

Predicting Fantasy Points Project Report

Jarett Smith, Max Thompson, Charlie Deaton, and Balin Allred

2023-12-14

Contents

| | |
|---|----------|
| Executive Summary | 2 |
| Results: | 2 |
| Motivation: | 2 |
| Value-Add: | 2 |
| Conclusion: | 2 |
| The Data | 3 |
| Player performance | 3 |
| Feature Engineering | 3 |
| Final Features: | 4 |
| Models | 5 |
| Train, Test, CV | 5 |
| XGBOOST Model | 5 |
| GBM Model | 5 |
| Random Forest Model | 6 |
| GLMNet Model | 6 |
| Neural Net Model | 7 |
| Appendix: | 8 |
| Libraries Used | 8 |
| Code to Produce the Final GBM Model | 8 |
| Code to Produce XGBOOST Model | 8 |
| Code to Produce Original GBM Model for Model Comparison | 9 |
| Code to Product Random Forest Model | 10 |
| Code to Produce GLMNet Model | 10 |
| Code to Produce NNET Model | 11 |

Executive Summary

Results:

Our model was able to better predict player performance compared to an industry standard, ESPN. We chose GBM model with an RMSE 3.14 points per game that beat ESPN's RMSE of 3.65 points per game. This represents a 14% beat of ESPN's model.

Motivation:

Fantasy sports allows the general public to essentially manage their own virtual franchise. Just like real-world sport franchises, fantasy managers draft their own players according to their belief in a player's ability to attempt to beat other teams in an imaginary "league". According to CNN, roughly 60 million people in North America play fantasy sports. Of the various types of sports, fantasy football is the most popular, with roughly 40 million users. With this large amount of participation, the fantasy sports industry as a whole is worth about \$20.3 billion and is expected to continue to grow at 14% annually. This monetary value stems from increased viewership of sporting events and the legalization of sports betting in several countries.

Value-Add:

With this fairly new industry rapidly growing, our goal is to gain a competitive advantage using data mining methods to better predict player performance. By outperforming industry giants like ESPN, there is real opportunity to fiscally benefit for those who are monetarily invested and opportunity to appeal to recreational participants. In short, this model is valuable to both financially invested individuals and to your "just-for-fun", casual participant.

Conclusion:

In summary, our model outperforms ESPN estimates of player performance by 14%. In doing so, we are positioned to provide services to a large audience as higher accuracy predictions have real monetary implications to both die-hard and recreational consumers.

The following table shows our results for the different models tested.

| Model | Train_RMSE | RMSE_SD | Test_RMSE |
|---------------|------------|-----------|-----------|
| XGB | 3.109668 | 0.0954447 | 3.114582 |
| GBM | 3.109920 | 0.0892284 | 3.091714 |
| Random Forest | 3.130287 | 0.1366360 | 3.143737 |
| GLM | 3.228058 | 0.1140531 | 3.180308 |
| Neural Net | 6.897551 | 0.1685324 | 6.513741 |

We selected the GBM model because it outperforms the simpler regularized regression model and had the best test RMSE based on our results. The following table shows the outputs for our final model after re-training our best model with all available data (seasons 2022 and before). The current season RMSE is how we are doing so far in the 2023 season using this model.

| Model | Train_RMSE | RMSE_SD | Current_Season_RMSE |
|-------|------------|-----------|---------------------|
| GBM | 3.321967 | 0.0869438 | 3.139346 |

| | Overall |
|-----------------|------------|
| points_pg_ly | 100.000000 |
| targets_pg_ly | 12.372682 |
| wopr_pg_ly | 14.106260 |
| pick | 16.828302 |
| air_yards_pg_ly | 2.652639 |

Note: These are the top 5 variables listed in order of importance.

The Data

Our data comes from play-by-play data of the National Football league dating back to 1999. nflfastR, along with the rest of the nflVerse packages, were used to scrape the data from ProFootballReference, ESPN, and other sites.

At its most granular level, a single row of data represents the performance of a single player, on a single team, for a single week in a single season.

