

swstr_modeling

Owen McGrattan

10/21/2018

```
# load in packages
```

```
library(dplyr)
```

```
##
```

```
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
```

```
##
```

```
## filter, lag
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
## intersect, setdiff, setequal, union
```

```
library(ggplot2)
```

```
library(baseballr)
```

```
library(readr)
```

```
library(randomForest)
```

```
## randomForest 4.6-14
```

```
## Type rfNews() to see new features/changes/bug fixes.
```

```
##
```

```
## Attaching package: 'randomForest'
```

```
## The following object is masked from 'package:ggplot2':
```

```
##
```

```
## margin
```

```
## The following object is masked from 'package:dplyr':
```

```
##
```

```
## combine
```

```
library(stats)
```

```
library(rfUtilities)
```

```
library(standardize)
```

```
# read in grouped df
```

```
df <- read_csv("swstr_pitchers.csv",  
  col_types = cols(X1 = col_skip()))
```

```
## Warning: Missing column names filled in: 'X1' [1]
```

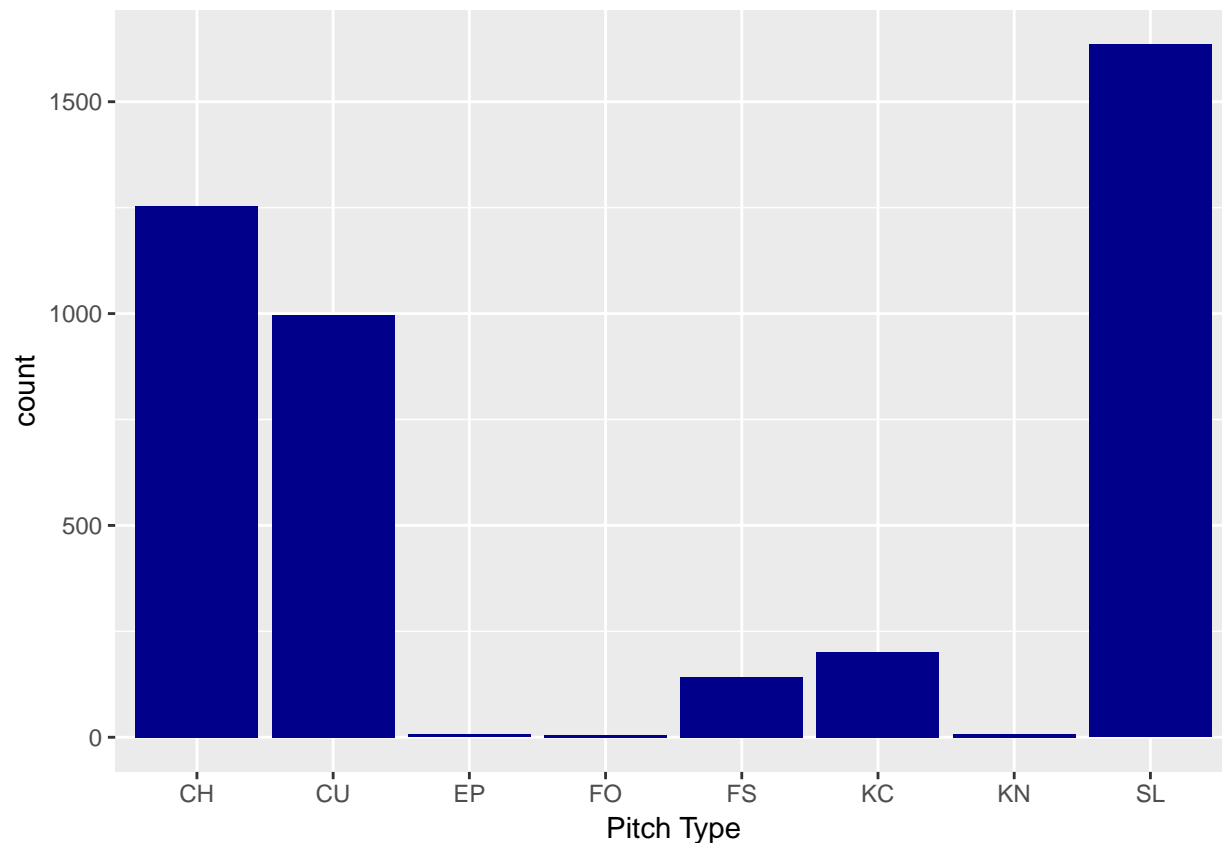
```
## Warning in read_tokens_(data, tokenizer, col_specs, col_names, locale_, :  
## length of NULL cannot be changed
```

```
## Warning in read_tokens_(data, tokenizer, col_specs, col_names, locale_, :  
## length of NULL cannot be changed
```

```
## Warning in read_tokens_(data, tokenizer, col_specs, col_names, locale_, :  
## length of NULL cannot be changed
```

