The histogram distribution of WLR (Win Loss Rate), our response variable, is roughly normally distributed with a mean WLR of .51, standard deviation of .097, a minimum value of .248 and a maximum value of .705.

The scatterplot matrix below shows the correlation between the variables. Using the VIF test we eliminated multicollinear variables.





# Data Exploration

From this point forward, we are making changes to 75% of our data (Training Data)

We used the stepwise procedure of variable selection to create our first model (baseball_1) that had only one variable, RunDiff.

To counter any possible nonlinear-pattern, we made a second model (baseball_2) that was identical to baseball_1, but with the inclusion of another predictor variable, RunDiff^2.
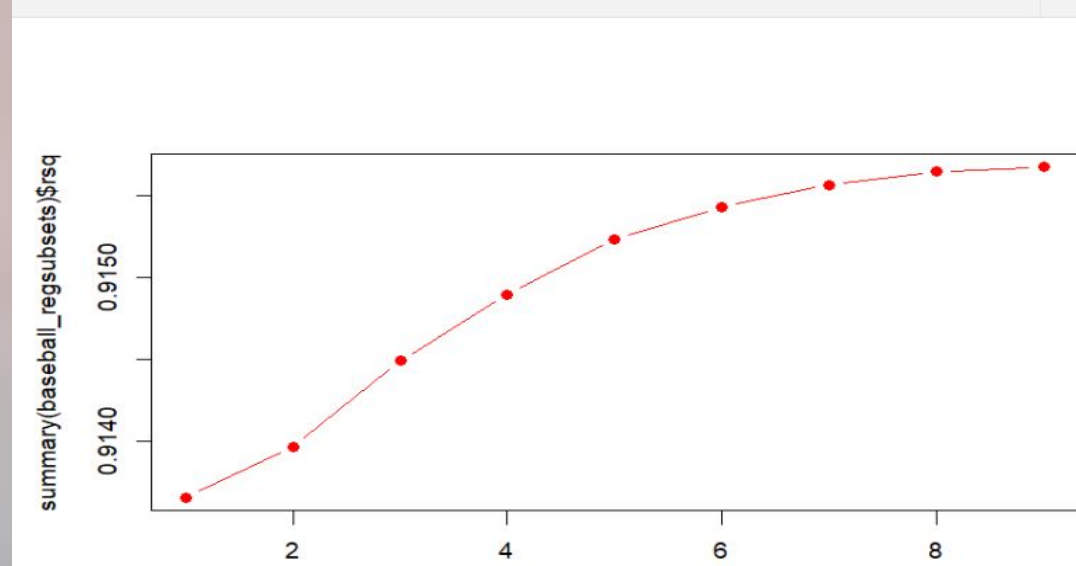
baseball_1:
Predicted WLR = .06181 + .4317(RunDiff)

