# Lab #6: Introduction to Linear Regression

Name

Date of lab session

---

## Lab report
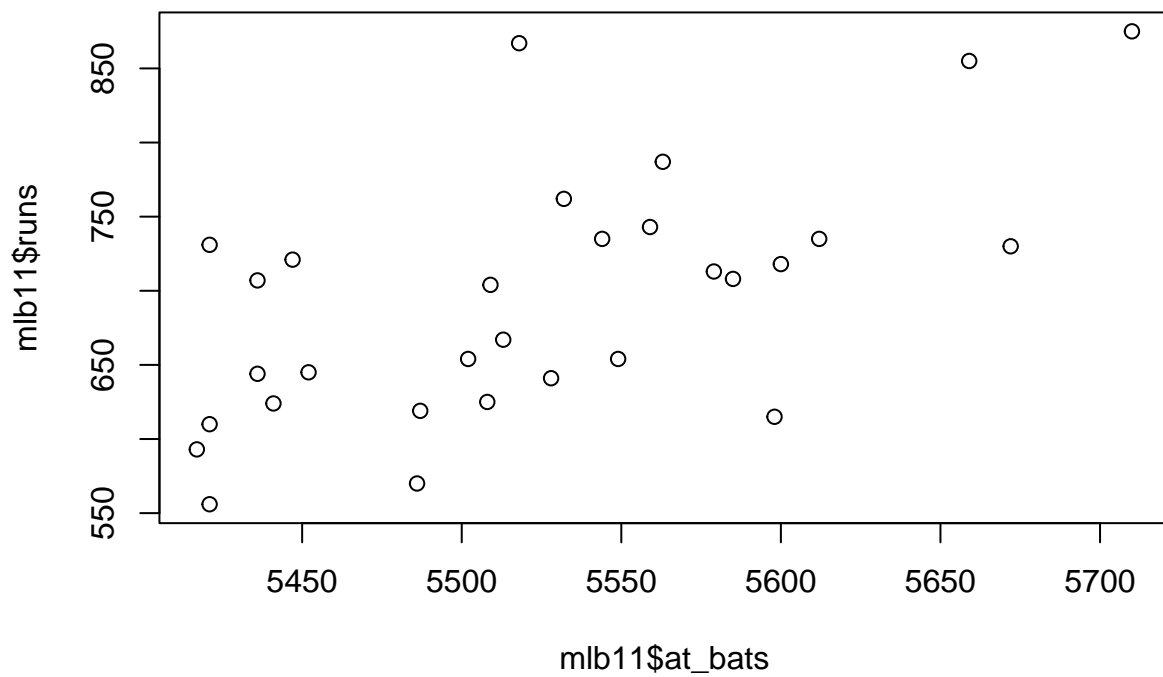
```
download.file("http://www.openintro.org/stat/data/mlb11.RData", destfile = "mlb11.RData")
load("mlb11.RData")
```
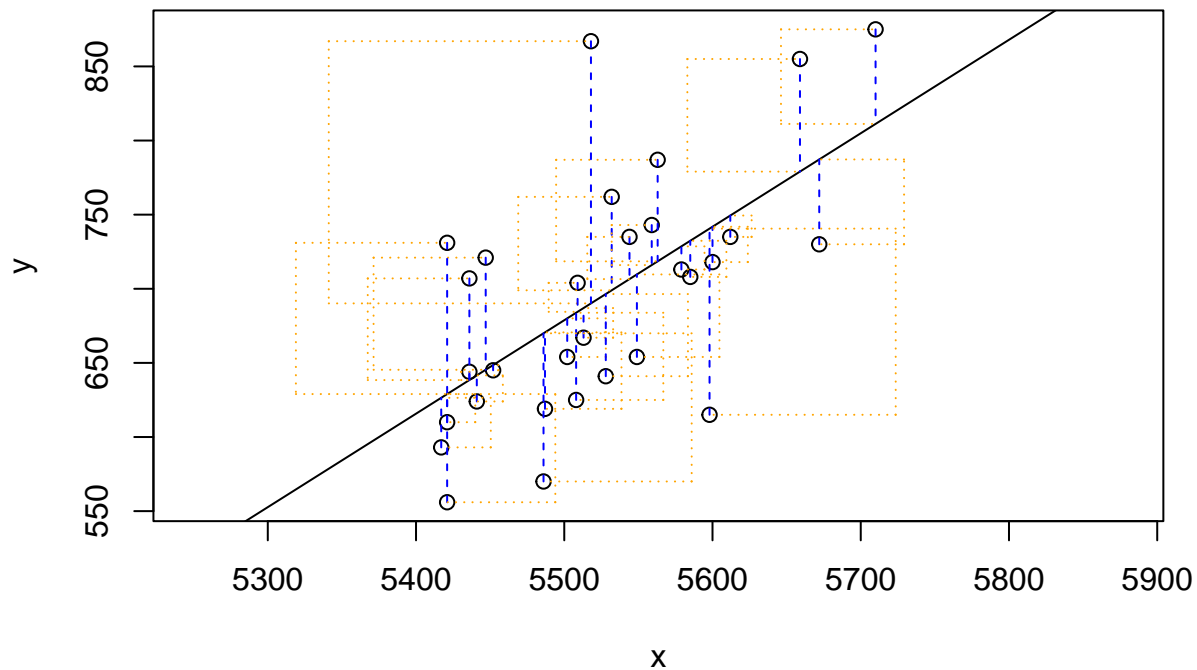
**Load data:**

