

data621_hw1_mia_wei

Wei Zhou, Mia Chen

2/29/2020

R Markdown

```
## -- Attaching packages ----- tidyverse 1.3.0 --
```

```
## v ggplot2 3.2.1    v purrr   0.3.3
## v tibble  2.1.3    v dplyr   0.8.3
## v tidyr   1.0.0    v stringr 1.4.0
## v readr   1.3.1    v forcats 0.4.0
```

```
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
```

```
# Read in the data
```

```
data <- read.csv("https://raw.githubusercontent.com/miachan410/DATA621/master/moneyball-training-data.csv")
glimpse(data)
```

```
## Observations: 2,276
```

```
## Variables: 17
```

```
## $ INDEX          <int> 1, 2, 3, 4, 5, 6, 7, 8, 11, 12, 13, 15, 16, 1...
## $ TARGET_WINS     <int> 39, 70, 86, 70, 82, 75, 80, 85, 86, 76, 78, 6...
## $ TEAM_BATTING_H   <int> 1445, 1339, 1377, 1387, 1297, 1279, 1244, 127...
## $ TEAM_BATTING_2B  <int> 194, 219, 232, 209, 186, 200, 179, 171, 197, ...
## $ TEAM_BATTING_3B  <int> 39, 22, 35, 38, 27, 36, 54, 37, 40, 18, 27, 3...
## $ TEAM_BATTING_HR  <int> 13, 190, 137, 96, 102, 92, 122, 115, 114, 96,...
## $ TEAM_BATTING_BB  <int> 143, 685, 602, 451, 472, 443, 525, 456, 447, ...
## $ TEAM_BATTING_SO  <int> 842, 1075, 917, 922, 920, 973, 1062, 1027, 92...
## $ TEAM_BASERUN_SB  <int> NA, 37, 46, 43, 49, 107, 80, 40, 69, 72, 60, ...
## $ TEAM_BASERUN_CS  <int> NA, 28, 27, 30, 39, 59, 54, 36, 27, 34, 39, 7...
## $ TEAM_BATTING_HBP <int> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, N...
## $ TEAM_PITCHING_H  <int> 9364, 1347, 1377, 1396, 1297, 1279, 1244, 128...
## $ TEAM_PITCHING_HR <int> 84, 191, 137, 97, 102, 92, 122, 116, 114, 96,...
## $ TEAM_PITCHING_BB <int> 927, 689, 602, 454, 472, 443, 525, 459, 447, ...
## $ TEAM_PITCHING_SO <int> 5456, 1082, 917, 928, 920, 973, 1062, 1033, 9...
## $ TEAM_FIELDING_E  <int> 1011, 193, 175, 164, 138, 123, 136, 112, 127,...
## $ TEAM_FIELDING_DP <int> NA, 155, 153, 156, 168, 149, 186, 136, 169, 1...
```

```
# Split data into training set and testing set by 7:3 ratio
```

```
set.seed(123)
```

```
train_ind = sample(seq_len(nrow(data)), size = nrow(data)*.7)
```

```
train = data[train_ind, ]
```

```
test = data[-train_ind, ]
```

```
glimpse(train)
```

```
## Observations: 1,593
## Variables: 17
## $ INDEX          <int> 2479, 582, 215, 2068, 1274, 1395, 1410, 1146,...
## $ TARGET_WINS    <int> 55, 76, 62, 60, 104, 71, 103, 72, 75, 54, 82,...
## $ TEAM_BATTING_H  <int> 1616, 1388, 1525, 1418, 1489, 1261, 1692, 154...
## $ TEAM_BATTING_2B <int> 268, 267, 277, 250, 244, 156, 269, 289, 204, ...
## $ TEAM_BATTING_3B <int> 145, 19, 62, 52, 47, 49, 76, 36, 88, 53, 35, ...
## $ TEAM_BATTING_HR <int> 72, 175, 58, 47, 213, 10, 60, 249, 16, 93, 13...
## $ TEAM_BATTING_BB <int> 757, 523, 469, 458, 714, 482, 572, 673, 379, ...
## $ TEAM_BATTING_SO <int> 460, 1090, 445, 773, 760, NA, 458, 1129, 539,...
## $ TEAM_BASERUN_SB <int> 306, 108, 96, 122, 84, 206, 144, 114, 164, 69...
## $ TEAM_BASERUN_CS <int> NA, 50, 70, 91, 59, NA, 94, 52, 149, 84, 48, ...
## $ TEAM_BATTING_HBP <int> NA, NA, NA, NA, NA, NA, NA, 83, NA, NA, NA, N...
## $ TEAM_PITCHING_H <int> 1983, 1388, 1615, 2088, 1577, 1353, 1815, 154...
## $ TEAM_PITCHING_HR <int> 88, 175, 61, 69, 226, 11, 64, 249, 17, 98, 13...
## $ TEAM_PITCHING_BB <int> 929, 523, 497, 675, 756, 517, 614, 673, 399, ...
## $ TEAM_PITCHING_SO <int> 565, 1090, 471, 1138, 805, NA, 491, 1129, 567...
## $ TEAM_FIELDING_E <int> 612, 104, 263, 150, 141, 283, 185, 133, 236, ...
## $ TEAM_FIELDING_DP <int> NA, 147, 176, 172, 165, 106, 178, 149, 156, 1...
```

1. DATA EXPLORATION (25 Points)

Describe the size and the variables in the moneyball training data set. Consider that too much detail will cause a manager to lose interest while too little detail will make the manager consider that you aren't doing your job. Some suggestions are given below. Please do NOT treat this as a check list of things to do to complete the assignment. You should have your own thoughts on what to tell the boss. These are just ideas.

```
# Cleaning the column names by removing TEAMS_
names(train) <- gsub("TEAM_", "", names(train))
names(test) <- gsub("TEAM_", "", names(test))
summary(train)
```