In other words, up to 17 rows of the same player for the same season would make up 1 year of performance for that player.

Player performance

Player performance is measured in fantasy points per game for a single player for a single season. These are generated by the following formula:

```
Fantasy Points = 1 * catches + 0.1 * receiving_yards + 6 * receiving_touchdowns +
0.1 * rushing_yards + 6 * rushing_touchdowns + 2 * fumbles_lost
```

In other words, pass catchers that catch more passes, for more yards and for more touchdowns will score more fantasy points than those who don't.

In order to account for injuries / incomplete seasons, points per game played is used instead of total points. This is calculated by dividing total points by the number of games played.

Feature Engineering

In order to convert the data from game-by-game to season-by-season, the data was grouped by player and season and aggregated for multiple different features that would later be used in the modeling process.

As a whole, there were four primary types of features that answer the following questions:

1. Player Performance – How good did this player do last year?
2. Player Attributes – What is this player's overall physical and historical background?
3. Team Performance – What is this player's surrounding team like?
4. Player Situation Attributes – What surrounding pieces around this player changed since last year?

Player Performance

In order to predict the performance of player X for season 2022, knowing any information about 2022 performance would be ‘cheating’ the model, so data was lagged by 1 season to account for this. If one was to predict Tyreek Hill’s 2022 points per game, they would only have information on his 2021 performance, not his 2022 performance.

Player Attributes

For a given year, there are attributes of a player that have nothing to do with past performance, but could still be relevant. Some examples include: * years of experience entering this given season * age entering this given season * when the player was drafted (if at all) * an assigned cluster based on physical attributes such as height, weight, and speed

Team Performance

If we know what team a player is starting a given season playing for, we can also know “past performance” of that team. This is done by aggregating the performance of all players on that team for the previous season. Some examples include: * total fantasy points scored by that offense last season * total passing yards / touchdowns scored by that offense last season * information on how good / bad the quarterback on that team played last season

Player Situation Attributes

Finally, even if we know a player’s attributes, their past performance, and how their team did last season, there are still changes season-to-season that could affect a player’s performance. Some examples include: * did the player change teams? * did this team add / lose any key players? * did this team change their head coach? * did this team invest in players that play the same position during the NFL draft over the offseason?

Final Features:

After considering all four categories, the following features were engineered to use in the modeling process:

Table 4: Fields used for Models

| | | |
|-----------------------------|-------------------------|--------------------------------|
| points_pg_ly | years_pro | is_returning_coach |
| targets_pg_ly | targets_added_this_year | hc_years_with_team |
| wopr_pg_ly | is_on_new_team | catch_rate_ly |
| pick | points_per_snap_ly | position |
| air_yards_pg_ly | targets_per_snap_ly | fp_dropoff |
| total_games_ly | wopr_per_snap_ly | starter_epa_passing_ly |
| pass_attempt_difference | total_snaps_ly | starter_epa_persnap_passing_ly |
| total_positional_investment | total_passing_tds_ly | combine_cluster |
| target_dropoff | total_passing_yds_ly | qbr_ly_bin |
| epa_pg_ly | total_passing_fp_ly | |

Models

Train, Test, CV

```
set.seed(2023)
train_rows <- sample(nrow(master), round(nrow(master)*0.7,1), replace = FALSE)

train <- master[train_rows,]
test <- master[-train_rows,]

ctrl <- trainControl(method = "cv", number = 5)
```

This sets up train and holdout (or test) data sets to train models on. By using train and holdout data sets, we can ensure that our model is not over-fitting to the training data and examine how well the model generalizes to new data.

XGBOOST Model

What XGBOOST models are and how they work:

XGBoost, short for Extreme Gradient Boosting, is an optimized and powerful implementation of the gradient boosting algorithm. It's widely used in machine learning for regression, classification, and ranking tasks due to its efficiency, speed, and high performance. It operates by building an ensemble of weak learners, commonly decision trees, to progressively minimize prediction errors. It showcases optimized tree construction methods, parallel processing, and regularization techniques like L1 and L2 regularization to prevent overfitting. XGBoost handles missing values, provides feature importance, and is known for its speed, scalability, and high performance in handling large datasets.