**Exercises:**

**Exercise 1:**  I would use a linear regression plot. The relationship between at_bats and runs seems linear. Although the model seems to be linear, the association does not seem to be strong, so the prediction might not be too accurate.

```
plot(mlb11$at_bats,mlb11$runs)
```

**Exercise 2:** The two variables seem to have a positive weak linear association. There are also a few outliers around 5600 at bats and 5520 at bats.

```
plot_ss(x = mlb11$at_bats, y = mlb11$runs, showSquares = TRUE)
```

```
## Click two points to make a line.
## Call:
## lm(formula = y ~ x, data = pts)
##
## Coefficients:
## (Intercept)              x
##  -2789.2429         0.6305
##
## Sum of Squares:  123721.9
```

**Exercise 3:** The smallest sum of squares I got was 132859.5, some of them with an opposite slope had sums 100s of times larger

**Exercise 4:** y = 415.2389 + 1.8345*homeruns For every homerun, there are are 1.8345 runs.
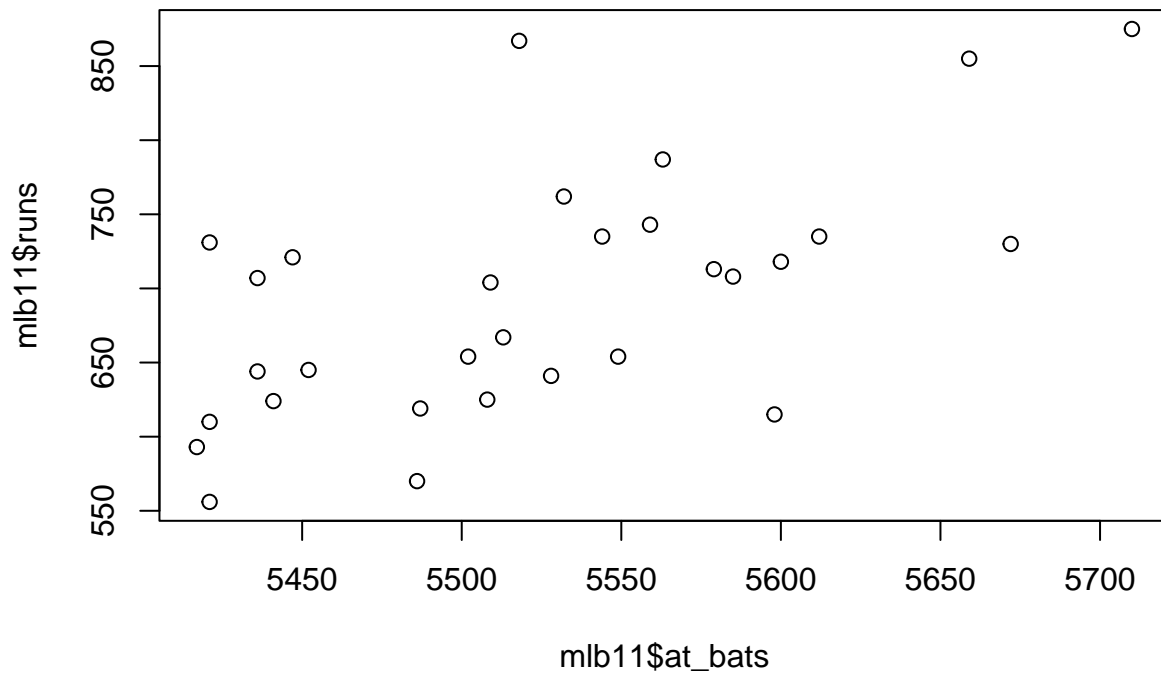
```
m1 <- lm(runs ~ homeruns, data = mlb11)
summary(m1)
```

```
##
## Call:
## lm(formula = runs ~ homeruns, data = mlb11)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
```

```
## -91.615 -33.410    3.231   24.292 104.631
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 415.2389    41.6779   9.963 1.04e-10 ***
## homeruns      1.8345     0.2677   6.854 1.90e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 51.29 on 28 degrees of freedom
## Multiple R-squared:  0.6266, Adjusted R-squared:  0.6132
## F-statistic: 46.98 on 1 and 28 DF,  p-value: 1.9e-07
```