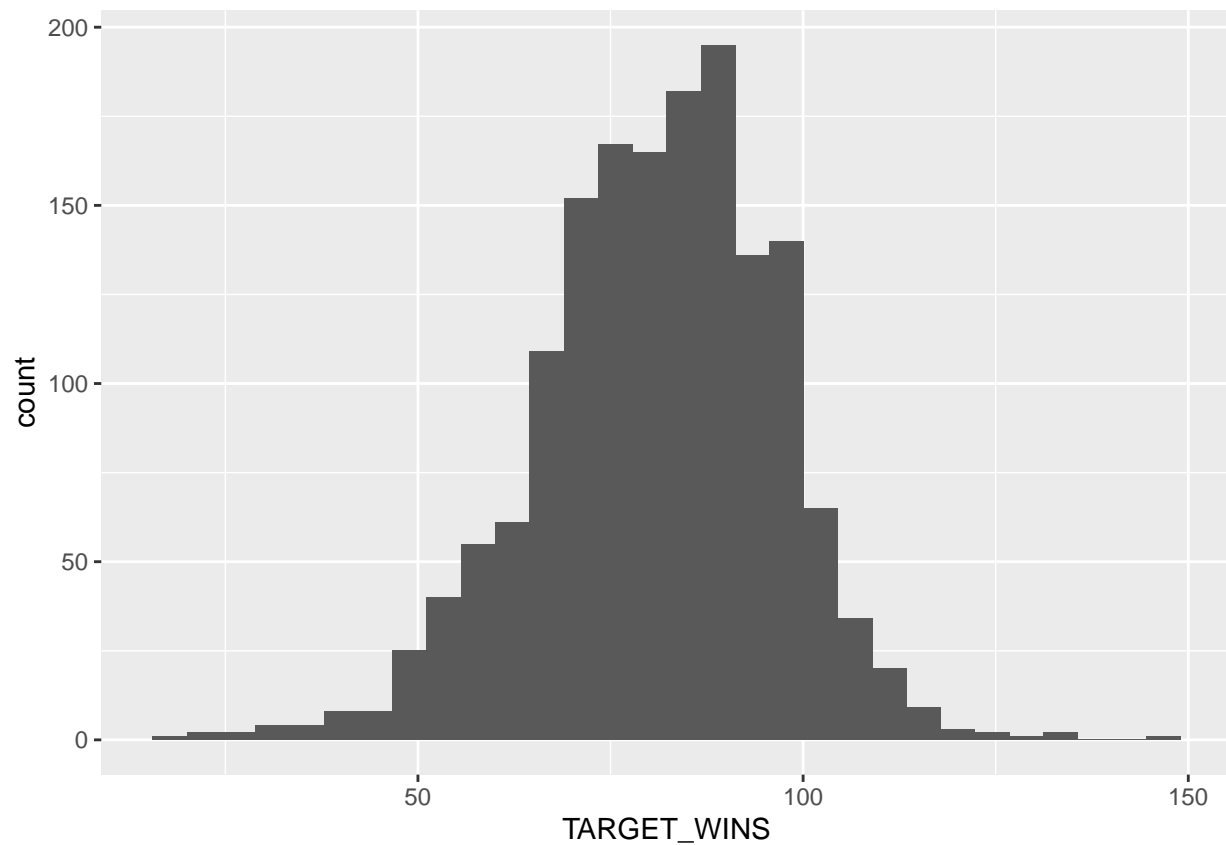
```
##      INDEX      TARGET_WINS      BATTING_H      BATTING_2B
## Min.   : 1      Min.   : 17.00      Min.   : 992      Min.   : 69.0
## 1st Qu.: 617      1st Qu.: 71.00      1st Qu.:1383      1st Qu.:209.0
## Median :1238      Median : 82.00      Median :1456      Median :239.0
## Mean   :1251      Mean   : 81.02      Mean   :1469      Mean   :242.1
## 3rd Qu.:1902      3rd Qu.: 92.00      3rd Qu.:1540      3rd Qu.:274.0
## Max.   :2535      Max.   :146.00      Max.   :2554      Max.   :458.0
##
##      BATTING_3B      BATTING_HR      BATTING_BB      BATTING_SO
## Min.   : 0.00      Min.   : 0.00      Min.   : 29.0      Min.   : 0.0
## 1st Qu.: 34.00      1st Qu.: 41.00      1st Qu.:452.0      1st Qu.: 546.8
## Median : 48.00      Median :101.00      Median :514.0      Median : 750.0
## Mean   : 55.08      Mean   : 99.28      Mean   :503.5      Mean   : 735.1
## 3rd Qu.: 72.00      3rd Qu.:147.00      3rd Qu.:581.0      3rd Qu.: 936.0
## Max.   :197.00      Max.   :264.00      Max.   :860.0      Max.   :1399.0
##
##                                     NA's   :81
##      BASERUN_SB      BASERUN_CS      BATTING_HBP      PITCHING_H
## Min.   : 0.0      Min.   : 12.00      Min.   :35.0      Min.   : 1168
```

```
## 1st Qu.: 67.0    1st Qu.: 39.00    1st Qu.:51.0    1st Qu.: 1419
## Median :102.0    Median : 50.00    Median :59.0    Median : 1520
## Mean   :125.6    Mean   : 53.43    Mean   :60.3    Mean   : 1757
## 3rd Qu.:159.0    3rd Qu.: 63.00    3rd Qu.:69.0    3rd Qu.: 1682
## Max.   :654.0    Max.   :200.00    Max.   :95.0    Max.   :30132
## NA's   :89      NA's   :539      NA's   :1464
## PITCHING_HR    PITCHING_BB    PITCHING_SO    FIELDING_E
## Min.    : 0    Min.    : 119.0    Min.    : 0.0    Min.    : 65.0
## 1st Qu.: 49    1st Qu.: 478.0    1st Qu.: 606.0    1st Qu.: 127.0
## Median :106    Median : 537.0    Median : 807.5    Median : 160.0
## Mean   :105    Mean   : 554.7    Mean   : 821.9    Mean   : 243.3
## 3rd Qu.:150    3rd Qu.: 613.0    3rd Qu.: 970.0    3rd Qu.: 245.0
## Max.   :343    Max.   :3645.0    Max.   :19278.0    Max.   :1898.0
##                                     NA's   :81
## FIELDING_DP
## Min.    : 52
## 1st Qu.:130
## Median :148
## Mean   :146
## 3rd Qu.:164
## Max.   :225
## NA's   :189
```

