

SEAHAWKS SHOULD HAVE RAN THE BALL

Artificial American Football Coach with Decision Trees

Benjamin Jones

February 28, 2020

Introduction

Have you ever dreamed of being an NFL coach? No? Me neither... Maybe you have a child on a pee wee football team in need of a coach? Also no? Same... But perhaps you like data? If *you're* reading this, chances are you answered 'yes' to at least one of those questions.

The following dataset consists of football situations and the correct play to call for each situation. Let's see if it is possible to create a decision tree which, given a situation, can output the correct play to call. We'll try a few different tree models. The dataset is titled *Football Strategy* and publicly available [here](#).

Setup

Load the necessary libraries as well as authenticate the [data.world](#) API:

```
library(data.world)
library(jsonlite)
library(stringr)
library(rpart)
library(rpart.plot)
library(ipred)
library(randomForest)
library(rpart)
library(caret)
setwd("~/Desktop")
dw_auth <- read_json("dw_cred.json")
saved_cfg <- save_config(dw_auth$token)
set_config(saved_cfg)
```

Query the data:

```
df <- query(qry_sql(paste0("SELECT football_scenarios_dfe_832307.antecedent as LABEL,
                             football_scenarios_dfe_832307.orig_antecedent as SITUATION ",
                             "FROM football_scenarios_dfe_832307))),
           dataset="https://data.world/wppdatacatalogue/football-strategy")
```

Preview the data:

```
head(df, 3)
```

```
## # A tibble: 3 x 2
##   LABEL          SITUATION
##   <chr>         <chr>
## 1 kick a field g~ It is first down and 10. The ball is on your opponent's 20 ya~
## 2 kick a field g~ It is second down and inches. The ball is on your opponent's ~
## 3 kick a field g~ It is second down and inches. The ball is on your opponent's ~
```

It's difficult to see, but an example of a full value of an entry in **SITUATION** is "It is first down and 10. The ball is on your opponent's 20 yardline. There is 3 seconds left in the second quarter. You are down by 3 points. Would you:"

Cleaning

The entries of the football situation is a messy, as it is a long string. Fortunately, the values are consistent and we can do use patterns to pull out useful quantitative and qualitative features. For example:

"It is **first down** and **10**. The ball is on **your opponent's 20** yardline. There is **3** seconds left in the **second quarter**. You are **down by 3** points. Would you:"

"It is **third down** and **inches**. The ball is on **your 15** yardline. There is **3** seconds left in the **fourth quarter**. You are **up by 7** points. Would you:"

Numerical values used will be: yards to go, yard line, seconds left, and point differential. When distance is 'inches', it will be converted to 0.

Categorical values used will be: down and quarter.

```
### Remove Cases that Don't Follow the Pattern
df <- df[grepl("It is", df$SITUATION),]
df <- df[grepl("The ball is", df$SITUATION),]
df <- df[grepl("There is", df$SITUATION),]
```

Down and Yards to Go

```
df$DOWN <- ""
df$YTG <- ""

for (i in 1:nrow(df)) {
  ss1 <- str_split_fixed(df[i,2], " down ", 2)[1,1]
  df$DOWN[i] <- str_split_fixed(ss1, "It is ", 2)[1,2]

  ss2 <- str_split_fixed(df[i,2], " down ", 2)[1,2]
  ss3 <- str_split_fixed(ss2, "and ", 2)
  df$YTG[i] <- str_split_fixed(ss3, "\\.", 2)[2,1]
}

df$YTG <- str_replace(df$YTG, "inches", "0")
df$YTG <- as.numeric(df$YTG)
```

Yardline

```
df$YL <- 0
for (i in 1:nrow(df)) {
  ss1 <- str_split_fixed(df[i,2], "The ball is on ", 2)[1,2]
  ss2 <- str_split_fixed(ss1, " yardline.", 2)[1,1]
  df$YL[i] <- as.numeric(gsub("[^0-9.-]", "", ss2))
  if (grepl("your opponent", ss1)) {
    df$YL[i] <- 50 + (50-df$YL[i])
  } else {
    df$YL[i] <- df$YL[i]
  }
}
```

Seconds and Quarter

```
df$SEC <- 0
df$Q <- ""
for (i in 1:nrow(df)) {
  ss1 <- str_split_fixed(df[i,2], "There is ", 2)[1,2]
  if (grepl("minute", ss1)) {
    df$SEC[i] <- as.numeric(str_split_fixed(ss1, " minute", 2)[1,1])*60
  } else{
    df$SEC[i] <- as.numeric(str_split_fixed(ss1, " second", 2)[1,1])
  }
  ss2 <- str_split_fixed(ss1, " quarter", 2)[1,1]
  df$Q[i] <- str_split_fixed(ss2, " in the ", 2)[1,2]
}
```