**Exercise 5:** Predicted runs for 5579 at bats is 728.3166 runs. This is a overestimate as the data point is at 713, so an overestimate by 15.3166. Which means the residual is also 15.3166.

```
plot(mlb11$runs ~ mlb11$at_bats)
abline(m1)
```
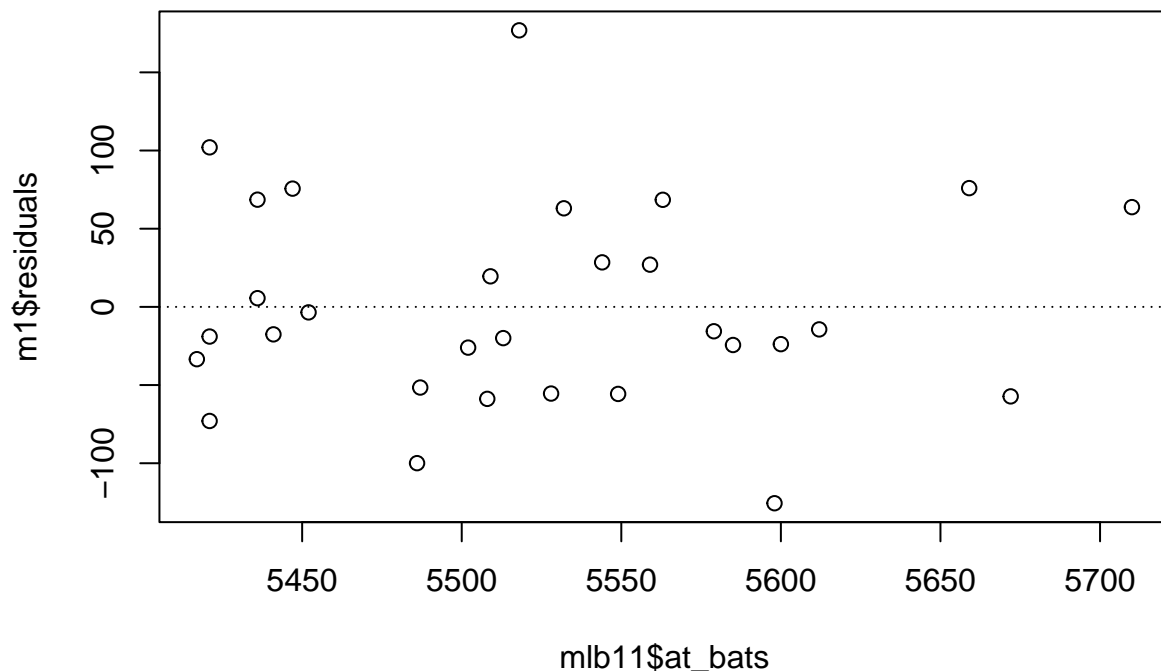


```
m1 <- lm(runs ~ at_bats, data = mlb11)
summary(m1)
```

```
##
## Call:
```

```
## lm(formula = runs ~ at_bats, data = mlb11)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -125.58  -47.05  -16.59   54.40  176.87
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept) -2789.2429   853.6957  -3.267 0.002871 **
