

---

# CS542 Final Project: Data Analysis of NBA Statistics

---

## Authors

Yichen Mu ycmu@bu.edu  
Fangxu Zhou fxzhou@bu.edu  
Jingxuan Guo krysguo@bu.edu  
Hao Fu haofu042@bu.edu

## Abstract

1 The goal of this paper was to analyze the NBA Statistics in some machine learning  
2 ways. By saying analyzing, we tried to predict NBA results of 2003 based on  
3 related statistics of players, coaches and results of season games of previous years.  
4 Also, we tried to figure out the outliers of players using the players' statistics.  
5 Besides, for specially exercising unsupervised learning technique, we clustered the  
6 players with statistics of their performance during season games. For the prediction  
7 part, we used techniques **KNN**, **SVM** and Decision Tree. We applied the **DBSCAN**  
8 clustering method to find the outliers. Then with the very common **Kmeans++** we  
9 cluster the players to explore this famous unsupervised learning technique.

## 10 ACKNOWLEDGEMENT

11 At the very beginning of this report, we would like to express our thanks to the professor **Peter Chin**  
12 and two teaching fellows **Peilun Dai** and **Andrew Wood** for giving us such an opportunity to try out  
13 the techniques we learned and improve ourselves.

## 14 1 Introduction

15 Basketball is all the time a fashionable sport all around the world. Machine learning is becoming  
16 more and more trendy nowadays. Therefore, that would be a nice try to combine these two elements  
17 to see if there are some findings fun and useful. As four novices to the domain of machine learning,  
18 we have to admit that we were somewhat overwhelmed by its complexity and interdisciplinarity.  
19 However, in spite of all the obstacles we've encountered, as well as the mistakes we may have made,  
20 we still found this science fascinating and were happy for having had chosen this **CS542**, which  
21 offered us this chance to have a brilliant embark on this venture of machine learning.

22 The goal of this paper was to exercise some ML techniques we learned from **CS542** as well as those  
23 we saw from the internet. We used techniques like **K-Nearest-Neighbors(KNN)**, **Support Vector**  
24 **Machine(SVM)** and Decision Tree for the prediction of outcomes of 2004 NBA games with data  
25 from previous years up to 1950s. For detecting the outliers of players and clustering the players, we  
26 tried **PCA** respectively with the unsupervised learning techniques **Kmeans++** and **DBSCAN**.

27 The data collected for this project was only use for the **CS542** project. All the methods defined below  
28 were used only for the **CS542** project. The code and the write-up for the **CS542** project are available  
29 at [https://github.com/FangxuZHOU/CS542\\_prj.git](https://github.com/FangxuZHOU/CS542_prj.git)

30 All the code for this paper is available at [https://github.com/FangxuZHOU/CS542\\_prj.git](https://github.com/FangxuZHOU/CS542_prj.git)  
31

32 All the pictures or tables showed on this paper are available at [https://github.com/FangxuZHOU/](https://github.com/FangxuZHOU/CS542_prj.git)  
33 [CS542\\_prj.git](https://github.com/FangxuZHOU/CS542_prj.git)

## 34 2 Related Work

35 I guess that there must be loads of similar work that have been done before, but we managed to find  
36 only one, with surely totally different dataset. This one comes from Nadir N who wrote a paper  
37 entitled *Assessing NBA player similarity with Machine Learning (R)* published on *towards data*  
38 *science*. The goal of that paper was to use Unsupervised Machine Learning and **PCA** to determine  
39 which NBA superstars are the most alike. For the clustering technique, they used **Kmeans** together  
40 with **PCA**. Using data from the 2017–18 NBA Regular season combined three data-sets from the  
41 websites Basketball-reference and NBAMining.com as input, they try to cluster the players to see  
42 who are most alike. For the prediction part, pitifully, we haven't found any similar work done. Maybe  
43 it was because we were not able to do much research due to the limitation of time. However, I'm  
44 pretty sure the plenitude of such work similar to ours.

## 45 3 Dataset and Features

46 The data we used in this project is of *Project F: NBA statistics data* from the <http://www.cs.cmu.edu/~ggordon/10601/projects.html>, which is a publicly available website specially for  
47 purposes of exercising techniques learnt from machine learning courses. Basically, the data provided  
48 from the website contains 2004-2005 NBA and ABA stats for:

- 50 • Player regular season stats
- 51 • Player regular season career totals
- 52 • Player playoff stats
- 53 • Player playoff career totals
- 54 • Player all-star game stats
- 55 • Team regular season stats
- 56 • coaches\_season.txt - nba coaching records by season
- 57 • coaches\_career.txt - nba career coaching records

58 What we've used are the 1<sup>st</sup>, 2<sup>nd</sup>, 4<sup>th</sup>, 6<sup>th</sup> and the last two items

### 59 3.1 Prediction of the outcomes of season games

60 For this part, we used datasets *Player regular season stats*, *coaches\_season.txt* and *Team regular*  
61 *season stats* which respectively contain the performance stats of every player, coach, team in their  
62 regular season matches starts from 1946. Since what we need to do is to use the data to fit learning  
63 models to predict the outcome, we have to combine and filter the above datasets to construct a  
64 comprehensive dataset. So after merging them, we observed the data and created some more  
65 representative feature columns. After that, this main dataset we use for this part has 72 feature  
66 columns, which mainly are some original features like WON, LOST, GP, MINUTES and some new  
67 features like T\_o\_Effic, T\_d\_3P\_rate, Coach\_Effic, Player\_Effic and others we created based on the  
68 original features.

### 69 3.2 Outliers detection for NBA players

70 In this part, we used the datasets *Player regular season career totals* and *Player playoff career totals*  
71 which contain respectively the performance stats of each players' whole career of regular season  
72 games and of playoff games start from 1946. There were 21 input features: 1. The first four show the  
73 basic information of a player, like his id, his league, his first name and last name. 2. Other features  
74 can be used to evaluate performance of a player, like OREB (stands for offensive rebounds), DREB  
75 (stands for defensive rebounds), etc. (Those are the features we used for analyzing the outliers)

### 76 3.3 Finding outstanding Players

77 This part is only for trying out the **Kmeans++** technique, also, because of the lack of time, we decided  
78 to only use the dataset of *Player regular season career totals*. Features are like what are listed at the  
79 previous section.

## 4 Methods

### 4.1 Prediction of the outcomes of season games

#### 4.1.1 Cluster the data to add labels

Since there are some congenital defects in our original datasets, including player, team, and coaches datasets. They all lack the label data. With that problem, after merging, filtering, calculating and get a comprehensive dataset, we still couldn't implement Supervised Learning algorithm on our dataset. So after preprocessing the data, in order to gain labels, we decide to use **Kmeans++** clustering algorithm(Unsupervised Learning) to cluster the data and get centroids, then we can set labels according to the sum of feature values of centroids. In this way, there will be label data in our training and testing dataset. After finishing these works, we can then use the supervised learning algorithms like **KNN** on the training data later.

