

## Problem Statement

### **Title:**

Predicting player performance and maximizing team performance in Fantasy Football

### **Problem:**

The inspiration for this project was work done by a Stanford University Computer Science student [here](#). Daily Fantasy Football websites allow users to select a new lineup of NFL players each week to compete against other fans. Each user has a fixed budget for their team, with players with high expected value drawing a high price tag. Users' teams then accrue points based on player performance. The goal of this project is to develop a regression algorithm to predict weekly performance for the players, then determine opportunities where predicted performance is significantly higher than cost to acquire the player. Likewise, I hope to provide an optimization algorithm that, given the fixed budget cap for a team, assembles the team with the maximum expected value. The previously cited project focused only on previous performance, while my aim is to expand this work to include more qualitative factors such as opponent, weather, week of season, etc. A successful model could be used by a Fantasy Football fan to outperform their peers.

### **Data:**

The main source of data is a public repo (<https://github.com/BurntSushi/nflldb>), containing a SQL database dump of raw NFL game data. This data includes, but is not limited to, game schedules, scores, rosters and play-by-play data for every preseason, regular season and postseason game dating back to 2009. I will attempt to merge in other data sources that may have an effect on player performance, such as weather data(<http://www.nflweather.com/>) or contract values and changes(<https://www.sportrac.com/nfl>). Finally, I will join in 'player cost' data from FanDuel, which determines the cost of adding a player to a team, with each team having a total budget cap.

### **Methodology:**

I'll begin by merging and flattening the full data set to produce a wide table where each row contains the 'fantasy score' (i.e. performance value) for a given player-game pair, with player and game attributes, such as opponent, age, salary. Next, I'll use recent historical fantasy scores to engineer features such as rolling average score and scoring variance. Also, I'll look to non-numeric columns to find certain player affinities, such as correlation between opponent and performance, as well as creating features from other data sources as outlined in the **Data** section. I'll apply a regression algorithm to predict the weekly performance for every offensive player. Given the data set covers the entire season from multiple years, I'll experiment with using previous years as training data, then using current year for test data as well as using first xx% of the season as training data, then testing the algorithm on the remainder of the season. After optimizing the algorithm, I'll acquire the 'player cost' data, and use the cost, predicted performance pairs to maximize my expected team performance, given the team salary cap. Performance of this project will be judged against the performance of others using the FanDuel platform to build teams, who presumably compile performance predictions from publicly available data, such as sports bloggers.

## Results:

The goal of this project is to provide a standalone repository with which others may replicate and use my prediction algorithm to help guide their own fantasy football season. Specifically, I will include a brief deck outlining my methodology, results, and suggestions for future improvement. Secondly, I will include the Jupyter notebooks (with instructions) used to wrangle and explore the data, as well as the code for the prediction algorithm.

## Project Goal Update:

Having attempted to apply and tune several different regression models to the data, with consistently high error, it's necessary to update the project's goal. Initially, I set out to accurately predict the fantasy for a given player, based on past behavior and other game metadata. Going forward, my goal will be to separate players, teams, and / or games into useful groups, using classification algorithms. Examples of these clusters could be 'low-scoring game overall' or 'advantageous to RBs'. Furthermore, I'll attempt to support or debunk common Fantasy Football techniques, such as never selecting a player in a Thursday game.

## Data Wrangling

### 1. Cleaning

1. Aggregations & Joins - The raw data included 4 tables (objects), each with at least one unique key, which could be used to join to the other tables. The tables are:
  1. Plays - what were the statistics for a given play and who participated in the play.
  2. Games - what were the total play statistics for the entire game and which teams participated in the game.
  3. Players - what are the summary attributes for a given player.
  4. Teams - Only contains full name, abbreviation, and home city for each team. Will be useful in joining to other data sets later.

I used GroupBy (by=player,game) to aggregate play-level data to game-player level data. Next I merged player- and game-level data into the previous game-player DataFrame. The result is aggregate statistics for each player, for each game, with metadata from the player (e.g. position) and the game (e.g. home team, away team) tables.

