

# **Global Networking:**

## **Standardizing International Basketball Statistics with a Network-Based Approach**

Alexander Beard

Advised by:

Michael J. Tomas III

Barbara Arel

## Abstract

European basketball is different than the NBA in that teams play much different schedules, due to the nature of their scheduling format. Because of this, it's difficult to analyze player statistics across leagues. By building an incomplete competition network from each player's actual schedule, and using supervised learning to complete it by predicting player performance in hypothetical matchups, the competition network can be completed, allowing for player statistics to be standardized across all leagues in the future.

## Introduction

The table below shows Jimmer Fredette's statistics from the 2014-2015 season (Jimmer Fredette NBA, RealGM), his last playing a significant amount of games in the NBA , and the 2016-2017 season (Jimmer Fredette International, RealGM), his first season playing in China after a brief stint on the Knicks.

Year	Team	MPG	FG%	PPG	RPG	APG
2014-2015	New Orleans	10.2	38.0	3.6	0.8	1.2
2016-2017	Shanghai Dongfang	40.5	47.4	37.6	8.2	4.2

How did Fredette's statistics improve so much? Did he develop that much in the 18 months in between? Did his five minutes in two games in New York, where he now has the franchise record for effective FG% after hitting his only shot and going 4/5 from the line (Fox Sports), inspire him to become a different player entirely? Unlikely.

What's really at play here is the difference in strength of schedule. Differences in strength of schedule are sometimes mentioned in the NBA, but the difference in skill level isn't that large (no matter what some tanking teams might make you think), and regardless – each team plays roughly the same schedule (NBA Stuffer). European basketball is likely quite different, for two reasons – there are approximately three times as many teams (so teams play less similar schedules), and the European schedule format is quite different than the NBA's (which will be discussed in greater detail in the Data Methodology section).

These differences make it difficult to take European statistics at face value. It would be extremely easy to compare statistics between players if they had all played every one of the 88 other teams in the league. To approximate this, I attempted to predict the player vs team matchup results that didn't

happen in a given year, which could then, when every player has played the “same” schedule, allows standardized statistics to be calculated.

## Literature Review

There are four international leagues that I have data for: The Greek Basketball League, the Spanish Liga ACB, the Italian Lega Basket Serie A, and the French LNB Pro A. Each of these four leagues is the highest level of competition in their respective countries.

In addition to games within leagues (which is usually referred to as club play), I have data for EuroCup and EuroLeague, which are tournaments where the participants are teams from different leagues. In international competition, EuroLeague is regarded as the top-tier tournament, and EuroCup as the second-tier (Fraschilla, 2017). Both tournaments have a regular season component, which happens throughout the year, in addition to club play. The top performing teams in the regular season advance to the playoffs (EuroCup also has an intermediate round robin component), where a winner is decided with a head-to-head bracket (7Days EuroCup Format, Turkish Airlines EuroLeague Basketball Format).

Previous papers and articles have discussed the challenges of analyzing international play. For example, Motomura (2014) examined the NBA’s history of drafting international player. This article is relevant because it illustrates the ultimate goal of this work, which is to look into the performance of international players. Their input is draft position, but in future work using standardizing international statistics, they could either be placed as a comparison to draft position (which one predicts NBA performance better), or as an input, with draft position as an output, to essentially test the accuracy of the standardized statistics, assuming that the draft positions are with perfect information (which is obviously not true, but they might be close enough to provide a decent approximation).

Other authors have looked more specifically about international statistics, and how they can be adjusted to be a more accurate representation. Vashro (2015) compares leagues to other leagues, by optimizing the different relative competition scores to minimize the disagreement, then standardizing them to a z-score (mean of 0, standard deviation of 1). This is actually something I tried to do, but on a player-game level, but I found it intractable.

In 2014, Vashro’s work on projecting international players also adjusts for strength of schedule, but on a player level. His series talks a lot about the challenges of strength of schedule, and about the differences between statistics in different leagues. For example, with his league-wise adjustments, he adjusts (for example) the French league being pass-happier). An extremely relevant quote that’s really why I tried to use a separate network-based approach for each statistic:

“Strength of competition does not necessarily impact all skill-sets in the same way, and different leagues may inflate one statistic while deflating another. These factors mean that information is lost in any simple SOS adjustment.”

