

INFO 4310 Final Project Design Doc

Abhimanyu Gupta (aag245), Natasha Armbrust (nka8)

Project Goals and Motivation

Our final project aims to answer the question whether Tom Brady, Quarterback of the New England Patriots, should have won the Most Valuable Player (MVP) award during the 2017 NFL season. We were inspired by our interest in football, as specifically one member of our team is a hardcore Patriots fan. The other member of our team enjoys exploring qualitative arguments by using data, hence the idea to prove if Tom Brady deserved the MVP. Through this project, we aim to explore how MVPs are chosen, and whether arguments made in favor of players are merely opinion or rooted in statistical backing. Currently, MVP's are chosen by members of the media by vote, which leads to a relatively arbitrary selection that can be based on many human biased reasons such as wisdom of the crowd, exposure to players, etc. We look forward to visualizing the data and to creating a better argument for this award rather than a vote.

In order to do so, we will be breaking down the 2017 NFL season and will attempt to educate the populous about football. Ultimately, we aim to make the data digestible in order to make an argument. We are inspired by [football visualizations](#) that are successful at displaying data in a way that readers can easily pinpoint important metrics. By similarly showcasing important dimensions of the large NFL dataset, we hope to provide numerical arguments which break through non-quantitative factors such as media hype or coastal bias (players in bigger cities receiving more coverage).

Intended Audiences and Insights

This project can be tailored to many different sets of people. By our analysis, we consider 6 major audiences who would be potentially be interested in our analysis:

1. Intense Football Fans

People who fall within this set have an extensive knowledge about football. Not only are they aware of players on their favorite team, but they are aware of all things NFL, including many players and teams, transactions, rumors, etc. Since these folks are inherently knowledgeable of the field, they have their own personal opinions about players and teams, and thus will also likely have an opinion on the MVP award. This article will serve as a more formal, less opinion based analysis of whether the award was awarded correctly. These users will gain the most insight out of the interactivity tools; they can adjust and check their biases by performing a deeper dive into the data such as comparing several quarterbacks who they believed should have won over Brady. We believe our visualization will allow intense football fans to come to their own conclusions based off prior knowledge. By engaging with our interactive elements to explore the data, these users might even uncover insights that Abhi or Natasha have not seen .

2. New England Patriots Fans

This set describes all fans of the Patriots football team, including some intense football fans. These fans will be most curious because the article will be about their quarterback, and will be interested in learning about whether he deserved the MVP designation. As most Pats fans are aware, many other

football fans dislike their team, and will be looking for insights and evidence about why their team is good, providing data for arguments they may engage in.

3. General Football Fans

Many people are aware about the game of football, its goals and rules, and follow the sport casually. It is likely that these fans occasionally or regularly follow a particular team. This set of fans is the most generic, and may include other sets of audiences mentioned in this section. Readers who fall into this group will be looking for insights about football in general, learning more about a sport they follow and figuring out which statistics are important for them to make their own opinions about players or teams.