##### 4.1.1.1 Data normalization

After having tried other Normalization method like `StandardScaler()`, we found `MinMaxScaler` is optimal. Since there are more than 60 features and the scales of data are pretty different, and what we want to do is to scale the features between a specified maximum and minimum value, which in this case is between 0 and 1, that's why we use this `MinMaxScaler` method.

##### 4.1.1.2 Implementation procedure

We defined two functions which are `elbow` function and `create_label` function. In `elbow` function, we used **Kmeans++** algorithm and `matplotlib` package to plot an image which is related to the optimal number for the clusters. As for the aim of `create_label` function, we use it to calculate the label for every row of data in both training and testing dataset. After implementing the `elbow` function, we can get images as following.

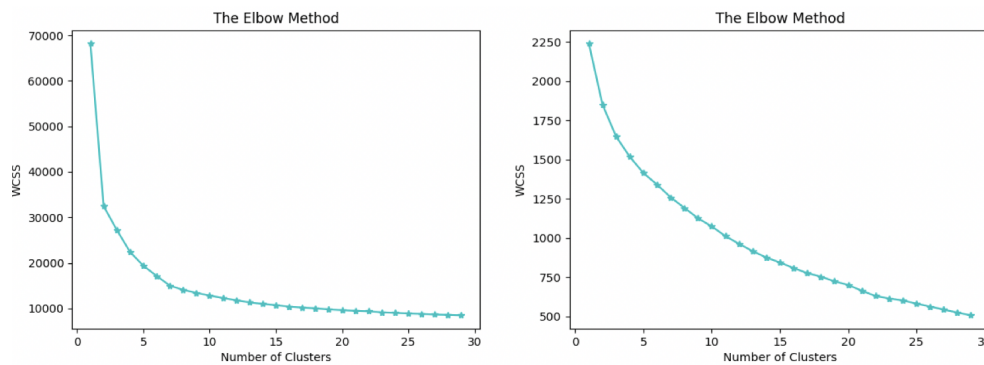


Figure 1: elbow

101

The left image is from the training data, and the right image is from the testing data. We can see that both two curves show an obvious inflection point when the Number of Clusters is equal to 5. So we set the 'n\_cluster' parameter in **Kmeans++** as 5 and implemented the **Kmeans++** algorithm on both two datasets. Then we can get the 5 centroids and then inverse and transform them into their original data. Then we can sum every centroid data and sort them in descending order to get 5 benchmarks(0 4). After that, we can set the label to every row of data based on the distance between the sum of every row of data and these 5 benchmarks. Then, we got the training dataset with labels and testing dataset with labels. As for the reason why we need to set labels to the testing dataset as well, since later we need to extract the labels in testing dataset as the control data to compare with the outcomes we will predict using the trained learning model to get the accuracy of the algorithm. Results which are training and testing dataset with labels will be shown in result section.

#### 4.1.2 Prediction of the outcomes of season games of 2003

Now we have preprocessed all the data and added labels to them, and we finally can implement the supervised learning algorithms on the data and see if we can predict the outcomes in 2003 correctly. As for this part, we used the `make_pipeline` method to construct a pipeline to deal with the data, the `make_pipeline` we used contains `minmaxscaler`, `Principal Component Analysis` and different learning algorithms. And as for the learning algorithms, we tried `K-NearestNeighbor`, `Support vector machine`, `decision tree` and `neural network`.

#### 4.1.2.1 Minmaxscaler

This has been mentioned in Data normalization in section 4.1.1.1

#### 4.1.2.2 PCA

This has been mentioned in section 4.2.2

#### 4.1.2.3 KNN

K-Nearest Neighbors is a method that classifies every record in a data set. The main idea of it is in the feature space, if most of the K most adjacent samples of a sample belong to a certain category, then the sample also belongs to that certain category and has the characteristics of sample in that category. And in determining classification decision, the method determines the category of sample according to the category of the closest one or several samples. The distance formula we used here is **minkowski**, which is

$$dist(X, Y) = \left( \sum_{i=1}^N |x_i - y_i|^p \right)^{1/p}$$

in which p is a variable. And we use sklearn package to implement this algorithm. The parameter is

#### 4.1.2.4 SVM

Support Vector Machine is a kind of generalized linear classifier that classifies data in a supervised learning way. Its decision boundary is the maximum-margin hyperplane solved for the learning sample. And **SVM** can implement the kernel function to do the non-linear classification. And the kernel function we used here is **RBF kernel**, which is And we use sklearn package to implement this

$$\kappa(\mathbf{X}_1, \mathbf{X}_2) = \exp \left( -\frac{\|\mathbf{X}_1 - \mathbf{X}_2\|^2}{2\sigma^2} \right)$$

algorithm. The parameters are kernel='rbf', class\_weight='balanced', C=3.0.

#### 4.1.2.5 decision tree

Decision tree algorithm is a method to approximate discrete function values. It uses the inductive algorithm to generate readable rules and decision trees, then these decisions are used to analyze the new data and classify the data through a series of rules. And the related parameters we used are *criterion = 'entropy', max\_depth = None, min\_samples\_split = 2, min\_samples\_leaf = 1, max\_features = None, max\_leaf\_nodes = None, min\_impurity\_decrease = 0*

#### 4.1.2.6 Neural Networks

Neural Networks is an algorithmic mathematical model that imitates the behavior characteristics of animal neural network and carries on distributed parallel information processing. It relies on the number of hidden layers in the network and the number of nodes on the hidden layers to adjust the complexity, and constructs an effective model for processing information by adjusting the interrelated relationship between a large number of nodes in the network. And in our project we implement the MLPClassifier in sklearn package. The main related parameters are solver='sgd', activation='relu', alpha=1e-4, hidden\_layer\_sizes=(300,300,300), random\_state=1, max\_iter=200, learning\_rate\_init=0.001, in which the 'sgd' means the sigmoid function which is

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

## 4.2 Outliers detection for players

This part, detecting the outliers of players, we removed first the players whose "minutes" column has value as 0, and take those players as outliers (We did so, because we thought it was quite reasonable to do so, a player cannot play 0 minute during his entire career and not have 0 as scores). We then divide each value of the features of the player by his total minutes. Normalizing by minutes gives us a fairer representation of each player's contributions. If we looked at game statistics instead, players who play more minutes would look better than players who play on elite teams and often don't play at the end of games that are blow-out wins. Then we normalized the features using **Z-Score**. We

continued by utilizing the **PCA** technique to reduce the dimensions. Subsequently, with **Silhouette Score**, we tried different combinations of minimum number of clusters and the value of epsilon to find the most appropriate one for the **DBSCAN** clustering method. For the next step, we clustered the players with **DBSCAN** and finally generated a list of outliers.

