



# NBA PLAYER VALUE

Using Machine Learning to Better Run an NBA Franchise: Milestone 1

# The Problem

- General Managers are tasked with the responsibility of constructing and paying a roster composed of high-level athletes—every one of which feels he ought to be paid what he is worth.
- All too often, these GMs succumb to external pressures: to sign the player the fans want or that the media has placed on a pedestal, to pay a player more than he is worth out of fear of missing out on a player, offering too little to a player out of pride or fear (or both) and missing out—the reasons go on and on.
- None of these pressures is based in anything objective. GMs need something on which to fall to make such costly and high stakes decisions, a means of defending the decisions they make.
- As we saw in Michael Lewis' "Moneyball," professional sports leagues can be ruled by forces other than logic, but the game always has potential to be affected by it.
- Using machine learning, we will look to give any GM an upper hand who is willing to take it.

# Initial Objectives



## **Identify Successful Teams**

Determine what makes them successful



## **Determine True Player Value**



## **Make Recommendations**



# The Data

- The data was acquired via webscraping a few different websites using the BeautifulSoup package.
- Data was drawn across multiple websites and many different urls within each website.
- Player data was a composite dataset of “per 100 possessions” stats, advanced stats, and shooting stats from [basketball-reference.com](https://www.basketball-reference.com), as well as salary stats from [hoopshype.com](https://www.hoopshype.com).
- Team data was also a composite dataset of win-loss data and advanced statistics pulled from [basketball-reference.com](https://www.basketball-reference.com) as well.
- All webscraping code can be found [here](#).