4. Football Novices

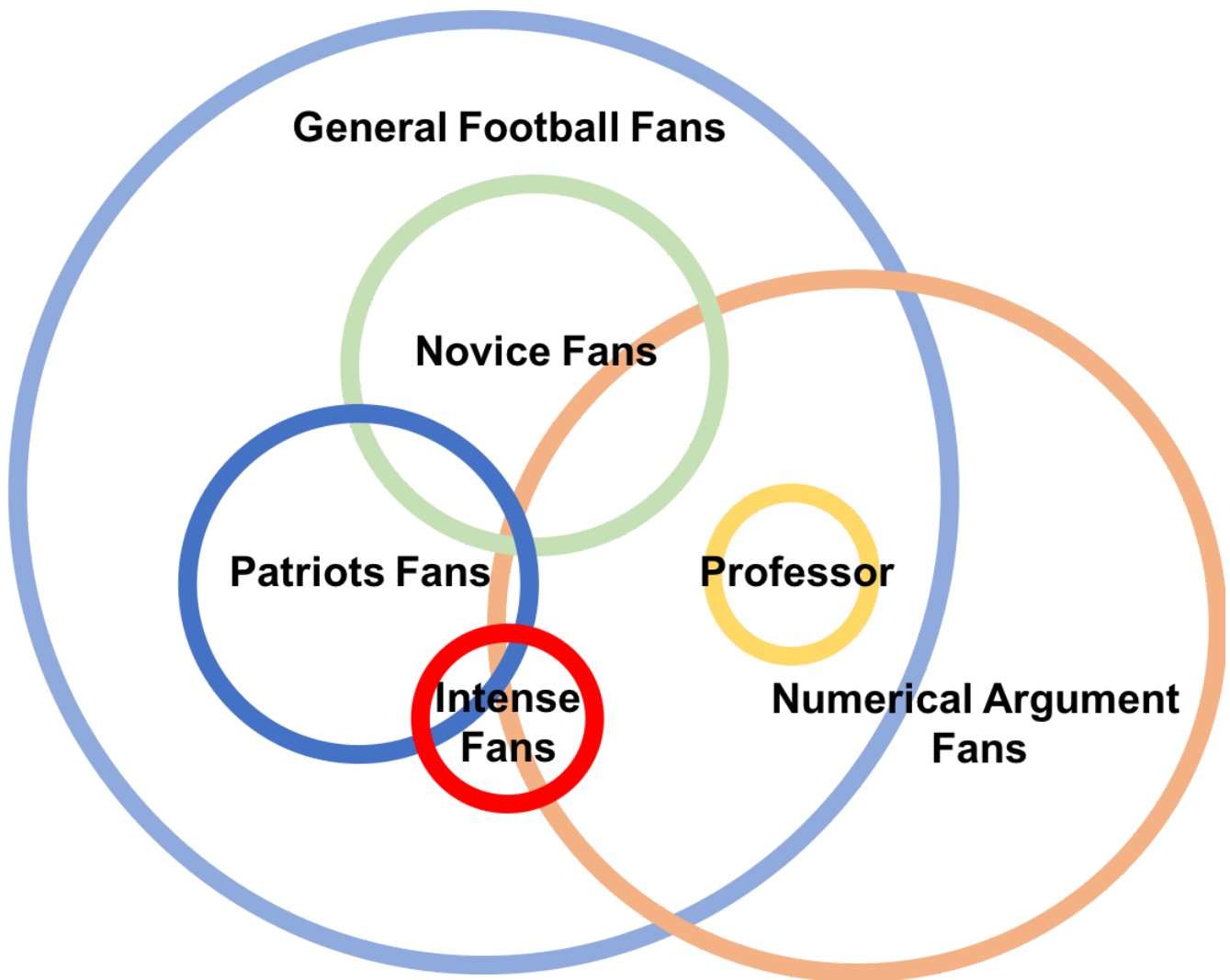
These fans have little to no prerequisite knowledge about the sport, but are aware of the existence of the sport. There is a chance they have attended a Super Bowl party and have caught parts of the game. These readers will be looking to learn more about the sport, but from more of a beginner's perspective. We plan to guide them through the article by pointing them to terms or stats which they are not aware of and are important.

5. Fans of Numerical Arguments

We do not assume that these readers have any knowledge or interest in football. It is possible that they may indeed be fans of football, but may simply be looking for a good statistical argument. Many readers of 538 do not read articles because they have knowledge of the given domain, but instead because they more easily digest information when presented in a numerical form, and thus may be reading our article for the same reasons.

6. Professor Jeff Rzeszutarski

The professor is a major audience who seems to fit under the groups of "General Football Fans" and "Fans of Numerical Arguments", and we find it fitting to include him in our intended audience listing as he is the most important of our audiences. We also make the assumption that the professor *definitely* does not fit under Audience 2, New England Patriots fan. We believe him to be a Pittsburgh Steelers fan, for which we are disappointed but do forgive him.