a. Mean / Standard Deviation / Median

```
ggplot(train, aes(x = TARGET_WINS)) +
  geom_histogram()
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



b. Bar Chart or Box Plot of the data

```
library(reshape)
```

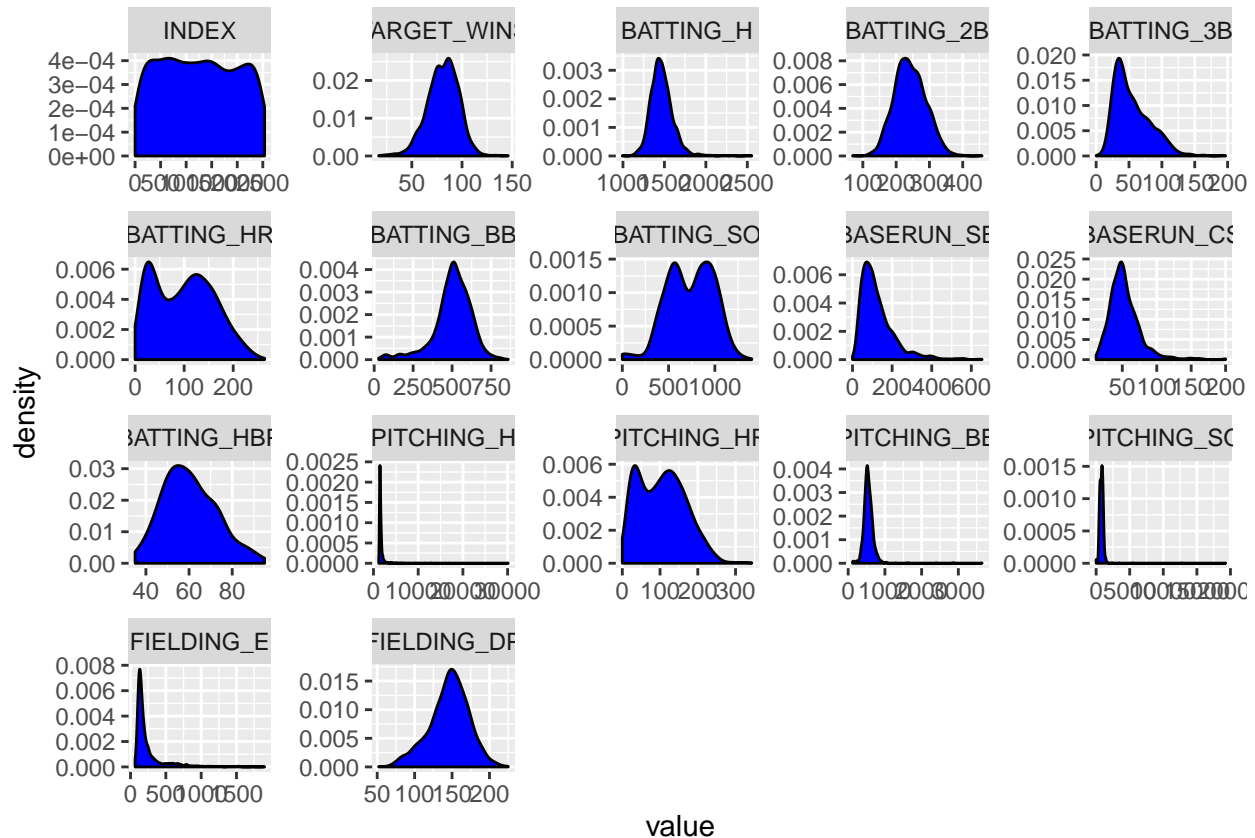
```
##  
## Attaching package: 'reshape'  
  
## The following object is masked from 'package:dplyr':  
##  
##   rename  
  
## The following objects are masked from 'package:tidyr':  
##  
##   expand, smiths
```

```
library(ggplot2)  
par(mfrow = c(3, 3))  
datasub = melt(train)
```

```
## Using   as id variables
```

```
ggplot(datasub, aes(x= value)) +
  geom_density(fill='blue') +
  facet_wrap(~variable, scales = 'free')
```