## at_bats         0.6305     0.1545   4.080 0.000339 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 66.47 on 28 degrees of freedom
## Multiple R-squared:  0.3729, Adjusted R-squared:  0.3505
## F-statistic: 16.65 on 1 and 28 DF,  p-value: 0.0003388
```
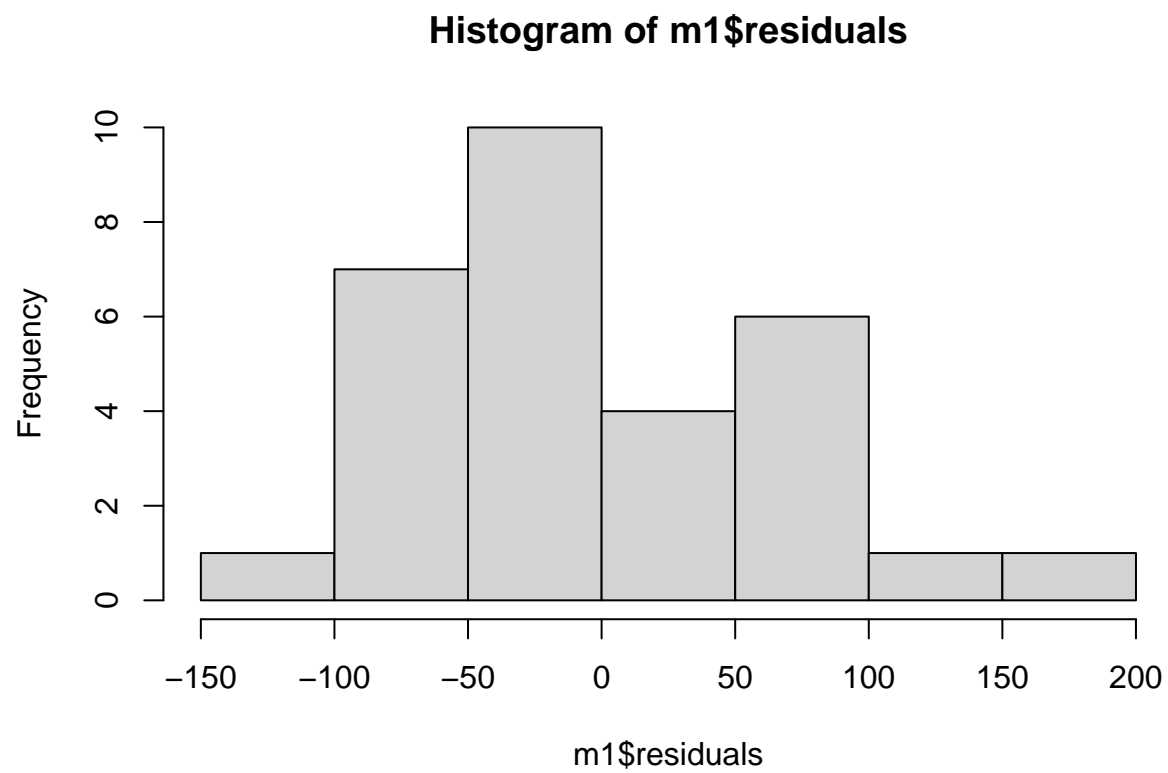
**Exercise 6:** Since there seems to be no apparent pattern, we can say that there is linearity.

```
plot(m1$residuals ~ mlb11$at_bats)
abline(h = 0, lty = 3)  # adds a horizontal dashed line at y = 0
```
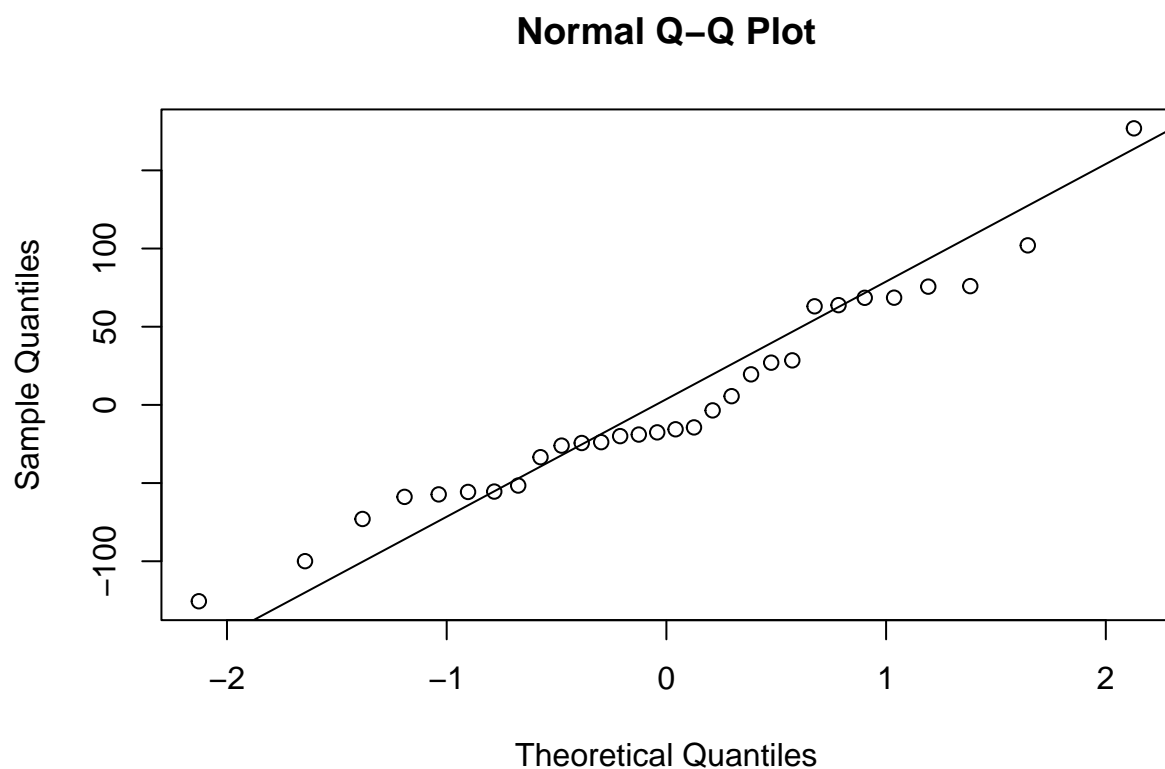