A set graph of the listed intended audiences, in order to understand how they may relate

By no means do we intend to cater to every listed audience, however we do aim to make our project *accessible* to all of these groups. In order to do so, we aim to make the article friendly to those who do not know the rules of football well, yet in a manner that does not distract from the premise of the article or turn off readers who are already in possession of this knowledge. In particular, we realize that there is heavy jargon in football, and we plan to allow users to learn these terms by providing a dynamic glossary that fits in the flow of the article. As mentioned by the professor in our design outline feedback, we will provide short descriptions of important words by highlighting the terms. These terms can be clicked to display a pop up box, which will have a short description of the word, and a link to an external resource which goes into further detail. A successful example of this occurs in the blog "[Wait But Why](#)", which we will mimic, as shown below.

its own food through photosynthesis. It takes CO_2 from the air and absorbs light energy (O_2). The plant keeps the carbon and oxygen stays in the plant as chemical energy of carbon and stored chemical energy.

What we're doing is reversing the photosynthesis in wood—that's why trees aren't constantly on fire.

A term can be highlighted or pointed to.

It takes CO_2 from the air and absorbs light energy from the sun.

The plant keeps the carbon and emits the oxygen.

What we're doing is reversing the photosynthesis. Normally, plants are constantly on fire.

That's why trees aren't constantly on fire.

What we're doing is reversing the photosynthesis. Normally, plants are constantly on fire.

That's why trees aren't constantly on fire.

What we're doing is reversing the photosynthesis. Normally, plants are constantly on fire.

That's why trees aren't constantly on fire.

Clicking the highlight will lead to more information

Related and Inspiring Materials

We have found many articles which are inspirational to this project, but we highlight 3 main articles that accurately describe our vision for our project.

1. 538's "[Sneaky Stats that could decide the Super Bowl](#)"

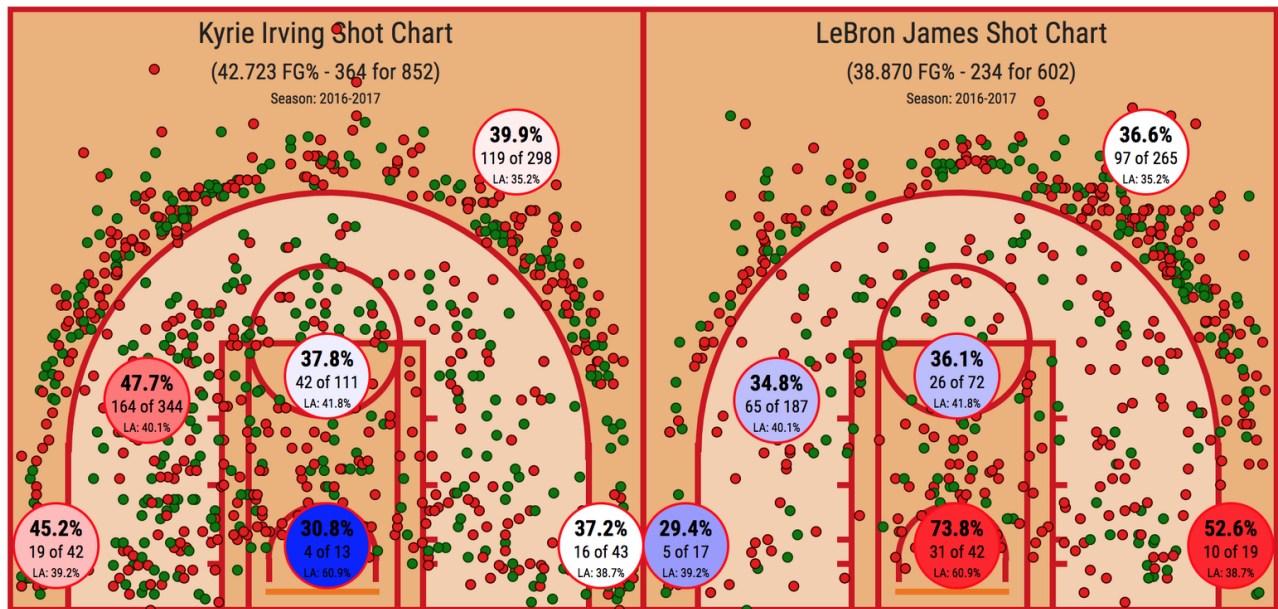
This article does a great job mixing writing with graphics. Although not interactive, the graphs shown in this article are easily digestible and are the center of the arguments being made by the author. All of the writing clearly explains key concepts succinctly, while also providing links to more detailed explanations. Statistics are used well to make an argument, and descriptions about their importance are also given. The writer also clearly highlights the data that they desire to use to focus on amongst relatively crowded graphics. If well executed, we envision our article having similar balance.