```
## Warning: Removed 2443 rows containing non-finite values (stat_density).
```



c. Is the data correlated to the target variable (or to other variables?)

Findings: 1. TEAM_BATTING_H exhibits the highest correlation to the response variable, 2. TEAM_FIELDING_E exhibits the lowest correlation 3. Both TEAM_PITCHING_HR and TEAM_PITCHING_BB exhibit positive correlations to the response variable 4. The correlation plot shows that TARGET_WINS is positively correlated with BATTING_H, BATTING_2B, BATTING_HR, BATTING_BB, PITCHING_H, PITCHING_HR, PITCHING_BB and negatively correlated with FIELDING_E. Thus we are going to construct our linear model by selecting from these attributes.

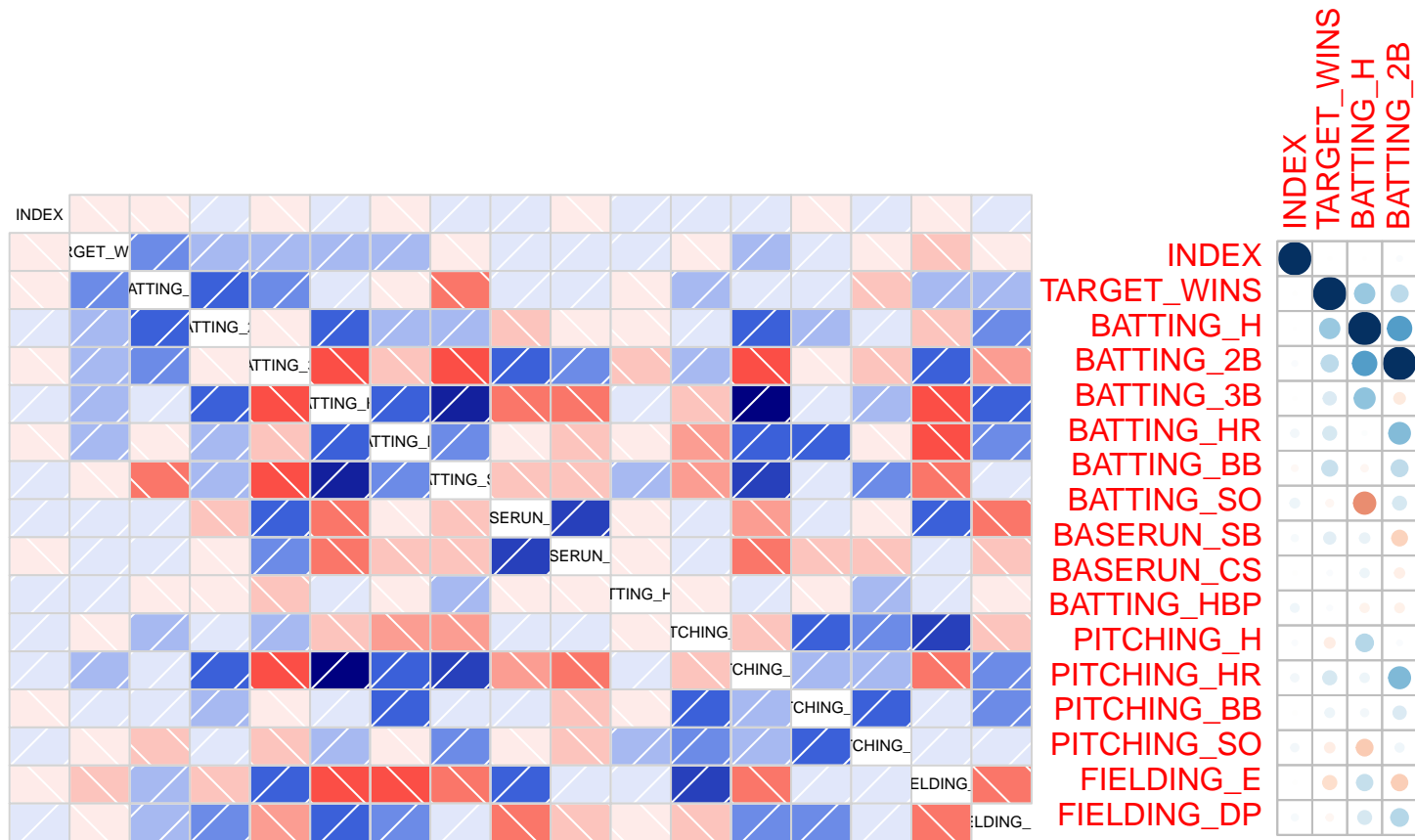
```
library(corrplot)
```

```
## corrplot 0.84 loaded
```

```
library(corrgram)
```

```
## Registered S3 method overwritten by 'seriation':
##   method      from
##   reorder.hclust gclus
```

```
corrplot(corrgram(train), method="circle")
```



d. Are any of the variables missing and need to be imputed “fixed”?

```
## Loading required package: colorspace
```

```
## Loading required package: grid
```

```
## Loading required package: data.table
```

```
##
```

```
## Attaching package: 'data.table'
```

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

```
##
```

```
## melt
```

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

```
##
```

```
## between, first, last
```

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