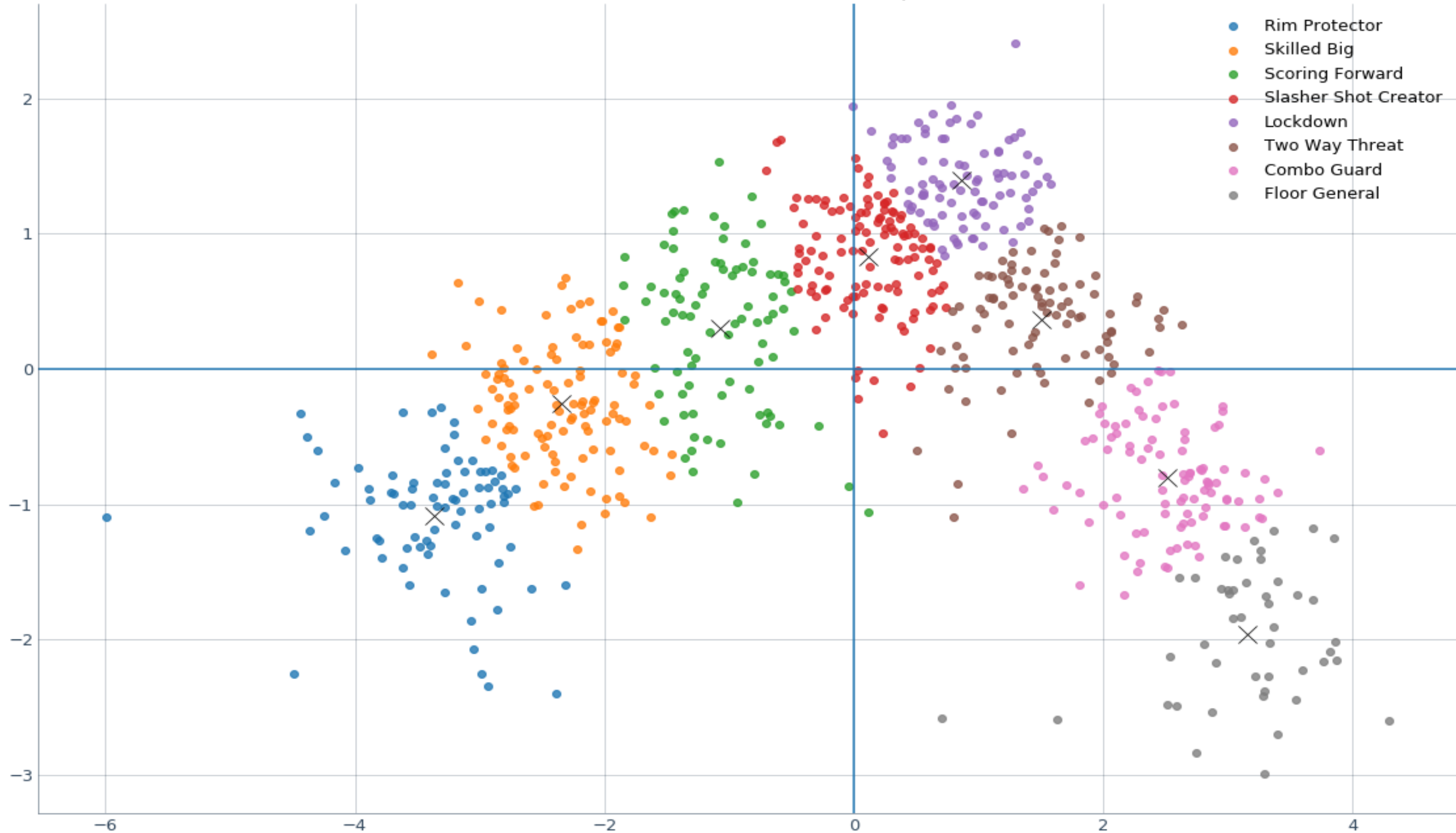
# Data Cleaning and Wrangling

- The bulk of the data wrangling needed for the player and team data was performed during the webscraping loops—each loop spit out list of pandas DataFrames that would ultimately be concatenated.
- After the webscraping, the majority of wrangling that remained for this data was simply structuring DataFrames such that they could be merged together using pandas merge.
- The data that needed the most cleaning was the salary data from [hoopshype.com](http://hoopshype.com). For example, a datapoint we might read as “\$20,000,000” was actually written in the HTML as “\n\t\t\t\$20,000,000\t\t\t\n”—a very messy string.
  - This required me to identify the characters that needed to be dropped before changing the datatype to something we could manipulate (i.e. string to float/integer).
- The data from [basketball-reference.com](http://basketball-reference.com) was for the most part very clean, handling null values was an important component of this project. Depending on the feature, some data was filled using the population mean, and for some data, a null value in the wrong feature meant dropping the record entirely.
- The bulk of the data wrangling can be found [here](#).

The background of the slide features a photograph of four basketball players in motion on a court. From left to right, the players are: a player in a dark jersey with the number 32, a player in a white jersey with the number 1, a player in a white Los Angeles Lakers jersey with the number 32, and a player in a dark Philadelphia 76ers jersey with the number 25. Above the players, height markers are visible: 6'4" above the first player, 6'7" above the second, 6'9" above the third, and 6'10" above the fourth. The title "New Player Positions" is overlaid in a large, white, serif font on the left side of the image.

# New Player Positions

- One of our primary objectives was to accurately determine player value. In order to best do this, it was necessary to distinguish differences between players who are listed as similar. Our hypothesis was that there is a more robust method for describing a player's position than just as one of five positions (point guard, shooting guard, small forward, power forward, or center).
- For example, compare the play styles of Ben Simmons and Steph Curry. Both are listed as Point Guards, but their contributions could not look more different.
- In order to do this, we used our extensive player-level data (per 100 possessions stats, advanced stats, and shooting stats) to cluster players into 8 distinct positions that better described their contributions to the team.