2. ESPN's "[Quarterback Carousel](#)"

In the 2018 NFL off-season, many different quarterbacks are going to be drafted or changing teams, with many different teams looking for a new quarterback. This visualization takes a very simple concept of pairing quarterbacks with teams and makes an incredibly easy and interesting visualization. Further, the visual elements make this article more interesting to engage with. The user's options are limited, in that they are only able to change the items on 2 axes, picking a quarterback and team, and clicking a "Pick this match" button to learn more about that pairing. We really liked this article and imagine the simplicity of this article to be reflected in our project.

3. NBA Savant's "[Shooting Chart Comparison](#)"

At its root, this site is purely a sports visualization site. You are able to compare basketball players and their shot charts. Here, we are able to see all the shots a player took across a season and what their accuracy is in a region of the court. We imagine that we will also be able to provide swaths of data, and a summarization statistic about the data that makes it much easier to understand. We have shown an example below.



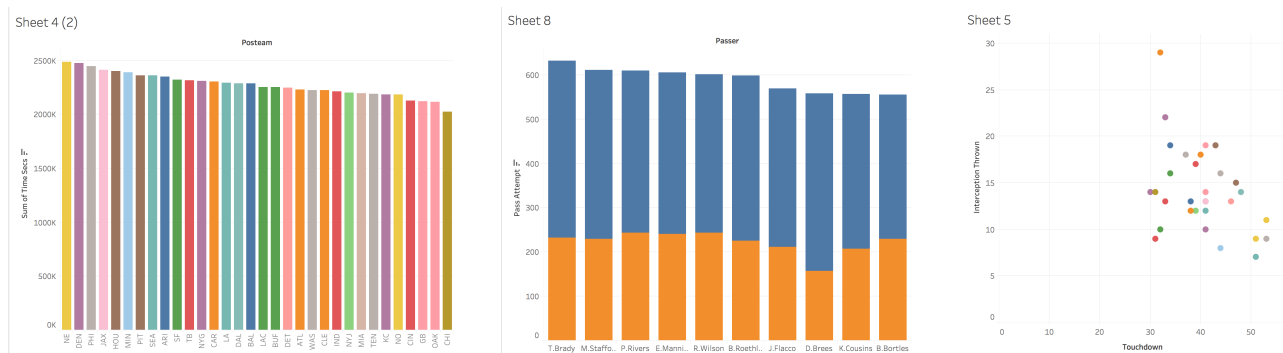
A shot chart comparison between Kyrie Irving and LeBron James from NBA Savant

Our Data Source

We are using the [nflscrapR-data](#) dataset, which contains a play by play documentation across the 2017 NFL Season. Abhi had previously worked on a Fantasy Football project and had exposure to the authors of this dataset. We were excited to see that the most recent season existed. The reason we chose this dataset is because it is heavily detailed and very comprehensive. The author originally scraped all the data off the [NFL API](#) and joined many tables together to form this one. The fields listed are well specified which makes it very easy for us to work with the dataset.

Design Considerations and Iterations

We first started off by exploring our dataset in Tableau. With some initial plots we realized the dataset was too large and multidimensional to get significant insights from Tableau plots. We realized we needed to think deeply about which statistics are relevant and important to explore. In addition, we did not want to bias our answer. It is easy to create visualizations showing Tom Brady deserved the MVP by exploring which statistics he did well in. An unbiased conclusion would need to come from significant football insight (Abhi's domain) and relevant data analysis (Natasha's domain).



Initial Design / Dataset Explorations in Tableau

We came upon three different themes of exploration:

How does Tom Brady at getting the ball down the field? Does this change under pressure?

How efficient is Tom Brady when he is on offense?

How does Tom Brady affect his team's ability to win when he gets onto the field?

Below we describe how we arrived at each and the tradeoffs associated.