The tuned hyper-parameters and output results are show below:

| | nrounds | max_depth | eta | gamma | colsample_bytree | min_child_weight | subsample |
|-----|---------|-----------|------|-------|------------------|------------------|-----------|
| 492 | 150 | 3 | 0.05 | 0 | 0.5 | 10 | 0.9 |

| Model | Train_RMSE | RMSE_SD | Test_RMSE |
|-------|------------|-----------|-----------|
| XGB | 3.109668 | 0.0954447 | 3.114582 |

GBM Model

What GBM models are and how they work:

GBM stands for Gradient Boosting Machine, it's a machine learning technique used for both regression and classification problems. It also happens to be an ensemble learning method just like Random Forests, but it operates differently. GBM's work by sequentially building multiple decision trees, each one correcting the errors made by the previous tree. It uses gradient descent optimization to minimize the overall error by fitting new trees to the residuals of the previous predictions. This process continues until the model converges, combining the predictions of these trees to make the final prediction. GBM excels in capturing complex relationships but requires hyperparameter tuning to prevent overfitting.

The tuned hyper-parameters and output results are show below:

| | n.trees | interaction.depth | shrinkage | n.minobsinnode |
|----|---------|-------------------|-----------|----------------|
| 24 | 150 | 7 | 0.03 | 9 |

| | Model | Train_RMSE | RMSE_SD | Test_RMSE |
|---|-------|------------|-----------|-----------|
| 2 | GBM | 3.10992 | 0.0892284 | 3.091714 |

Random Forest Model

What Random Forest models are and how they work:

Random Forest is a versatile machine learning algorithm used for both classification and regression tasks. It's an ensemble learning method combining multiple models to make more accurate predictions compared to a single model. It constructs multiple decision trees using random subsets of the data and features. It combines these trees to make predictions by averaging for regression. By leveraging the diversity of these trees and their collective decision-making, Random Forest reduces overfitting, handles large datasets well, and provides insights into feature importance.

The tuned hyper-parameters and output results are show below:

| | n.trees | interaction.depth | shrinkage | n.minobsinnode |
|----|---------|-------------------|-----------|----------------|
| 24 | 150 | 7 | 0.03 | 9 |

| | Model | Train_RMSE | RMSE_SD | Test_RMSE |
|---|---------------|------------|----------|-----------|
| 3 | Random Forest | 3.130287 | 0.136636 | 3.143737 |

GLMNet Model

What GLMNet models are and how they work:

GLMNet models use regularized regression to control the complexity of the model and prevent over fitting. These models are particularly useful for high-dimensional data sets where the number of predictors is large relative to the number of observations. The models use traditional linear regression in addition to penalty terms, alpha and lambda, to perform variable selection and handle correlated predictors. The tuning parameter, lambda, controls the strength of regularization, with larger lambda values leading to simpler models. Similarly, larger alpha levels emphasizes lasso regression which also leads to simpler models.

Optimal alpha and lambda levels were found using 5-fold cross-validation. We found the optimal tuning parameters shown below. Essentially, these are the tuning parameters that minimized the RMSE of the model.

The tuned hyper-parameters and output results are show below:

| | alpha | lambda |
|----|-------|--------|
| 77 | 1 | 0 |

| | Model | Train_RMSE | RMSE_SD | Test_RMSE |
|---|-------|------------|-----------|-----------|
| 4 | GLM | 3.228058 | 0.1140531 | 3.180308 |

Neural Net Model

What Neural Net models are and how they work:

Neural network models, often referred to as artificial neural networks (ANNs), are a class of machine learning models inspired by the structure and functioning of the human brain’s neural networks. It consists of interconnected nodes organized in layers. They process data through these layers, adjusting connection weights iteratively during training to make accurate predictions. Using activation functions and backpropagation, neural networks learn patterns and relationships within data, enabling them to handle complex tasks such as image recognition, language processing, and more.