#### 4.2.1 Z-Score (standard score)

With Z-Score, we normalized our dataset, which means that for each value, we are going to subtract the mean of the feature-range from that value and divide it by the standard deviation of the feature. This ensures that features with high-value ranges do not have a greater impact on our overall similarity comparison than features with low-value ranges.

$$z = \frac{x - \mu}{\sigma}$$

where:  $\mu$  is the mean of the population.  $\sigma$  is the standard deviation of the population.

#### 4.2.2 Principal component analysis (PCA)

We decided to apply the **PCA** technique to do the dimension reduction. A sklearn library was used (Scikit-learn: Machine Learning in Python) to help us to deal with **PCA**. **PCA** reduce dimensions by Projecting each data point onto only the first few principal components to obtain lower-dimensional data while preserving as much of the data's variation as possible

#### 4.2.3 Silhouette Score

Silhouette Score is a metric used to calculate the goodness of a clustering technique. Its value ranges from -1 to 1. Here, we used it to find the most appropriate combination of minimum samples and epsilon value for **DBSCAN**. We used a method of the sklearn.metrics library to carry it out.

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i) - b(i)\}}$$

where:  $a(i)$  can be interpreted as a measure of how well the point  $i$  is assigned to its cluster (the smaller the value, the better the assignment).  $b(i)$  is the smallest mean distance of  $i$  to all points in any other cluster.

Its value ranges from -1 to 1:

- 1: Means clusters are well apart from each other and clearly distinguished.
- 0: Means clusters are indifferent, or we can say that the distance between clusters is not significant.
- -1: Means clusters are assigned in the wrong way.

#### 4.2.4 DBSCAN clustering

With a library of sklearn.cluster, we carried out this clustering technique. **DBSCAN**, whose full name is Density-based spatial clustering of applications with noise, given a set of points in some space, it groups together points that are closely packed together (points with many nearby neighbors), marking as outliers points that lie alone in low-density regions (whose nearest neighbors are too far away).

Unlike **Kmeans++** clustering, we don't have to pre-decide the number of clusters while using **DBSCAN**, the algorithm itself would take care of it.

But we do have to choose two parameters appropriately first, they are: *minPts*, the minimum number of points required to form a dense region, and  $\epsilon$  who specifies how close points should be to each other to be considered a part of a cluster.

### 4.3 Finding outstanding players

For this part, we used **Kmeans++**, since we had already tried **DBSCAN** at the previous part, we wanted to try a new one. And we are told by the internet that **Kmeans++** clustering technique is a very commonly used one. For the data pre-processing phase, it's almost the same as what we did for the "Outliers detection" part. Except that we removed the outliers we found at the outlier detection part. Then with the elbow method we tried to find the most appropriate number of clusters, which will be used later to cluster the points with **Kmeans++** technique. Introduction of these two methods, elbow method and **Kmeans++**, are given at the section 4.1

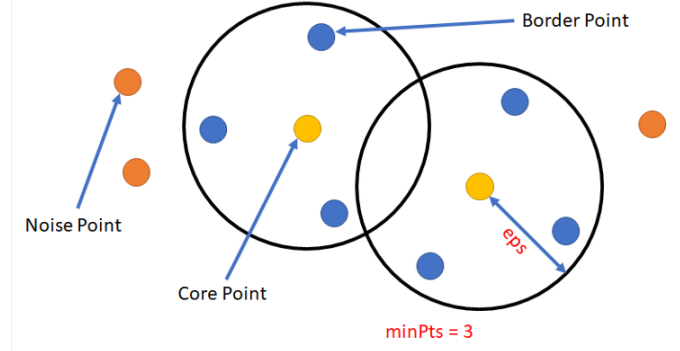


Figure 2: DBSCAN

## 5 Results

### 5.1 Prediction of the outcomes of season games

#### 5.1.1 Results of clusters and labels

After we implemented the **Kmeans++** clustering algorithms in section 4.1.1, we got five centroids and five benchmarks and we sorted them in descending orders. After that, we set label to every row of data based on the distance between the sum of every row of data and these 5 benchmarks. Then, we get the training dataset with labels and testing dataset with labels. Here due to the space limit, I just place the training data here as an example, as shown in Figure 3. And the testing data have the same structures. As for the meaning of labels, since there are 5 different labels which are

year	a_fgm	a_fga	a_ftm	a_fta	a_oreb	a_dreb	a_reb	a_ast	a_pf	a_stl	a_to	a_blk	a_3pm	a_3pa	a_pts	d_fgm	d_fga	d_ftm	d_fta	d_oreb	d_dreb	d_ast	d_pf	d_stl	d_to	d_blk	d_3pm	d_3pa	d_pts	pace	won	lost	yr_order
season_win	season_loss	gp	minutes	pts	oreb	dreb	reb	asts	stl	blk	turnover	pf	fga	fgm	fta	ftm	tpa	tpm	T_o_Effic	T_o_FG_rate	T_o_FT_rate	T_o_3P_rate	T_d_Effic	T_d_FG_rate	T_d_FT_rate	T_d_3P_rate	T_WIN_LOSE_rate	Coach_Effic	P_Season_Effic	P_Season_FG_rate	P_Season_FT_rate	P_Season_PTS_rate	
0	1946	1397	5133	811	1375	0	0	0	470	1202	0	0	0	0	0	3605	0	0	0	0	0	0	0	0	0	0	0	0	0	3605	0	0	0
	0	0	0	0	0	0	0	0	0	0	0	0	3900	0.000000	22	38	1	-225															
	22	38	6	0	13	0	0	1	0	0	0	0	15	22	5	4	3	0															
	0.272161	0.589818	0.000000	3900	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.366667	0.733333																				
	-0.666667	0.227273	0.750000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
1	1946	1397	5133	811	1375	0	0	0	470	1202	0	0	0	0	0	3605	0	0	0	0	0	0	0	0	0	0	0	0	0	3605	0	0	0
	0	0	0	0	0	0	0	0	0	0	0	0	3900	0.000000	22	38	1	-225															
	22	38	58	0	567	0	0	60	0	0	0	0	115	870	223	193	121	0															
	0.272161	0.589818	0.000000	3900	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.366667	0.733333																			
	-1.586207	0.256322	0.626943	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
	Name: Label																																

Figure 3: training dataset with labels

from 0 to 4, and we sorted them in descending order, we can define the label 0 as having a very high probability of winning and we define the remaining labels as the winning rates which are in descending order.