```
##
```

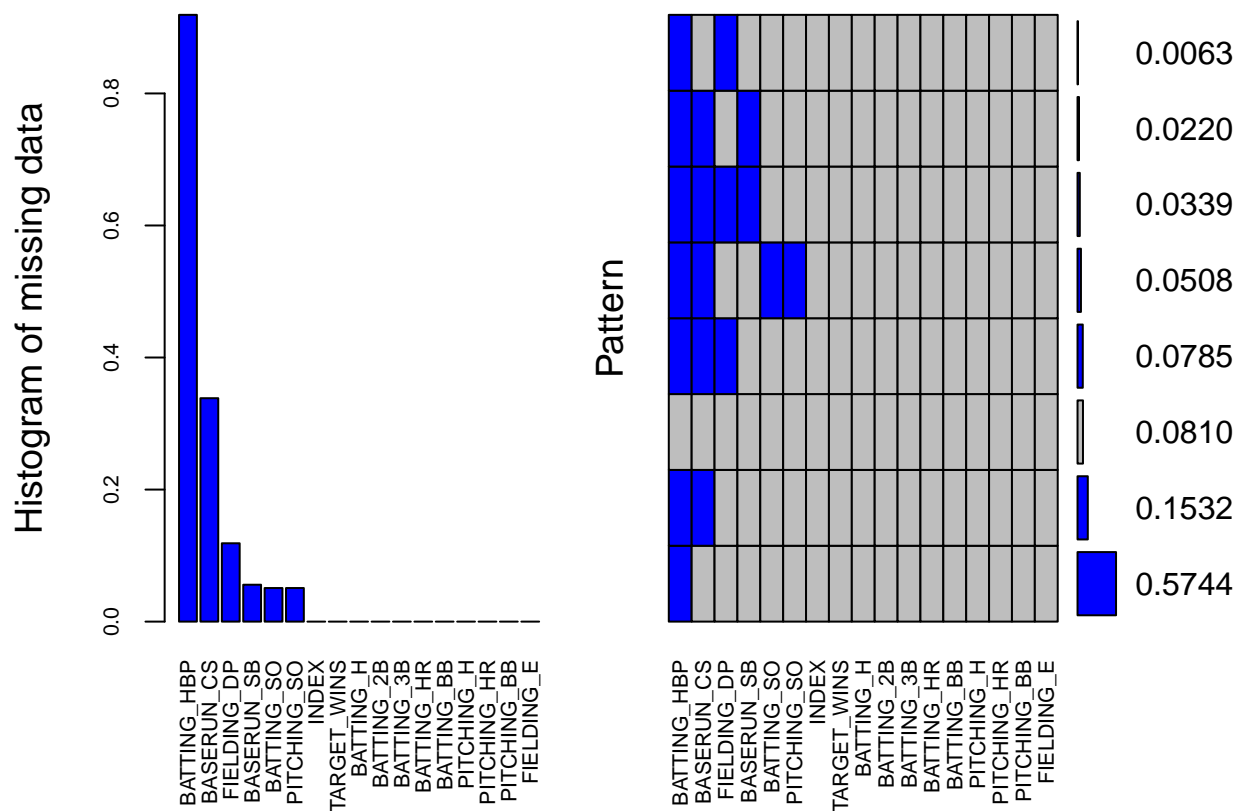
```
## transpose
```

```
## VIM is ready to use.
## Since version 4.0.0 the GUI is in its own package VIMGUI.
##
## Please use the package to use the new (and old) GUI.

## Suggestions and bug-reports can be submitted at: https://github.com/alexkova/VIM/issues

##
## Attaching package: 'VIM'

## The following object is masked from 'package:datasets':
##
## sleep
```



```
##
## Variables sorted by number of missings:
## Variable Count
## BATTING_HBP 0.91902072
## BASERUN_CS 0.33835530
## FIELDING_DP 0.11864407
## BASERUN_SB 0.05586943
## BATTING_SO 0.05084746
## PITCHING_SO 0.05084746
## INDEX 0.00000000
## TARGET_WINS 0.00000000
## BATTING_H 0.00000000
```

```
## BATTING_2B 0.00000000
## BATTING_3B 0.00000000
## BATTING_HR 0.00000000
## BATTING_BB 0.00000000
## PITCHING_H 0.00000000
## PITCHING_HR 0.00000000
## PITCHING_BB 0.00000000
## FIELDING_E 0.00000000
```

2. DATA PREPARATION (25 Points)

Describe how you have transformed the data by changing the original variables or creating new variables. If you did transform the data or create new variables, discuss why you did this. Here are some possible transformations. a. Fix missing values (maybe with a Mean or Median value) b. Create flags to suggest if a variable was missing c. Transform data by putting it into buckets d. Mathematical transforms such as log or square root (or use Box-Cox) e. Combine variables (such as ratios or adding or multiplying) to create new variables

Missing imputation

Considering some columns has outliers, we'll fill in the missing values using their respective median values.

```
train_clean = train %>% mutate(
  PITCHING_SO = ifelse(is.na(train$PITCHING_SO), median(train$PITCHING_SO, na.rm = TRUE), train$PITCHING_SO),
  BATTING_SO = ifelse(is.na(train$BATTING_SO), median(train$BATTING_SO, na.rm = TRUE), train$BATTING_SO),
  BASERUN_SB = ifelse(is.na(train$BASERUN_SB), median(train$BASERUN_SB, na.rm = TRUE), train$BASERUN_SB),
  BASERUN_CS = ifelse(is.na(train$BASERUN_CS), median(train$BASERUN_CS, na.rm = TRUE), train$BASERUN_CS),
  FIELDING_DP = ifelse(is.na(train$FIELDING_DP), median(train$FIELDING_DP, na.rm = TRUE), train$FIELDING_DP))
```

Feature engineering

We'll add a new variable BATTING_HBP_YN that is 1 when the TEAM_BATTING_HBP exists and 0 when it does not.

```
train_clean = train_clean %>% mutate(BATTING_HBP_YN = ifelse(is.na(BATTING_HBP), 0, 1),
  BATTING_1B = BATTING_H - BATTING_2B - BATTING_3B - BATTING_HR)
```

Create ratios: TARGET_WINS_Ratio = TARGET_WINS / 162 (i.e. the percentage of wins)
 TEAM_H_Ratio = (TEAM_BATTING_1B + TEAM_BATTING_2B + TEAM_BATTING_3B + TEAM_BATTING_HR) / TEAM_PITCHING_H (i.e. the ratio of hits earned to hits allowed)
 TEAM_BASERUN_Ratio = TEAM_BASERUN_SB / TEAM_BASERUN_CS (i.e. the ratio of successful steals to unsuccessful ones)
 TEAM_HR_SO_Ratio = TEAM_BATTING_HR / TEAM_BATTING_SO (i.e. the ratio of home runs to strikeouts)

