📕 **tejseth** / **nfl-r-tutorials**    Public

<> **Code**    ⊙ Issues    ⥮ Pull requests    ▷ Actions    ▥ Projects    📖 Wiki    ⊘ Security    ⭘ Insights    ⚙ Settings

ᛘ **master** ⌄                                                             ···

**nfl-r-tutorials** / **mfans-rf-xgboost.R**

| 👤 **tejseth** random forest and xgboost tutorial done | 🕘 History |

👥 **1 contributor**

203 lines (165 sloc)    6.39 KB                                                        ···

```
 1    # Tutorial 1: Introduction to nflfastR: https://www.youtube.com/watch?v=uw6PH1YiCtI&t=155s
 2    # Tutorial 2: Data Viz With nflfastr: https://www.youtube.com/watch?v=tGaXtAAt6-4&t=255s
 3    # Tutorial 3: Modeling I: https://www.youtube.com/watch?v=J4p8ZfYW5Oo&t=200s
 4    # Tutorial 4: Modeling II (Today): https://github.com/tejseth/nfl-r-tutorials/blob/master/mfans-rf-xgboost.R
 5
 6    # Load packages
 7    library(tidyverse)
 8    library(nflfastR)
 9    library(ggthemes)
10    library(ranger)
11    library(vip)
12    library(caret)
13    library(xgboost)
14    library(ggimage)
15    options(scipen = 9999)
16
17    # Make our custom package
```

```r
18   theme_reach <- function() {
19     theme_fivethirtyeight() +
20       theme(
21         legend.position = "none",
22         plot.title = element_text(size = 22, hjust = 0.5, face = "bold"),
23         plot.subtitle = element_text(size = 18, hjust = 0.5),
24         plot.caption = element_text(size = 16),
25         axis.title.x = element_text(size=18),
26         axis.title.y = element_text(size=18),
27         axis.text = element_text(size = 14),
28         strip.text = element_text(size = 16, face = "bold"),
29         legend.text = element_text(size = 14)
30       )
31   }
32
33   # Load in play-by-play data
34   pbp <- load_pbp(2014:2021)
35
36   # Check what type of plays happen on 4th down
37   pbp %>%
38     filter(down == 4) %>%
39     group_by(play_type) %>%
40     tally(sort = T)
41
42   # Get 4th downs
43   fourth_downs <- pbp %>%
44     filter(down == 4, !play_type %in% c("no_play", "qb_kneel", NA)) %>%
45     mutate(went_for_it = ifelse(play_type %in% c("pass", "run"), 1, 0)) %>%
46     select(posteam, defteam, home_team, season, week, game_id, play_id, desc,
47            play_type, down, yardline_100, ydstogo, half_seconds_remaining, wp,
48            wpa, score_differential, ep, epa, temp, wind, went_for_it) %>%
49     filter(!is.na(epa))
50
51   # Check for NA's
52   colSums(is.na(fourth_downs))
53
54   fourth_downs <- fourth_downs %>%
```

```r
55      mutate(temp = ifelse(is.na(temp), 70, temp),
56             wind = ifelse(is.na(wind), 0, wind))
57
58  # Select the data we want for the model
59  model_data <- fourth_downs %>%
60    select(went_for_it, yardline_100, ydstogo, half_seconds_remaining, wp,
61           score_differential, ep, temp, wind, season)
62
63  # Build our random forest
64  rf_4th <- ranger(went_for_it ~ ., data = model_data,
65                   num.trees = 100, importance = "impurity")
66
67  # Check variable importance of random forest
68  vip(rf_4th) + theme_reach()
69
70  # Make a grid for tuning
71  dim(model_data)
72  rf_grid <- expand.grid(mtry = seq(2, 8, by = 1),
73                 splitrule = "variance") # For classification
74
75  set.seed(2014) # go lions
76
77  # Use the tuning grid
78  rf_4th_tune <-
79    train(went_for_it ~ ., data = model_data,
80          method = "ranger", num.trees = 100,
81          trControl = trainControl(method = "cv", number = 5),
82          tuneGrid = rf_grid)
83
84  # Get the results from the best tune
85  rf_4th_tune$bestTune
86
87  # Remake random forest with tuning parameters
88  rf_4th_best <- ranger(went_for_it ~ ., data = model_data,
89                   num.trees = 100, importance = "impurity",
90                   mtry = 5)
91
```

```r
 92   # Get predictions
 93   rf_preds <- data.frame(predict(rf_4th_best, data.frame(model_data))$predictions)
 94
 95   names(rf_preds)
 96
 97   rf_preds <- rf_preds %>%
 98     rename(exp_go = predict.rf_4th_best..data.frame.model_data...predictions)