**Exercise 7:** Based on these two plots, it does not seem to be following a normal distribution.

```
hist(m1$residuals)
```

## Histogram of m1$residuals



```
qqnorm(m1$residuals)
qqline(m1$residuals)
```
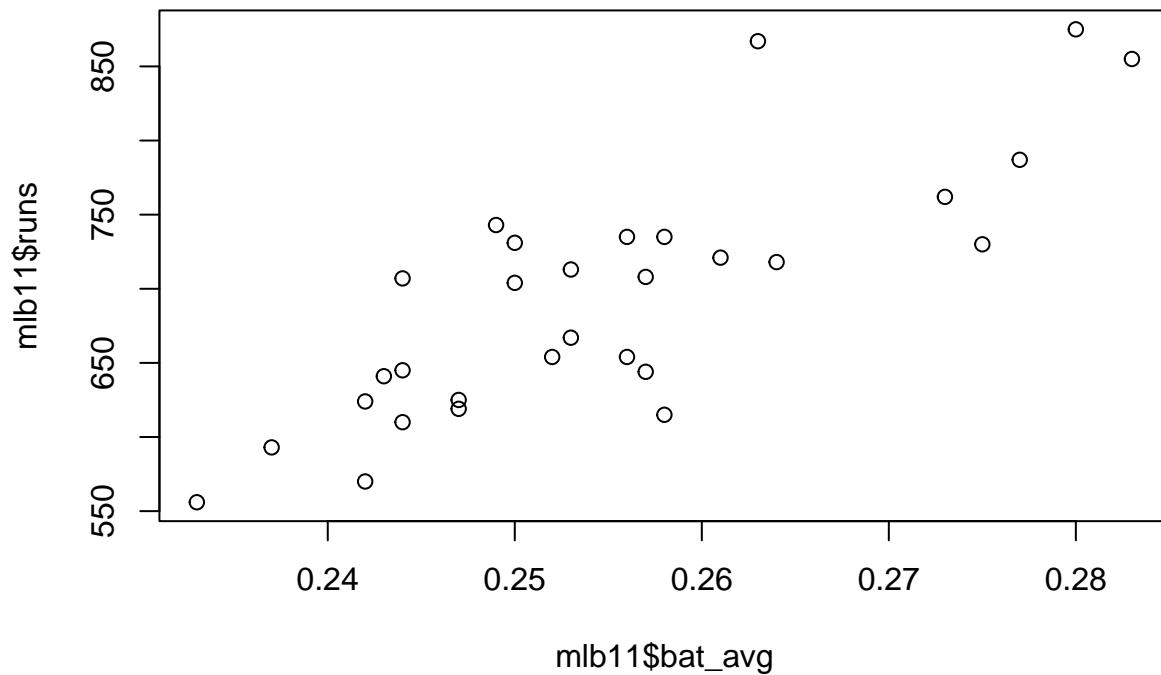
## Normal Q–Q Plot