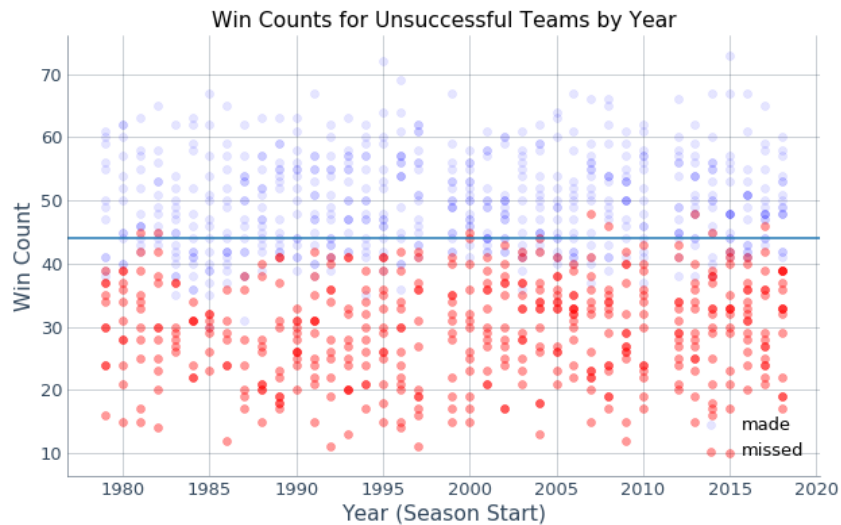
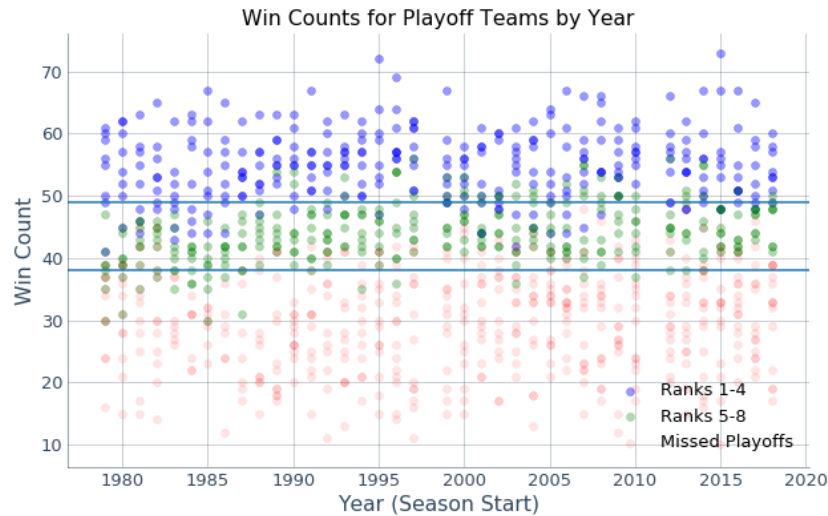
# Dimensionality Reduction

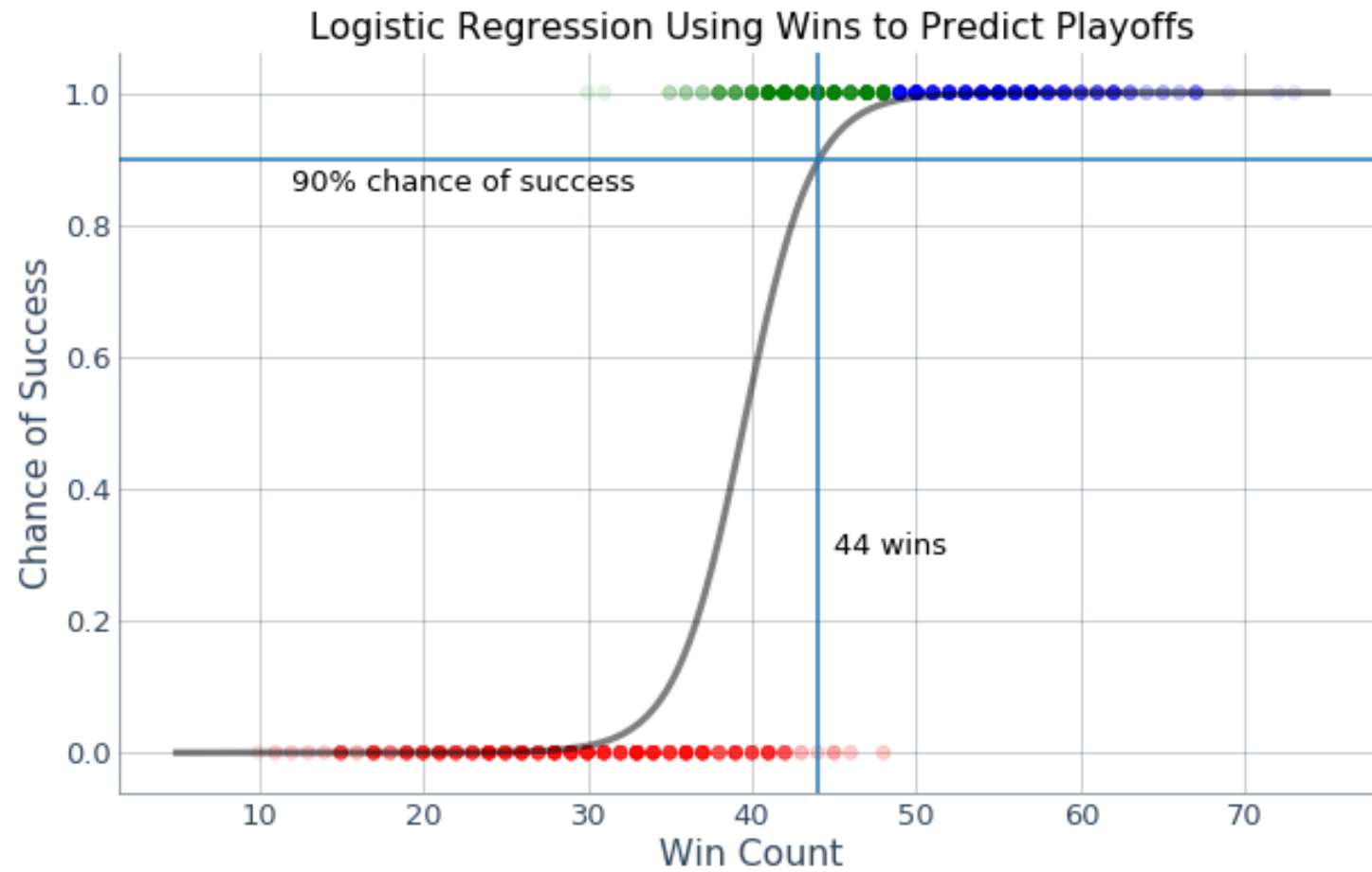
- Principal Component Analysis (PCA)
- Linear Discriminant Analysis (LDA)
- K-means clustering
- Once new labels were merged with our player-level data, we were able to add another layer to our EDA and inferential statistics as we could then consider a player's cluster when asking questions of value.