 99
100   # Bind the original dataset and predictions together
101   fourth_downs_rf_projs <- cbind(fourth_downs, rf_preds)
102
103   fourth_downs_rf_projs <- fourth_downs_rf_projs %>%
104     mutate(go_over_expected = went_for_it - exp_go)
105
106   # Check 2021 stats
107   rf_team_stats <- fourth_downs_rf_projs %>%
108     filter(season == 2021) %>%
109     group_by(posteam) %>%
110     summarize(avg_gooe = 100*mean(go_over_expected),
111               wpa = 100*mean(wpa)) %>%
112     left_join(teams_colors_logos, by = c("posteam" = "team_abbr"))
113
114   # Make graph
115   rf_team_stats %>%
116     ggplot(aes(x = avg_gooe, y = wpa)) +
117     geom_image(aes(image = team_logo_espn), asp = 16/9, size = 0.05) +
118     theme_reach() +
119     labs(x = "Go For It Rate Over Expected",
120         y = "WPA Added on 4th Downs",
121         title = "Go For It Rate Over Expected and WPA on 4th Downs in 2021",
122         subtitle = "WPA = Win Probability Added",
123         caption = "By Tej Seth | @tejfbanalytics | M-FANS")
124   ggsave('go-oe.png', width = 15, height = 10, dpi = "retina")
125
126   #############################################################################
127
128   # Make xgboost grid
```

```r
129    xgboost_tune_grid <- expand.grid(nrounds = seq(from = 20, to = 200, by = 20),
130                                     eta = c(0.025, 0.05, 0.1, 0.3), gamma = 0,
131                                     max_depth = c(1, 2, 3, 4), colsample_bytree = 1,
132                                     min_child_weight = 1, subsample = 1)
133    xgboost_tune_control <- trainControl(method = "cv", number = 5, verboseIter = FALSE)
134    set.seed(2011) # go lions
135
136    # Tune xgboost grid
137    xgb_tune <- train(x = as.matrix(dplyr::select(model_data, -went_for_it)),
138                      y = model_data$went_for_it, trControl = xgboost_tune_control,
139                      tuneGrid = xgboost_tune_grid,
140                      objective = "reg:squarederror", method = "xgbTree",
141                      verbose = TRUE)
142
143    # Get the best tune
144    xgb_tune$bestTune
145
146    # Set xgboost parameters
147    nrounds <- 100
148    params <-
149      list(
150        booster = "gbtree",
151        objective = "binary:logistic",
152        eval_metric = c("logloss"),
153        eta = 0.025,
154        gamma = 5,
155        subsample = 0.8,
156        colsample_bytree = 0.8,
157        max_depth = 4,
158        min_child_weight = 6,
159        base_score = mean(model_data$went_for_it)
160      )
161
162    seasons <- seq(2014, 2021, 1)
163
164    # Make the holdout model
165    cv_results <- map_dfr(seasons, function(x) {
```

```r
166        test_data <- model_data %>%
167          filter(season == x) %>%
168          select(-season)
169        train_data <- model_data %>%
170          filter(season != x) %>%
171          select(-season)
172
173        full_train <- xgboost::xgb.DMatrix(model.matrix(~ . + 0, data = train_data %>% select(-went_for_it)),
174                                           label = train_data$went_for_it
175        )
176        xg_4th <- xgboost::xgboost(params = params, data = full_train, nrounds = nrounds, verbose = 2)
177
178        preds <- as.data.frame(
179          matrix(predict(xg_4th, as.matrix(test_data %>% select(-went_for_it))))
180        ) %>%
181          dplyr::rename(exp_go = V1)
182
183        cv_data <- bind_cols(test_data, preds) %>% mutate(season = x)
184        return(cv_data)
185  })
186
187  # Get the predictions
188  xg_preds <- cv_results %>% select(exp_go)
189
190  # Put it all together
191  fourth_downs_xg_projs <- cbind(fourth_downs, xg_preds)
192
193  fourth_downs_xg_projs <- fourth_downs_xg_projs %>%
194    mutate(go_oe = went_for_it - exp_go)
195
196  # Check the team stats from xgboost
197  xg_team_stats <- fourth_downs_xg_projs %>%
198    filter(season == 2021) %>%
199    group_by(posteam) %>%
200    summarize(avg_gooe = 100*mean(go_oe),
201              wpa = 100*mean(wpa))
202
```

203