```
train_clean = train_clean %>%
  mutate(H_Ratio = (BATTING_1B + BATTING_2B + BATTING_3B + BATTING_HR) / PITCHING_H,
    BASERUN_Ratio = BASERUN_SB / BASERUN_CS,
    HR_SO_Ratio = BATTING_HR / ifelse(BATTING_SO == 0, median(BATTING_SO), BATTING_SO))
```

3. BUILD MODELS (25 Points) Using the training data set, build at least three different multiple linear regression models, using different variables (or the same variables with different transformations).

Since we have not yet covered automated variable selection methods, you should select the variables manually (unless you previously learned Forward or Stepwise selection, etc.). Since you manually selected a variable for inclusion into the model or exclusion into the model, indicate why this was done. Discuss the coefficients in the models, do they make sense? For example, if a team hits a lot of Home Runs, it would be reasonably expected that such a team would win more games. However, if the coefficient is negative (suggesting that the team would lose more games), then that needs to be discussed. Are you keeping the model even though it is counter intuitive? Why? The boss needs to know.

Model 1: Simple linear regression using all features in training dataset

```
library(caret)

## Loading required package: lattice

##
## Attaching package: 'lattice'

## The following object is masked from 'package:corrgram':
##
##   panel.fill

##
## Attaching package: 'caret'

## The following object is masked from 'package:purrr':
##
##   lift

train_model1 = train_clean
train_model1 = train_model1 %>% select(-INDEX, -BATTING_HBP)
model1 = train(TARGET_WINS ~ ., data = train_model1, method = 'lm', na.action=na.exclude)

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading
```

```
## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading

## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading
```

```
## Warning in predict.lm(modelFit, newdata): prediction from a rank-deficient
## fit may be misleading
```

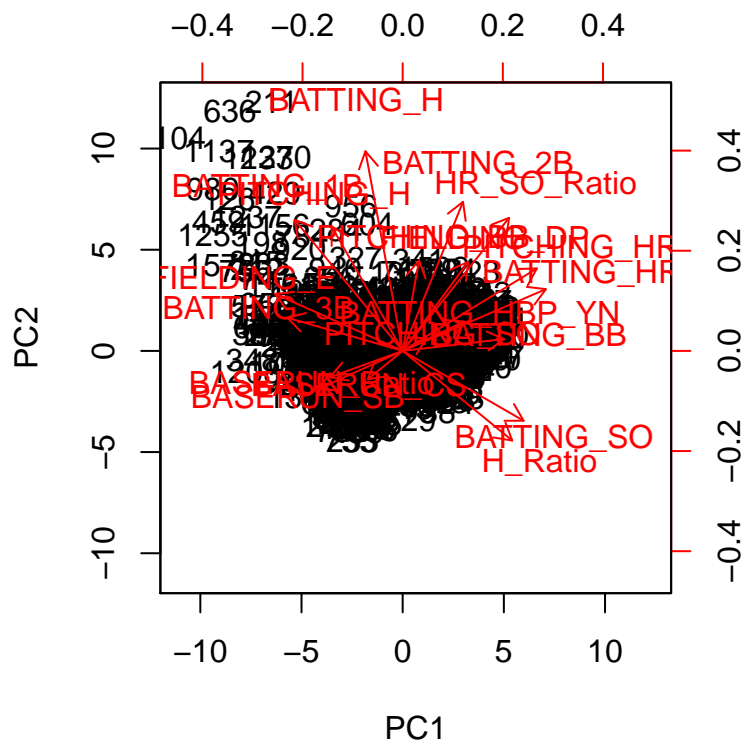
```
summary(model1)
```