#### 5.1.2 Results of predictions of the outcomes of 2003

After trying different machine learning models to fit the data and then use them to predict the outcomes of 2003, finally we compare the result we predicted and the true results in the original testing dataset, we got the different accuracies as shown in Figure 4. So eventually we can tell that the Neural

When PCA (n\_components=0.98, svd\_solver='auto'),  
 KNN: RMSE on testing set = 1.3280510201777347  
 Decision Tree: RMSE on testing set = 1.3043080587787235  
 SVM: RMSE on testing set = 1.86429873038029

Figure 4: performances of different learning models



Figure 5: accuracy of NNs and outcome it predicts

217 Networks has the highest accuracy and it will be the optimal learning model to predict the outcome  
 218 of 2003. The accuracy of Neural Networks and outcome it predicted are shown in Figure 5.

## 219 5.2 Outliers detection for players

220 We tried a long list of combinations of *minPts* and  $\epsilon$ , we computed respectively the silhouette score  
 221 of **DBSCAN** clustering using each of the combinations. Then we chose the one with the largest score,  
 which is shown in the Figure 6.



Figure 6: Silhouette score of the most appropriate *minPts* and  $\epsilon$

222

223 With the combination of *minPts* and  $\epsilon$  gained, we clustered with **DBSCAN**, and got a plot as shown  
 224 in Figure 7.

225 Then, as presented in Figure 8, we came up a list of players who correspond to the black circles in  
 226 the plot, and who therefore are considered to be outliers. We haven't yet analyzed why these players  
 227 are set to be outliers since we ran out of time. But at a very first glance of the list, we noticed that  
 228 those players, compared to others, they have very small values as their "minutes" features.

229 So we made a bold guess, maybe it can be one of the reasons: A player is taken as an outlier maybe  
 230 it's partly because he played too less minutes during his entire career. Surely, there would be other  
 231 reasons, but for now, we've figured only this one out.

## 232 5.3 Finding outstanding players

233 As what we have said before, this part is done with **Kmeans++**. With **Elbow** method, showed in  
 234 Figure 9, we found the appropriate number of clusters is 10. Then we clustered and got a plot like  
 235 what you can see in Figure 10. Just to let you know, there is a plot with names of players attached in  
 236 the Annex section. Still, we don't have time to finish this part, we haven't done the analysis of the  
 237 main characteristics of each cluster, so that we won't be able to figure out players of which cluster  
 238 can be seen as outstanding.

239 But we made a guess, the cluster that is the farthest away from the others may be the most outstanding  
 240 one, or the least. We can continue this part as future work later.

