

# DATA 605 - Discussion 11

Omer Ozeren

## Table of Contents

For this discussion, I will look at Kaggle's Powerlifting Database dataset.....	1
Visualize the data (EDA) .....	2
Residual Analysis.....	3

Using R, build a regression model for data that interests you. Conduct residual analysis.

Was the linear model appropriate? Why or why not?

### For this discussion, I will look at Kaggle's Powerlifting Database dataset.

- It's a dataset containing competitor results in powerlifting from the OpenPowerlifting Database and do residual analysis.

I will build a simple linear regression model of body weights vs best bench press for seniors to see if a linear relation exists between them.

- Dataset can be found here: <https://www.kaggle.com/open-powerlifting/powerlifting-database>
- Get the data and examine a preview

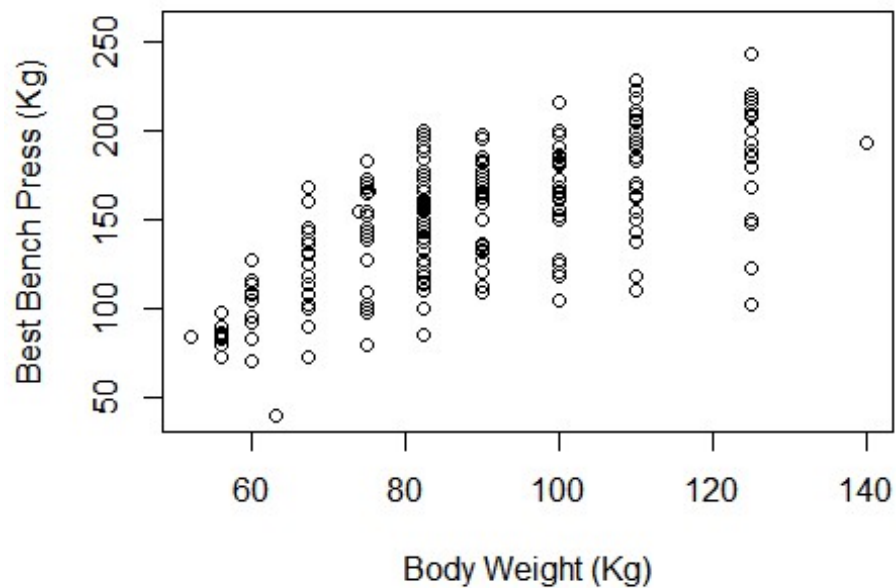
```
library(data.table)

## Warning: package 'data.table' was built under R version 3.5.3

powerlift <-
read.csv('C:\\Users\\OMERO\\Documents\\GitHub\\DATA605\\openpowerlifting.csv')
head(powerlift, n=5)

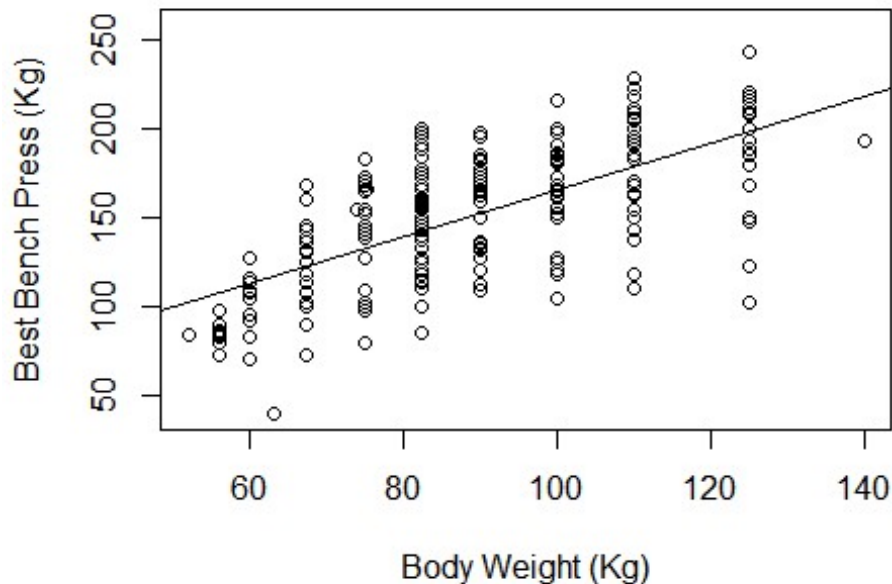
##           Name Sex Event Equipment Age AgeClass Division BodyweightKg
## 1  Abbie Murphy  F  SBD  Wraps    29   24-34   F-OR           59.8
## 2  Abbie Tuong   F  SBD  Wraps    29   24-34   F-OR           58.5
## 3 Ainslee Hooper F    B    Raw    40   40-44   F-OR           55.4
## 4 Amy Moldenhauer F  SBD  Wraps    23   20-23   F-OR           60.0
## 5  Andrea Rowan  F  SBD  Wraps    45   45-49   F-OR          104.0
##  WeightClassKg Squat1Kg Squat2Kg Squat3Kg Squat4Kg Best3SquatKg Bench1Kg
## 1             60      80      92.5      105      NA           105      45.0
## 2             60     100     110.0      120      NA           120      55.0
## 3             56      NA      NA      NA      NA           NA       27.5
## 4             60     -105    -105.0      105      NA           105      67.5
```