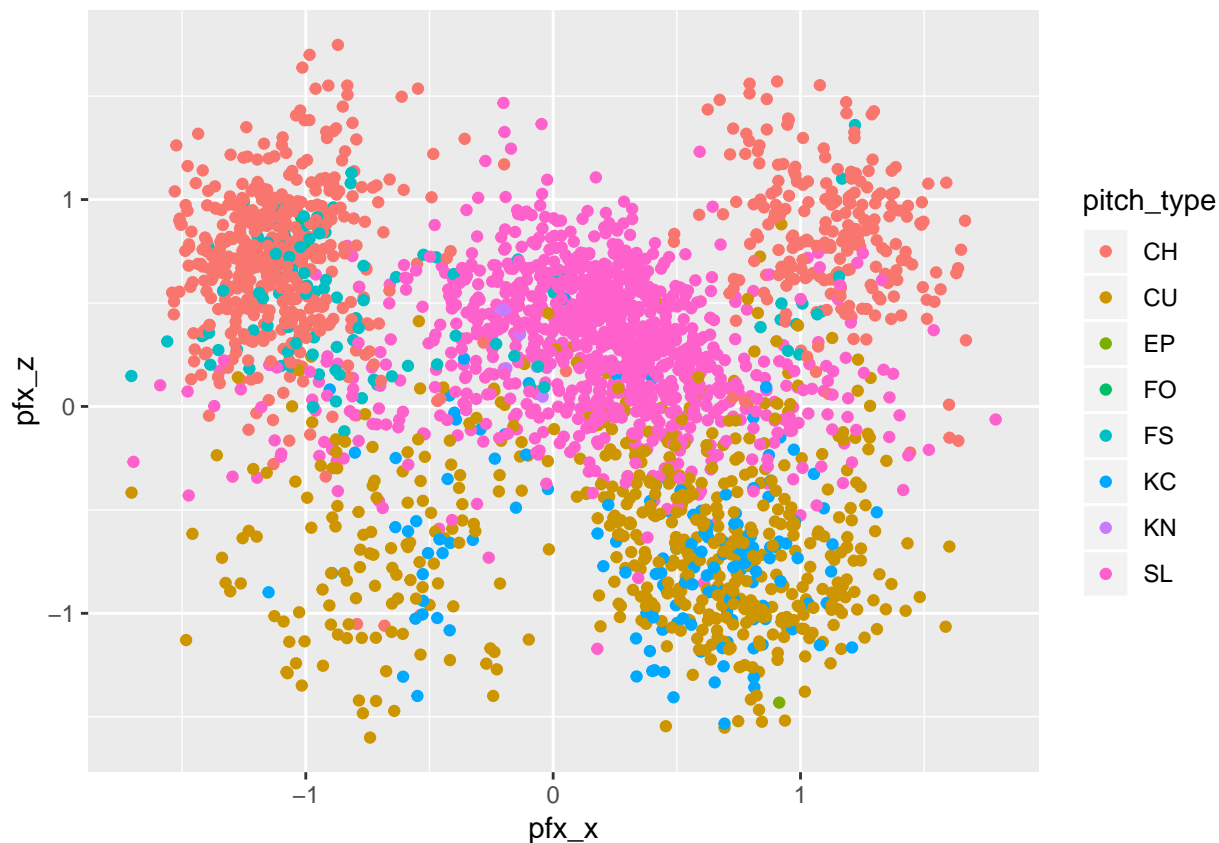
```
# rename our fastball velo column appropriately
colnames(df)[which(names(df) == "fb_diff")] <- "fb_velo"

ggplot(data.frame(df$pitch_type), aes(df$pitch_type)) + geom_bar(fill = "blue4") + labs(x = "Pitch Type", y = "count")
```