# Exploratory Data Analysis

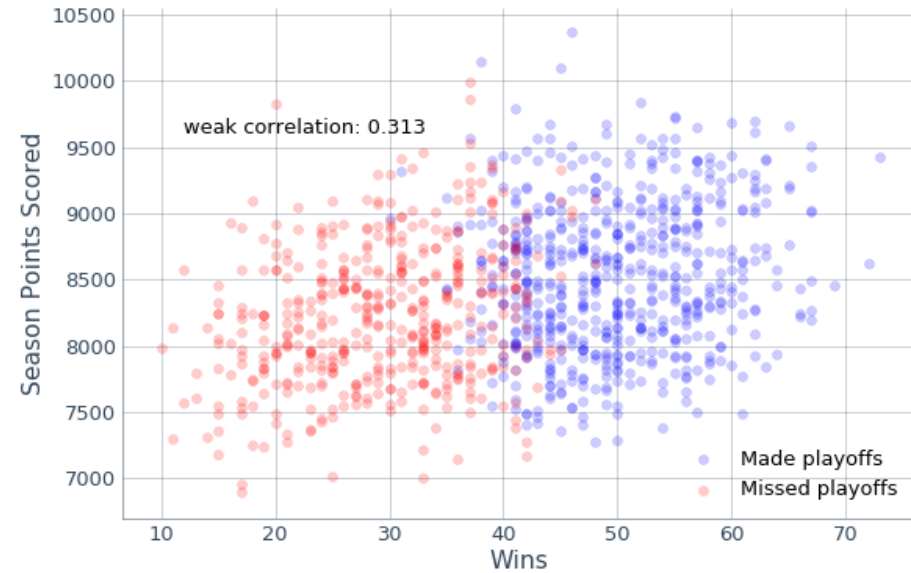
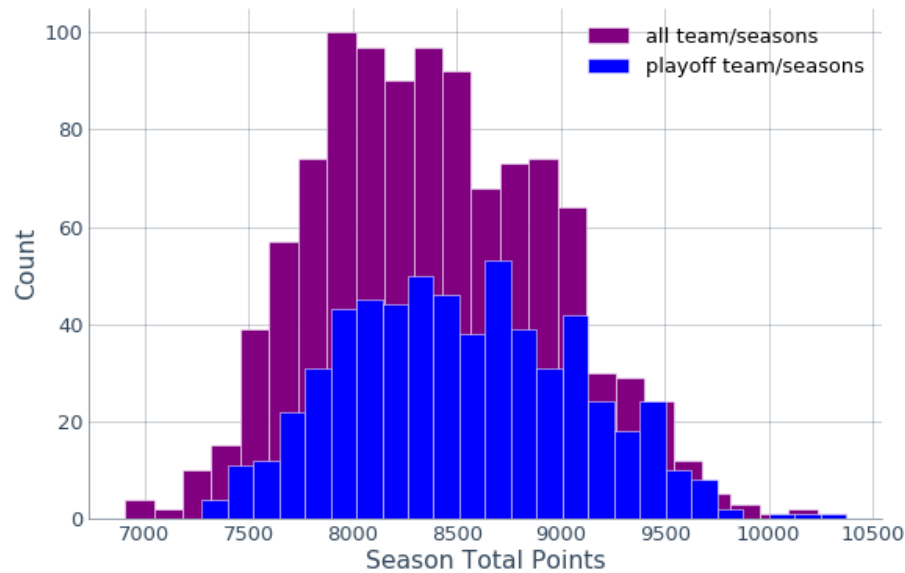
- On the team-level, we wanted to identify success, so that required of us that we define our terms.
  - “Success” for the scope of this project equated to making playoffs, or even better, making the playoffs with a good seeding (1<sup>st</sup> through 4<sup>th</sup>).
- Left are two similar plots that show roughly how many wins generated successful season.
  - At 38 wins, a team has a good chance of making playoffs, but at 44 wins, it's highly unlikely that a team would miss playoffs.
  - At 49 wins, a team will very likely have a good playoff seed.





