

Capstone 1 Milestone Report

Problem Statement

The aim of this project is to take a mathematical approach to betting on NFL contests. We want to predict the total number of points scored between two teams playing against each other in any particular NFL game. We will be using regression techniques for numeric predictions, and will be comparing our predictions to the game totals that professional sports books put out for the public to view and bet on. Comparing our predictions to the betting market will paint a clear picture on which games present us with an opportunity to profit. Once we have identified a game or games with glaring values in comparison to the betting market, we can recommend these plays to a potential client, friend, or even play them ourselves. Additionally, we can recommend that fantasy football players use football players in their lineups from games that we project to go “over”, and recommend to stay away from players from games that we project to go “under” the total.

Attribute Selection

The attributes that will be used for this project will NFL team matchup statistics such as the offensive team’s passing matchup, running matchup, and overall skill matchup. Additionally, we will be using team grades from profootballfocus.com, venue data, home and away data, points scored, and a defense’s points allowed. Another important piece of information is data on game totals, which will tell us what values the sports books set the betting market for each game, and can serve as a marker from which to compare our predictions and make recommendations.

Obtaining the Data

For all of the attributes described, there are many sources from which to obtain the data. One source is pro-football-reference.com, which has a database that contains nearly every relevant football statistic for every NFL game played dating back to 1920. Another source for obtaining the needed data will be profootballfocus, which has defensive and offensive grades for every team in every contest dating back to the early 2000's.

Data Collection

- Built five web data scrapers using BeautifulSoup4 to collect all data for this project
 - Pro-football-reference.com referee and game data from 2006 onward
 - Profootballfocus.com game grades for every team dating back to 2006
 - Footballoutsiders.com DVOA metrics for every team and every game since 2006
- For missing values, I explored why they were missing and found a pattern that led me to changing the code in two of my scrapers.
 - Games that went into overtime had additional columns of data that I hadn't accounted for
- Put all data into .csv files and Pandas dataframes

Data Cleaning

- Used Pandas to import data from .csv files into dataframes using .read_csv()
- Performed dataframe column operations to create new variables
- Used string methods on dataframe columns to clean string variables
- Changed date columns from strings to datetime objects
- Replaced old team names with up-to-date team names for data consistency among all tables

- Used a dictionary that mapped all old team names to their current team name and applied the `.map()` function to the ‘team’ columns to transform the team column to a new, updated team column for every table
- Used joins and merges on tables of data to prepare data for analysis in the next part of this project

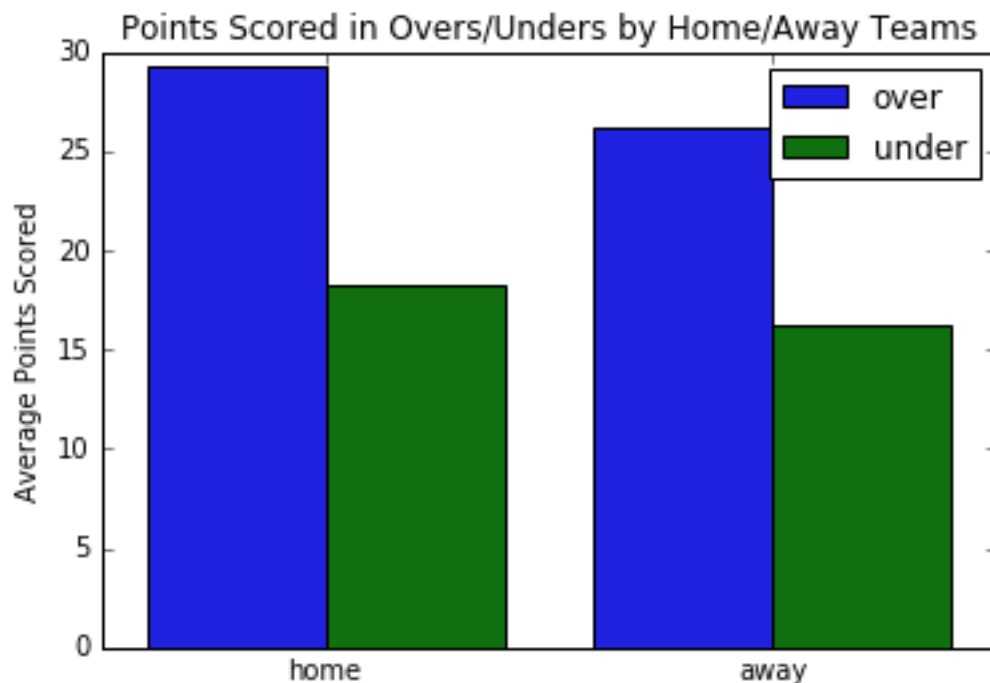
Data Storytelling

With the steps of data wrangling completed, it is now time to explore the data further. Our main goal is to identify trends and correlations, as well as to communicate insights found within our dataset. We want to ask questions of our data that will drive us toward our main goal of this project, which is to profit from predicting over-under outcomes of NFL games. Remember that sportsbooks set a total number of points to be scored for every NFL game, and bettors can either bet the “over” if they think the teams will combine for more points than the set total, or “under” otherwise. With this in mind, let’s proceed with exploring our data in such a way that we are pushed toward our goal of predicting NFL outcomes.

Firstly, we look at the outcomes of games within our dataset. We are using NFL games dating back to 2007 for this project. We can see within our jupyter notebook that 1,350 games went over the total set by the sportsbook, 1,341 went under, and 45 were a “push”, or scored as many points as projected. These numbers suggest that the sportsbooks are incredibly good at setting totals for games, as almost exactly 50% of games hit the over, and almost exactly 50% of the games hit the under.