We considered many different designs across the visualizations we were thinking of. First, we decided that *completion percentage* would be the metric we would use to gauge a quarterback's accuracy, because it is the predominant metric used when quantitatively accessing a QB's completed passes. Completion percentage is calculated as $\frac{\text{\#passes completed}}{\text{\#passes attempted}}$. We are unable to show total number of passes attempted with this metric (a tradeoff). Completion percentage might not tell the whole story, such as the comparison between a 100% completion percentage for 1 pass attempt versus a 65% completion percentage for 100 pass attempts. We acknowledge this tradeoff and aim to alleviate it by filtering quarterbacks by a minimum number of passes attempted.

We next considered how we would bucketize the data we wanted to show in this graphic. Since we aim to show accuracy by distance downfield, we believed it would be good to bin the data by magnitudes (0-10, 11-20, etc.). In doing so, we explored the idea of scaling the height of each bin according to the distance downfield from the starting line. 100% accuracy in the 0-10 bin would reach a maximum height that corresponded to 10 yards on the football field background. Likewise, 100% accuracy in the 11-20 bin would have a max height at the 20 yard line. Although this shows distance well, it skews a user's impression of the importance of a bin by the height of the bin. It is not true that a pass 60 yards downfield is more important than a pass 10 yards downfield, and thus we did not want to enforce such a scale as each bin is as important as the others. We also wanted to allow for easy comparison across bins, and this would not lend itself well to that goal. Given these considerations, we decided to nix the idea of making the height of a bin correspond to the maximum yardage in that bin.

We also realized it was important to analyze the time a team spent on offense to understand the impact of a quarterback on his team. However, after a closer look, we determined that a simple analysis of time on offense was not a sufficient metric. A team may spend more time on offense and not score, or a team may spend a sparse amount of time on offense and score every time. This lead to a realization that we care about *efficiency* rather than an objective measure of time. As a result, we decided to analyze several dimensions about a team's time on offense to compare the effectiveness of the quarterback.

This thought process forced us to think about both time on offense and defense. We knew that there was some insight we could extract from the two. Eventually, we determined that we would analyze a QB's impact on his team by looking at how his team's win probability changes when he steps onto the field. Thus, we decided on a slope graph to measure the change in win probability during changes from defense to offense in a game.

Finally, we explored creating our own statistic that was similar to Baseball's WAR, or Wins Above Replacement, to quantify how valuable a player is to his team. However, after further consideration, we realized many issues would arise by creating our own statistic. Namely we would need to explain and defend our statistic in regards to why it makes sense, is better than existing statistics, and unbiased. This

would force our analysis to diverge from the focal point of whether Brady deserved the MVP, and would not allow for much visualization during discussion about the statistic.

Poster Session Insights

1. Classmates mentioned the desire for [previous MVP's](#) in our analyses. This is a good suggestion to help users not only compare against players from the current season, but players who won the MVP in previous years to understand what metric ballpark is required to win the award.
2. The slope graph between offense and defense win probabilities received strong support by classmates once the idea was explained. It will definitely be in our final project.
3. In our completion percentage graphs, we were originally going to allow for 2 quarterbacks and the NFL Average, with Tom Brady and the NFL Average being fixed items. However, many viewers believed that it would be better if we did not fix any item except Brady so that Brady could be compared to more than 1 quarterback. In our final visualization, we will allow NFL Average to be toggled.
4. As mentioned earlier, there were discussions about football novices would access the football terminology we will be using. To alleviate this, we will allow for explanations in pop up boxes.
5. Preset defaults would help immensely in telling a story. We can show readers a story by orienting interactive graphs with a default setting, say comparing Tom Brady with Ben Roethlisberger. As a result, we will carefully curate our graph defaults.
6. We will be using some static visualizations, or fixed graphs, to make our argument. This will help us convey our reasoning to novice fans without overwhelming them with too many options they may not know how to use.
7. We will begin our writeup with an exploration into how the NFL currently chooses its MVP, to provide perspective for readers.

Final Design Description

We have mentioned how the feedback received has helped us modify or extend our plan in the previous section, such that we can have an effective final design.