| | size | decay |
|---|------|-------|
| 5 | 2 | 1e-04 |

| | Model | Train_RMSE | RMSE_SD | Test_RMSE |
|---|------------|------------|-----------|-----------|
| 5 | Neural Net | 6.897551 | 0.1685324 | 6.513741 |

Appendix:

Libraries Used

```
library(dplyr)
library(xgboost)
library(DALEX)
library(glmnet)
library(caret)
library(gbm)
library(tidyverse)
library(nflverse)
library(zoo)
library(plotly)
library(visdat)
library(groupdata2)
library(ggplot2)
library(knitr)
```

Code to Produce the Final GBM Model

```
GBM_FULLL <- train(points_per_game ~ points_pg_ly + targets_pg_ly + wopr_pg_ly +
  pick + air_yards_pg_ly + total_games_ly +
  pass_attempt_difference + total_positional_investment +
  target_dropoff + epa_pg_ly + years_pro +
  targets_added_this_year + is_on_new_team +
  points_per_snap_ly + targets_per_snap_ly + wopr_per_snap_ly
  + total_snaps_ly + total_passing_tds_ly +
  total_passing_yds_ly + total_passing_fp_ly +
  is_returning_coach + hc_years_with_team + catch_rate_ly +
  position + fp_dropoff + starter_epa_passing_ly +
  starter_epa_persnap_passing_ly + combine_cluster +
  qbr_ly_bin,
  data=train,
  method="gbm",
  tuneGrid=expand.grid(n.trees = 150,
    interaction.depth = 7,
    n.minobsinnode = 9,
    shrinkage = 0.03),
  metric="RMSE",
  verbose=FALSE,
  trControl=ctrl)
```

Code to Produce XGBOOST Model

```
nrounds <- seq(50, 150, by = 50)
max_depth <- seq(2, 3, by = 1)
eta <- seq(0, 0.1, by = 0.05)
```



```

gamma <- seq(0, 1, by = 0.5)
colsample_bytree <- seq(0.5, 1, by = 0.5)
min_child_weight <- seq(10, 12, by = 1)
subsample <- seq(0.8, 1, by = 0.1)

XGB <- train(points_per_game ~ points_pg_ly + targets_pg_ly + wopr_pg_ly +
  pick + air_yards_pg_ly + total_games_ly +
  pass_attempt_difference + total_positional_investment +
  target_dropoff + epa_pg_ly + years_pro +
  targets_added_this_year + is_on_new_team +
  points_per_snap_ly + targets_per_snap_ly + wopr_per_snap_ly
  + total_snaps_ly + total_passing_tds_ly +
  total_passing_yds_ly + total_passing_fp_ly +
  is_returning_coach + hc_years_with_team + catch_rate_ly +
  position + fp_dropoff + starter_epa_passing_ly +
  starter_epa_persnap_passing_ly + combine_cluster +
  qbr_ly_bin,
  data=train,
  method="xgbTree",
  tuneGrid=expand.grid(nrounds = nrounds,
    max_depth = max_depth,
    eta = eta,
    gamma = gamma,
    colsample_bytree = colsample_bytree,
    min_child_weight = min_child_weight,
    subsample = subsample),
  metric="RMSE",
  verbose=FALSE,
  trControl=ctrl)

```

Code to Produce Original GBM Model for Model Comparison

```

trees <- seq(from=50,to=150,by=50)
interaction_d <- seq(from=5,to=7,by=1)
minobsinnode <- seq(from=8,to=10,by=1)
shrinkage <- seq(from=0.03,to=0.05,by=0.01)

GBM <- train(points_per_game ~ points_pg_ly + targets_pg_ly + wopr_pg_ly +
  pick + air_yards_pg_ly + total_games_ly +
  pass_attempt_difference + total_positional_investment +
  target_dropoff + epa_pg_ly + years_pro +
  targets_added_this_year + is_on_new_team +
  points_per_snap_ly + targets_per_snap_ly + wopr_per_snap_ly
  + total_snaps_ly + total_passing_tds_ly +
  total_passing_yds_ly + total_passing_fp_ly +
  is_returning_coach + hc_years_with_team + catch_rate_ly +
  position + fp_dropoff + starter_epa_passing_ly +
  starter_epa_persnap_passing_ly + combine_cluster +
  qbr_ly_bin,
  data=train,