baseball_2:
Predicted WLR = −.1082+.76741(RunDiff)−.159(RunDiff)^2

# Creating Models

# Regsubset



```{r}
plot(1:9,summary(baseball_regsubsets)$rsq,
    type="b",
    col="red",
    pch=16)
```

```{r}
plot(1:9,summary(baseball_regsubsets)$adjr2,
    type="b",
    col="blue",
    pch=16)
```
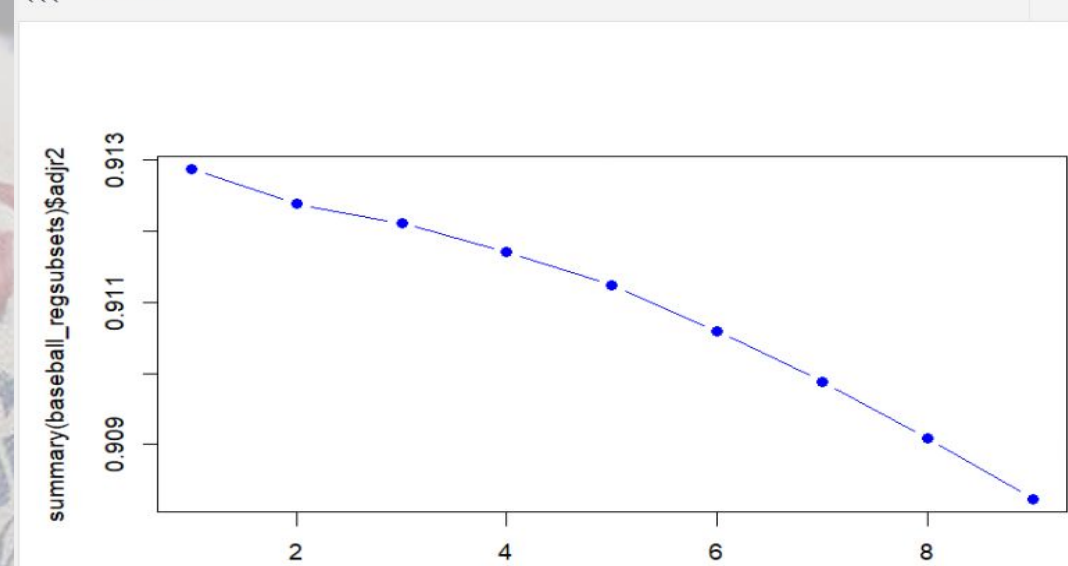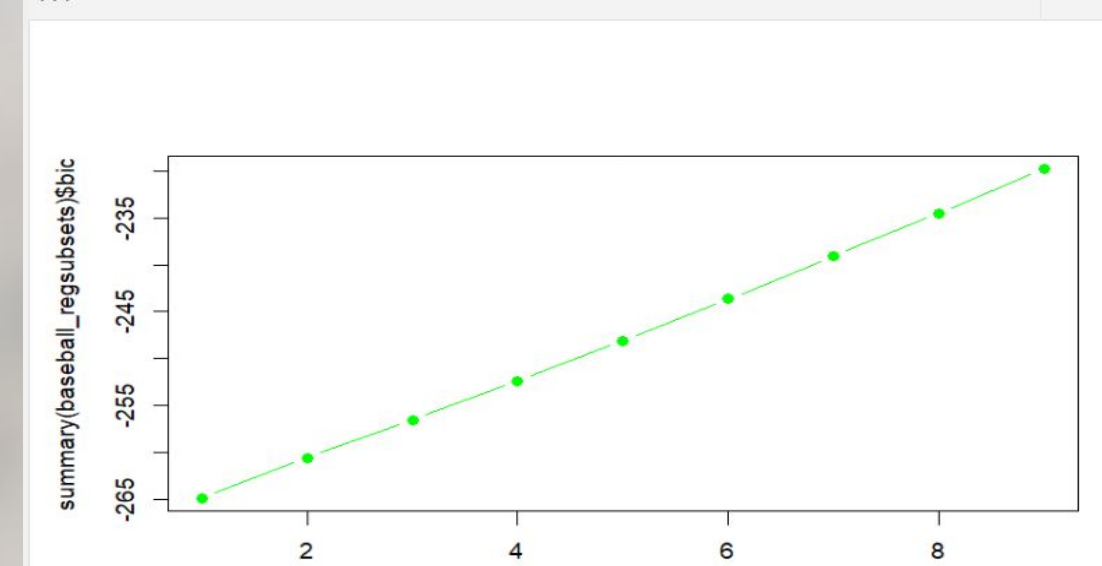
```{r}
plot(1:9,summary(baseball_regsubsets)$bic,
    type="b",
    col="green",
    pch=16)
```

Looking at the $R^2$ values, the model indicates it would be best to use all 9 predictor variables.

Looking at the Adjusted $R^2$ values, the model indicates it would be best to use only 1 predictor variable, which is RunDiff. .

Looking at the BIC values, the model indicates it would be best to use only 1 predictor variable, which is RunDiff.
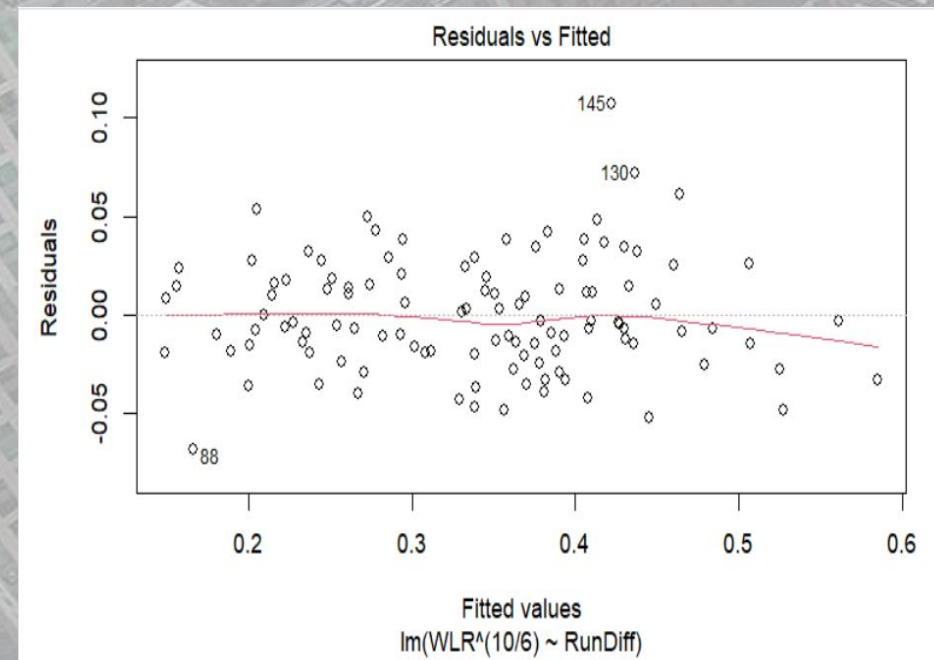
# Creating Models

