



## Introduction & Motivation

National Football League (NFL) teams spend a large amount of time and resources studying their opponents to gain insights into their tendencies. One such characteristic is a team's propensity to run or pass the ball in a given situation. For a defense, having a sense of whether an opposing offense will run or pass informs decision-making about play-calling, personnel groupings to deploy, and physical positioning on the field. Using a combination of NFL play-by-play data, information on offensive formation, and metrics of player quality for each position group, we hope to build a model to predict whether any given offensive play will be a run or a pass.

## Data



- Football Outsiders
  - Proprietary NFL play-by-play data (includes formation data)
  - Snap counts for every NFL player
- Publicly available Madden video game ratings downloaded from [maddenratings.weebly.com](http://maddenratings.weebly.com).

| YEAR | WEEK | OFFENSE | DEFENSE | ACCURACY | PASS PROPORTION |
|------|------|---------|---------|----------|-----------------|
| 2012 | 15   |         |         | 0.938    | 0.446           |
| 2012 | 12   |         |         | 0.932    | 0.865           |
| 2013 | 8    |         |         | 0.929    | 0.829           |
| ...  | ...  | ...     | ...     | ...      | ...             |
| 2014 | 3    |         |         | 0.545    | 0.519           |
| 2012 | 6    |         |         | 0.536    | 0.464           |
| 2013 | 6    |         |         | 0.486    | 0.667           |

## Features

### Raw Input

- Score Difference
- Current Quarter
- Time Remaining in Quarter
- Current Down
- Distance to First Down
- Number of Offensive Players per Position
- Number of Defensive Players per Position
- Offensive Formation (e.g. shotgun, no huddle)
- Indicator of an offensive player out of position
- Turnovers
- Indicator of whether offense is at home

### Derived Features

- Proportion of pass plays over week, season, last 50 plays
- Proportion of passes faced by defensive team over week, season, last 50 plays
- Indicator of a team being up or down by more than once score
- QB pass completion rate over week, season, last 25 plays
- Weighted Madden rating of each offensive/defensive position group.

We utilized a total of 17 features capturing the context behind a given play as well as overall tendencies and strengths of each team, which are generally critical factors in play-calling decisions.

## Models

- Logistic Regression** - Classifies training examples through the logistic function  $\frac{e^{\beta_0 + \beta_1 x^{(i)}}}{1 + e^{\beta_0 + \beta_1 x^{(i)}}}$  where  $\beta_0$  and  $\beta_1$  are fit via maximum likelihood
- Linear Discriminant Analysis** - Dimensionality reduction technique and classifier that uses Bayes' Theorem to make a linear classification.
- Random Forests** - Uses a collection of bootstrapped training sets to train decision trees to make a classification. To reduce high variance among trees, each split of the decision tree chooses from a subset of all features (of size  $\sqrt{n}$ ).
- Gradient Boosting Machine** - An ensemble method that combines several weak-learning decision trees into a strong classifier. In our model, we found that 300 weak learners achieved the best results.
- Mixed Model** - Weighted average of the probabilities from both the random forest (40%) and the gradient boosting machine (60%) to derive a classifier that is better than each individually.

## Results

| Model                        | Training Accuracy | Test Accuracy |
|------------------------------|-------------------|---------------|
| Logistic Regression          | 0.728             | 0.727         |
| Linear Discriminant Analysis | 0.721             | 0.714         |
| Random Forest                | 1.000             | 0.737         |
| Gradient Boosting Machine    | 0.750             | 0.744         |
| <b>Mixed</b>                 | <b>0.906</b>      | <b>0.746</b>  |

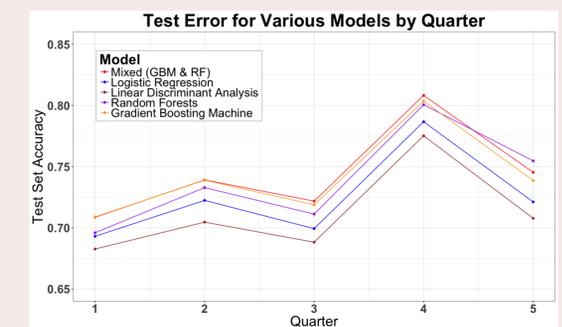
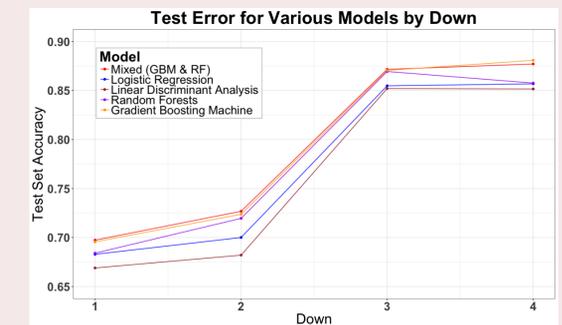
### Best and Worst Games for Prediction Accuracy

| Year | Week | Offense | Defense | Accuracy | Pass Proportion |
|------|------|---------|---------|----------|-----------------|
| 2012 | 15   |         |         | 0.938    | 0.446           |
| 2012 | 12   |         |         | 0.932    | 0.865           |
| 2013 | 8    |         |         | 0.929    | 0.829           |
| ...  | ...  | ...     | ...     | ...      | ...             |
| 2014 | 3    |         |         | 0.545    | 0.519           |
| 2012 | 6    |         |         | 0.536    | 0.464           |
| 2013 | 6    |         |         | 0.486    | 0.667           |

### Best and Worst Team-Seasons for Prediction Accuracy

| Year | Offense | Prediction Accuracy | Pass Proportion |
|------|---------|---------------------|-----------------|
| 2014 |         | 0.850               | 0.516           |
| 2013 |         | 0.813               | 0.699           |
| 2012 |         | 0.812               | 0.665           |
| ...  | ...     | ...                 | ...             |
| 2014 |         | 0.655               | 0.502           |
| 2013 |         | 0.667               | 0.633           |
| 2013 |         | 0.670               | 0.499           |

## Plots



## Discussion

- Our plots for down and quarter align with our intuition - teams are easier to predict with fewer yards to go, which correlates with down. 2nd and 4th quarter prediction accuracies are higher because the ends of these quarters have outside effects on the outcome of the game. In particular, 4th quarter accuracy is highest because score margin and time remaining often directly dictate play-calling in end-of-game scenarios.
- We hypothesize that our models tended to do worse with mobile QBs because signal-callers with the ability to scramble often turn designed pass plays into runs.

## Future Work

- Our dataset included a timeout feature, but failed to include which team called the timeout. We suspect that knowing how many timeouts a team has remaining would help prediction of end-of-half scenarios, and, in the future, we hope to obtain this data.
- The dataset indicated the direction of the offensive play. We could explore predicting a team's next play as well as the direction of the play.

## References

- [1] W. Burton and M. Dickey (2015). NFL play predictions. In JSM Proceedings, Statistical Computing Section.
- [2] Hastie, T., Tibshirani, R., & Friedman, J. H. (2009). The elements of statistical learning: Data mining, inference, and prediction. New York: Springer.
- [3] Booz Allen Hamilton (2016). Assessing the Predictability of an NFL Offense. MIT Sloan Sports Analytics Conference.