The next step that I take is not assuming that within leagues, the inflation/deflation for each statistic is the same. If there are significant differences between defenses within a league, that’s something my model should capture. Also, some players may have unique challenges or advantages against certain defenses, which is why looking at player performance as a individual player-game level competition network could improve the predictive power.

My network-based follows a similar approach to Park and Yook (2014). Their peer-reviewed paper is looking at a slightly different problem (inferring the ranks of nodes in a competition network, rather than the specific adjustments), but the general concept is similar. As was described in the introduction, we start out with an incomplete network, because not every player has played every team. To infer the hypothetical complete network, they infer the missing edges (as can be seen in Figure 1b of their paper). The main differences are what is trying to be predicted in the end (as mentioned), the fact that my competition network is a bipartite graph (players do not play other players specifically, they play defenses), and, of course, the specific application: While their main application is college football, mine is international basketball.

## Data Methodology

### Summary – before any pruning

Year	# Players	# Teams	# Matchups that happened	# Potential Matchups	Network Completion %
2014-2015	1563	104	34315	162552	21.11
2015-2016	1542	102	34746	157284	22.09
2016-2017	1378	89	30541	122642	24.90

### Pruning

The tournament games mentioned in the previous section are extremely important, because they allow players in different leagues to be compared, by providing connections between teams in different leagues Without those connections, there would be no way to project player performance against teams not in their league.

Figure 1 shows the full dataset, before any pruning. The gray dots in the graphs are teams that are not in any of the four leagues that Basketball Reference has club play boxscore data for, and they appear in

this dataset because their EuroLeague or EuroCup play is tracked. For this analysis, because teams with no club play data will have about half as much overall data, those teams will not be considered. Figure 2 shows the competition network with those teams removed.

### Summary – Teams with no club play data removed

Year	# Players	# Teams	# Matchups that happened	# Potential Matchups	Network Completion %
2014-2015	1042	66	24509	68772	35.64
2015-2016	1050	66	25061	69300	36.16
2016-2017	1032	65	22972	67080	34.25

One special case, as can be seen in Figures 1 and 2, is that the French teams are isolated from the rest of the teams in 2016-2017. In 2016, the French Basketball Federation made the decision to no longer participate in EuroCup or EuroLeague, and instead be a part of the FIBA Champions League (Stroggylakis, 2016). Because of this, the 2016-2017 French players and teams were discarded. Figure 3 shows the effect of removing those teams on the 2016-2017 competition network.

### Summary – 2016-2017 French teams removed

Year	# Players	# Teams	# Matchups that happened	# Potential Matchups	Network Completion %
2014-2015	1042	66	24509	68772	35.64
2015-2016	1050	66	25061	69300	36.16
<b>2016-2017</b>	<b>765</b>	<b>47</b>	<b>16916</b>	<b>35955</b>	<b>47.05</b>

There are some hypothetical player vs defense matchups where there's not enough information to predict the result. A minimum number of games played could be used, but there could be situations where the distribution of games played skews the amount of information available – it would be very difficult to predict hypothetical matchups for a player who has played all of his team's club games, but none of their EuroCup/EuroLeague games. To fix this, matchups were only predicted if they had at least 10 triangles (which will be explained later in this section). As you can see in the summary below, very few matchups had to be removed (around 50 each year), so there shouldn't be much of an effect on the end result.

## Summary – Player-defense matchups with under 10 triangles removed

Year	# Players	# Teams	# Matchups that happened	# Potential Matchups	Network Completion %
2014-2015	983	66	24449	64878	37.68
2015-2016	1012	66	25022	66792	37.46
2016-2017	716	47	16864	33652	50.11

## Scraping

All of this data is scraped from Basketball Reference. Matt Goldberg wrote a python scraping library for the NBA and NFL (Goldberg, Github). I forked it, and added international basketball support (Beard, Github). Figure 4 shows an example of a source webpage, and Figure 5 shows the resulting pandas DataFrame and the python snippet that scrapes it.