It is important to note that we have not yet computed the graphs that we will be searching through, and thus the metrics we have chosen to determine an MVP are not biased in any form towards Tom Brady and do not lean us towards an answer.

Tom Brady for MVP?

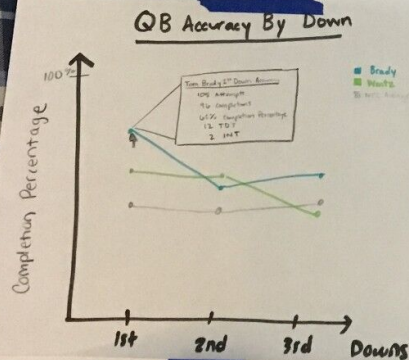
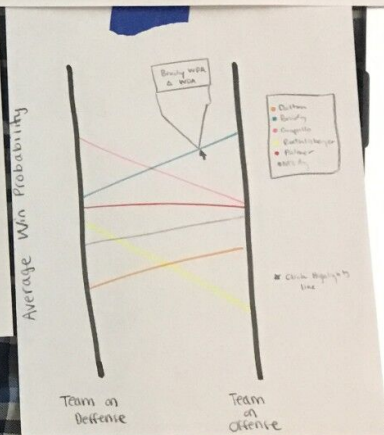
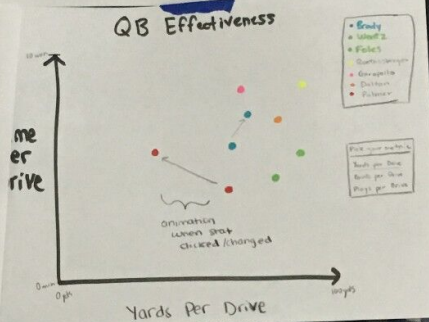
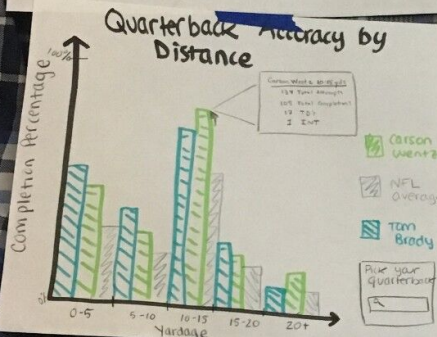
A case by #s.

Abhi Gupta, aag295
Natasha Ambrust, nka8

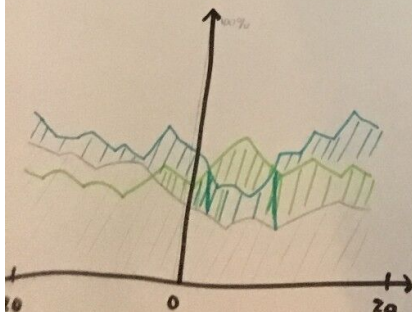
Is Tom Brady more effective at getting the ball down the field?

How effective is Brady when he has the ball?

How Valuable is Tom Brady to the Patriots?



QB Accuracy By Point Spread



Final poster framework for the final project

As seen, we have our general overview of the writeup, the major questions we ask, and how we measure these fields. We go into detail about these graphics below.

MVP Breakdown / Introduction

We answer introductory questions such as “What is an MVP?”, describe baseline knowledge that primes the reader with information they should know, such as how the NFL selects the MVP, which players were in contention for the award, and who won. We then set the stage for the controversy around the award, and propose the central question, “Did Tom Brady deserve the NFL MVP award in 2017”?