The logistic regression to the left reinforce the claims of the previously slide.

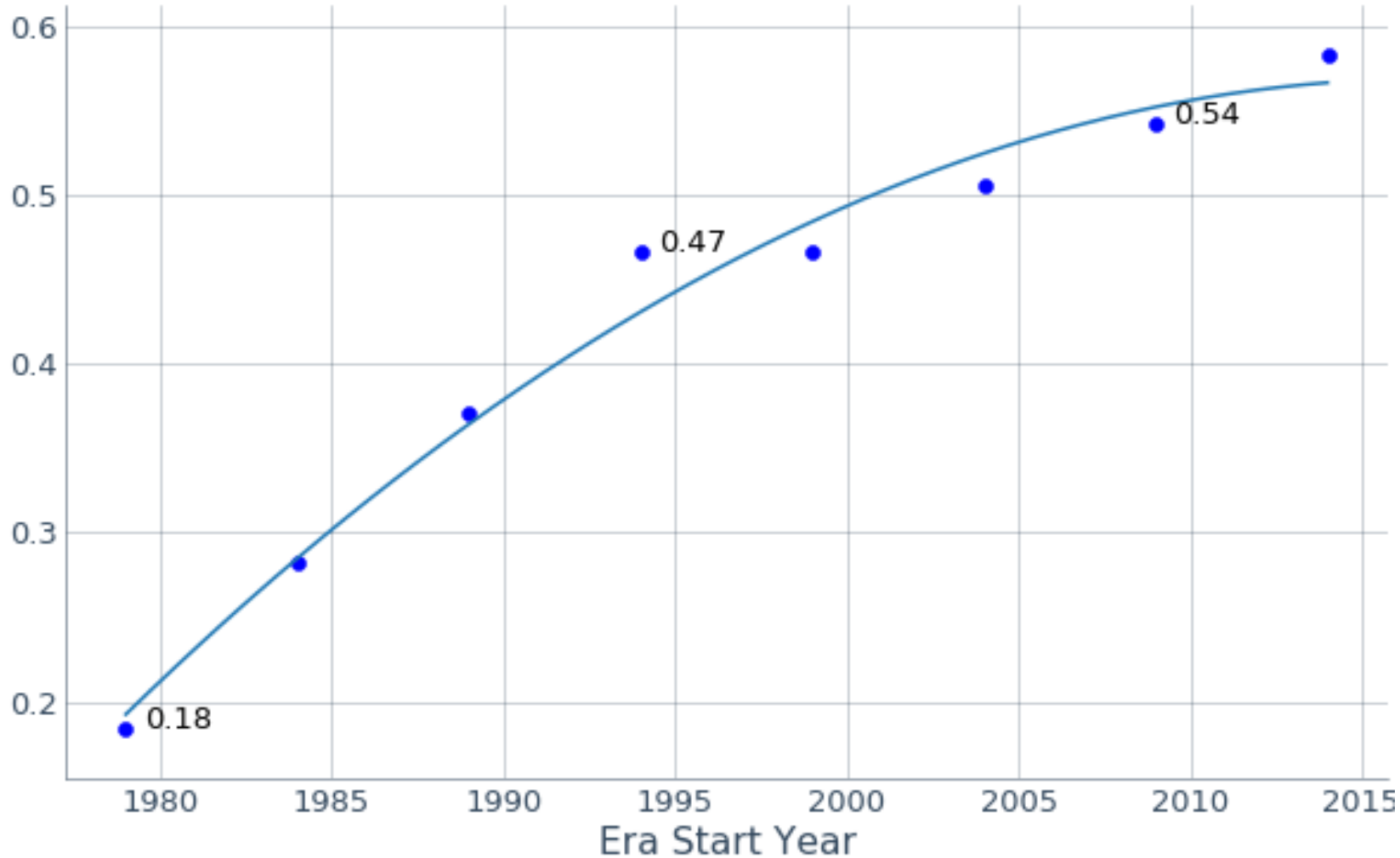
## Point Totals for Playoff Team/Seasons