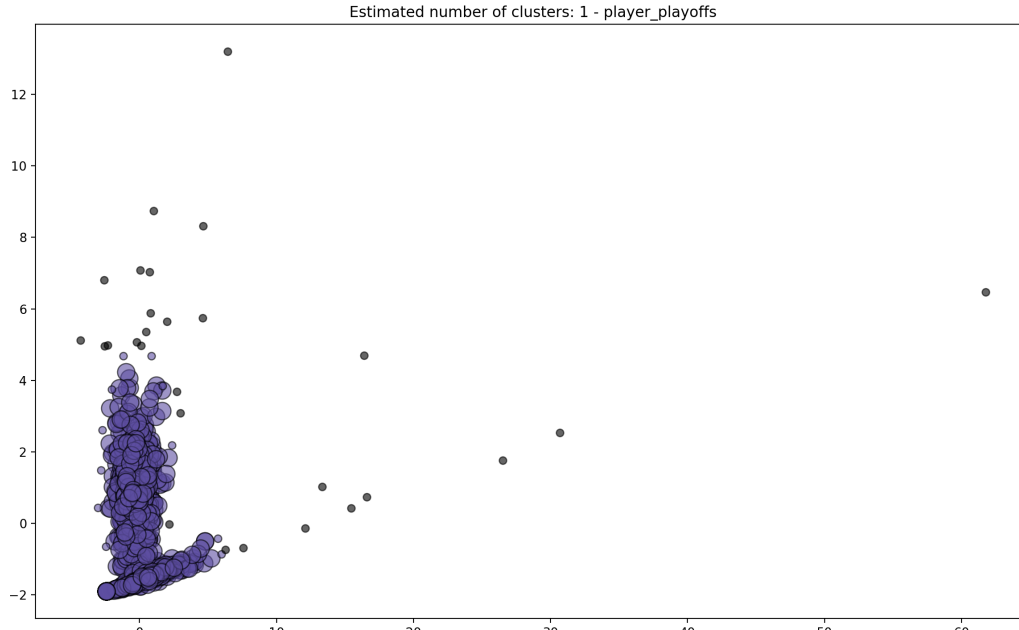


Figure 7: Plot of DBSCAN clustering

lbid	firstname	lastname	gp	minutes	pts	dreb	oreb	reb	asts	stl	blk	turnover	pf	fga	fgm	fta	ftm	tpa	tpm		
77	BREWEJA01	Jamison	Brewer	2	6	-0.68038635	3.54665183	6.44700639	3.88675404	0.83867720	-0.65438045	-0.36398182	3.20388489	2.69381024	-0.11531909	-1.18969088	1.44455845	1	-0.41238581	-0.2403	
78	BRITTM01	Mike	Brittain	1	2	1.76083703	-0.81729117	9.99981290	2.52086339	-0.82263293	-0.65438045	-0.36398182	1.91887527	2.69096736	-0.81354132	0	-0.41238581	-0.2403	0	-0.41238581	-0.2403
79	BROOKKE01	Kevin	Brooks	2	5	1.76083703	1.80107462	3.60476117	1.70132899	-0.82263293	-0.65438045	-0.36398182	-0.77151527	-1.02339191	3.13939189	1.91483571	1.89617841	1	2.83202264	-0.2403	
80	CERVIALD1	Al	Cervi	27	116	4.28624054	-0.81729117	-0.65860664	1.46112063	8.28592951	-0.65438045	-0.36398182	-0.77151527	5.51375671	3.04995748	2.15570415	7.01194240	116	-0.41238581	-0.2403	
81	CLOSSB01	Bill	Closs	11	21	13.1997123	-0.81729117	-0.65860664	4.66726299	10.0945479	-0.65438045	-0.36398182	-0.77151527	10.30523411	16.73943421	12.1154231	11.4447446	31	-0.41238581	-0.2403	
82	COOPERDU1	Duane	Cooper	2	4	-1.16863103	-0.81729117	9.99981290	2.52086339	1.66933226	-0.65438045	-0.36398182	-0.77151527	-1.02339191	3.44452105	-1.18969088	-0.81354132	0	11.7541459	-0.2403	
83	EDMONKE01	Keith	Edmonson	1	2	1.76083703	-0.81729117	9.99981290	2.52086339	4.16129746	-0.65438045	-0.36398182	-0.77151527	-1.02339191	3.9322949	2.69096736	-0.81354132	0	-0.41238581	-0.2403	
84	HAYEST01	Steve	Hayes	1	1	4.69030510	-0.81729117	-0.65860664	6.61853535	-0.82263293	-0.65438045	-0.36398182	-0.77151527	-1.02339191	1.91887527	6.57162561	-0.81354132	0	-0.41238581	-0.2403	
85	HEALSH01	Shane	Heal	2	3	4.69030510	-0.81729117	-0.65860664	-1.57680857	2.49998733	-0.65438045	-0.36398182	-0.77151527	-1.02339191	0.90177808	3.98452011	-0.81354132	0	10.4023090	19.7999	
86	HENRYCA01	Carl	Henry	1	2	3.22557106	-0.81729117	-0.65860664	-1.57680857	-0.82263293	-0.65438045	-0.36398182	-0.77151527	-1.02339191	0.39322949	2.69096736	-0.81354132	0	7.69863533	14.7898	
87	HOLLAD01	Joe	Holland	21	37	6.66967542	0.24420847	0.49365492	1.18861971	7.25941635	4.25279051	-0.36398182	1.80712267	8.42031086	7.77405636	8.03998278	3.39750962	11	-0.41238581	-0.2403	
88	HOLZMR01	Red	Holzman	24	79	2.61372014	-0.81729117	-0.65860664	1.12039322	2.07940248	-0.65438045	-0.36398182	-0.77151527	1.23516129	2.92309780	1.54360059	3.13098740	26	-0.41238581	-0.2403	
89	JOHNSAR01	Arnie	Johnson	22	166	2.64320693	-0.81729117	-0.65860664	6.17420948	3.32063450	-0.65438045	-0.36398182	-0.77151527	3.72387831	1.55124978	1.70911407	4.16516060	92	-0.41238581	-0.2403	
90	JONESWA02	Wallace	Jones	6	8	26.2951321	-0.81729117	-0.65860664	-1.57680857	26.5898842	-0.65438045	-0.36398182	-0.77151527	24.9970232	27.8548534	21.1240940	27.9772309	29	-0.41238581	-0.2403	
91	JUOKIED1	Jeff	Judkins	7	10	1.46789022	0.49189172	5.73644508	1.70132899	-0.82263293	2.97950758	-0.36398182	-0.77151527	-1.02339191	1.30861696	1.91483571	-0.81354132	0	4.45422687	2.76571	
92	LAUDPR01	Priest	lauderdale	3	7	-1.16863103	1.05297011	2.38665608	0.76471826	-0.82263293	-0.65438045	-0.36398182	9.45094231	0.03866584	0.17528009	-1.18969088	-0.81354132	0	-0.41238581	-0.2403	
93	MCINHO01	Horace	McKinney	7	20	5.56914552	-0.81729117	-0.65860664	4.97946657	5.65647658	-0.65438045	-0.36398182	-0.77151527	9.75649434	11.1200076	1.8355978	2.57360834	8	-0.41238581	-0.2403	
94	MEYERLO1	Loren	Meyer	3	14	-1.16863103	2.92323140	2.38665608	1.93548168	-0.11064287	-0.65438045	-0.36398182	3.91087912	4.33971351	-1.02339191	-0.26061869	-1.18969088	-0.81354132	0	0.74633149	-0.2403
95	MOSLEG01	Glenn	Mosley	3	6	1.27259235	-0.81729117	-0.65860664	-0.21091792	0.83867720	-0.65438045	-0.36398182	-0.77151527	-1.02339191	3.9322949	1.39741461	2.57360834	1	-0.41238581	-0.2403	
96	MURRAKE01	Ken	Murray	6	15	6.64328381	-0.81729117	-0.65860664	6.07217906	5.82260759	-0.65438045	-0.36398182	-0.77151527	5.41975849	10.66659110	7.60646781	2.34779836	6	-0.41238581	-0.2403	
97	RANDAMA01	Mark	Randall	2	6	-1.16863103	7.91059484	2.89419987	5.25264470	-0.82263293	-0.65438045	-0.36398182	3.20388489	0.21567547	-0.62386768	-1.18969088	-0.81354132	0	-0.41238581	-0.2403	
98	RATKOG01	George	Ratkovic	24	59	11.6912372	-0.81729117	-0.65860664	10.23004285	5.59734520	-0.65438045	-0.36398182	-0.77151527	10.8212522	8.95236088	9.33412809	16.17961634	99	-0.41238581	-0.2403	
99	RAUTLED1	Leo	Rautins	3	5	0.58904980	-0.81729117	7.86812899	1.70132899	1.17091922	6.61339561	-0.36398182	-0.77151527	1.95036981	0.69835865	0.36257241	-0.81354132	6	0.0743110	5.71775	
100	SCHFEED01	Tom	Scheffer	3	10	0.88199661	1.80107462	5.73644508	2.52086339	-0.82263293	2.97950758	-0.36398182	1.61372482	-1.02339191	-0.21702881	0.36257241	1.89617841	3	-0.41238581	-0.2403	
101	STROED01	John	Stroeder	1	1	1.71977317	-0.81729117	-0.65860664	-1.57680857	-0.82263293	-0.65438045	-0.36398182	-0.77151527	-1.02339191	1.91887527	6.57162561	-0.81354132	0	5.8096564	29.8200	
102	VAUGHDA02	David	Vaughn	1	1	-1.16863103	-0.81729117	-0.65860664	-1.57680857	-0.82263293	-0.65438045	-0.36398182	-0.77151527	-1.02339191	1.13241628	-1.18969088	-0.81354132	0	-0.41238581	-0.2403	
103	WARDH01	Henry	Ward	1	1	10.5492412	-0.81729117	-0.65860664	-1.57680857	-0.82263293	-0.65438045	-0.36398182	-0.77151527	-1.02339191	8.02145838	14.3329421	-0.81354132	0	-0.41238581	-0.2403	
104	WHEATED01	Dejuan	Wheat	1	3	0.78434767	3.54665183	-0.65860664	1.54972731	-0.82263293	11.4585796	-0.36398182	-0.77151527	-1.02339191	0.90177808	1.39741461	-0.81354132	0	-0.41238581	-0.2403	

Figure 8: List of outliers

## 6 Conclusion and Future Work

During this project, we tried to do the outliers detection, prediction of outcomes, with techniques like DBSCAN, Kmeans++, KNN, SVM, Neural Networks, etc. And eventually we got a pretty good result. We are really happy that we, every member of our team, never gave up and kept trying and learning new things though they are kind of hard to us. And meanwhile, during this project, we learned so many techniques and gained much project experience. For the future work, as what we have said, we can try to finish the analysis part of our unsupervised learning part. We can continue to try to figure out why the outliers are outliers, which cluster of players among all the clusters can be seen as outstanding or not that outstanding. Besides, in the future, we can try to collect more data like stats which are about the win-loss situation between every pair of two teams. With these data, we can set them as labels and I guess maybe we can get more valuable and concrete predictions. Although there are still some details to retouch, we've gained a lot during the process and we are proud of what we have done. We think we had a pretty good start in machine learning.