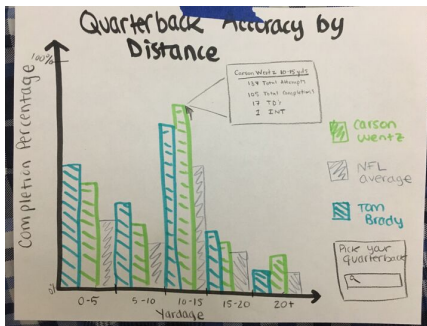
Who are we comparing

Our comparison is only amongst quarterbacks, and the reasoning will be described in detail in this section. We aim to highlight the QB’s high level of impact on a game relative to other positions, and potentially link other sources which highlight this fact. Namely, QB’s touch the ball on every play, and make the most important decisions that influence his team’s success (audibles, checking the defense and modifying the play, getting the ball from the snapper to the running back or receiver). We further want to compare metrics which Tom Brady is relevant to, which allows for QB comparisons to be the most applicable. We will also compare previous MVP’s to Brady, and looking at the previous 30 years, 22.5 MVP awards have gone to Quarterbacks (before 2003, Co-MVP’s were allowed). Thus, we deem that QB’s are the most likely to win the award, and thus we will analyze them in our article.

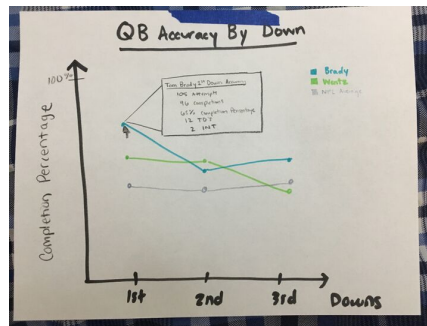
We will be analyzing Brady on three different dimensions.

1. Is Tom Brady more effective at getting the ball downfield?

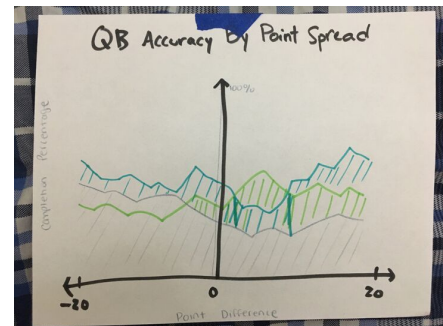
This dimension analyzes QB accuracy by looking at completion percentage. We look at effectiveness by pass distance, accuracy by down, and accuracy by point spread in the game. Each of these graphs highlight the QB’s effectiveness based on difficulty of pass (distance), difficulty of situation (down), and overall pressure (point spread).



QB Accuracy by Distance



QB Accuracy by Down

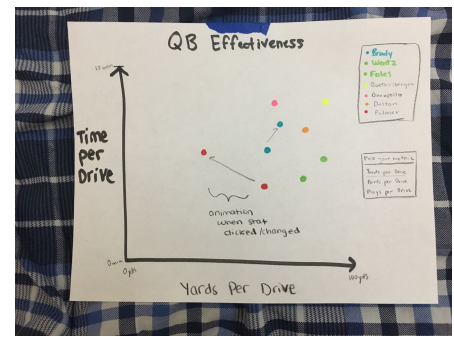


QB Accuracy by Point Spread

Our interactions are as follows: for each graph users will be able to plot up to 2 quarterback selections (including NFL average) next to Tom Brady. Upon hovering over a mark, we will display more information on the quarterback and the specific value that is being plotted for that quarterback. Upon selecting quarterbacks we hope to have an animation feature to make the corresponding UI changes in the graph engaging.

2. How efficient is Tom Brady?

This metric looks at the efficiency of a QB based on a series of metrics including time per possession, yards per drive, points per drive, and plays per drive. This graph is effectively three graphs in one due to the ability to vary the x-axis. Readers will be able to select (yards per drive, points per drive, and plays per drive) to plot on the x-axis against the y-axis of time per drive. We maintain the y-axis as constant (time per drive) to display the effectiveness and efficiency of a QB. After our analysis, we will determine which metric tells the most important story and we will set the graph to show this metric on the x-axis as default.



