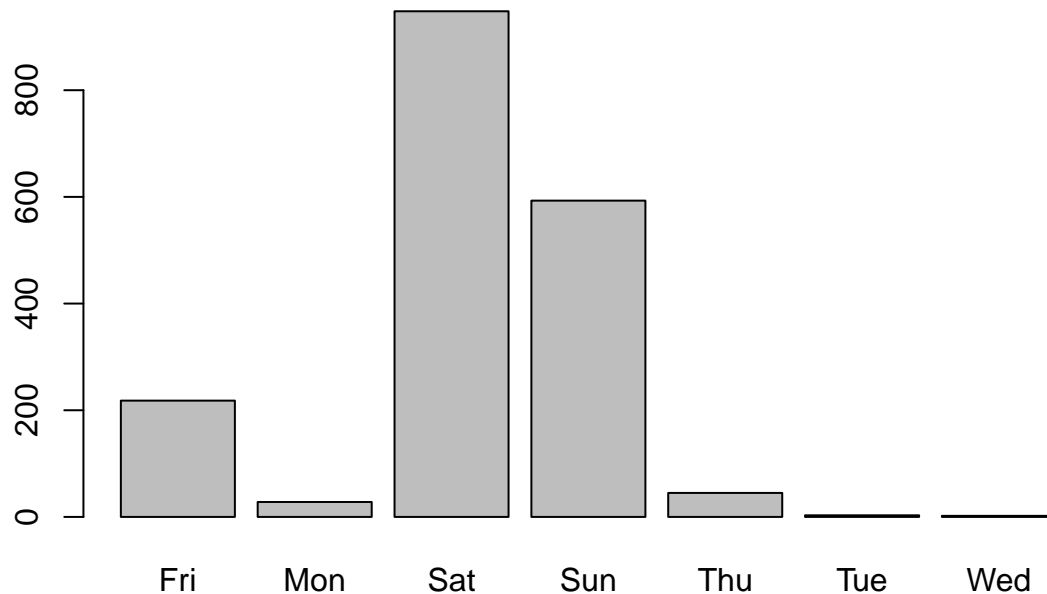


## Basic\_Analysis

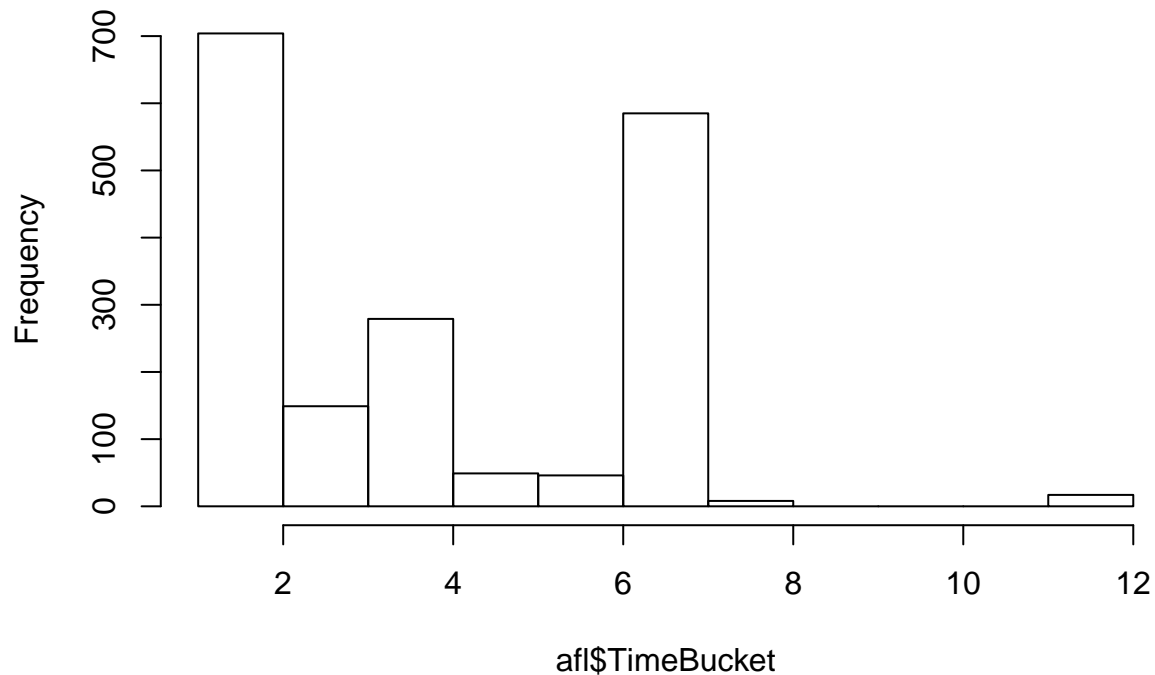
```
plot(afl$Day)
```



We see that the majority of games are on Saturdays and Sundays. We can perform further analysis to determine how day of week affects attendance.

```
hist(afl$TimeBucket)
```

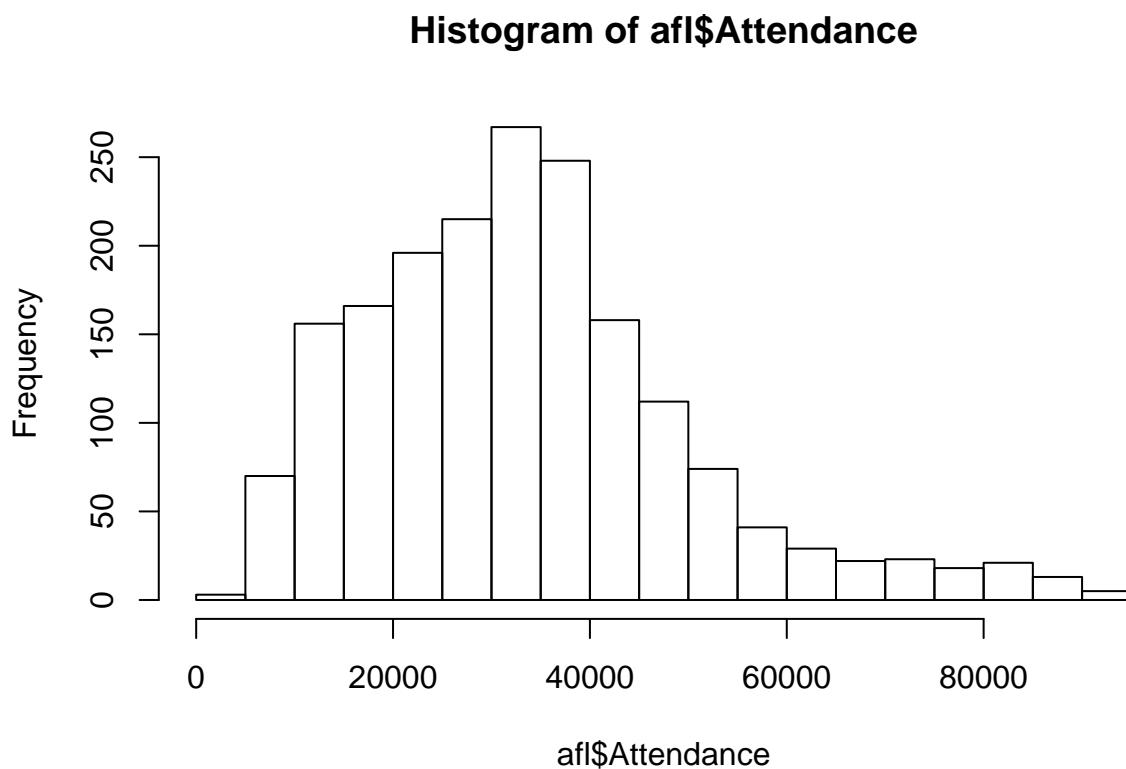
### Histogram of afl\$TimeBucket



bucket allows us to split the times into more easily workable data. We can see here that most of the games occur in the 1 o'clock and the 7 o'clock hours (start between 1 and 2 and between 7 and 8). We can perform