## Team-level Win Correlates

- The next step was to look for correlates of win count. “What did a team that won a bunch of games do well that season?”
- On first pass, we checked the statistic of “points scored”. This, interestingly had a very weak correlation with win count—just because a team is scoring a lot of points, doesn’t mean they’re winning games. Same could be said with the statistics of “points against”.

3P% Correlation with Win Count against Era Start Year



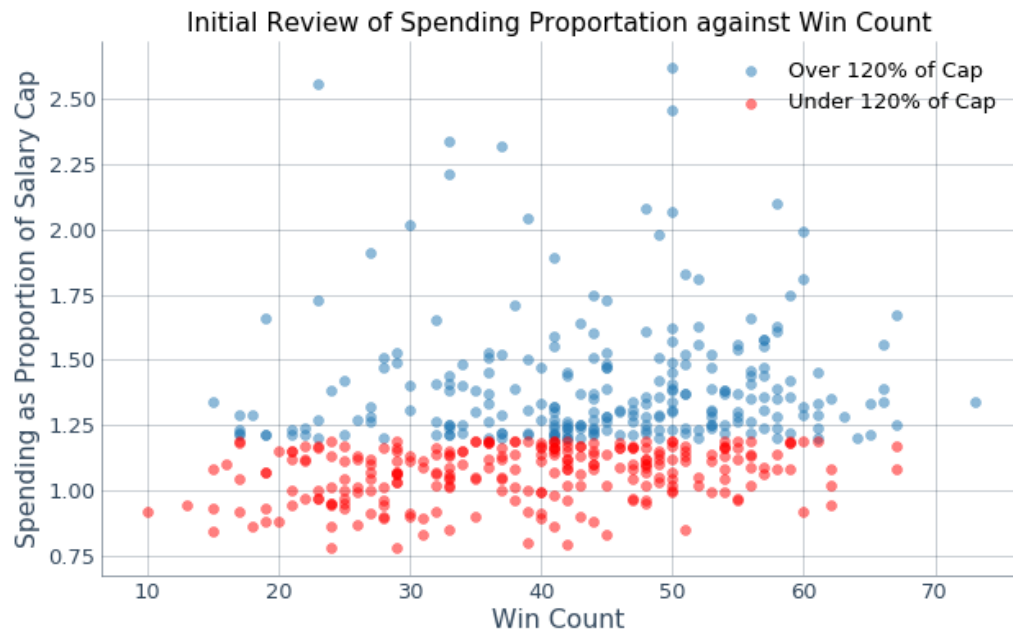
- Where we did see our most notable strong correlation was with the feature “3P%”, or 3-point percentage. Teams that were more accurate from the 3-point line won more games.
  - When we first looked at this correlate, we saw a coefficient of 0.18, but this included team data dating back to 1979.
  - When we limited our data to just team-seasons dating back 5 years (2014-2019) instead, we saw a correlation of 0.58.