```

```

method="gbm",
tuneGrid=expand.grid(n.trees = trees,
                      interaction.depth = interaction_d,
                      n.minobsinnode = minobsinnode,
                      shrinkage = shrinkage),
metric="RMSE",
verbose=FALSE,
trControl=ctrl)

```

Code to Product Random Forest Model

```

mtry <- seq(2, 10, by=1)

RF <- train(points_per_game ~ points_pg_ly + targets_pg_ly + wopr_pg_ly +
            pick + air_yards_pg_ly + total_games_ly +
            pass_attempt_difference + total_positional_investment +
            target_dropoff + epa_pg_ly + years_pro +
            targets_added_this_year + is_on_new_team +
            points_per_snap_ly + targets_per_snap_ly + wopr_per_snap_ly
            + total_snaps_ly + total_passing_tds_ly +
            total_passing_yds_ly + total_passing_fp_ly +
            is_returning_coach + hc_years_with_team + catch_rate_ly +
            position + fp_dropoff + starter_epa_passing_ly +
            starter_epa_persnap_passing_ly + combine_cluster +
            qbr_ly_bin,
            data = train,
            method = "rf",
            tuneGrid = expand.grid(mtry = mtry),
            metric = "RMSE",
            verbose = FALSE,
            trControl = ctrl)

```

Code to Produce GLMNet Model

```

alpha <- seq(0,1,by=0.1)
lambda <- seq(-3,0,by=0.5)

GLM <- train(points_per_game ~ points_pg_ly + targets_pg_ly + wopr_pg_ly +
            pick + air_yards_pg_ly + total_games_ly +
            pass_attempt_difference + total_positional_investment +
            target_dropoff + epa_pg_ly + years_pro +
            targets_added_this_year + is_on_new_team +
            points_per_snap_ly + targets_per_snap_ly + wopr_per_snap_ly
            + total_snaps_ly + total_passing_tds_ly +
            total_passing_yds_ly + total_passing_fp_ly +
            is_returning_coach + hc_years_with_team + catch_rate_ly +
            position + fp_dropoff + starter_epa_passing_ly +
            starter_epa_persnap_passing_ly + combine_cluster +
            qbr_ly_bin,

```

```

data=train,
method="glmnet",
tuneGrid=expand.grid(alpha = alpha,
                      lambda = lambda),
metric="RMSE",
verbose=FALSE,
trControl=ctrl)

```

Code to Produce NNET Model

```

size <- seq(1, 5, by = 1)
decay <- c(0.1,0.01,0.001,0.0001)

NNET <- train(points_per_game ~ points_pg_ly + targets_pg_ly + wopr_pg_ly +
              pick + air_yards_pg_ly + total_games_ly +
              pass_attempt_difference + total_positional_investment +
              target_dropoff + epa_pg_ly + years_pro +
              targets_added_this_year + is_on_new_team +
              points_per_snap_ly + targets_per_snap_ly + wopr_per_snap_ly
              + total_snaps_ly + total_passing_tds_ly +
              total_passing_yds_ly + total_passing_fp_ly +
              is_returning_coach + hc_years_with_team + catch_rate_ly +
              position + fp_dropoff + starter_epa_passing_ly +
              starter_epa_persnap_passing_ly + combine_cluster +
              qbr_ly_bin,
              data = train,
              method = "nnet",
              tuneGrid = expand.grid(size = size,
                                    decay = decay),
              metric = "RMSE",
              trace = FALSE,
              trControl = ctrl)

```