The articles in the Literature Review use season-level statistics, because for a long time, that's all that was available. However, Basketball Reference has game-level data available from the 2014-2015 season on, as can be seen in Figure 6.

By going through the the full schedule for each league, and pulling out each boxscore and concatenating it,, we obtain a list of each player's performances against each defense. Figure 7 shows the format of this dataset.

## Network Representation

This allows a weighted bipartite network to be constructed, where the two disjoint groups of nodes are players and defenses, and the edge weights are the resulting statistic for that matchup. For example, if Figure 8 was the dataset of player performances against defenses, Figure 9 would show the resulting bipartite network.

## Triangles

The main projection building block used is a triangle. For example, if we were trying to complete the network shown in Figures 8 and 9, the hypothetical amount of points that Player Y would score against Defense A would need to be calculated. While Player Y hasn't played Defense A, they have played Defense B, and scored 25 points. Therefore, we need some way to compare Defense A to Defense B. Luckily, Player X has played both Defense A and Defense B (see Figure 10).

Because Player X scored 22 points against Defense B, and 20 against Defense A, we assume that Defense A will allow  $(20/22) = 0.91 = 91\%$  of the points Defense B does to the same player. So, we would expect Player Y to score  $25 \times 0.91 = 22.75$  points against Defense A.

This prediction by itself would probably be very noisy. However, for each hypothetical matchup, there are multiple triangles that can be used to predict, which can be seen in the larger example in Figure 11. As you can see in the below table, to predict the result if Player Y faced Defense B (if it hadn't have actually happened), there would be four different triangles that could be used, and therefore, four different potential predictions:

### Player Y vs Defense B predictors

Triangle Components	Prediction
Defense A, Player W	5.75
Defense C, Player W	35.11
Defense A, Player X	5.71
Defense C, Player X	29.00

Because each individual prediction is so noisy, many predictions will have to be used.

## Triangle Selection/Ordering Methods

The obstacle once the triangle predictions are calculated is that to use traditional supervised learning methods, we need a properly shaped, rectangular matrix, but because each different matchups have different numbers of triangles, we actually have a jagged matrix, as can be seen in Figure 12. Also, we need some way to order the triangles, so they can then be made into a properly shaped matrix.

The ordering method used was number of games played. As can be seen in the Triangles section, a triangle is made up of three matchups – in the example shown in Figure 8 and 9, the three matchups that make up the prediction are Player X vs Defense A, Player X vs Defense B, and Player Y vs Defense B. If Player X has played Defense A three times, Player X has played Defense B two times, and Player Y has played Defense B one time, the total number of games played for this triangle is  $3 + 2 + 1 = 6$  games. This ordering is used because as there is more data used to make a prediction, the prediction should be less noisy.

Once the triangles were ordered, they need to be turned into a consistent set of predictions or aggregate predictions. Again operating on the assumption that more data is better, instead of only taking the top 5 or 10 predictions and leaving out the rest, the set of triangles is split into either 5 or 10 buckets. For example, if there were 100 triangles, split into 5 buckets, the first bucket would be the top 20 triangles by number of games played, the second bucket would be triangles 21-40, and so on. Because, in the

Pruning section, matchups with less than 10 triangles are removed, either method will work for any hypothetical matchup.

Once the triangles were split into buckets, either the median or mean of the predictions from the triangles in that bucket were calculated. Combining the two decisions, there were four possible treatments:

- 5 buckets, mean
- 5 buckets, median
- 10 buckets, mean
- 10 buckets, median

In addition to these 5 or 10 data points, the player's mean and median overall average was also considered (the matrix structure for the 5 buckets/mean treatment is shown in Figure 13). The null model, for comparison on if the additional data helped at all, was just the player's mean and median overall average.