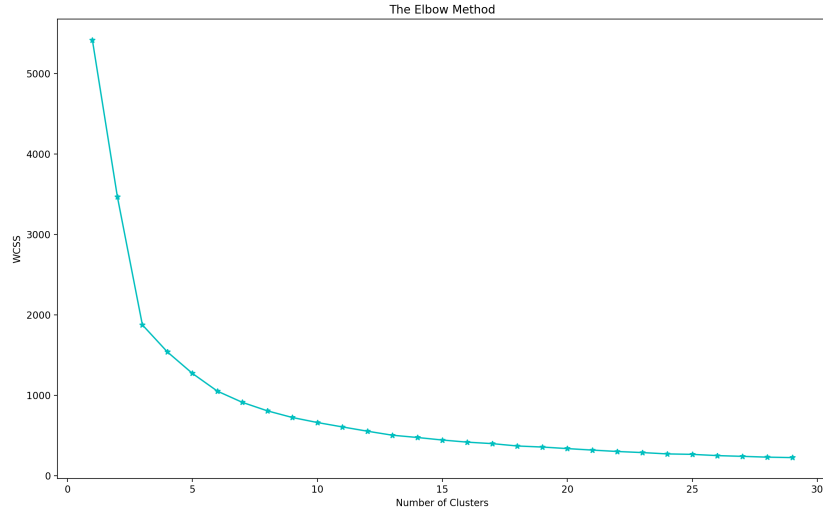


Figure 9: Elbow for finding the appropriate number of clusters

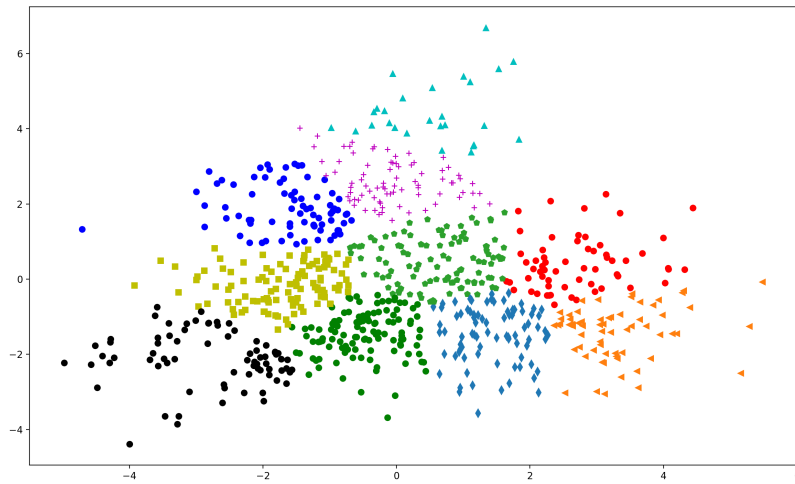


Figure 10: Clustering of players

## 7 Contributions

All of the work described above was done by our group members: Yichen Mu, Fangxu Zhou, Fu Hao, Jingxuan Guo.

## References

- [1] Nadir, N. (2018) *Assessing NBA player similarity with Machine Learning (R)*
- [2] Sebastian, R. (2018) *An overview of proxy-label approaches for semi-supervised learning*

