10.1080/01605682.2017.1409158.
- Wright, C, (2012). *Statistical Predictors of March Madness: An Examination of the NCAA Men's' Basketball Championship*. USA, Pomona College Economics Department.
- Xu, X., Zhang, M., and Yi, Q., (2022). "Clustering Performances in Elite Basketball Matches According to the Anthropometric Features of the Line-ups Based on Big Data Technology", *Frontiers in Psychology*, 11, Sec. Movement Science and Sport Psychology <https://doi.org/10.3389/fpsyg.2022.955292>.
- Yang, C.H., Lin, H.Y., and, Chen, C.P., (2014). "Measuring the efficiency of NBA teams: Additive efficiency decomposition in two-stage DEA", *Annals of Operations Research*, 217 (1), p. 565-589, DOI: 10.1007/s10479-014-1536-3.
- Yang, S.Y. (2015). *Predicting Regular Season Results of NBA Teams Based on Regression Analysis of Common Basketball Statistics*. University of California at Berkeley (A thesis submitted in fulfillment of the requirement for the degree of honors in Statistics).
- Yuan, L.H., Liu, A., Yeh, A., Kaufman, A., Reece, A., Bull, P., Franks, A., Wang, S., Illushin, D., and Bornn, L., (2015). "A mixture-of-modelers approach to forecasting NCAA tournament outcomes", *Journal of Quantitative Analysis in Sports*, 11(1), p. 13-27.
- Zhang, D.S., (2019). *Modelling and Simulation in Game Performances of Basketball Players and Teams in the National Basketball Association*, Madrid, Facultad de Ciencias de la Actividad Física y del Deporte (INEF) (UPM), DOI: <https://doi.org/10.20868/UPM.thesis.55641> (PhD Thesis).
- Zhang, X. (2019). *Modeling of NBA Game Data and their Correlation Structure*. University of Houston (Degree of Doctor of Philosophy in Mathematics).

- Zuccolotto, P., and Manisera, M., (2020). *Basketball Data Science With Applications in R*. UK, Chapman & Hall.
- Zuccolotto, P., Sandri, M., and Manisera, M., (2021). “Spatial Performance Indicators and Graphs in Basketball”, *Social Indicators Research* 156(2):1-14, DOI:10.1007/s11205-019-02237-2.

## Electronic References

- “A detailed guide for developing player ratings for WNBA (and other leagues)”, *THE ATHLYTICS BLOG*, (2022), <https://412sportsanalytics.wordpress.com/2022/08/18/a-detailed-guide-for-developing-player-ratings-for-wnba-and-other-leagues/> .
- “What is the best advanced statistic for basketball? NBA executives weigh in”, (2021), *Klutch Basket*, <https://klutchbasket.com/nba-news/what-is-the-best-advanced-statistic-for-basketball-nba-executives-weigh-in/> .
- “Introducing DARKO: An NBA playoffs game projection and betting guide”, *The Athletic blog* (2021), <https://theathletic.com/2613015/2021/05/26/introducing-darko-an-nba-playoffs-game-projection-and-betting-guide/> .
- “LEARN A STAT: BOX PLUS MINUS AND VORP”, *Hack a Stat blog*, <https://hackastat.eu/en/learn-a-stat-box-plus-minus-and-vorp/> .
- Cheema, A., (2021). “Calculating Regularized Adjusted Plus-Minus for 25 Years of NBA Basketball”, *The Spax blog*, <https://www.thespax.com/nba/calculating-regularized-adjusted-plus-minus-for-25-years-of-nba-basketball/> .
- Clemens, A., (2015). “Nylon Calculus 101: Plus-Minus and Adjusted Plus-Minus”, *FANSIDED blog*, <https://fansided.com/2014/09/25/glossary-plus-minus-adjusted-plus-minus/> .
- Erler, M., (2014). “The problem with ESPN's Real Plus-Minus”, *SBNATION blog*, <https://www.poundingtherock.com/2014/4/8/5594238/problem-with-real-plus-minus> .
- Goldstein, J., (2018). “Nylon Calculus: Introducing Player Impact Plus-Minus”, *FANSIDED blog*, <https://fansided.com/2018/01/11/nylon-calculus-introducing-player-impact-plus-minus/> .
- Goldstein, J., (2018). “PLAYER IMPACT PLUS-MINUS”, *BBALL-INDEX blog*, <https://www.bball-index.com/player-impact-plus-minus/> .
- Hamilton H., (2010). “Some starting documents on the plus/minus problem”, *SOCCEMTRICS RESEARCH blog*, <https://www.soccermetrics.net/player-performance/some-starting-documents-on-the-plusminus-problem-2> .
- Hamilton H., (2014). “Adjusted Plus/Minus in football – why it’s hard, and why it’s probably useless”, *SOCCEMTRICS RESEARCH blog*, <https://www.soccermetrics.net/player-performance/adjusted-plus-minus-deep-analysis> .
- Ilardi, S., (2014). “The next big thing: Real Plus-Minus”, *ABC News blog*, <https://abcnews.go.com/Sports/big-thing-real-minus/story?id=23226009> .
- Ilardi, S., and Barzilai, A., (2007). “Adjusted Plus-Minus: An Idea Whose Time Has Come”, *82games*, <http://www.82games.com/ilardi2.htm> .