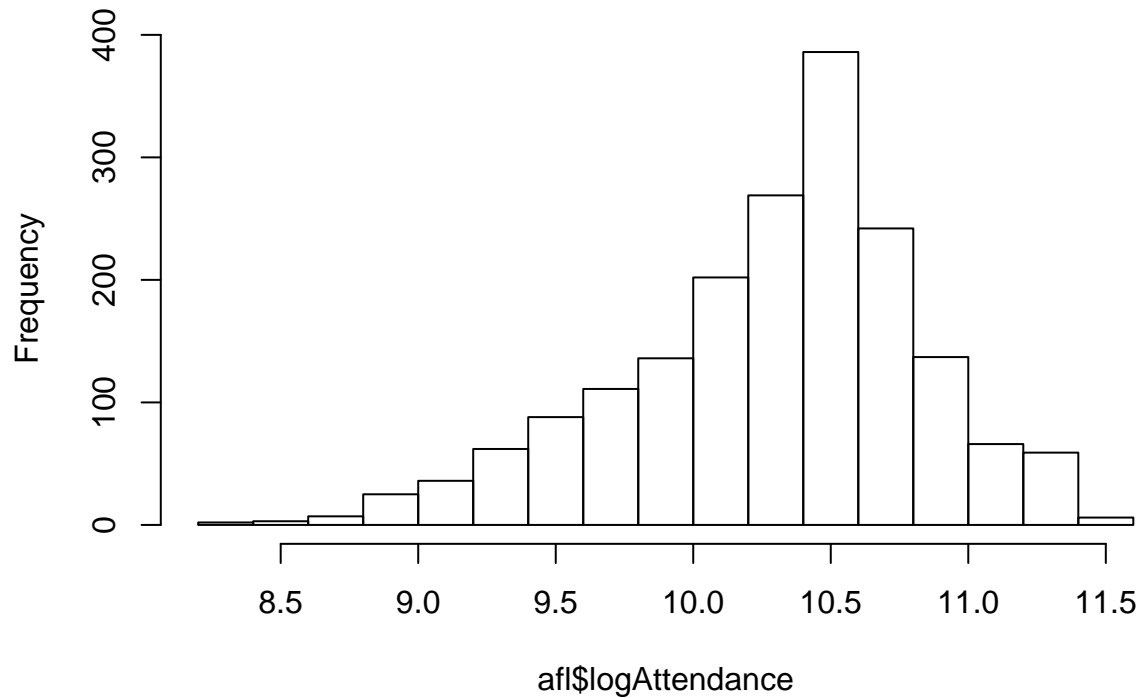
further analysis to determine if time of day affects attendance. Possible interaction between time of day and day of week.

```
hist(afl$Attendance, breaks = 30)
```



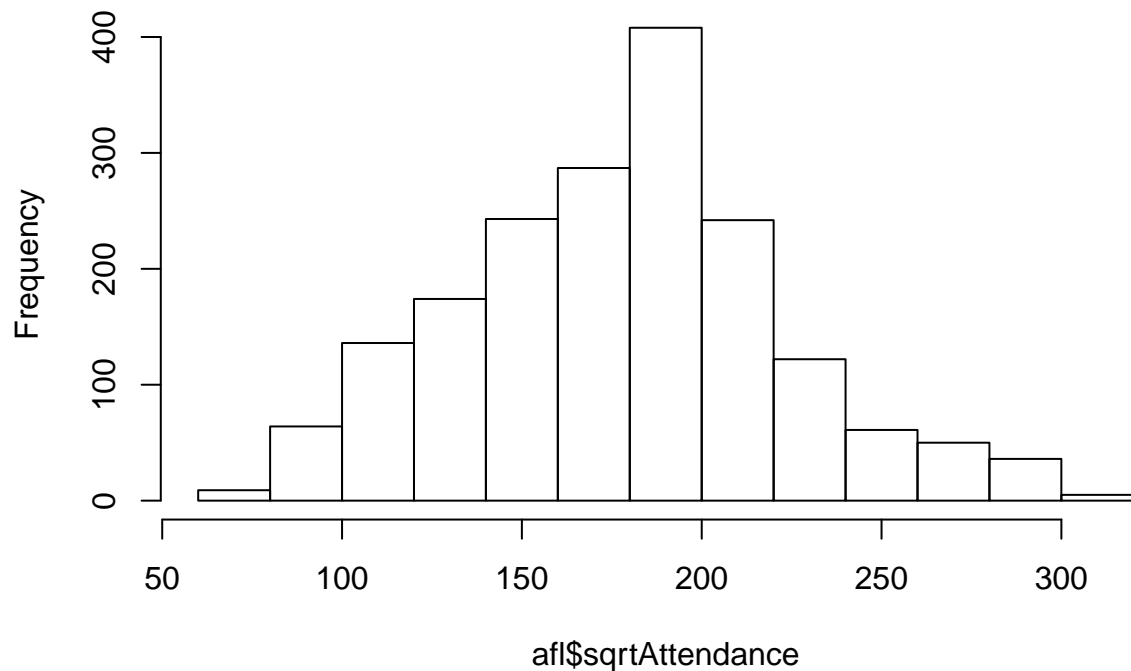
```
afl$logAttendance = log(afl$Attendance)
hist(afl$logAttendance)
```