```
##
## Call:
## lm(formula = .outcome ~ ., data = dat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -38.635  -8.443  -0.144   8.062  64.447
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   9.201e+01  1.009e+01   9.122 < 2e-16 ***
## BATTING_H      4.288e-02  4.570e-03   9.382 < 2e-16 ***
## BATTING_2B     -1.553e-02  1.094e-02  -1.419  0.15600
## BATTING_3B      1.009e-01  2.011e-02   5.018 5.81e-07 ***
## BATTING_HR      3.175e-01  4.440e-02   7.152 1.31e-12 ***
## BATTING_BB      2.327e-02  7.218e-03   3.224  0.00129 **
## BATTING_SO      1.816e-03  3.906e-03   0.465  0.64203
## BASERUN_SB      5.935e-02  2.790e-02   2.127  0.03358 *
## BASERUN_CS     -5.014e-02  4.743e-02  -1.057  0.29061
## PITCHING_H     -7.145e-04  5.072e-04  -1.409  0.15911
## PITCHING_HR    -2.799e-01  4.426e-02  -6.324 3.32e-10 ***
## PITCHING_BB    -9.831e-03  5.155e-03  -1.907  0.05671 .
## PITCHING_SO     1.227e-03  9.767e-04   1.256  0.20918
## FIELDING_E     -4.506e-02  4.237e-03 -10.637 < 2e-16 ***
## FIELDING_DP    -1.131e-01  1.502e-02  -7.530 8.51e-14 ***
## BATTING_HBP_YN -4.349e+00  1.461e+00  -2.975  0.00297 **
## BATTING_1B      NA         NA         NA         NA
## H_Ratio        -7.033e+01  8.500e+00  -8.275 2.71e-16 ***
## BASERUN_Ratio  -1.077e+00  1.427e+00  -0.754  0.45066
## HR_SO_Ratio     3.806e+01  1.557e+01   2.444  0.01465 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 12.67 on 1574 degrees of freedom
## Multiple R-squared:  0.3393, Adjusted R-squared:  0.3317
## F-statistic: 44.9 on 18 and 1574 DF, p-value: < 2.2e-16
```

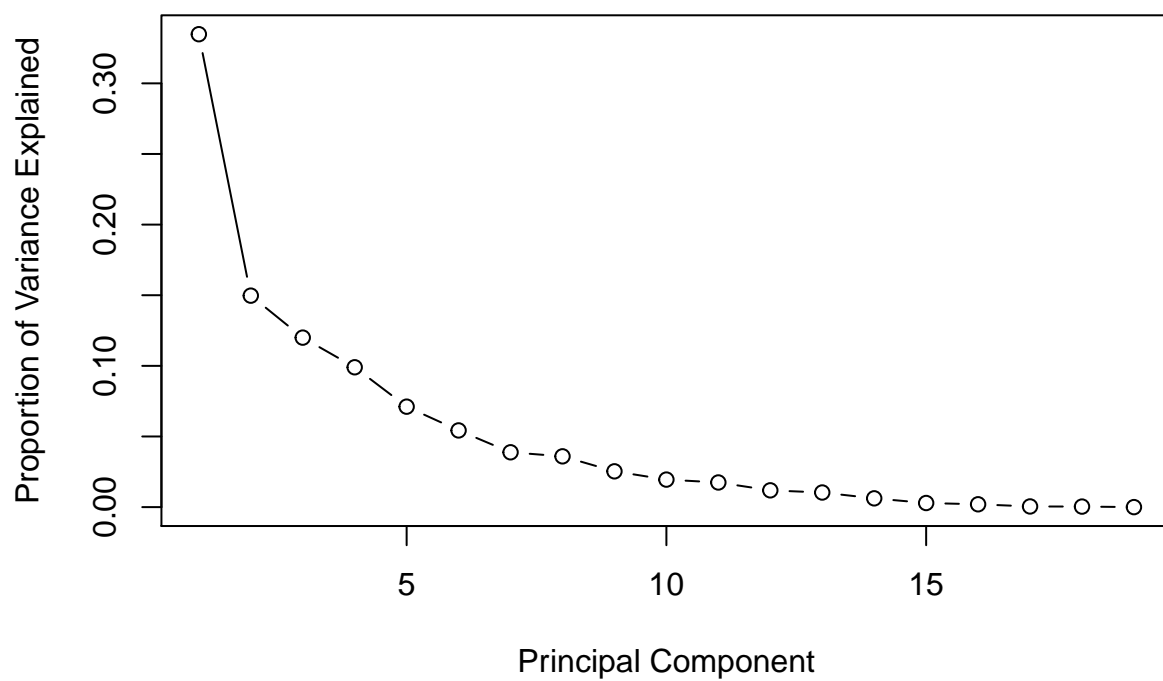
###Model2 Principal Component Analysis Given there is strong multicollinearity among variable, it is better to conduct principal component analysis on dataset in order to eliminate the colinearity.

Welcome! Want to learn more? See two factoextra-related books at <https://goo.gl/ve3WBa>

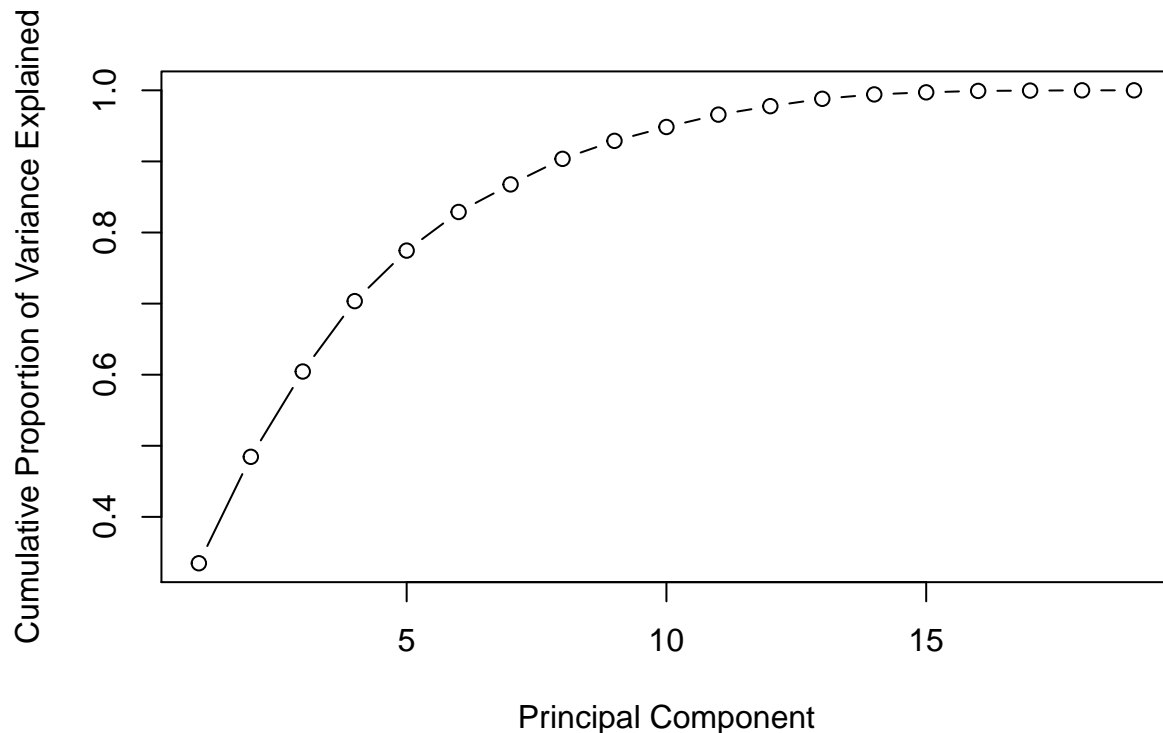
```
biplot(prin_comp, scale = 0)
```



```
std_dev <- prin_comp$sdev
pr_var <- std_dev^2
prop_varex <- pr_var/sum(pr_var)
plot(prop_varex, xlab = "Principal Component",
      ylab = "Proportion of Variance Explained",
      type = "b")
```



```
plot(cumsum(prop_varex), xlab = "Principal Component",
     ylab = "Cumulative Proportion of Variance Explained",
     type = "b")
```



This plot shows that 15 components results in variance close to ~ 98%. Therefore, in this case, we'll select number of components as 15 [PC1 to PC15] and proceed to the modeling stage. This completes the steps to implement PCA on train data. For modeling, we'll use these 15 components as predictor variables and follow the normal procedures.