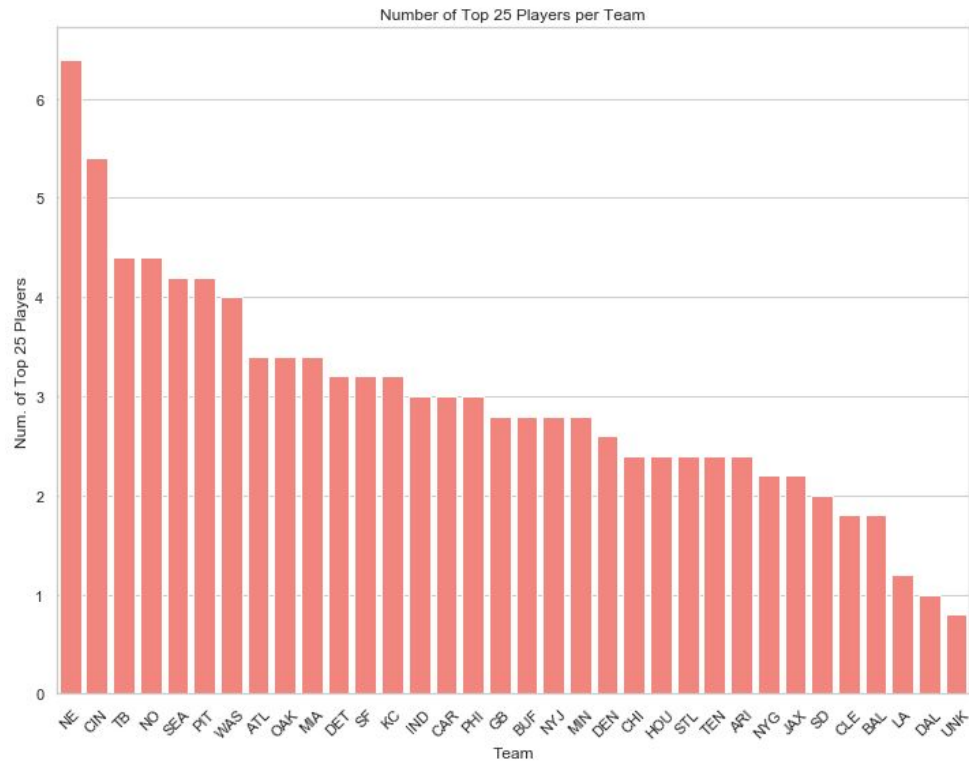
2. Incomplete Data - There were four columns for name, inconsistently filled throughout the data. I wrote a function to check the columns in order, return the first value found, or combine columns if necessary, and write to a new name column.

3. Removing Irrelevant data - I removed rows for defensive players, which is not in scope for this project, and removed columns that either have no impact on players' scores or will not be used in features.
  4. Parsing non-numeric data - Weather data provided the temperature and a description of the weather. With something like 65 unique weather descriptions, and no clear pattern for RegEx, I manually categorized the descriptions into Good Weather, Wet Weather, Windy (but dry) Weather, and Winter Weather, and engineered a new column: desc\_simple.
2. Missing Values
    1. Missing weather data - ~16% of the games were missing weather data, with all of 2010 missing weather data. With that outlier removed, 6% at most of any other season is missing data. Later on, I may have to remove 2010 from the model or remove weather as a feature. Filling with the mean does not seem appropriate in this case.
    2. Bye weeks - Each player will have one bye week each season where they will not play. Bye weeks are not included in the data set, so there won't be an observation of 0 point performance due to bye weeks. I won't fill the week with a value, but will just skip over it in any rolling calculations.
  3. Outliers
    1. Outliers aren't an issue with this project. In fact, finding outliers with abnormally high scores relative to cost would be an ideal outcome for the project. However, this data set contains data for all players in the selected positions, and given the team salary cap constraint, it is highly likely that I will be selecting only from the top 25 players in any position. The bottom performers will likely score 0 points on most weeks, considering that they won't play at all, so I'll remove all players outside of the top 25 per position.

## Exploratory Data Analysis

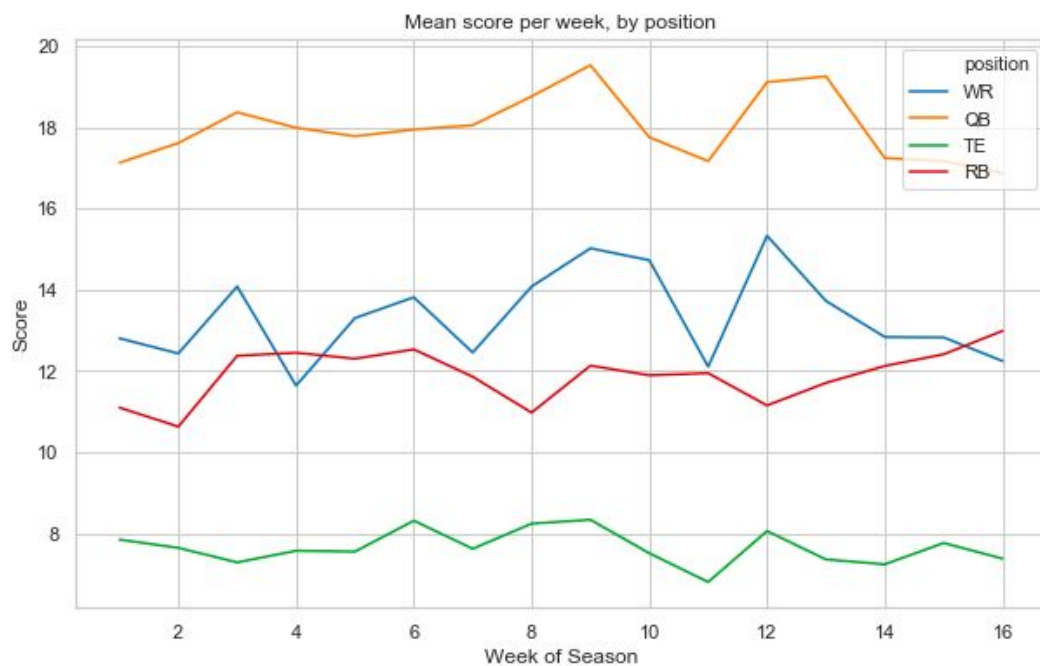
### **Top Talent Distribution:**

The top 25 players at each position over the last 5 years is not evenly distributed amongst all teams. In fact, the most talent dense team, New England Patriots have twice the number of top players than the median and more than 6 times more than the least talent rich team with at least one top 25 player.