### Histogram of afl\$logAttendance



```
afl$sqrtAttendance = sqrt(afl$Attendance)
hist(afl$sqrtAttendance)
```

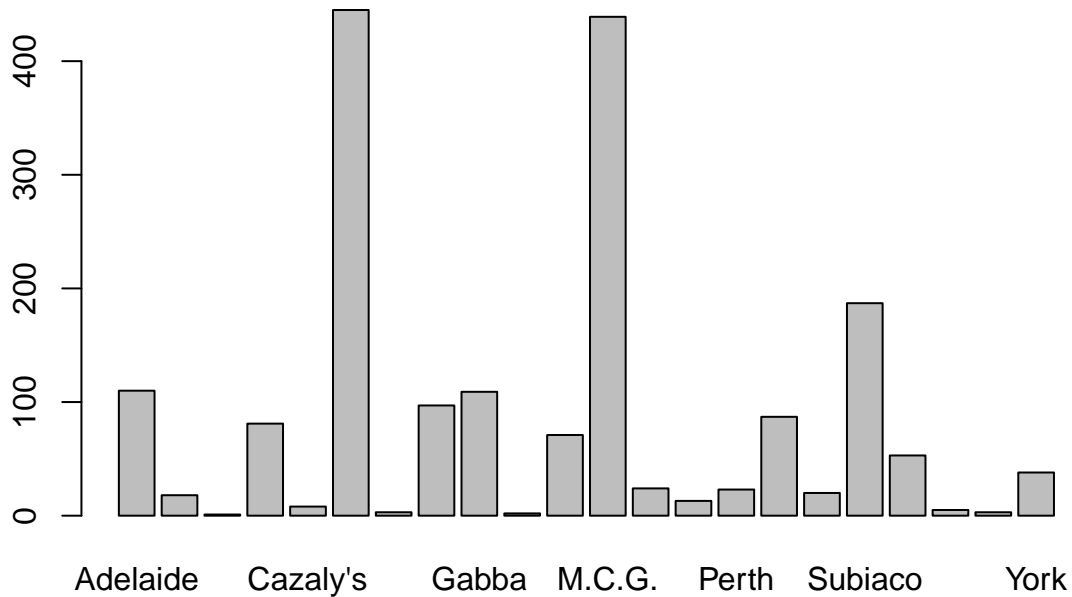
### Histogram of afl\$sqrtAttendance



Attendance may be our response variable as we are trying to determine the factors which influence fans to attend games. The histogram shows a right skew, so I tried an exponential transformation. This results in a left-skewed histogram. As such, I tried a square root transformation which gives us the most normal

distribution of the three.

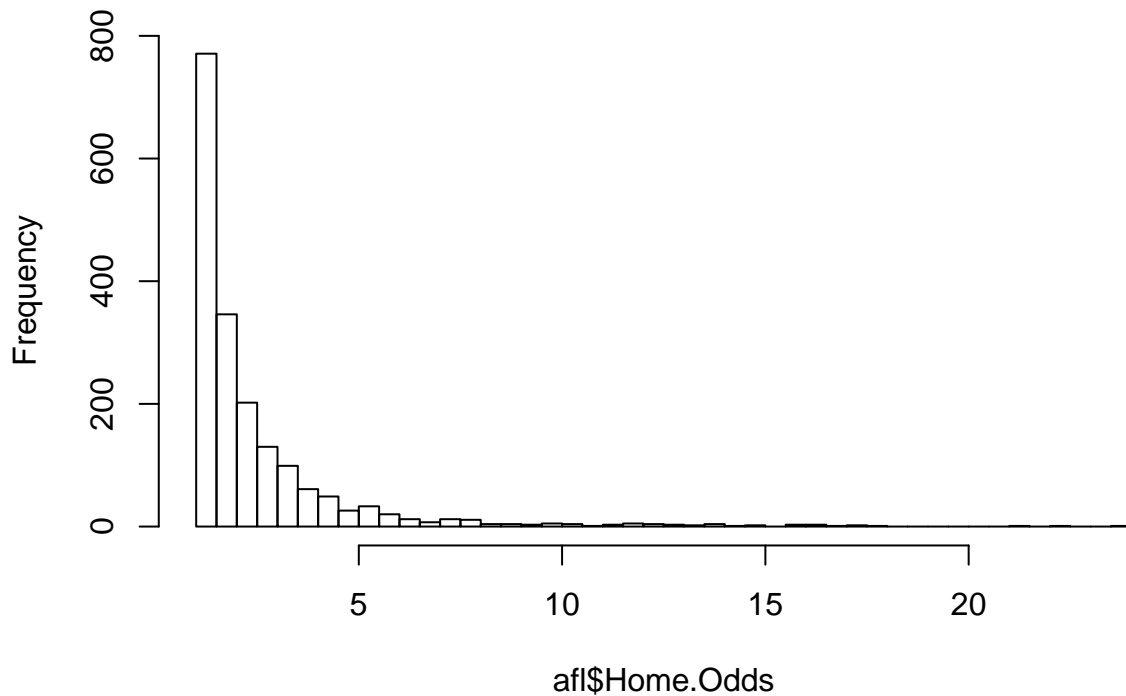
```
plot(afl$Venue)
```



This plot simply shows which stadiums host the most games. We see that a majority of the games are played at M.C.G. and Docklands.

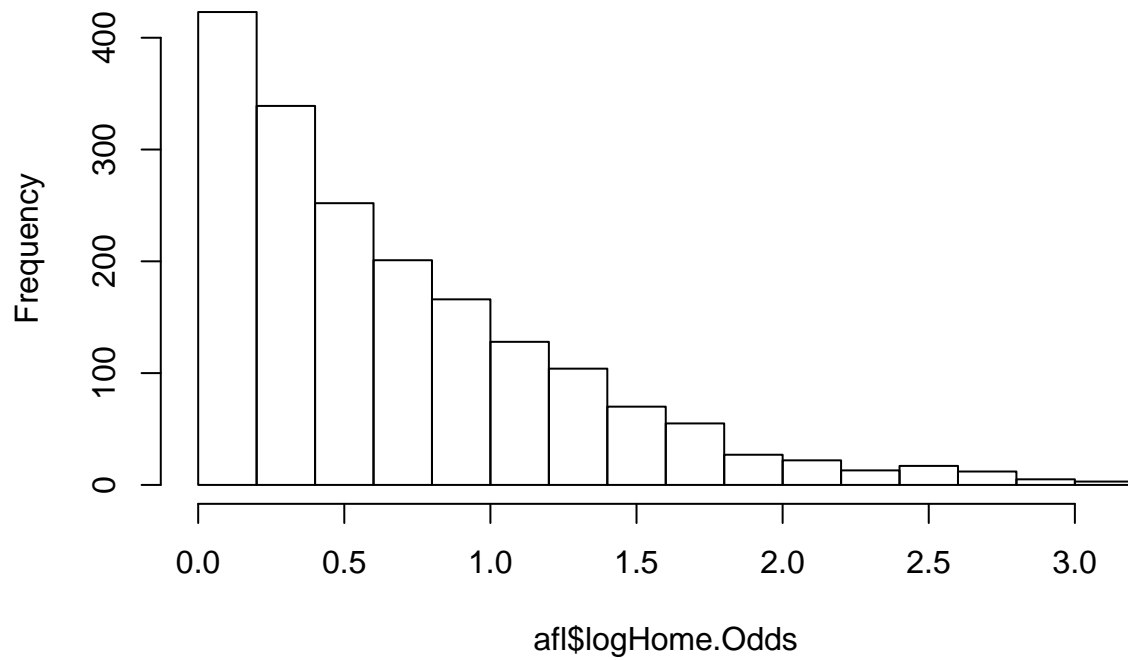
```
hist(afl$Home.Odds, breaks = 40)
```