### Residual Analysis

```
lm_powerlift_senior <- lm(Best3BenchKg ~ BodyweightKg, data =
powerlift_senior)
with(powerlift_senior, plot(BodyweightKg, Best3BenchKg,xlab = "Body Weight
(Kg)",
                           ylab = "Best Bench Press (Kg)"))
abline(lm_powerlift_senior)
```



```
summary(lm_powerlift_senior)
```

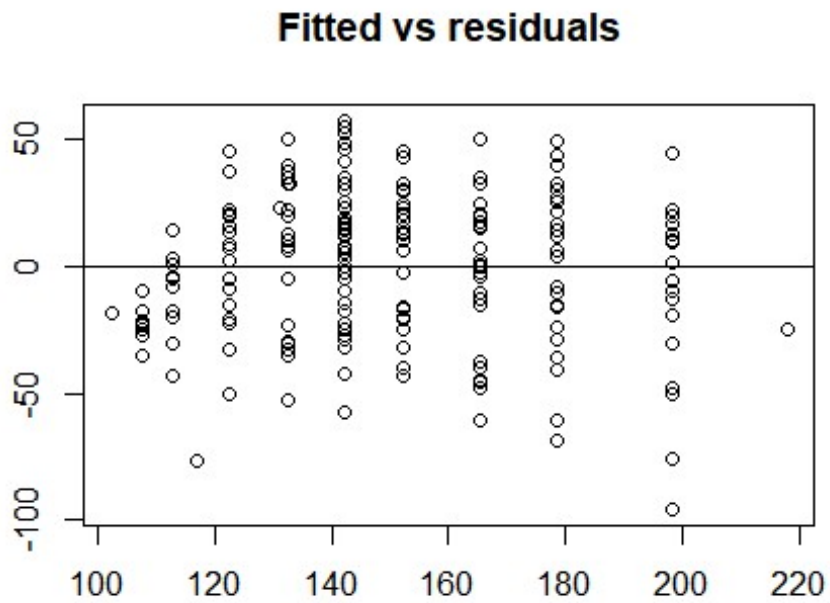
```
##
## Call:
## lm(formula = Best3BenchKg ~ BodyweightKg, data = powerlift_senior)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -95.697 -20.860   4.566  20.359  57.206
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   34.0119     9.2492   3.677 0.000297 ***
## BodyweightKg   1.3135     0.1023  12.838 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 28.77 on 218 degrees of freedom
## (12 observations deleted due to missingness)
## Multiple R-squared:  0.4305, Adjusted R-squared:  0.4279
## F-statistic: 164.8 on 1 and 218 DF,  p-value: < 2.2e-16
```