A logical follow-up would be to explore what happened in games that went over/under. In games that went over the total, the home team averaged 29.3 points per game (ppg), and the away team averaged 26.2 ppg. In games that went under, the home team averaged 18.2 ppg, and

the away team averaged 16.2. Therefore, in games that go over, the average number of total points scored is 55.5 ppg, and for unders it is 34.4 ppg. That is a massive difference of about 21 points, or the football equivalent of three touchdowns/extra points. Fantasy football players can see the immediate value of being able to accurately predict whether or not the game will go over or under, as they will want to choose players from games that are predicted to go over the total set by the sportsbook, since touchdowns are incredibly valuable in fantasy football.



Next, we take a look at games in which the total is set as a high score or low score. These games are interesting to investigate because in these games, the betting public usually has more of a consensus opinion on what will happen in the game. For example, if two teams with very bad offenses are playing against each other, the total for the game might be set at a number like 40, or even lower. This is considered a very low point total for a sports book to set. The public might lean toward betting the under on this game, since the two teams' offenses have been poor during the course of that season. We will use below or equal to 40 as our "low point" threshold,

and above or equal to 48.5 as our “high point” threshold. Looking at games with totals below or equal to 40 points, we can see that the over hit 260 times, the under hit 207 times, and push hit 10 times. The over hits about 55.7% of the time in this case, which exceeds the 52.4% needed to beat the rake of the sportsbook. The average score for these games is 21.2 ppg for the home team, and 18.6 ppg for the away team, for a total of 39.8 points, on average. For the fantasy football player, this suggests that although we may predict one of these “low total” games to go “over”, it doesn’t mean that we should be targeting players from these games. Looking at the “high total” games with totals set at or above 48.5, we find that the under hit 254 times, the over hit 236 times, and push hit 7 times. The average score for these games is 27.4 ppg for the home team, and 24.1 ppg for the away team, for a total of 51.5 points, on average.

Lastly, we look into how one of our variables, ‘offMatchup’ is related to points scored for a team. offMatchup represents how good or bad a particular offense’s matchup is going into that game. The variable combines how good the team’s offense is with how bad that team’s opponent’s defense is. A high offMatchup score is good for a team’s offensive outlook for that game. Looking at the initial scatterplot, we can sense a positive correlation between offMatchup and points for (pf). Using the corrcoef function in Python’s numpy library, we can see that the correlation coefficient between the two (for home teams) is 0.236. Investigating further, we can see that home teams who have an offMatchup score of 10 or greater average 27 ppg, and away teams average 24.4 ppg. On the other hand, home teams who have an offMatchup score of -2 or lower average 19.9 ppg, and away teams average 17.5. Both sports bettors and fantasy players can benefit from using offMatchup scores when attempting to predict how an offense will perform.

One final piece of information to note is that when an away team is facing a team that has an average coverage score of 73 or higher, the under hits at a 63% rate. This would be an extremely profitable trend to follow as a sports bettor or daily fantasy player.

We can see that digging into the data is a worthwhile endeavor when it comes to placing bets, or attempting to determine how a game will play out. Someone who wishes to profit from playing fantasy sports or betting on games would be wise to take advantage of some of the information gathered from digging deeper into the data we collected and wrangled.

Inferential Statistics

The next step in our process of predicting over/unders for NFL games is to use inferential techniques to explore the data further. Using an iPython notebook, we break our data down into two tables: homeMatchup, and awayMatchup. The homeMatchup table contains instances of games for the home team. Each instance has variables that contain information about a particular home team's matchup for that game. One example includes the ptsMatchup variable, which is calculated by taking the home team's average points scored over the past seven games, and adding that number to the home team's opponent's average points allowed over the past seven games. There are ten independent variables for each table, and the dependent variable for each table is points scored (pf).

Looking at the correlation heatmap matrix for both tables, we can see some positive correlations between the independent variables and pf variable for both tables. Additionally, there are some notably strong correlations between independent variables in both tables as well (passMatchup and offMatchup have a pearson correlation of about .77 in both tables). We should make sure to keep this in mind when we begin the model-building process, as the issue of multicollinearity may arise.

There are also some positive correlations between the independent variables, and the dependent variable for both home teams and away teams. The strongest correlations for home teams with respect to points scored are ptsMatchup, totalDvoaMatchup, offMatchup, passMatchup, and ovrMatchup, with $r = .2853, .2693, .2360, .2341,$ and $.2334$, respectively. The strongest correlations for away teams with respect to points scored are ptsMatchup, totalDvoaMatchup, offMatchup, passMatchup, and ovrMatchup, with $r = .2925, .2804, .2490, .2567,$ and $.2385$, respectively.

Using critical values for Pearson's Correlation Coefficient with $\alpha = 0.05$, we have enough samples, and thus degrees of freedom, from both the home and away tables to put a baseline correlation of 0.06 as our mark of whether or not a variable's correlation is significant. Looking at the correlation heatmap, we see that for the homeMatchup table, all dependent variables have significant correlations to the target variable (pf). In the awayMatchup table, all independent variables except for stDvoa and pblkMatchup have significant correlations.

For the home table, we can say that there was a significant positive relationship between all individual independent variables, and the dependent variable ($p < 0.05$). For the away table, we can say that there was a significant positive relationship between all individual independent variables ($p < 0.05$) except for stDvoa and pblkMatchup ($p > 0.05$), and the dependent variable.

Pictured on the next page is the scatterplot matrix, which shows how all of the variables for home teams correlate with each other, as well as a histogram for each variable. It is a reasonable assumption to deem each histogram as “normal” for each variable, based on the visualizations. The dependent variable, pf, is the very last row and column.