### Histogram of afl\$Home.Odds



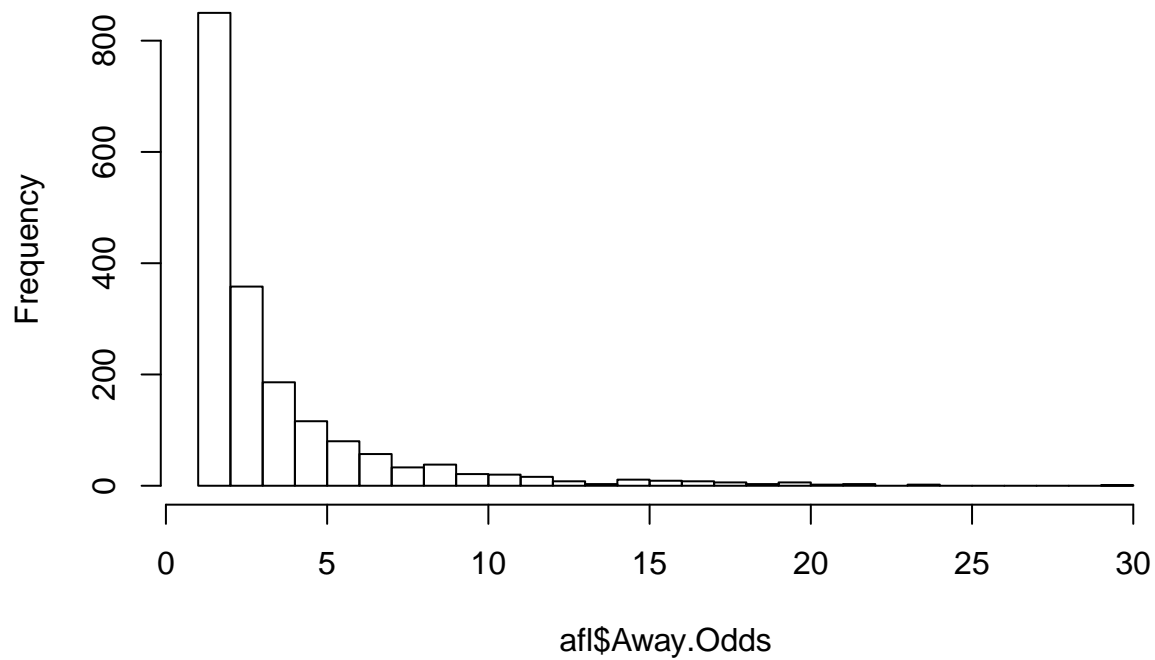
```
afl$logHome.Odds = log(afl$Home.Odds)  
hist(afl$logHome.Odds)
```

**Histogram of afl\$logHome.Odds**



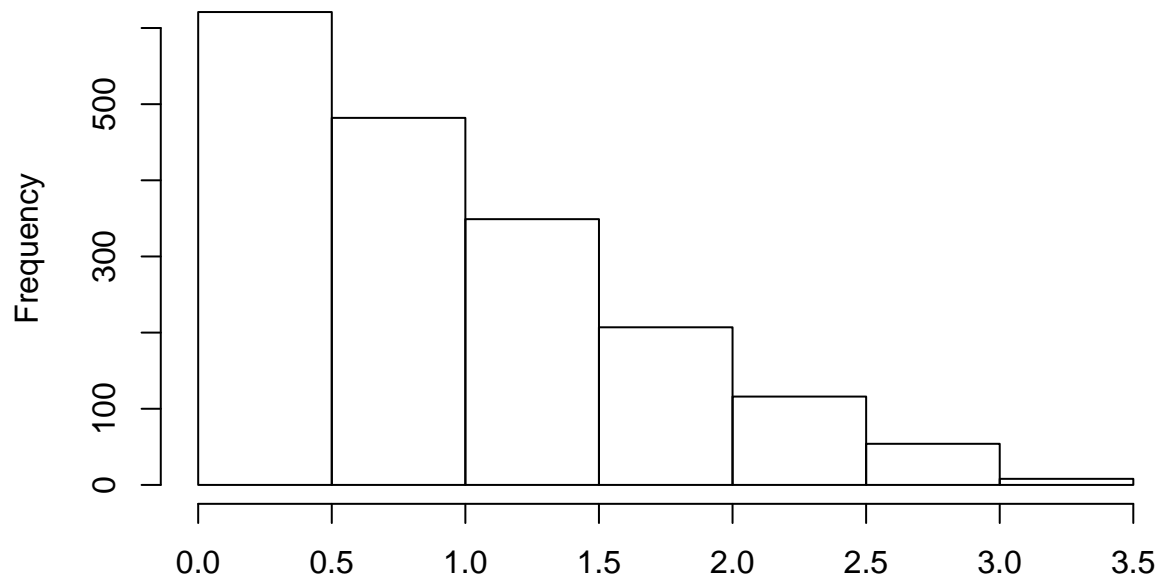
```
hist(afl$Away.Odds, breaks = 40)
```

**Histogram of afl\$Away.Odds**



```
afl$logAway.Odds = log(afl$Away.Odds)  
hist(afl$logAway.Odds)
```

## Histogram of afl\$logAway.Odds



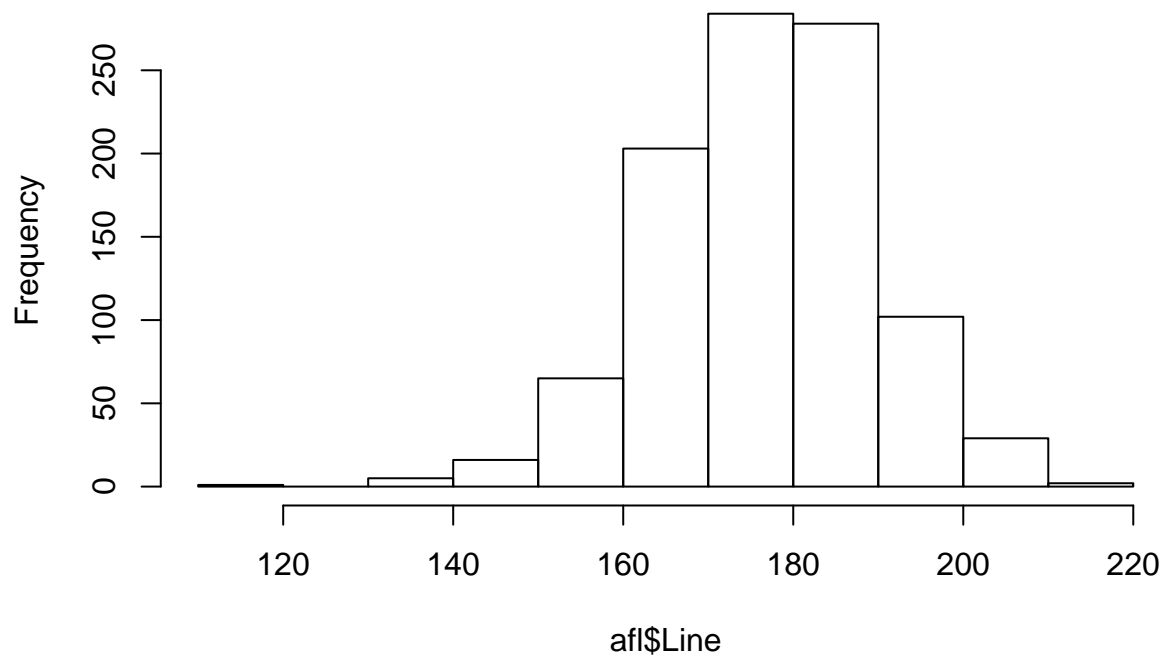
afl\$logAway.Odds

Both

home odds and away odds are extremely skewed. Taking the log of the log creates a much more difficult result to interpret and still doesn't give us normally distributed data. These variables will be two of the most important of our data, but as they correlated we can only use one. I show other ways we can incorporate this data below.

```
hist(afl$Line)
```

