

# moneyball

9/6/2021

```
# load data from Doug's github repo
df <- read.csv('https://raw.githubusercontent.com/douglasbarley/DATA621/main/Homework1/moneyball-training.csv')
summary(df)
```

```
##      INDEX      TARGET_WINS      TEAM_BATTING_H TEAM_BATTING_2B
## Min.   : 1.0    Min.   : 0.00    Min.   : 891    Min.   : 69.0
## 1st Qu.: 630.8  1st Qu.: 71.00    1st Qu.:1383    1st Qu.:208.0
## Median :1270.5  Median : 82.00    Median :1454    Median :238.0
## Mean   :1268.5  Mean   : 80.79    Mean   :1469    Mean   :241.2
## 3rd Qu.:1915.5  3rd Qu.: 92.00    3rd Qu.:1537    3rd Qu.:273.0
## Max.   :2535.0  Max.   :146.00    Max.   :2554    Max.   :458.0
##
## TEAM_BATTING_3B TEAM_BATTING_HR TEAM_BATTING_BB TEAM_BATTING_SO
## Min.   : 0.00    Min.   : 0.00    Min.   : 0.0    Min.   : 0.0
## 1st Qu.: 34.00    1st Qu.: 42.00    1st Qu.:451.0    1st Qu.: 548.0
## Median : 47.00    Median :102.00    Median :512.0    Median : 750.0
## Mean   : 55.25    Mean   : 99.61    Mean   :501.6    Mean   : 735.6
## 3rd Qu.: 72.00    3rd Qu.:147.00    3rd Qu.:580.0    3rd Qu.: 930.0
## Max.   :223.00    Max.   :264.00    Max.   :878.0    Max.   :1399.0
##                                     NA's   :102
## TEAM_BASERUN_SB TEAM_BASERUN_CS TEAM_BATTING_HBP TEAM_PITCHING_H
## Min.   : 0.0    Min.   : 0.0    Min.   :29.00    Min.   : 1137
## 1st Qu.: 66.0    1st Qu.: 38.0    1st Qu.:50.50    1st Qu.: 1419
## Median :101.0    Median : 49.0    Median :58.00    Median : 1518
## Mean   :124.8    Mean   : 52.8    Mean   :59.36    Mean   : 1779
## 3rd Qu.:156.0    3rd Qu.: 62.0    3rd Qu.:67.00    3rd Qu.: 1682
## Max.   :697.0    Max.   :201.0    Max.   :95.00    Max.   :30132
## NA's   :131     NA's   :772     NA's   :2085
## TEAM_PITCHING_HR TEAM_PITCHING_BB TEAM_PITCHING_SO TEAM_FIELDING_E
## Min.   : 0.0    Min.   : 0.0    Min.   : 0.0    Min.   : 65.0
## 1st Qu.: 50.0    1st Qu.: 476.0    1st Qu.: 615.0    1st Qu.: 127.0
## Median :107.0    Median : 536.5    Median : 813.5    Median : 159.0
## Mean   :105.7    Mean   : 553.0    Mean   : 817.7    Mean   : 246.5
## 3rd Qu.:150.0    3rd Qu.: 611.0    3rd Qu.: 968.0    3rd Qu.: 249.2
## Max.   :343.0    Max.   :3645.0    Max.   :19278.0    Max.   :1898.0
##                                     NA's   :102
## TEAM_FIELDING_DP
## Min.   : 52.0
## 1st Qu.:131.0
## Median :149.0
## Mean   :146.4
## 3rd Qu.:164.0
## Max.   :228.0
## NA's   :286
```

```
str(df)
```

```
## 'data.frame':    2276 obs. of  17 variables:
## $ INDEX          : int  1 2 3 4 5 6 7 8 11 12 ...
## $ TARGET_WINS     : int  39 70 86 70 82 75 80 85 86 76 ...
## $ TEAM_BATTING_H   : int  1445 1339 1377 1387 1297 1279 1244 1273 1391 1271 ...
## $ TEAM_BATTING_2B  : int  194 219 232 209 186 200 179 171 197 213 ...
## $ TEAM_BATTING_3B  : int  39 22 35 38 27 36 54 37 40 18 ...
## $ TEAM_BATTING_HR  : int  13 190 137 96 102 92 122 115 114 96 ...
## $ TEAM_BATTING_BB  : int  143 685 602 451 472 443 525 456 447 441 ...
## $ TEAM_BATTING_SO  : int  842 1075 917 922 920 973 1062 1027 922 827 ...
## $ TEAM_BASERUN_SB  : int  NA 37 46 43 49 107 80 40 69 72 ...
## $ TEAM_BASERUN_CS  : int  NA 28 27 30 39 59 54 36 27 34 ...
## $ TEAM_BATTING_HBP : int  NA NA NA NA NA NA NA NA NA NA ...
## $ TEAM_PITCHING_H  : int  9364 1347 1377 1396 1297 1279 1244 1281 1391 1271 ...
## $ TEAM_PITCHING_HR : int  84 191 137 97 102 92 122 116 114 96 ...
## $ TEAM_PITCHING_BB : int  927 689 602 454 472 443 525 459 447 441 ...
## $ TEAM_PITCHING_SO : int  5456 1082 917 928 920 973 1062 1033 922 827 ...
## $ TEAM_FIELDING_E  : int  1011 193 175 164 138 123 136 112 127 131 ...
## $ TEAM_FIELDING_DP : int  NA 155 153 156 168 149 186 136 169 159 ...
```

```
kable(head(df))
```