**Exercise 8:**   Since there is no apparent pattern in the residual plot, we can see the variability condition is met.

---

## On your own:

**1:**   There seems to be a linear association.

```
plot(mlb11$runs ~ mlb11$bat_avg)
```

**2:** Based on the two R Squareds it seems that batting average is a better fit for the data than at bats as there is a higher R Squared value.

```
m1 <- lm(runs ~ at_bats, data = mlb11)
summary(m1)
```

```
##
## Call:
## lm(formula = runs ~ at_bats, data = mlb11)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -125.58  -47.05  -16.59   54.40  176.87
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -2789.2429   853.6957  -3.267 0.002871 **
## at_bats         0.6305     0.1545   4.080 0.000339 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 66.47 on 28 degrees of freedom
## Multiple R-squared:  0.3729, Adjusted R-squared:  0.3505
## F-statistic: 16.65 on 1 and 28 DF,  p-value: 0.0003388
```

```
m2 <- lm(runs ~ bat_avg, data = mlb11)
summary(m2)
```
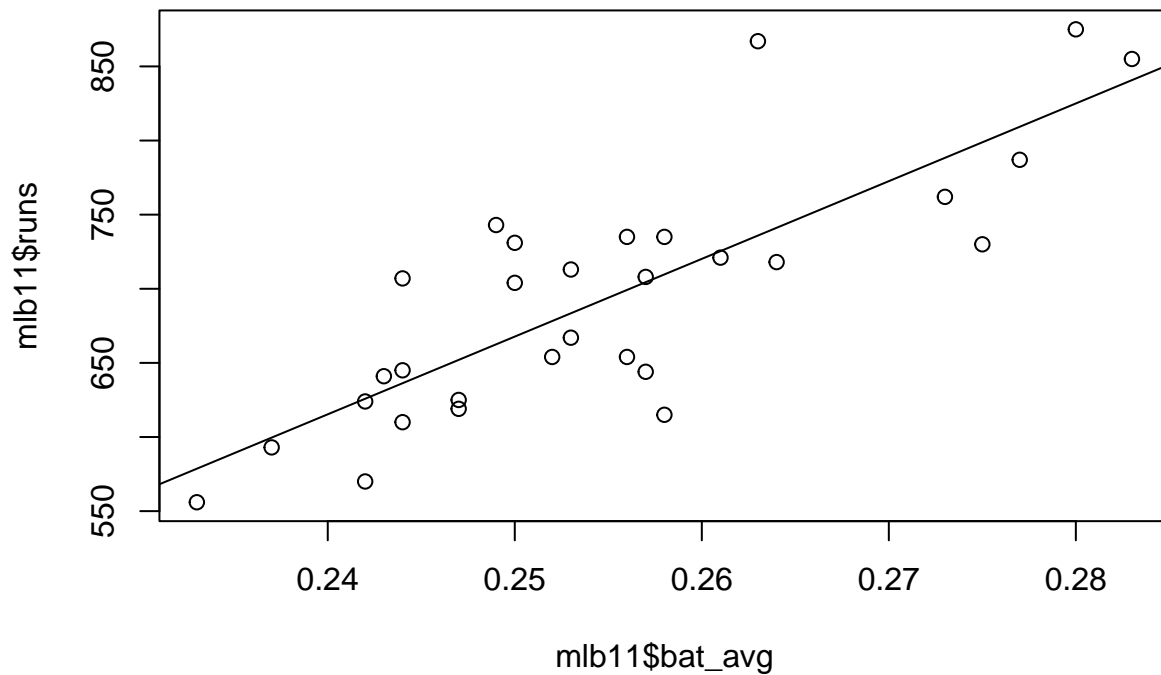
```
##
## Call:
## lm(formula = runs ~ bat_avg, data = mlb11)
##
## Residuals:
##     Min     1Q  Median     3Q    Max
## -94.676 -26.303  -5.496  28.482 131.113
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -642.8      183.1  -3.511  0.00153 **
## bat_avg       5242.2      717.3   7.308 5.88e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 49.23 on 28 degrees of freedom
## Multiple R-squared:  0.6561, Adjusted R-squared:  0.6438
## F-statistic: 53.41 on 1 and 28 DF,  p-value: 5.877e-08
```

**3:**  Batting average seems to be the best predictor for runs as it has the highest R squared at 0.64.

```
m2 <- lm(runs ~ bat_avg, data = mlb11)
summary(m2)
```

```
##
## Call:
## lm(formula = runs ~ bat_avg, data = mlb11)
##
## Residuals:
##     Min     1Q  Median     3Q    Max
## -94.676 -26.303  -5.496  28.482 131.113
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -642.8      183.1  -3.511  0.00153 **
## bat_avg       5242.2      717.3   7.308 5.88e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 49.23 on 28 degrees of freedom
## Multiple R-squared:  0.6561, Adjusted R-squared:  0.6438
## F-statistic: 53.41 on 1 and 28 DF,  p-value: 5.877e-08
```

```
plot(mlb11$runs ~ mlb11$bat_avg)
abline(m2)
```