## Histogram of afl\$Line

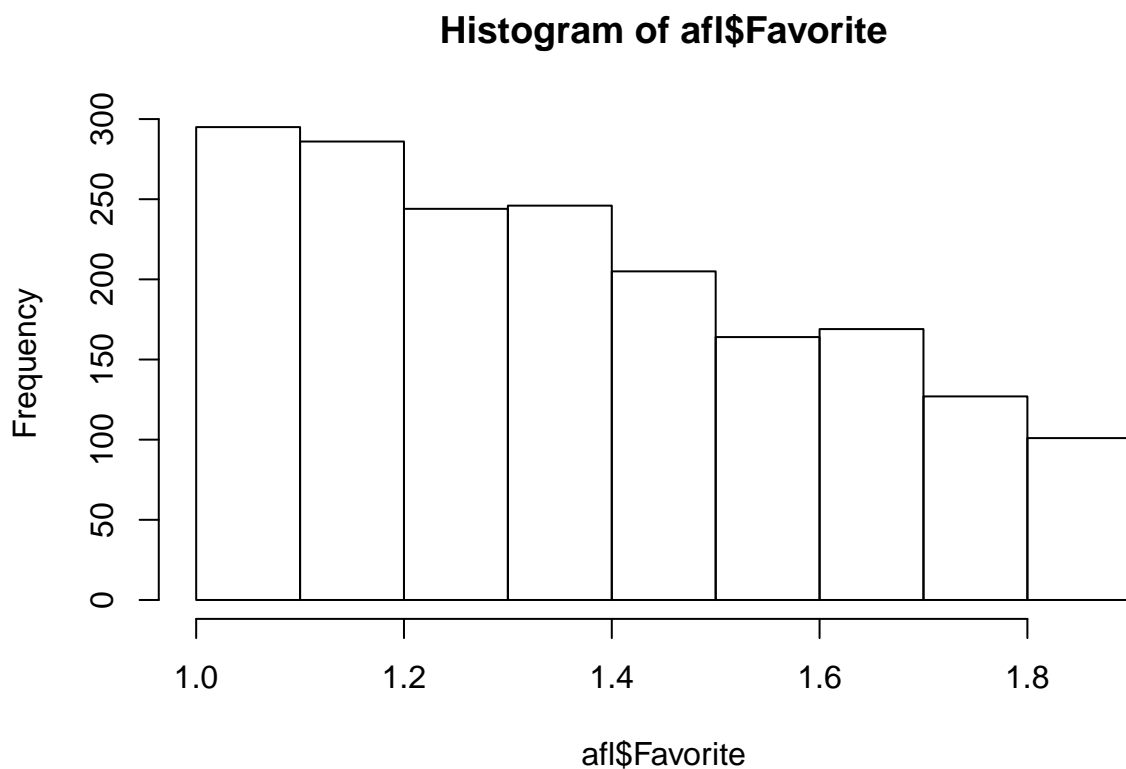


afl\$Line

The

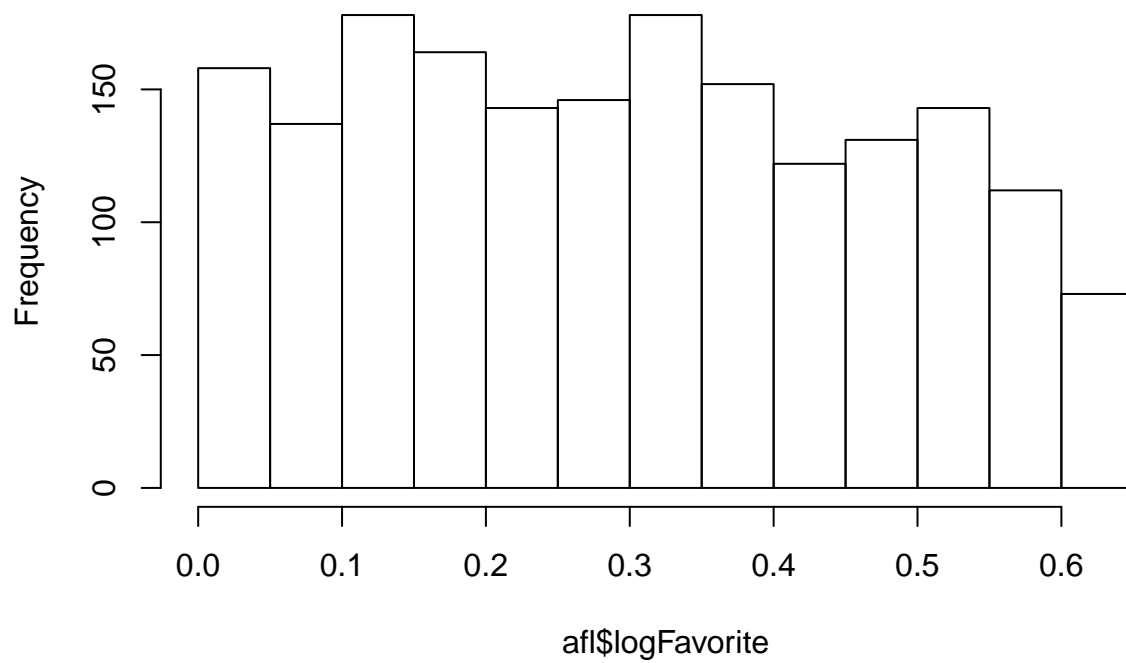
data for the over-under of the games is relatively normally distributed. However, only games from 2014 and later have over-under data.

```
hist(afl$Favorite, breaks = 10)
```