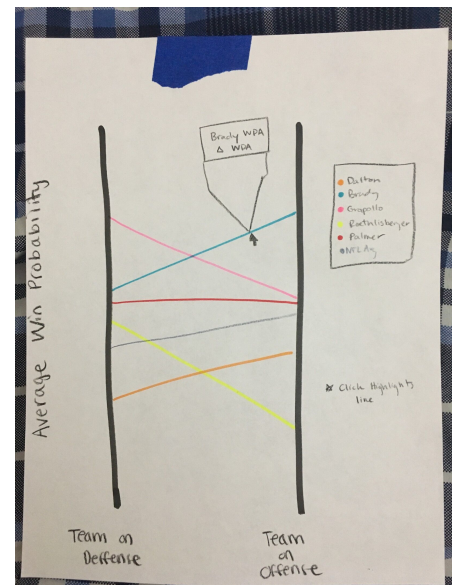
QB Efficiency by Metric Selection

Our interactions are as follows: a user will be able to pick which metric they wish to view on the x-axis: yards per drive, points per drive, and plays per drive. Upon metric selection, our graphic will animate each player mark in the filtered by the current metric to their corresponding mark in the new metric. Since time per drive will not change, a user will effectively see only an x-axis change, which we hope will reduce the confusingness of many points moving in the animation. Users will be able to hover over the points for more information including exact numbers for each of the metrics.

3. How Valuable is Tom Brady to his team?

We look for player value by looking at his impact on his team's win percentage probability, utilizing a slope graph. By comparing offense and defense WPP, the slope from defense to offense will indicate how much that given player impacts their team.

Our interactions are as follows: users will be able to select which quarterbacks to plot with Tom Brady. Users will be able to hover over the lines for more information including exact slope of the graph. Since we are worried that the graph will become too cluttered, we plan to have a "Show All" feature which shows all the lines and a filtering feature in which users can filter specific quarterbacks to plot.



QB Value by WPP Change

Time permitting, we also aim to look at Red Zone efficiency to look at the ability of a QB to guide his team towards points when he is close to scoring (within 20 yards), and another graphic which shows the QB's DVOA (Defense-adjusted Value Over Average), or the predominant "WAR" statistic in Football. We also do not want to overload our reader with information.

It will be hard for a novice reader to get to this point in our article. So depending on how easy our graphs are to interpret, we may or may not reduce the number of graphics we have in our final article. We will determine this through user testing. We will conduct user testing with each of our audiences described above.

Conclusion

Finally, we will summarize our analysis and reach a conclusion that answers our question. Again, since we have not preprocessed the data yet, our final design is not biased towards or against Tom Brady (Natasha will keep Abhi's inherent biases toward Brady in check 😊). We will form our conclusion based off the insights from our visualizations.

Project Implementation Steps

Since all of our data has already been collected in our dataset, we do not need to scrape any sites.

In order to clean our data, we plan to cut the data into separate datasets for each graphic. This is because there will be a high load time for our project if we compute each of the given dimensions each time we load the page. By generating separate CSV's, we will be able to preprocess, perform all of our calculations ahead of time, and ensure we have the data in the form we require.

Once we have our data separated, we can manipulate our data within graphics by applying filters and using sorting logic (to bucketize, or list smallest to largest). Most of the heavy manipulation will be completed when we generate our CSV's in the clean step.

We will be using d3 to implement our final project designs. We gained more experience using animations in d3 in homework 3. Although frustrating at times, d3 will give us the necessary flexibility to be able to create the graphics we desire. We also plan to use frameworks such as bootstraps to help with some of the web development.

Identify how each team member will contribute to the final project

At the time of reading this article, we have several tasks in front of us to complete. We see these steps and contribution assignments as follows:

1. Creation of CSV's

The dataset we are working with is massive, and performing in browser computations for each graph will be very expensive and slow. Instead, we aim to generate the modified datasets for each graphic beforehand and increase our load time. Both members will help with this task.

2. Writing the Article

Our article will require a decent amount of explanation and football related knowledge, and thus Abhi will take charge of this section. Natasha will be providing proofreading for this task.

3. Visualization

Within the article, we will be making our interactive graphics. Natasha will be taking lead on this section, as she has more extensive knowledge of web development tools. Abhi will be helping develop as well.

4. User Testing

Once the article is complete, we will be performing some user tests in order to gauge the effectiveness of our graphics. This will mainly help us learn how users understand our tool, and determine if there are any changes we need to make. Both members will contribute.

5. Final Touches

There are many small final touches we expect to perform towards the end of this project. We anticipate both members modifying the final product as they see changes that need to be improved.

6. Video

The final task for this project will be to create a video explaining our final product. Both members will help create this.