Score

```
df$PTS <-0
for (i in 1:nrow(df)) {
  ss1 <- str_split_fixed(df[i,2], "You are ", 2)[1,2]
  ss2 <- str_split_fixed(ss1, "by ", 2)[1,2]
  df$PTS[i] <- as.numeric(str_split_fixed(ss2, " points", 2)[1,1])
  if (grepl("down", ss1)) {
    df$PTS[i] <- df$PTS[i]*-1
  }
}
```

Final Dataset

Due to the fact there were only 9 cases of “Don’t Know/it depends”, and that it’s not helpful, remove those cases.

```
df$SITUATION <- NULL
df <- subset(df, df$LABEL != "Don't know / it depends")
df$LABEL <- as.factor(df$LABEL)
df$DOWN <- as.factor(df$DOWN)
df$Q <- as.factor(df$Q)
head(df, 5)
```

```
## # A tibble: 5 x 7
##   LABEL          DOWN    YTG    YL    SEC Q      PTS
##   <fct>         <fct> <dbl> <dbl> <dbl> <fct> <dbl>
## 1 kick a field goal first     10    80     3 second    -3
## 2 kick a field goal second     0    95     3 second    -3
## 3 kick a field goal second     0    80     3 second    -3
## 4 kick a field goal second     0    95     3 fourth    -3
## 5 kick a field goal second     0    80     3 fourth    -3
```

Split into Training/Validation

While not shown here, distribution of labels between training and validation sets was assumed to be similar and was verified.

```
set.seed(74)
split_size <- floor(0.75*nrow(df))
split_index <- sample(seq_len(nrow(df)), size = split_size)

train <- df[split_index, ]
val <- df[-split_index, ]
```

Decision Trees

Bootstrap Aggregation aka *Belicheck Yourself*

```
tree1 <- bagging(LABEL~., data=train)
results1 <- table(predict(tree1, val[,2:7]), val$LABEL)
confusionMatrix(results1)
```

```
## Confusion Matrix and Statistics
```

```
##
```

```
##
```

```
##               kick a field goal kneel down pass punt run
##  kick a field goal           112         1    8    6    6
##  kneel down                   0         9    1    1    0
##  pass                         9         3  219    0  23
##  punt                         2         2    8   52   1
##  run                          3         2   28    4 151
```

```
##
```

```
## Overall Statistics
```

```
##
```

```
##               Accuracy : 0.8341
```

```
##               95% CI : (0.8033, 0.8619)
```

```
##      No Information Rate : 0.4055
```

```
##      P-Value [Acc > NIR] : < 2e-16
```

```
##
```

```
##               Kappa : 0.7669
```

```
##
```

```
##  McNemar's Test P-Value : 0.06056
```

```
##
```

```
## Statistics by Class:
```

```
##
```

```
##               Class: kick a field goal Class: kneel down Class: pass
## Sensitivity                0.8889                0.52941        0.8295
## Specificity                0.9600                0.99685        0.9096
## Pos Pred Value              0.8421                0.81818        0.8622
## Neg Pred Value              0.9730                0.98750        0.8866
## Prevalence                  0.1935                0.02611        0.4055
## Detection Rate              0.1720                0.01382        0.3364
## Detection Prevalence        0.2043                0.01690        0.3902
## Balanced Accuracy           0.9244                0.76313        0.8696
```

```
##               Class: punt Class: run
```

```
## Sensitivity                0.82540        0.8343
## Specificity                0.97789        0.9213
## Pos Pred Value              0.80000        0.8032
## Neg Pred Value              0.98123        0.9352
## Prevalence                  0.09677        0.2780
## Detection Rate              0.07988        0.2320
## Detection Prevalence        0.09985        0.2888
## Balanced Accuracy           0.90164        0.8778
```

```
#varImp(tree1)
```