- If we're fitting a model it may just be best to leave the forkballs, knuckles, and ephuses out.
- How do these pitches look on a movement spectrum?

```
ggplot(filter(df, len > 150)) + geom_point(aes(x = pfx_x, y = pfx_z, color = pitch_type))
```



include closest comp pitcher, weighted for distance from pitcher

```
pitch_comp <- function(player, pitch_tpy, year) {
```

```
  dat <- filter(df, pitch_type == pitch_tpy, game_year == year)
```

```
  scaled <- select(dat, velo, release_pos_x, release_pos_z, pfx_x, pfx_z, pct)
```

```
  scaled <- as.data.frame(scale(scaled))
```

```
  scaled$player_name <- dat$player_name
```

```
  pitch <- filter(scaled, player_name == player)
```

```
  dat$dist <- sqrt((pitch$velo - scaled$velo)^2 +
    (pitch$release_pos_x - scaled$release_pos_x)^2 +
    (pitch$release_pos_z - scaled$release_pos_z)^2 +
    (pitch$pfx_x - scaled$pfx_x)^2 +
    (pitch$pfx_z - scaled$pfx_z)^2 +
    (pitch$pct - scaled$pct)^2)
```

```
  dat <- dat[order(dat$dist),]
```

```
  return(mean(head(dat$swstr, n = 6)[2:6]))
```

```
}
```

```
df$near_comp <- mapply(pitch_comp, df$player_name, df$pitch_type, df$game_year)
```

```
df$appch_diff <- (df$fb_appch - df$appch_angle)^2
```

```
df$dist_from_fb <- sqrt((df$plate_x - df$fb_plate_x)^2 + (df$plate_z - df$fb_plate_z)^2)
```

```
df$velo_diff <- df$velo - df$fb_velo
```

shrink our sample down to those who have thrown at least 150 of an offspeed pitch, none of the more o

```
samp <- df %>% filter(len >= 150, pitch_type != "KN", pitch_type != "FO", pitch_type != "EP")
```

```

samp$pitch_type <- as.factor(samp$pitch_type)
samp$handedness <- as.factor(samp$handedness)

set.seed(121)
# create test and train data
samp <- samp %>% mutate(id = row_number())

train <- na.omit(samp %>% sample_frac(.80))
test <- na.omit(anti_join(samp, train, by = "id"))
train <- select(train, -c(id, handedness, velo, fb_appch, fb_platex, fb_platez, near_comp, appch_angle))
test <- select(test, -c(id, handedness, velo, fb_appch, fb_platex, fb_platez, near_comp, appch_angle))

# random forest
set.seed(123)
Swstr_predict <- randomForest(formula = swstr ~ .,
                              data = train[, -c(1, 3)], na.action = na.omit, ntree = 500, replace = T, corr.bias = T,
                              print(Swstr_predict))

##
## Call:
## randomForest(formula = swstr ~ ., data = train[, -c(1, 3)], ntree = 500,      replace = T, corr.bias
##              Type of random forest: regression
##              Number of trees: 500
## No. of variables tried at each split: 6
##
##              Mean of squared residuals: 0.001377431
##              % Var explained: 49.58
##              Test set MSE: 0
##              % Var explained: 44.44
## Bias correction applied:
## Intercept: -0.0002227149
## Slope: 1.160082

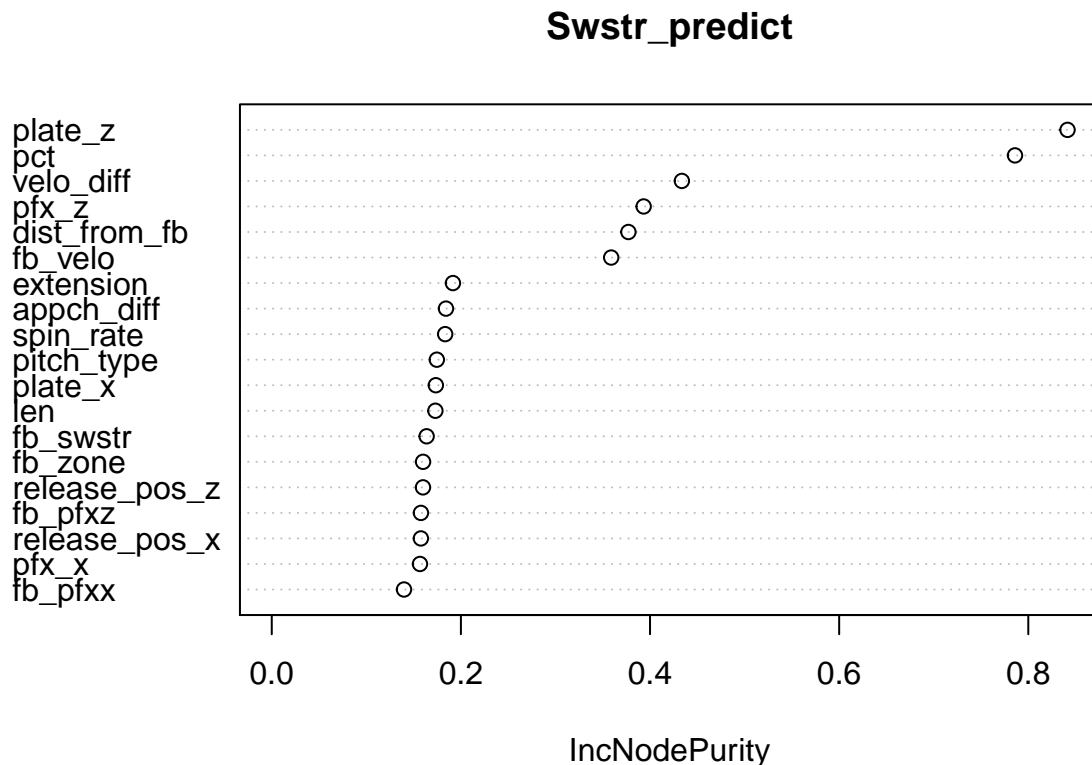
# just a linear model for reference
a <- lm(formula = swstr ~ .,
        data = train[, -c(1, 3)], na.action = na.omit)
summary(a)

##
## Call:
## lm(formula = swstr ~ ., data = train[, -c(1, 3)], na.action = na.omit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.12330 -0.02720 -0.00280  0.02422  0.16739
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -2.383e-01  4.246e-02  -5.613  2.27e-08 ***
## pitch_typeCU -2.562e-02  4.996e-03  -5.128  3.21e-07 ***
## pitch_typeFS  6.816e-03  4.875e-03   1.398  0.16216
## pitch_typeKC -2.678e-02  5.738e-03  -4.667  3.26e-06 ***
## pitch_typeSL -8.860e-03  3.157e-03  -2.807  0.00505 **
## pfx_x        -1.371e-03  1.718e-03  -0.798  0.42519
## pfx_z         1.533e-02  2.898e-03   5.289  1.36e-07 ***