INDEX	TEAM	BATTING_H	BATTING_2B	BATTING_3B	BATTING_HR	BATTING_BB	BATTING_SO	BASERUN_SB	BASERUN_CS	BATTING_HBP	PITCHING_H	PITCHING_HR	PITCHING_BB	PITCHING_SO	FIELDING_E	FIELDING_DP
1	39	1445	194	39	13	143	842	NA	NA	NA	9364	84	927	5456	1011	NA
2	70	1339	219	22	190	685	1075	37	28	NA	1347	191	689	1082	193	155
3	86	1377	232	35	137	602	917	46	27	NA	1377	137	602	917	175	153
4	70	1387	209	38	96	451	922	43	30	NA	1396	97	454	928	164	156
5	82	1297	186	27	102	472	920	49	39	NA	1297	102	472	920	138	168
6	75	1279	200	36	92	443	973	107	59	NA	1279	92	443	973	123	149

Clipping outliers and imputing missing values:

We should look at using median for outliers instead of clipping

```
# Impute median for these
```

```
df$TEAM_BASERUN_CS[is.na(df$TEAM_BASERUN_CS)] = median(df$TEAM_BASERUN_CS, na.rm=T)
df$TEAM_BASERUN_SB[is.na(df$TEAM_BASERUN_SB)] = median(df$TEAM_BASERUN_SB, na.rm=T)
df$TEAM_BATTING_SO[is.na(df$TEAM_BATTING_SO)] = median(df$TEAM_BATTING_SO, na.rm=T)
df$TEAM_PITCHING_SO[is.na(df$TEAM_PITCHING_SO)] = median(df$TEAM_PITCHING_SO, na.rm=T)
df$TEAM_FIELDING_DP[is.na(df$TEAM_FIELDING_DP)] = median(df$TEAM_FIELDING_DP, na.rm=T)
```

```
# Correlation Between All Variables
```

```
mydata.cor = cor(df$TARGET_WINS, df)
kable(mydata.cor)
```

INDEX	TEAM_WINS	TEAM_BATTING_H	TEAM_BATTING_2B	TEAM_BATTING_3B	TEAM_BATTING_HR	TEAM_BATTING_BB	TEAM_BATTING_SO	TEAM_BASERUN_SB	TEAM_BASERUN_CS	TEAM_BATTING_HBP	TEAM_PITCHING_H	TEAM_PITCHING_HR	TEAM_PITCHING_BB	TEAM_PITCHING_SO	TEAM_FIELDING_E	TEAM_FIELDING_DP
-	1	0.388762591036426084761532325599	-	0.1236109159598A	-	0.1890137241745	-	-	-	-	-	-	-	-	-	-
0.0210564						0.0305814				0.1099371			0.07579917648418300863			

```
corrplot(mydata.cor)
```