- Equation of line is

$$\text{bestbenchpress} = 0.1376 + 1.2804 * \text{bodyweight}$$

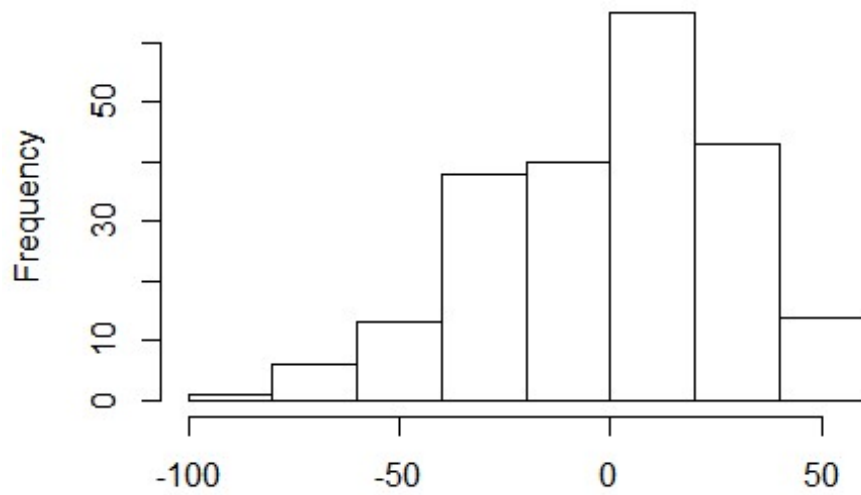
```
#### Residual plots
```

```
plot(fitted(lm_powerlift_senior), resid(lm_powerlift_senior),  
     main = "Fitted vs residuals", xlab = "", ylab = "")  
abline(h = 0)
```



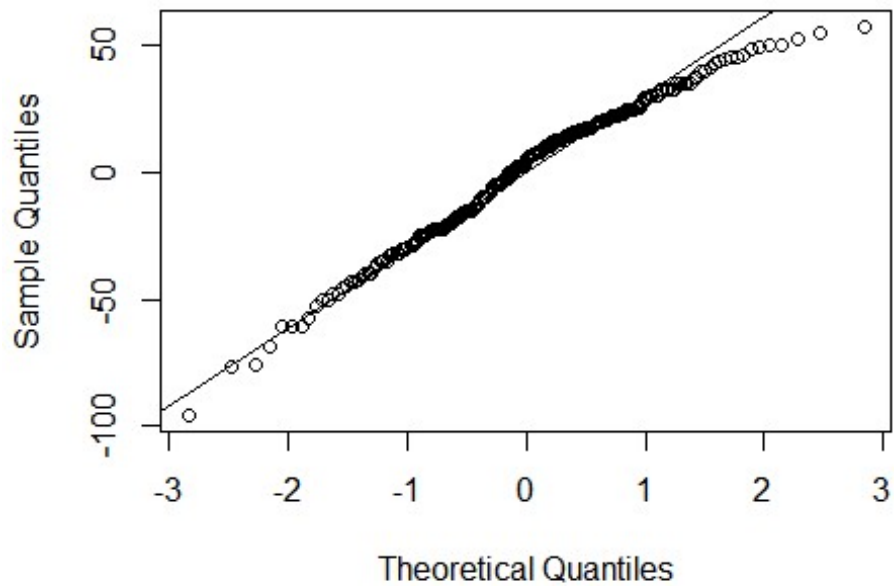
```
hist(resid(lm_powerlift_senior), xlab = "", main = "Histogram of Residuals")
```

### Histogram of Residuals



```
qqnorm(resid(lm_powerlift_senior))  
qqline(resid(lm_powerlift_senior))
```

### Normal Q-Q Plot



### Summary

- We see that a linear model based on one explanatory variable doesn't explain the data well. The  $R^2$  value is quite low which shows that the fitted model doesn't accurately predict the values of Senior divisions competitors bench press best based on their weight.