```

```
## plate_x      5.225e-03  3.397e-03  1.538  0.12420
## plate_z     -8.438e-02  5.241e-03 -16.099 < 2e-16 ***
## release_pos_x 1.026e-03  1.104e-03  0.929  0.35290
## release_pos_z -4.270e-04  2.313e-03 -0.185  0.85358
## spin_rate     4.651e-06  3.702e-06  1.256  0.20916
## extension     6.457e-03  2.209e-03  2.923  0.00351 **
## len          1.300e-05  4.689e-06  2.772  0.00562 **
## pct          1.538e-01  9.587e-03 16.037 < 2e-16 ***
## fb_velo       4.736e-03  4.254e-04 11.134 < 2e-16 ***
## fb_pfx       5.362e-04  2.603e-03  0.206  0.83685
## fb_pfxz      -1.361e-02  3.442e-03 -3.954 7.95e-05 ***
## fb_swstr     -4.043e-02  1.431e-02 -2.825  0.00477 **
## fb_zone       1.403e-02  2.177e-02  0.644  0.51943
## appch_diff   -1.397e-03  3.525e-04 -3.963 7.65e-05 ***
## dist_from_fb  3.502e-02  6.247e-03  5.606 2.35e-08 ***
## velo_diff    -2.418e-03  4.187e-04 -5.777 8.82e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.03933 on 1997 degrees of freedom
## Multiple R-squared:  0.4403, Adjusted R-squared:  0.4341
## F-statistic: 71.39 on 22 and 1997 DF,  p-value: < 2.2e-16
```

```
varImpPlot(Swstr_predict)
```



- Taking a general overview of our model, an R^2 of ~ 0.44 tells us a strong deal given just pure “stuff” and not taking the context of each individual pitch. Of course there’s a good deal of variation that we should not expect to be explained in the data I fed into the model. While I do give pct usage and total num of pitches thrown, I do not give any information about usage by count, sequencing, batter handedness, leverage, so on and so forth.

- We've put together a decent model but it has a finite application. Given that $swstr\%$ is something that'll stabilize fairly quickly it doesn't necessarily help identify current MLB candidates that will see a jump in performance, we can simply look at MLB pitchers who generate a strong amount of whiffs on one or more misused offspeed pitches. This prediction is something that better serves identifying the minor league talent (in an ideal scenario where we have movement data on these pitches) who project to produce above average whiffs and therefore strikeouts. It'd greatly benefit our model to include $swstr\%$'s for given players in our data set who have pitched in either the AA or AAA level. Unfortunately, while we can grab $swstr\%$ for pitchers, we do not have publically available pitch-by-pitch data with each pitch type tracked.