## Annexe

RkId	first_name	last_name	gp	minutes	pts	reb	orb	reb	asts	stl	blk	turnover	pf	fga	fgm	fta	ftm	tpa	tpm	
BRUYLAND	Jason	Brenner	2	0	0.000000	5.34663383	0.44700000	0.88675406	0.00677201	0.0543804	-0.3038182	0.20084897	0.00600200	-0.1331359	1.0849008	0.44450805	1	-0.41238081	-0.2403266119363793	
BRITTMOL	Mike	Brittain	1	2	1.76081703	-0.81729117	5.99981290	2.52066390	0.82261293	0.0543804	0.145803146	-0.77151527	-0.02339191	0.184807272	0.69067363	-0.81354132	0	-0.41238081	-0.2403266119363793	
BROOKERKE	Kevin	Brooks	2	5	1.76081703	1.80376426	1.60476176	1.70132896	-0.82261293	0.0543804	-0.3038182	-0.77151527	-0.02339191	0.130931884	0.54403754	0.896178412	1	1.83022646	-0.2403266119363793	
CERVALLO	Al	Cone	27	116	4.28624044	-0.81729117	-0.03800664	0.461120058	0.80202516	0.0543804	-0.3038182	-0.77151527	0.113704713	0.040951487	0.155704517	0.01244505	116	-0.41238081	-0.2403266119363793	
GLOSBOIS	Bill	Coss	11	21	15.10971235	-0.81729117	-0.03800664	4.46726290	10.09454793	0.0543804	-0.3038182	-0.77151527	10.30524215	16.73945425	12.11542311	11.44474464	31	-0.41238081	-0.2403266119363793	
COOPERUD	Duane	Cosper	2	4	1.16601109	-0.81729117	5.99981290	2.52066390	1.69313267	0.0543804	-0.3038182	-0.77151527	-0.02339191	0.04412105	1.0849008	-0.81354132	0	0.17441491	-0.2403266119363793	
EDMONDSON	Keith	Edmondson	1	2	1.76081703	-0.81729117	5.99981290	2.52066390	0.161297464	0.0543804	-0.3038182	-0.77151527	-0.02339191	0.393228495	0.69067363	-0.81354132	0	-0.41238081	-0.2403266119363793	
HAVESTOT	Steve	Hayes	1	1	4.690305104	12.74031984	-0.03800664	6.03831016	0.82261293	0.0543804	-0.3038182	-0.77151527	-0.02339191	0.181807572	0.57621612	-0.81354132	0	-0.41238081	-0.2403266119363793	
HEALDRE	Shane	Heal	2	3	4.690305104	-0.81729117	-0.03800664	1.57640807	7.499667312	0.0543804	-0.3038182	-0.77151527	-0.02339191	0.06177808	0.94620111	-0.81354132	0	0.40220902	13.7992820799356	
HENYCARD	Carl	Henry	1	2	3.22571049	-0.81729117	-0.03800664	-1.57640807	0.82261293	0.0543804	-0.3038182	-0.77151527	-0.02339191	0.393228495	0.69067363	-0.81354132	0	0.40803138	14.7898034949676	
HOLLADAY	Jim	Holland	21	97	6.609175412	0.442404705	0.933640182	1.0861317	7.239413554	0.91677905	-0.3038182	0.180712874	0.40203889	7.74050463	0.09967793	0.39705045	17	-0.41238081	-0.2403266119363793	
HOLZMIRE	Red	Holzman	24	79	2.61370141	-0.81729117	-0.03800664	1.120392232	0.79402486	0.0543804	-0.3038182	-0.77151527	1.235161296	2.02097809	0.54000591	1.13087408	26	-0.41238081	-0.2403266119363793	
JOHNSARD	Arnie	Johnson	22	166	2.64420094	-0.81729117	-0.03800664	0.17420942	1.02634058	0.0543804	-0.3038182	-0.77151527	3.73877811	1.51149784	1.70714074	1.16310064	92	-0.41238081	-0.2403266119363793	
JONESWARD	Walter	Jones	6	8	0.2913212	-0.81729117	-0.03800664	1.57640807	70.18894614	0.0543804	-0.3038182	-0.77151527	24.90702327	27.84451482	1.12490460	27.97733091	29	-0.41238081	-0.2403266119363793	
JONESJEE	Jeff	Judkins	7	10	1.40780126	0.401891727	1.78440281	1.70132896	-0.82261293	0.0543804	-0.3038182	-0.77151527	-0.02339191	1.108616962	0.54403754	-0.81354132	0	0.45422877	2.7671449499444	
LAUDERPEL	Pratt	Lauderdale	3	7	1.16601109	1.020701143	1.886569816	0.74718205	-0.82261293	0.0543804	-0.3038182	0.16004313	0.03866404	0.17338099	1.0849008	-0.81354132	0	-0.41238081	-0.2403266119363793	
McKINNOE	Norace	McKinney	7	20	5.56914525	-0.81729117	-0.03800664	4.97946571	5.656476583	0.0543804	-0.3038182	-0.77151527	9.756484348	7.411200076	0.18355977	2.573608347	8	-0.41238081	-0.2403266119363793	
MEYERUD	Loren	Meyer	3	14	1.16601109	2.32312402	0.36655081	0.19346186	-0.1164387	0.0543804	-0.3038182	0.318978744	0.0971518	-0.3061869	1.0849008	-0.81354132	0	0.74621046	-0.2403266119363793	
MOSLEGILL	Glen	Medley	3	6	1.27259255	-0.81729117	-0.03800664	-0.21091792	0.838677201	0.0543804	-0.3038182	-0.77151527	-0.02339191	0.393228495	1.07414612	2.573608347	1	-0.41238081	-0.2403266119363793	
MURKARD	Ken	Murray	6	15	0.64928818	-0.81729117	-0.03800664	6.07217096	5.622607396	0.0543804	-0.3038182	-0.77151527	5.419758495	10.65391106	0.06467812	2.34778869	6	-0.41238081	-0.2403266119363793	
RANDAMARK	Mark	Randell	2	6	1.16601109	7.920504862	0.841098715	0.25644705	-0.82261293	0.0543804	-0.3038182	0.423550463	0.20848973	0.216174747	-0.01280706	1.0849008	-0.81354132	0	-0.41238081	-0.2403266119363793
RATKOGEZ	George	Ratkoicz	24	59	11.6913277	-0.81729117	-0.03800664	10.22004285	5.97748205	0.0543804	-0.3038182	-0.77151527	10.81215225	8.952368889	0.34128094	16.17981636	99	-0.41238081	-0.2403266119363793	
RAUTLEIGH	Lee	Rautins	9	5	0.58949865	-0.81729117	0.80239971	7.0232896	1.70393217	0.0543804	-0.3038182	-0.77151527	1.95038916	0.68838602	0.162372454	-0.81354132	0	0.67042105	0.7173265610522	
SCHETTEL	Tom	Scheffler	3	10	0.86196612	1.81074625	0.78445087	2.52066390	-0.82261293	0.0543804	-0.3038182	0.161374026	-0.02339191	-0.21702861	0.26272414	0.896178412	3	-0.41238081	-0.2403266119363793	
STROUDJO	John	Shneider	1	1	7.61977374	-0.81729117	-0.03800664	-1.57640807	-0.82261293	0.0543804	-0.3038182	-0.77151527	-0.02339191	0.181807572	0.57621612	-0.81354132	0	0.80961649	29.420051449897157	
VALUHEMD	David	VanHeine	1	1	1.16601109	-0.81729117	-0.03800664	-1.57640807	-0.82261293	0.0543804	-0.3038182	-0.77151527	-0.02339191	1.12414128	1.0849008	-0.81354132	0	-0.41238081	-0.2403266119363793	
WEATHERS	Deway	Ward	1	1	10.54924124	-0.81729117	-0.03800664	-1.57640807	-0.82261293	0.0543804	-0.3038182	-0.77151527	-0.02339191	8.021458381	14.33294211	-0.81354132	0	-0.41238081	-0.2403266119363793	
WARDHEID	Dwain	Wheat	1	3	0.784474707	0.346631833	-0.03800664	1.50871731	-0.82261293	0.0543804	-0.3038182	-0.77151527	-0.02339191	0.961378888	1.07414612	-0.81354132	0	-0.41238081	-0.2403266119363793	

Table 1: Full list of outliers ("Outliers detection" part)

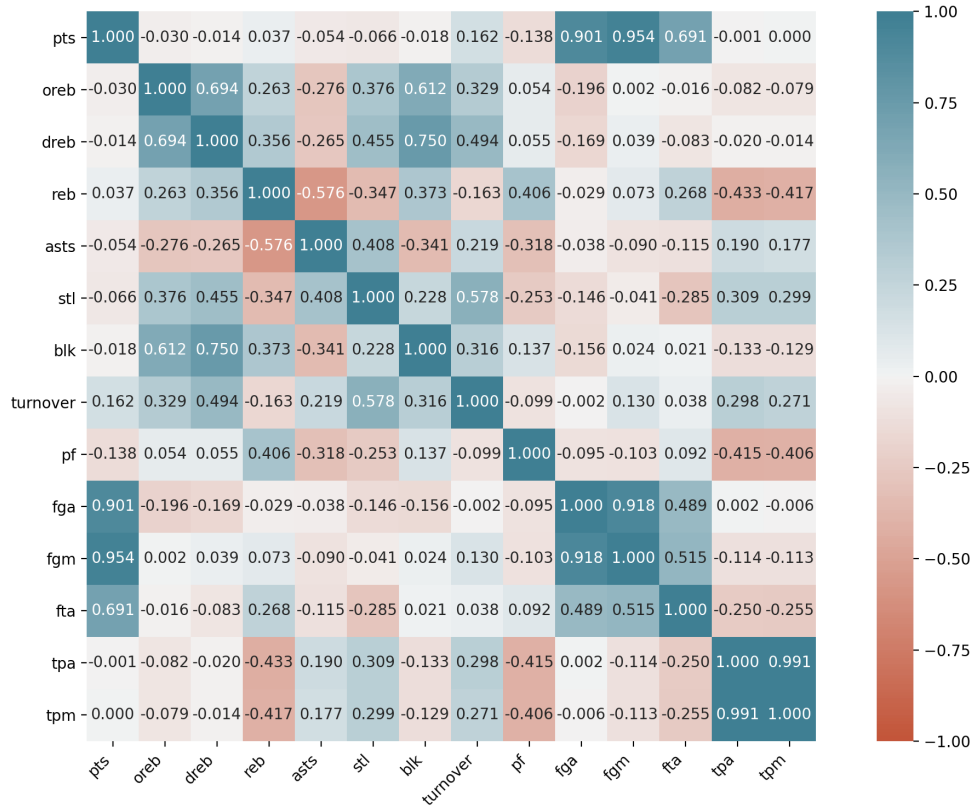


Figure 11: Correlation between normalized features("Outliers detection" part)