Random Forest aka *Andy Feed*

```
tree2 <- randomForest(LABEL~., data=train)
results2 <- table(predict(tree2, val[,2:7]), val$LABEL)
confusionMatrix(results2)
```

```
## Confusion Matrix and Statistics
```

```
##
```

```
##
```

```
##               kick a field goal kneel down pass punt run
## kick a field goal           118           2   10   10   7
## kneel down                   0           8    0    0    0
## pass                         6           4  232    0  25
## punt                        2           2   8   51   1
## run                         0           1  14    2 148
```

```
##
```

```
## Overall Statistics
```

```
##
```

```
##               Accuracy : 0.8556
```

```
##               95% CI : (0.8262, 0.8817)
```

```
##       No Information Rate : 0.4055
```

```
##       P-Value [Acc > NIR] : < 2.2e-16
```

```
##
```

```
##               Kappa : 0.7965
```

```
##
```

```
## Mcnemar's Test P-Value : 0.0002021
```

```
##
```

```
## Statistics by Class:
```

```
##
```

```
##               Class: kick a field goal Class: kneel down Class: pass
## Sensitivity                0.9365                0.47059        0.8788
## Specificity                0.9448                1.00000        0.9096
## Pos Pred Value             0.8027                1.00000        0.8689
## Neg Pred Value             0.9841                0.98600        0.9167
## Prevalence                 0.1935                0.02611        0.4055
## Detection Rate             0.1813                0.01229        0.3564
## Detection Prevalence       0.2258                0.01229        0.4101
## Balanced Accuracy          0.9406                0.73529        0.8942
```

```
##               Class: punt Class: run
```

```
## Sensitivity                0.80952                0.8177
## Specificity                0.97789                0.9638
## Pos Pred Value             0.79688                0.8970
## Neg Pred Value             0.97956                0.9321
## Prevalence                 0.09677                0.2780
## Detection Rate             0.07834                0.2273
## Detection Prevalence       0.09831                0.2535
## Balanced Accuracy          0.89371                0.8908
```

```
#varImp(tree2)
```

Classification and Regression aka *Mike Shanahan-me-down*

```
tree3 <- rpart(LABEL~., data=train)
results3 <- table(predict(tree3, val[,2:7], type="class"), val$LABEL)
confusionMatrix(results3)
```

```
## Confusion Matrix and Statistics
```

```
##
```

```
##
```

```
##               kick a field goal kneel down pass punt run
## kick a field goal           113           3   20   13  18
## kneel down                   0           0   0    0   0
## pass                         13          11  183   0  25
## punt                         0           3   25   50   2
## run                          0           0   36   0 136
```

```
##
```

```
## Overall Statistics
```

```
##
```

```
##               Accuracy : 0.7404
```

```
##               95% CI : (0.7049, 0.7737)
```

```
##       No Information Rate : 0.4055
```

```
##       P-Value [Acc > NIR] : < 2.2e-16
```

```
##
```

```
##               Kappa : 0.6397
```

```
##
```

```
## McNemar's Test P-Value : NA
```

```
##
```

```
## Statistics by Class:
```

```
##
```

```
##               Class: kick a field goal Class: kneel down Class: pass
## Sensitivity                0.8968                0.00000    0.6932
## Specificity                0.8971                1.00000    0.8734
## Pos Pred Value              0.6766                  NaN    0.7888
## Neg Pred Value              0.9731                0.97389    0.8067
## Prevalence                  0.1935                0.02611    0.4055
## Detection Rate              0.1736                0.00000    0.2811
## Detection Prevalence        0.2565                0.00000    0.3564
## Balanced Accuracy           0.8970                0.50000    0.7833
```

```
##               Class: punt Class: run
```

```
## Sensitivity                0.79365    0.7514
## Specificity                0.94898    0.9234
## Pos Pred Value              0.62500    0.7907
## Neg Pred Value              0.97723    0.9061
## Prevalence                  0.09677    0.2780
## Detection Rate              0.07680    0.2089
## Detection Prevalence        0.12289    0.2642
## Balanced Accuracy           0.87132    0.8374
```

```
#varImp(tree3)
```

Conclusion

Feature Importance can also be evaluated. Similar to *Guess Who?*, feature importance describes the most efficient path of questioning to get to an answer. In this context, importance reveals which factors of the current play/game will lead to a play calling decision in the least amount of steps. In this experiment, the features ranked in descending order of importance are:

1. What down it is
2. Yards to Go
3. Yardline (position on field)
4. Score differential
5. Seconds remaining in the quarter
6. Which quarter it is

What I'm trying to say is... wait let me verify...

```
i <- nrow(val)+1
val[i, 2] <- as.factor('second') #second down
val[i, 3] <- as.numeric(1) #1 yard to go
val[i, 4] <- as.numeric(99) #99 yard line
val[i, 5] <- as.numeric(26) #26s on the clock
val[i, 6] <- as.factor('fourth') #fourth quater
val[i, 7] <- as.numeric(-4) #down by 4
```

```
predict(tree1, val[i,2:7])
```

```
## [1] kick a field goal
## Levels: kick a field goal
```

```
predict(tree2, val[i,2:7])
```

```
##                1
## kick a field goal
## Levels: kick a field goal kneel down pass punt run
```

```
predict(tree3, val[i,2:7], type="class")
```

```
##                1
## kick a field goal
## Levels: kick a field goal kneel down pass punt run
```

The Seahawks should have *kicked* the ball? Well nobody's perfect...