- Ilardi, S., and Barzilai, A., (2008). “Adjusted Plus-Minus Ratings: New and Improved for 2007-2008”, 82games, <http://www.82games.com/ilardi2.htm>.
- Jacobs, J., (2017). “Deep Dive on Regularized Adjusted Plus Minus II: Basic Application to 2017 NBA Data with R”, *Squared Statistics: Understanding Basketball Analytics* blog, <https://squared2020.com/2017/09/18/deep-dive-on-regularized-adjusted-plus-minus-ii-basic-application-to-2017-nba-data-with-r/>.
- Jacobs, J., (2017). “Deep Dive on Regularized Adjusted Plus-Minus I: Introductory Example”, *Squared Statistics: Understanding Basketball Analytics* blog, <https://squared2020.com/2017/09/18/deep-dive-on-regularized-adjusted-plus-minus-i-introductory-example/>.
- Jacobs, J., (2018). “Regularized Adjusted Plus-Minus Part III: What Had Really Happened Was...”, *Squared Statistics: Understanding Basketball Analytics* blog, <https://squared2020.com/2018/12/24/regularized-adjusted-plus-minus-part-iii-what-had-really-happened-was/>.
- Johnson, A., (2014). “Introducing Player Tracking Plus Minus”, *Counting The Baskets* blog, <https://counting-the-baskets.typepad.com/my-blog/2014/09/introducing-player-tracking-plus-minus.html>.
- Myers, D., (2014). “Introducing Box Plus/Minus (BPM)”, *Sports Reference* blog, <https://www.sports-reference.com/blog/2014/10/introducing-box-plusminus-bpm-2/>.
- Myers, D., (2020). “A Review of Adjusted Plus/Minus and Stabilization”, *Chosen Stats (Carefully Chosen Sports Stats)* blog, <http://godismyjudgeok.com/DStats/2011/nba-stats/a-review-of-adjusted-plusminus-and-stabilization/>.
- Myers, D., (2020). “About Box Plus/Minus (BPM)”, *Basketball Reference* blog, <https://www.basketball-reference.com/about/bpm2.html>.
- Myers, D., (2020). “ASPM and VORP”, *Chosen Stats (Carefully Chosen Sports Stats)* blog, <http://godismyjudgeok.com/DStats/aspm-and-vorp/>.
- Narsu, K., and Cranjjs, T., (2021). “LEBRON: THE MAN, THE MYTH, THE METRIC?”, *BBALL-INDEX* blog, <https://www.bball-index.com/lebron-introduction/>.
- Paine, N., (2009). “Predicting with Statistical Plus/Minus”, *Basketball Reference* blog, <https://www.basketball-reference.com/blog/index48d3.html?p=1351>.
- Silly, A., (2011). “Deconstructing a Model (Advanced Plus Minus)”, *Arturo's Silly Little Stats v2.0* blog, <https://arturogalletti.wordpress.com/2011/03/04/deconstructing-a-model/>.
- Silver, N., (2015). “We’re Predicting The Career Of Every NBA Player. Here’s How.”, *FiveThirtyEight* blog, <https://fivethirtyeight.com/features/how-were-predicting-nba-player-career/>.
- Silver, N., (2019). “How Our RAPTOR Metric Works”, *FiveThirtyEight* blog, <https://fivethirtyeight.com/features/how-our-raptor-metric-works/>.
- Silver, N., (2019). “Introducing RAPTOR, Our New Metric For The Modern NBA”, *FiveThirtyEight* blog, <https://fivethirtyeight.com/features/introducing-raptor-our-new-metric-for-the-modern-nba/>.