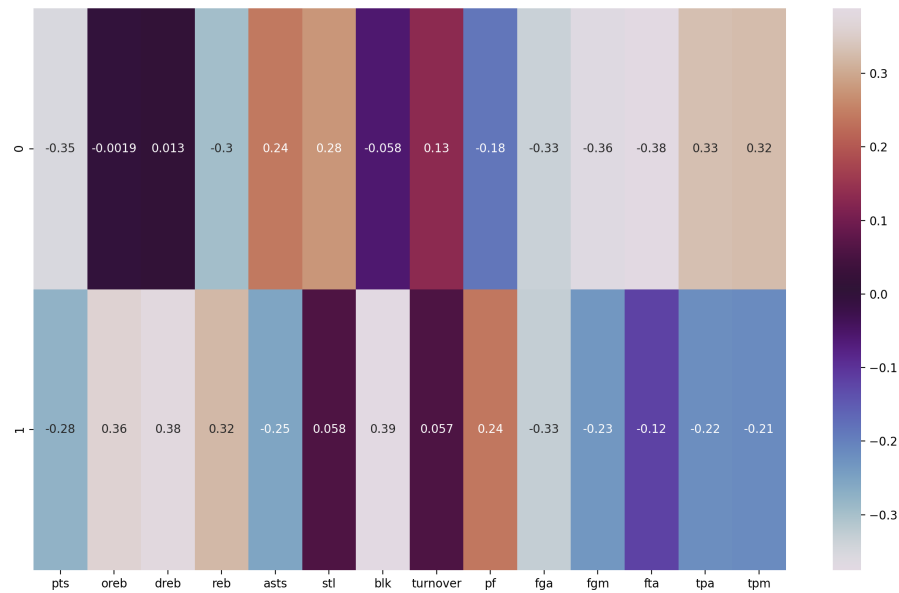


Figure 12: Correlation between features and principal components ("Outliers detection" part)

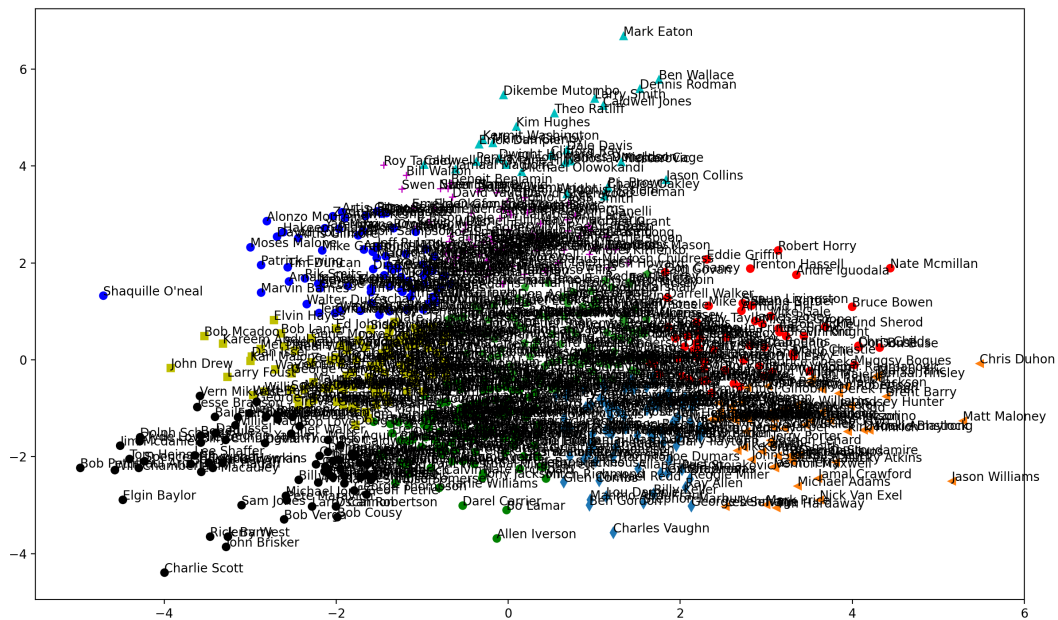


Figure 13: Plot of clustering with names attached ("Finding outstanding players" part)

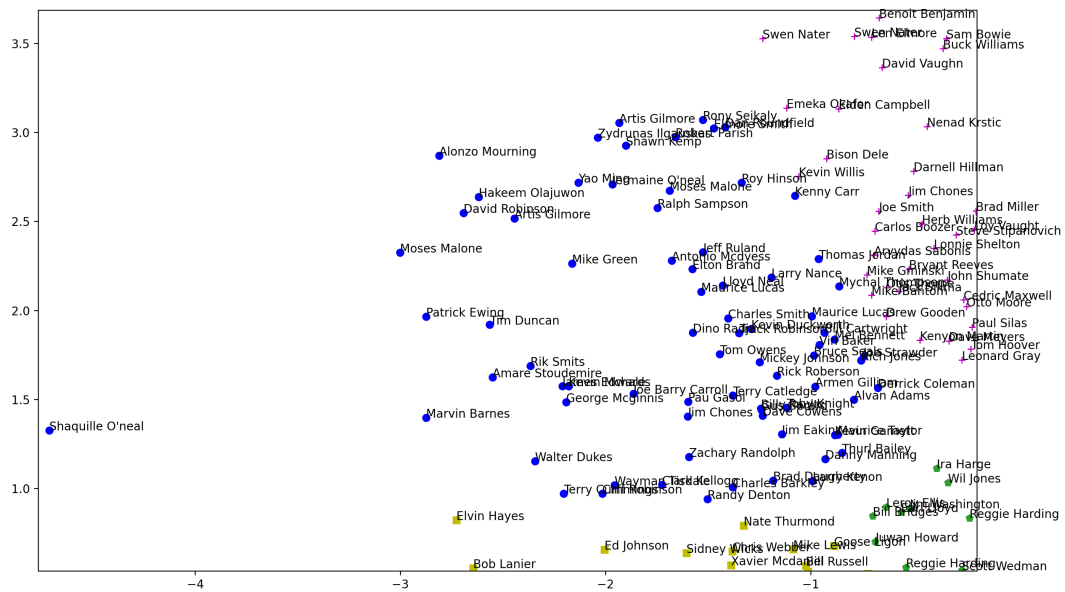


Figure 14: Plot of clustering with names attached - zoomed in ("Finding outstanding players" part)