We have 3 models to consider using this point forward: baseball_1, baseball_2, and baseball_3

Baseball_3 is the most complex and has the lowest Adjusted R^2, so we will not use it.

Baseball_2 is slightly more complex than baseball_1, and has a slightly better Adjusted R^2. It isn't clear which of these two models (baseball_1 and baseball_2) may be better, so we will continue using both for the time being.

# Model Consensus
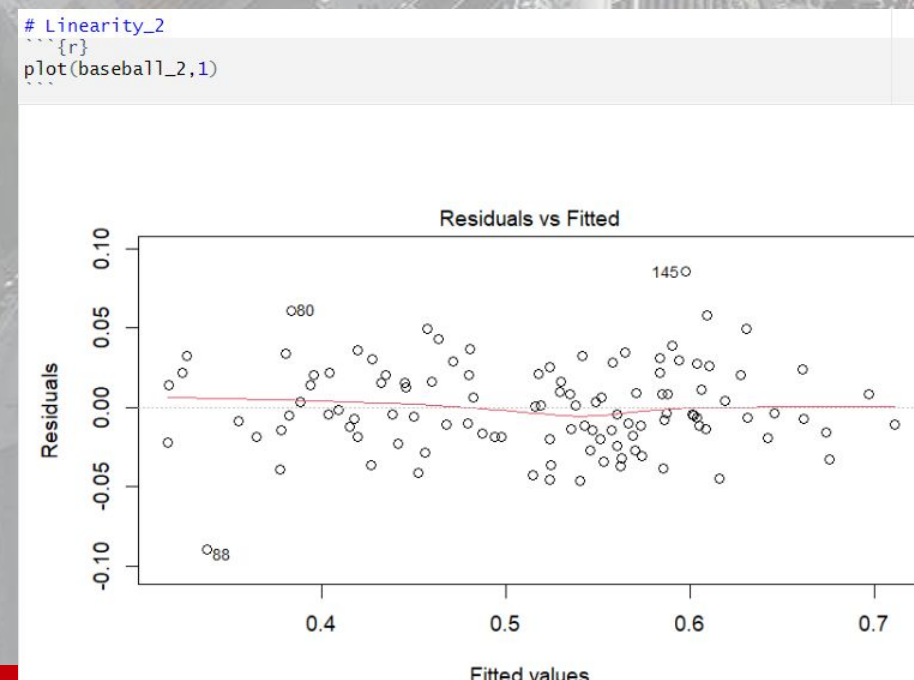
# Model Conditions

## Baseball 1



baseball_1 had violated the linearity condition, so it was refit using box-cox, then met linearity and constant variance. It had a high shapiro-wilk (.13) and Durbin Watson (.3)p-value, satisfying the normality and independence conditions. 91.91% of the variation in WLR is explained in baseball_1.

## Baseball 2



baseball_2 satisfied the linearity and constant variance conditions. It also had a high shapiro-wilk (.4)and Durbin Watson (.32) p-value, satisfying the normality and independence conditions. 92.16% of the variation in WLR is explained in baseball_2.

Both Models were tested and found to be significant.

# Mean

# Square

# Error

Comparing the mean square error between both models calculated against the 25% of testing data (0.0087 for baseball_2 vs 0.1712 for baseball_1), we concluded that model two (baseball_2) is the best model to fit our data to in order to predict Win-Loss Rate for the Atlanta Braves.

$$E(y) = -.1082 + .7674(RunDiff) - .159(RunDiff)^2$$

# Conclusions