```
model2_pca.data <- data.frame(TARGET_WINS = train_model2$TARGET_WINS, prin_comp$x)
model2_pca.data = model2_pca.data[1:16]
model2 = train(TARGET_WINS ~ ., data = model2_pca.data, method = 'lm', na.action=na.exclude)
summary(model2)
```

```
##
## Call:
## lm(formula = .outcome ~ ., data = dat)
##
## Residuals:
```

	Min	1Q	Median	3Q	Max
	-48.040	-8.500	0.101	8.221	58.626

```
##
## Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	81.0195	0.3237	250.314	< 2e-16 ***
PC1	0.4573	0.1284	3.562	0.000379 ***
PC2	2.4262	0.1920	12.638	< 2e-16 ***
PC3	-2.7176	0.2144	-12.673	< 2e-16 ***
PC4	-2.8697	0.2361	-12.155	< 2e-16 ***
PC5	0.1303	0.2786	0.468	0.639989

```
## PC6          1.0814      0.3188    3.392 0.000710 ***
## PC7          0.1545      0.3771    0.410 0.682085
## PC8         -4.5923      0.3918   -11.721 < 2e-16 ***
## PC9          1.4375      0.4667    3.080 0.002105 **
## PC10         -1.4451      0.5322   -2.715 0.006694 **
## PC11         -2.1563      0.5626   -3.833 0.000132 ***
## PC12          1.3906      0.6816    2.040 0.041493 *
## PC13          1.4084      0.7304    1.928 0.054009 .
## PC14         -4.0765      0.9436   -4.320 1.66e-05 ***
## PC15         -7.0788      1.3792   -5.133 3.21e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 12.92 on 1577 degrees of freedom
## Multiple R-squared:  0.3118, Adjusted R-squared:  0.3053
## F-statistic: 47.63 on 15 and 1577 DF,  p-value: < 2.2e-16
```

```
model3 <- lm(TARGET_WINS ~ BATTING_H+BATTING_2B+BATTING_3B+BATTING_HR+BATTING_BB+BATTING_HBP-BATTING_SO
summary(model3)
```

```
##
## Call:
## lm(formula = TARGET_WINS ~ BATTING_H + BATTING_2B + BATTING_3B +
##     BATTING_HR + BATTING_BB + BATTING_HBP - BATTING_SO + BASERUN_SB -
##     BASERUN_CS - FIELDING_E + FIELDING_DP - PITCHING_BB - PITCHING_H -
##     PITCHING_HR + PITCHING_SO, data = train_clean)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -21.7659  -5.6300  -0.0169   5.0174  20.5997
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  38.400045  23.380068   1.642   0.1031
## BATTING_H     0.032959   0.016736   1.969   0.0512 .
## BATTING_2B    0.001462   0.039598   0.037   0.9706
## BATTING_3B   -0.168639   0.092540  -1.822   0.0709 .
## BATTING_HR    0.036132   0.032383   1.116   0.2668
## BATTING_BB    0.063442   0.012614   5.030 1.76e-06 ***
## BATTING_HBP   0.098657   0.065525   1.506   0.1348
## BASERUN_SB    0.042778   0.029610   1.445   0.1512
## FIELDING_DP  -0.091334   0.044775  -2.040   0.0436 *
## PITCHING_SO  -0.036425   0.008461  -4.305 3.45e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.64 on 119 degrees of freedom
## (1464 observations deleted due to missingness)
## Multiple R-squared:  0.4588, Adjusted R-squared:  0.4178
## F-statistic: 11.21 on 9 and 119 DF,  p-value: 1.523e-12
```

Compare the RMSE(Root Mean Squared Error) among the 3 models - model 3 appears to have the lowest RMSE

```
fit1 <- fitted.values(model1)
error1 <- fit1 - test$TARGET_WINS
```

```
## Warning in fit1 - test$TARGET_WINS: longer object length is not a multiple
## of shorter object length
```

```
rmse1 <- sqrt(mean(error1^2))
rmse1
```

```
## [1] 18.14997
```

```
fit2 <- fitted.values(model2)
error2 <- fit2 - test$TARGET_WINS
```

```
## Warning in fit2 - test$TARGET_WINS: longer object length is not a multiple
## of shorter object length
```

```
rmse2 <- sqrt(mean(error2^2))
rmse2
```

```
## [1] 18.01896
```

```
fit3 <- fitted.values(model3)
error3 <- fit3 - test$TARGET_WINS
```

```
## Warning in fit3 - test$TARGET_WINS: longer object length is not a multiple
## of shorter object length
```

```
rmse3 <- sqrt(mean(error3^2))
rmse3
```

```
## [1] 17.94966
```

Model selection rationale

As discussed above, we selected Model2, which was based on principal component analysis, followed by removal of any highly collinear variables. Although Model2 did not have the lowest RMSE, it was the most stable (little collinearity between variables).

Inference and regression diagnostics

For our inferences to be valid, we need to perform some regression diagnostics and validate some assumptions:

Independence of errors: Based on the residual plot below, the residuals appear random over the index values

Outliers and leverage: Based on the leverage plots below, there do not appear to be any data points exerting undue leverage on the regression

Normality: Based on the qq-plot below, the residuals are fairly normally distributed, although there are some outliers in the tails

Constant variance: Based on the spread-level plot below, variance appears relatively constant, although again with a few outliers