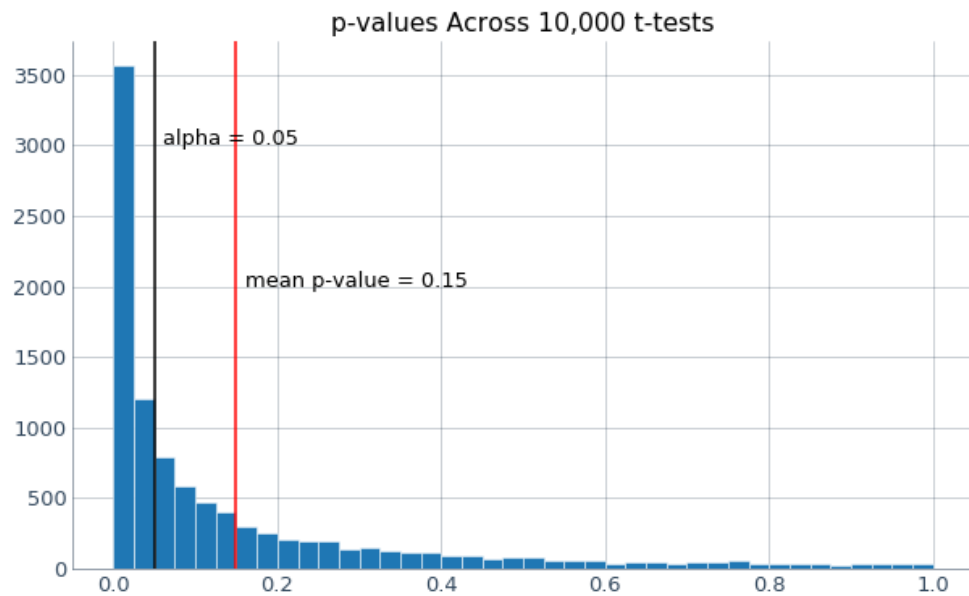
# INFERENTIAL STATISTICS

© Andy Marlin-USA TODAY Sports



# Spend to Win?

- With our team-level inferential statistics, two sample t-tests of independence showed us that money doesn't buy wins--at least not at a statistically significant level. The differences in mean Win count for teams paying over 120% of the salary cap for any given year, and teams paying under 120% of the salary cap for any given year varied greatly--44 and 38, respectively.
- However, when we sampled from these groups and ran a few t-tests (read: ten thousand t-tests), we saw that that this difference was not in fact significant at even a 0.10 confidence. The resultant p-values were plotted in the histogram (bottom).



# Inferential Statistics cont.

The following slide demonstrates how python allowed us to scan for statistically significant differences in 6 different features across 8 subsets of the data, using 10,000 t-tests each time (48,000 in total) in about 2.5 minutes.

- From this scan, we learned the following:
  - the "Lockdown" cluster was flagged for statistically lower mean Player Efficiency Rating (PER)
  - "Rim Protectors" were flagged for having statistically lower mean Usage Rate (USG%)
  - "Combo Guards" were conversely flagged for statistically higher mean Usage Rate
  - "Rim Protectors" were the only cluster with a significant difference (lower) in Points Scored (PTS)

```

clusterlist = players.cluster.unique()
value_stats = ['PER', 'WS', 'BPM', 'VORP', 'USG%', 'PTS']
confidence = 0.05 * 2

for h in value_stats:

    pop_mean = np.mean(players[h])

    for i in tqdm(clusterlist):

        players_T_list = []
        players_p_list = []
        cluster = list(players[players.cluster == i][h])

        for j in range(10000):

            cluster_samp = sample(cluster, 30)
            T, p = stats.ttest_1samp(cluster_samp, pop_mean)
            players_T_list.append(T)
            players_p_list.append(p)

        T = np.mean(players_T_list)
        p = np.mean(players_p_list)

        if p <= confidence:
            print('Cluster: {}; \t Value Stat: {}'.format(i, h))
            print('Sample mean {}'.format(h), round(np.mean(cluster_samp),
                2), 'Population mean {}'.format(h), round(pop_mean, 2))
            print('Test statistic (T):', round(T, 2), 'p-value:', round(p, 3))
            print('Reject the null. Significant difference in mean {} between population and the {} cluster.'.format(h, i))
            if T < 0:
                print('For the {} cluster, mean {} is lower than the league average.'.format(i, h))
            else:
                print('For the {} cluster, mean {} is above the league average.'.format(i, h))

```



# Conclusions and Next Steps

- Certain clusters are probably more valuable than others, objectively, meaning our LDA and clustering was at least worth it in that regard.
- Inferential statistics showed us that, at least at the 120% of the salary cap mark, spending more doesn't result in significantly more wins than spending less.
- Making recommendations will greatly on a deep understanding of our features; a few factors have emerged that are slowing us down:
  - Early regression analysis has proven fruitless—salary is not simple to predict based purely on the composite of many features.
  - There is not a consistent value system currently in place from team to team, and certainly not from year to year.
- See the full code in [this GitHub repository](#).