### Season Fatigue:

Towards the end of the season, the average score for both WRs and QBs tends to decrease, while RBs actually see an increase in average score in the last few weeks.



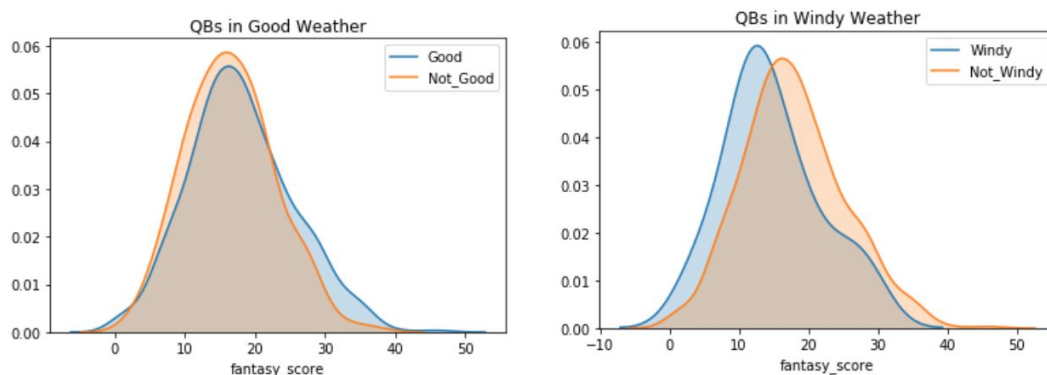
### Weather:

*Assumption: Quarterbacks can't throw accurately into the wind.*

Does windy weather significantly decrease a quarterback's fantasy score? And conversely, does good weather (or playing in a dome) improve performance?

Steps:

1. Filter data set of player-game data to only quarterbacks and split the set into games with `windy` weather vs. non-`windy` weather.
2. Plot distribution of fantasy score for both sets.



3. Perform 2 sample, one-sided z-test, with alpha = 0.05 where:  
 $H(0): \mu_{windy} = \mu_{not\ windy}$   
 $H(a): \mu_{windy} < \mu_{not\ windy}$
4. Determine significance and repeat with `good` vs. not `good` weather data.
5. We find a significant positive ( $p = 0.0027$ ) for `good` weather and a significant negative ( $p = 0.0476$ ) impact of wind on quarterback performance.

### Thursday Night Curse:

*Assumption: Teams that play on Thursday have had a shorter week to recover and are less likely to play with intensity.*

Is the average fantasy score lower than the average score the rest of the week? In terms of actual score (not fantasy points), are Thursday games lower scoring on average?

Steps:

#### A. Average Fantasy Score

- a. Split player-game data set into Thursday and non-Thursday games.
- b. Perform 2 sample, one-sided z-test, with alpha = 0.05 where:

$$H(0): \mu_{Thurs} = \mu_{non-Thurs}$$

$$H(a): \mu_{Thurs} < \mu_{non-Thurs}$$

- c. Determine significance ( $p=0.303$ ) and fail to reject the null hypothesis. No significant difference in players' fantasy scores.

#### B. Average Total Game Score

- a. Compute total game score as home\_score + away\_score
- b. Subset data into Thursday, non-Thursday

- c. Apply bootstrapping by 1) computing difference in means between the two data sets 2) permuting the samples and 3) counting bootstrapped samples where difference of means is at least as extreme as our original samples ( $p=0.1948$ ).
- d. Determine significance - again, no significant difference in total game score between Thurs. and non-Thurs. games.

### Share of Targets:

*Assumption: Teams force the ball to their best players. Getting more of the targets will yield them a higher fantasy score.*

What's the correlation between the relevant 'opportunity' metrics and fantasy score at each position?

Using `pearsonr()`:

Metric	Pearson R (Correlation Coefficient)
QB - Passing Attempts	0.3178
QB - Rushing Attempts	0.2052
RB - Rushing Attempts	0.598
RB - Targets	0.373
WR - Targets	0.5575
TE - Targets	0.6845

There's a moderate positive correlation between opportunity and fantasy score for RBs, WRs, and TEs. Strangely enough, the position handing out the football, the QB, has a low correlation with passes or rushing attempts with fantasy score. This may raise additional questions around efficiency of player e.g. pass completion percentage.