**4:** It seems that on base plus slugging is the best as R squared is 0.93.

```
m2 <- lm(runs ~ new_onbase, data = mlb11)
summary(m2)
```

```
##
## Call:
## lm(formula = runs ~ new_onbase, data = mlb11)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -58.270 -18.335   3.249  19.520  69.002
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -1118.4      144.5  -7.741 1.97e-08 ***
## new_onbase    5654.3      450.5  12.552 5.12e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 32.61 on 28 degrees of freedom
## Multiple R-squared:  0.8491, Adjusted R-squared:  0.8437
## F-statistic: 157.6 on 1 and 28 DF,  p-value: 5.116e-13
```
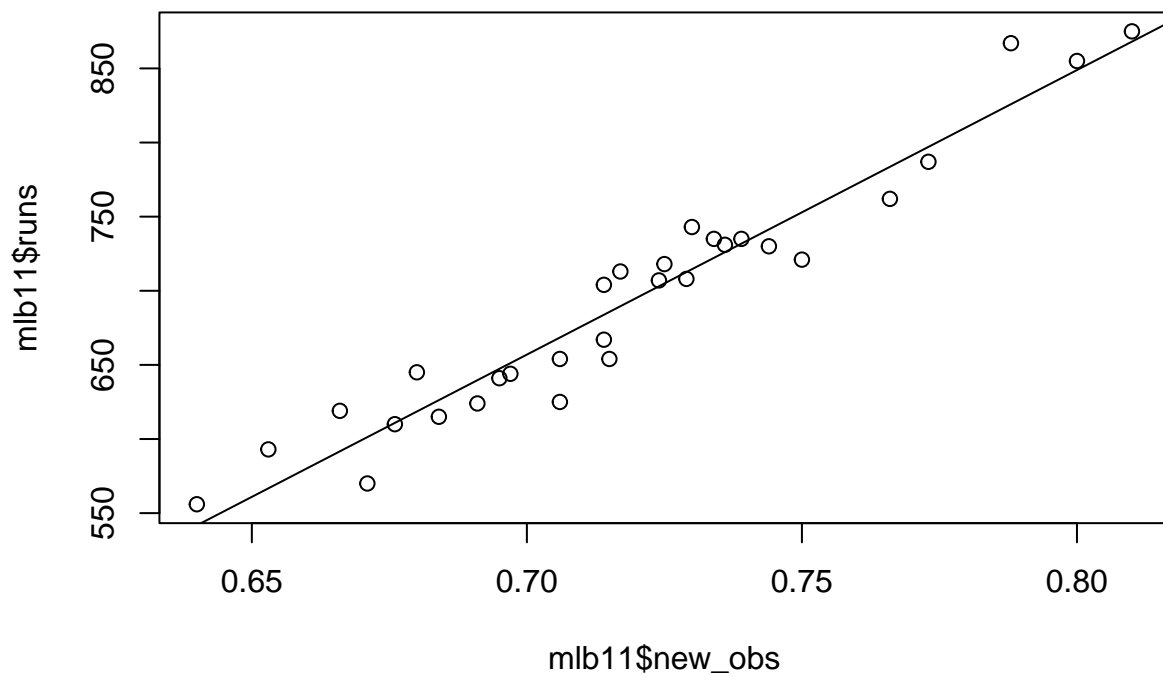
```
m3 <- lm(runs ~ new_slug, data = mlb11)
summary(m3)
```

```
##
## Call:
## lm(formula = runs ~ new_slug, data = mlb11)
##
## Residuals:
##     Min     1Q Median     3Q    Max
## -45.41 -18.66  -0.91  16.29  52.29
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -375.80      68.71   -5.47 7.70e-06 ***
## new_slug     2681.33     171.83   15.61 2.42e-15 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 26.96 on 28 degrees of freedom
## Multiple R-squared:  0.8969, Adjusted R-squared:  0.8932
## F-statistic: 243.5 on 1 and 28 DF,  p-value: 2.42e-15
```

```
m4 <- lm(runs ~ new_obs, data = mlb11)
summary(m4)
```

```
##
## Call:
## lm(formula = runs ~ new_obs, data = mlb11)
##
## Residuals:
##      Min      1Q  Median      3Q     Max
## -43.456 -13.690   1.165  13.935  41.156
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -686.61      68.93  -9.962 1.05e-10 ***
## new_obs      1919.36      95.70  20.057  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 21.41 on 28 degrees of freedom
## Multiple R-squared:  0.9349, Adjusted R-squared:  0.9326
## F-statistic: 402.3 on 1 and 28 DF,  p-value: < 2.2e-16
```
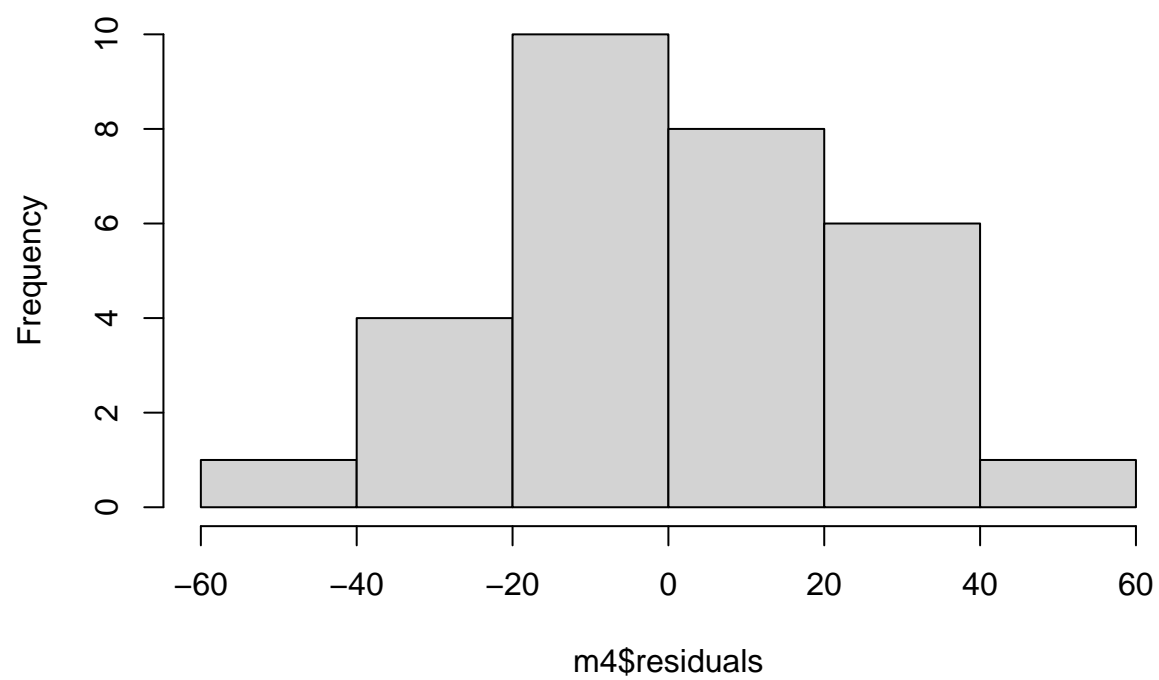
```
plot(mlb11$runs ~ mlb11$new_obs)
abline(m4)
```