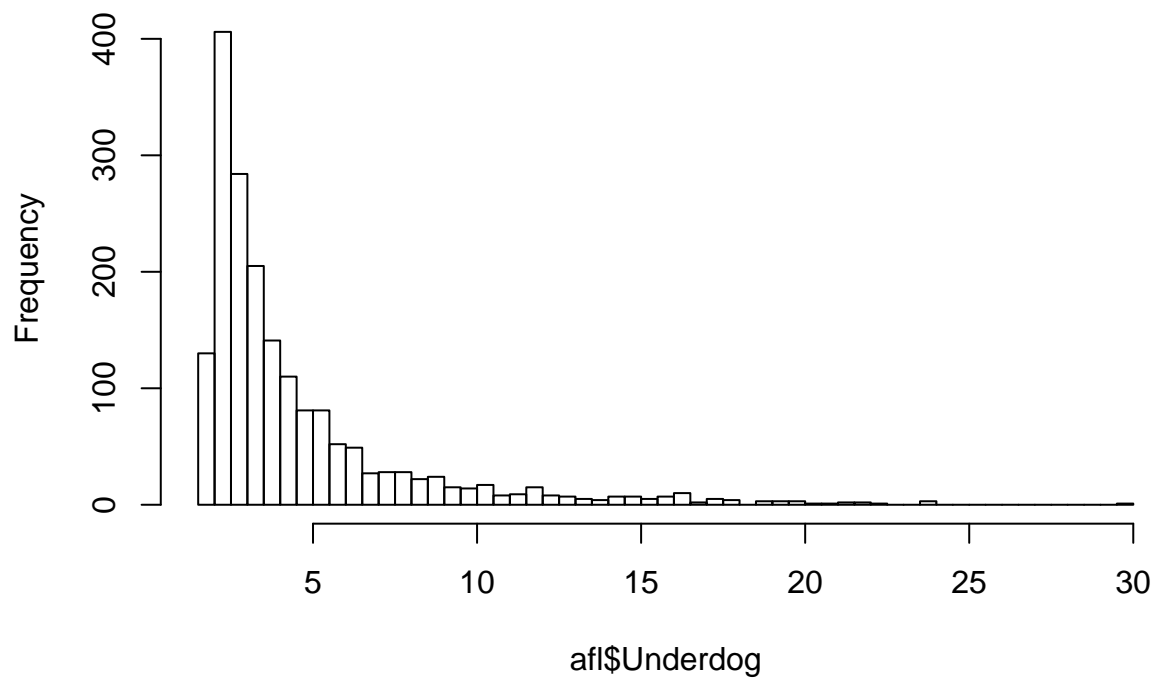
```
afl$logFavorite = log(afl$Favorite)  
hist(afl$logFavorite)
```

**Histogram of afl\$logFavorite**



```
hist(afl$logUnderdog, breaks = 40)
```

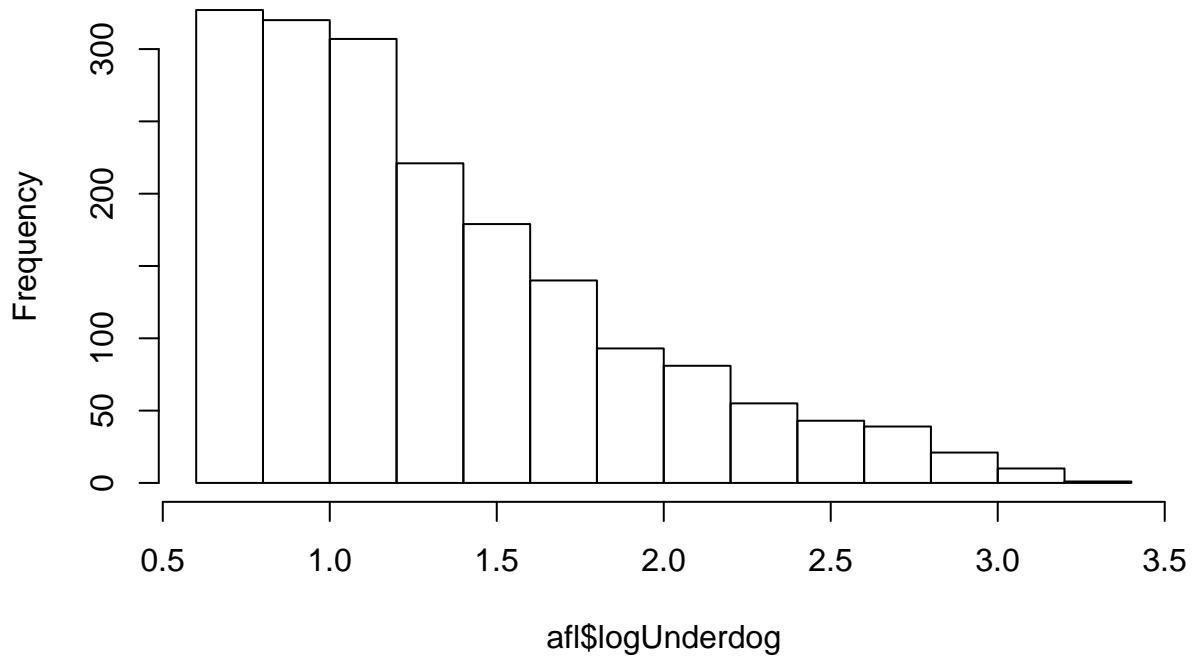
**Histogram of afl\$Underdog**



```
afl$logUnderdog = log(afl$Underdog)  
hist(afl$logUnderdog)
```



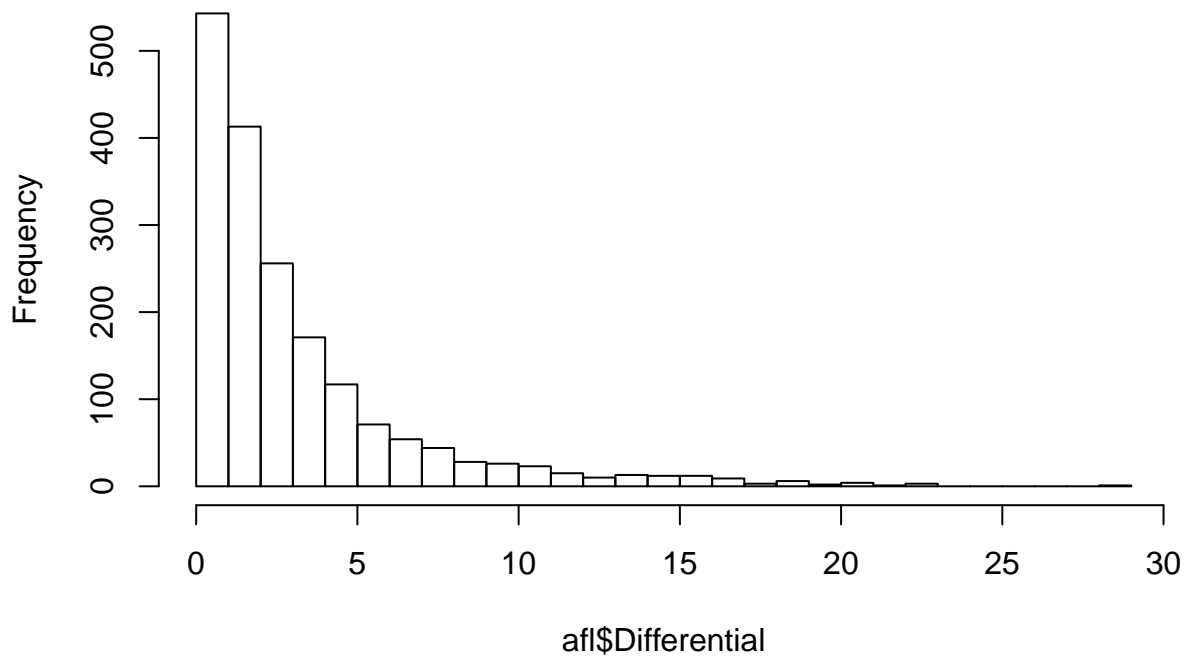
## Histogram of afl\$logUnderdog



The favorite odds and the underdog odds also aren't great to work with as they have huge amounts of right skewness as well. These variables will also require further transformation if we are to use them.