## **Models**

The two prediction models I used were Linear Regression, and Random Forest Regressor, both from the scikit-learn package (Pedregosa et al, 2011). As you can see in Figure 14, a Linear Regression fits a line to the y-data given, using the features given in the x-matrix. It's best for data with linear relationships between the features and the output, and it also is very transparent, which is why it was chosen.

The other model chosen is a good contrast for the Linear Regression – Random Forests are not very transparent, but generally have higher performance. An example of the Random Forest Regressor can be seen in Figure 15.

## **Statistics**

The five statistics I used were points, total rebounds (offensive + defensive), assists, steals, and blocks, all adjusted to be per minute. These five were chosen instead of a more holistic, overall performance statistic like Game Score because, as spoken to in the Literature Review section, defenses could be stingy in regards to rebounds, but allow many assists. Because all triangles were being considered, and therefore players across archetypes were being used to compare, this separation is important. Therefore, for each year, and statistic, there is a separate network.



## Cross-Validation

For each of the three iterations, a random selection of 5% of the data was “hidden” for the model training. A new network was built *without* this data. The matrices used for training the model were then built, with each row being a matchup that actually happened that had not been hidden. The triangles were generated from the new network without the hidden data. After the model was trained with this training data, the model was then used to predict the 5% of matchups that had been hidden.

All of the statistics in the results are calculated from the predictions of the hidden data vs the actual values. The main statistics used were mean absolute error (MAE), and root mean squared error (RMSE). Because the different statistics can have very different averages (for example, points per minute vs blocks per minute), the MAE and RMSE were also presented scaled by the average actual value among the hidden set. To make the tables easier to read, these statistics were averaged across the three iterations – as can be seen in the two examples in Figure 16, most of the difference between iterations is the distribution of the actual values (which is variation that we want to average out)- the slope of the residuals vs actual values doesn’t change much.

## Results

As can be seen in Figures 18 through 21, looking just at the linear regression variations, there’s very little difference between the null model and the other variants with additional information. So little that it’s almost definitely statistically insignificant. What’s interesting is that in almost every case, every linear regression variation performs better than almost every random forest variation.

Looking specifically at the random forest variations, as you can see in both:

null → median 5 → median 10

and: null → mean 5 → mean 10

It almost always improves with more data (since it technically contains the same data, this could be also stated as more precise data). Also, for the most part, random forest methods median variations outperform random forest mean methods, which speaks to how noisy the data is.

## Conclusion

While the predictions are definitely not good, they’re somewhat encouraging – given how noisy a particular matchup is, especially with players and teams that don’t play each other often (or even have

many of the same opponents), predicting points per minute within approximately 40% (which is what a 0.4 scaled MAE is equivalent to) could be a lot worse.

Even though the random forest variations are less accurate, the fact that they consistently get better with more information suggests there's potential. Increasing the number of maximum features (considered at each split of the tree that makes up the forest), or feeding it more data somehow, looks like it would continue to improve performance. One possible way to have more data to train on would be to combine years, in conjunction with some kind of age curve to adjust for players aging into or out of their prime.

Part of the noise might be from the fact, mentioned in the Data Methodology section, that all archetypes are thrown together. One way to add more data would be to add on more features that were more of a drill down into similar players: buckets ordered by closest player heights, same player position. The same could also be done on the defensive side – similar teams, either by league, win/loss record, playoff performance, etc. could be weighted higher, or at least considered separately so that the model could decide whether they should be weighted higher.

## Possible Extensions

The clear next step would be to use the best-performing model to predict the matchups that actually haven't happened, actually doing what was simulated with cross-validation. With a completed competition network, or at least mostly complete (accounting for the pruning), a player's performances across all matchups (both real and projected) could be averaged, giving them standardized per minute statistics.

With these standardized statistics in hand, there are many avenues for future analysis. Which player archetypes are over/under drafted compared to their statistics, standardized statistics vs draft position as an indicator for NBA success, which statistics translate best to the NBA... the list could keep going on. Because the predictions were so inaccurate, it seems like there's no real point in standardizing the statistics now, but with more work on the prediction side, it's definitely a possibility that they could be accurate enough to perform future analysis with.