**5:** Seems to be relatively normal, but there seems to be no linearity or constant variability as there is an apparent pattern in the residual graph.
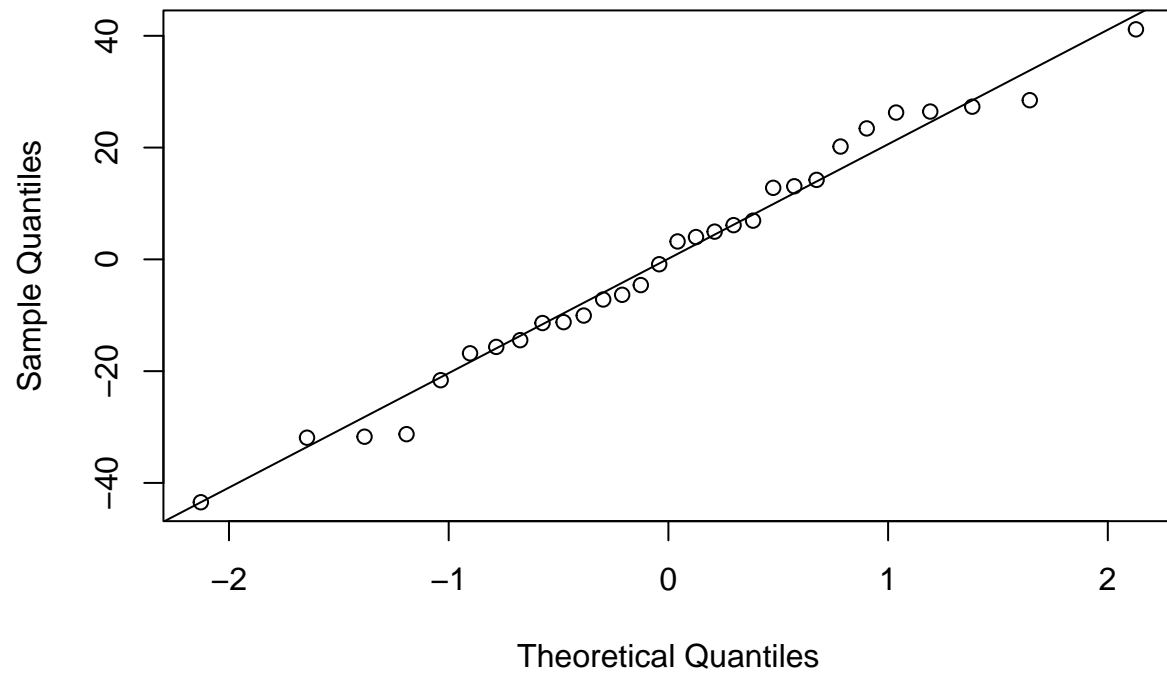
```
m4 <- lm(runs ~ new_obs, data = mlb11)
hist(m4$residuals)
```

## Histogram of m4$residuals



```r
qqnorm(m4$residuals)
qqline(m4$residuals)
```

## Normal Q–Q Plot



```r
plot(m1$residuals ~ mlb11$new_obs)
abline(h = 0, lty = 3)  # adds a horizontal dashed line at y = 0
```