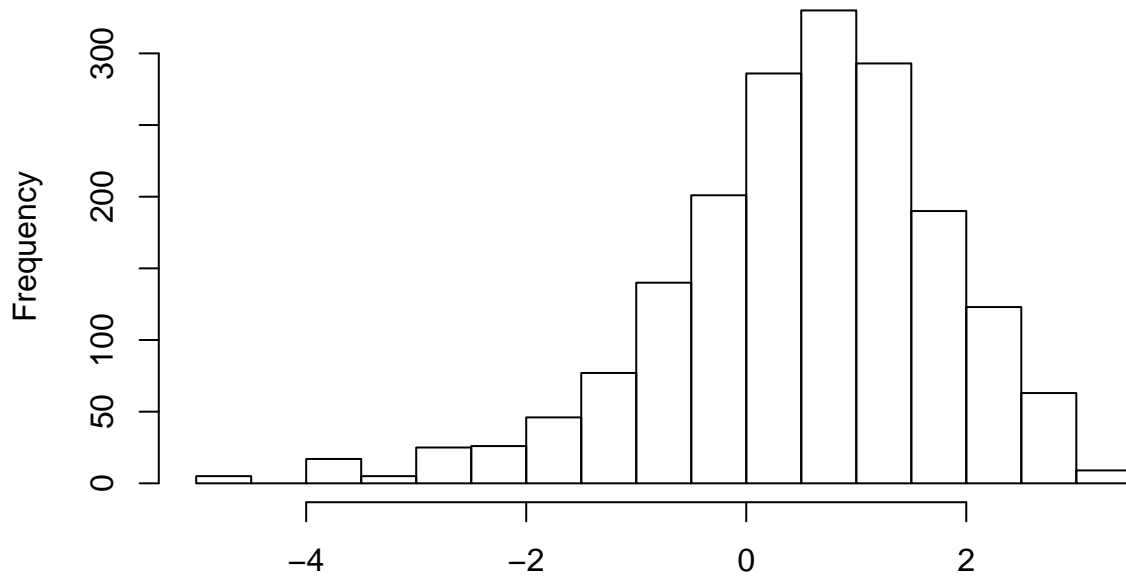
```
hist(afl$Differential, breaks = 25)
```

## Histogram of afl\$Differential



```
afl$logDifferential = log(afl$Differential)
hist(afl$logDifferential)
```

## Histogram of afl\$logDifferential

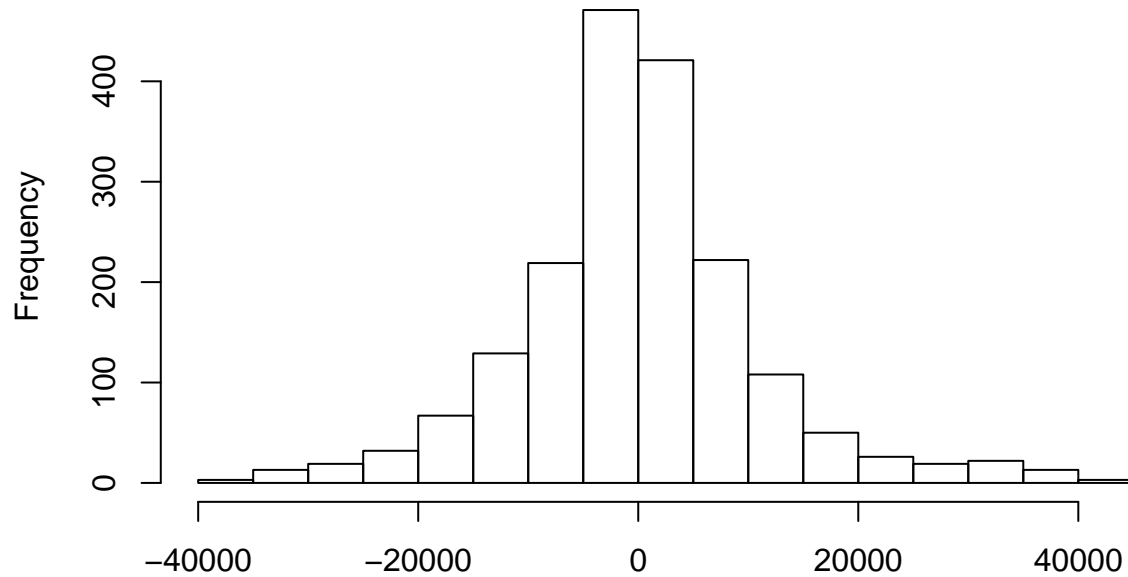


afl\$logDifferential

This is my personal favorite of the odds variables. It shows the difference between the odds of the favorite and the odds of the underdog. It also has a severe right skew but the log transformation works to perfection and creates an almost perfectly normal histogram (with slight left skew). I believe this variable provides the most pertinent information regarding how far apart the teams are in chances of winning the game. It could potentially be used in conjunction with home odds to create a model as the two may be able to provide different information.

```
hist(afl$Attendance_Differential)
```

## Histogram of afl\$Attendance\_Differential



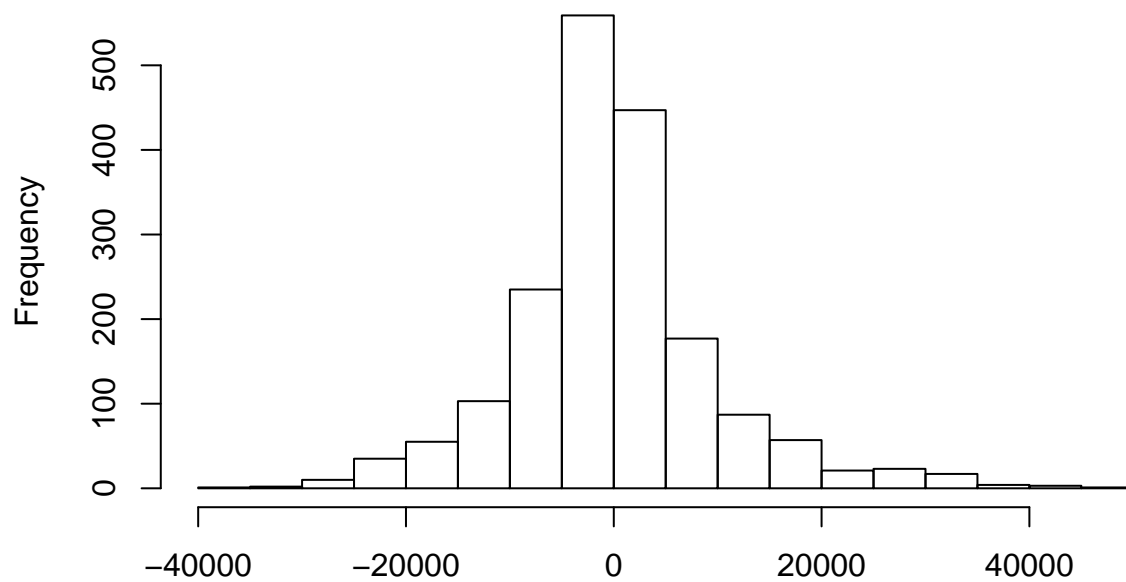
afl\$Attendance\_Differential

Atten-

dance differential shows the difference in attendance between the actual versus the average attendance (for that venue) for each game. The data are relatively normal although the tails are slightly long.

```
hist(afl$HomeDifferential, breaks = 14)
```

## Histogram of afl\$HomeDifferential



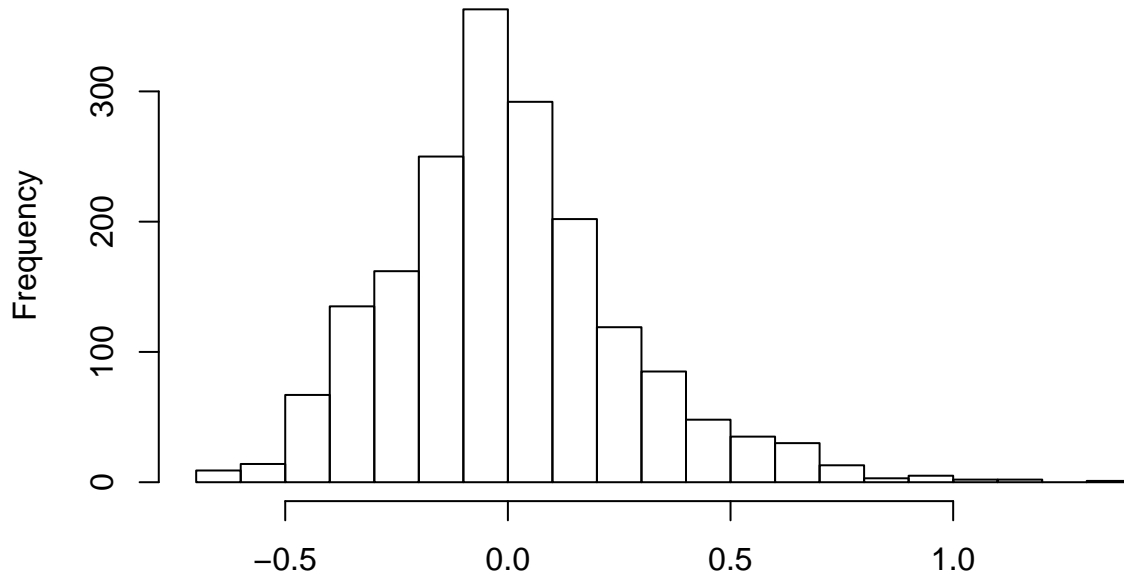
afl\$HomeDifferential

This is the

same as differential but it shows the difference between the actual versus average attendance at the venue but only for the home team's games (teases out variance due to home team).

```
hist(afl$HomeDifferentialPct, breaks = 20)
```

## Histogram of afl\$HomeDifferentialPct



afl\$HomeDifferentialPct

This is my

favorite attendance statistic because it shows the percentage difference between the team's average home attendance at a specific venue and the actual attendance for every individual game. This gives better information than strictly numerical values because it scales the change based on how many people generally attend that team's home games at that stadium.

```
fit <- lm(HomeDifferentialPct ~ Day + TimeBucket + Home.Odds + Differential, data = afl)
library(car)
```

```
## Loading required package: carData
```

```
vif(fit)
```

```
##              GVIF Df GVIF^(1/(2*Df))
## Day          1.404252 6      1.028696
## TimeBucket   1.412345 1      1.188421
## Home.Odds     1.136975 1      1.066290
## Differential  1.162453 1      1.078171
```

```
fit1 <- step(fit)
```

```
## Start:  AIC=-5204.66
## HomeDifferentialPct ~ Day + TimeBucket + Home.Odds + Differential
##
##              Df Sum of Sq  RSS    AIC
## - Home.Odds    1    0.0198 106.91 -5206.3
## <none>          0    106.89 -5204.7
## - TimeBucket    1    0.8460 107.73 -5192.2
## - Differential    1    2.9537 109.84 -5156.6
## - Day            6   12.4280 119.31 -5014.6
##
```

```
## Step: AIC=-5206.32
## HomeDifferentialPct ~ Day + TimeBucket + Differential
##
##           Df Sum of Sq    RSS    AIC
## <none>                106.91 -5206.3
## - TimeBucket      1     0.8330 107.74 -5194.1
## - Differential    1     3.5324 110.44 -5148.6
## - Day             6    12.4082 119.31 -5016.6
```

```
summary(fit1)
```

```
##
## Call:
## lm(formula = HomeDifferentialPct ~ Day + TimeBucket + Differential,
##     data = afl)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.81804 -0.15519 -0.01185  0.13456  1.01578
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   0.104799   0.025188   4.161 3.32e-05 ***
## DayMon         0.374247   0.049216   7.604 4.56e-14 ***
## DaySat        -0.119941   0.019547  -6.136 1.04e-09 ***
## DaySun        -0.164349   0.022196  -7.404 2.00e-13 ***
## DayThu         0.004777   0.039602   0.121 0.904005
## DayTue         0.588658   0.140825   4.180 3.05e-05 ***
## DayWed         0.436806   0.172229   2.536 0.011289 *
## TimeBucket     0.010092   0.002674   3.774 0.000166 ***
## Differential -0.012150   0.001563  -7.772 1.28e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2418 on 1828 degrees of freedom
## Multiple R-squared:  0.1756, Adjusted R-squared:  0.172
## F-statistic: 48.68 on 8 and 1828 DF,  p-value: < 2.2e-16
```

```
fit <- lm(HomeDifferentialPct ~ Day + TimeBucket + Differential, data = afl)
summary(fit)
```

```
##
## Call:
## lm(formula = HomeDifferentialPct ~ Day + TimeBucket + Differential,
##     data = afl)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.81804 -0.15519 -0.01185  0.13456  1.01578
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   0.104799   0.025188   4.161 3.32e-05 ***
## DayMon         0.374247   0.049216   7.604 4.56e-14 ***
## DaySat        -0.119941   0.019547  -6.136 1.04e-09 ***
```

```

## DaySun      -0.164349    0.022196   -7.404 2.00e-13 ***
## DayThu       0.004777    0.039602    0.121 0.904005
## DayTue       0.588658    0.140825    4.180 3.05e-05 ***
## DayWed       0.436806    0.172229    2.536 0.011289 *
## TimeBucket   0.010092    0.002674    3.774 0.000166 ***
## Differential -0.012150    0.001563   -7.772 1.28e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2418 on 1828 degrees of freedom
## Multiple R-squared:  0.1756, Adjusted R-squared:  0.172
## F-statistic: 48.68 on 8 and 1828 DF,  p-value: < 2.2e-16

```