

DAS - Poster

anushapanil

28/06/2021

```
knitr::opts_chunk$set(  
  eval = TRUE,  
  echo = FALSE,  
  message = FALSE,  
  warning = FALSE  
)  
library(dplyr)
```

```
##  
## Attaching package: 'dplyr'  
  
## The following objects are masked from 'package:stats':  
##  
##   filter, lag  
  
## The following objects are masked from 'package:base':  
##  
##   intersect, setdiff, setequal, union
```

```
library(ggplot2)  
library(janitor)
```

```
##  
## Attaching package: 'janitor'  
  
## The following objects are masked from 'package:stats':  
##  
##   chisq.test, fisher.test
```

```
library(moderndiver)  
library(infer)  
library(broom)  
library(kableExtra)
```

```
##  
## Attaching package: 'kableExtra'  
  
## The following object is masked from 'package:dplyr':  
##  
##   group_rows
```

```
library(GGally)
```

```
## Registered S3 method overwritten by 'GGally':
##   method from
##   +.gg      ggplot2
```

```
library(skimr)
library(knitr)
library(gridExtra)
```

```
##
## Attaching package: 'gridExtra'
```

```
## The following object is masked from 'package:dplyr':
##
##   combine
```

```
library(readr)
library(kableExtra)
library(olsrr)
```

```
##
## Attaching package: 'olsrr'
```

```
## The following object is masked from 'package:datasets':
##
##   rivers
```

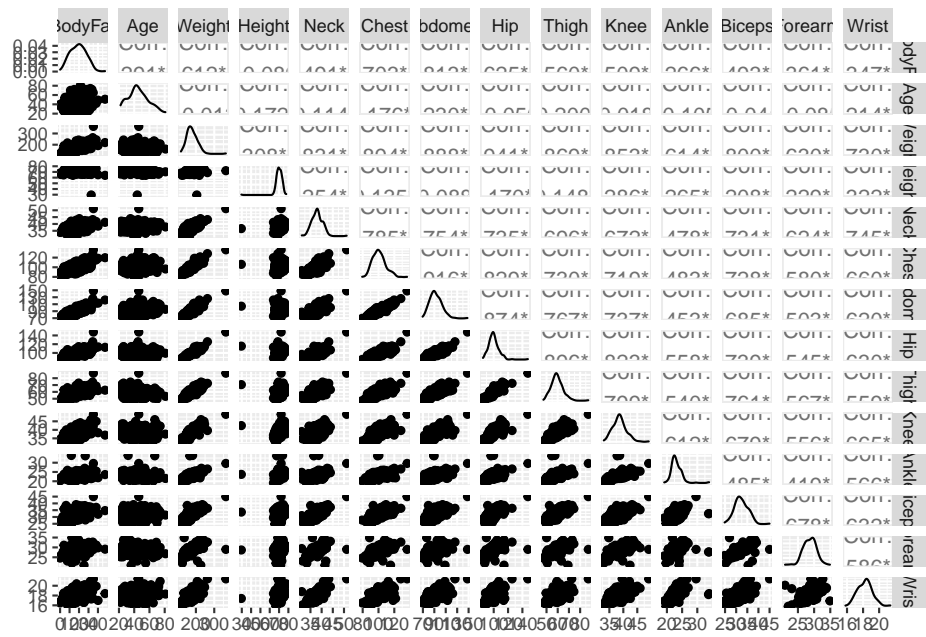


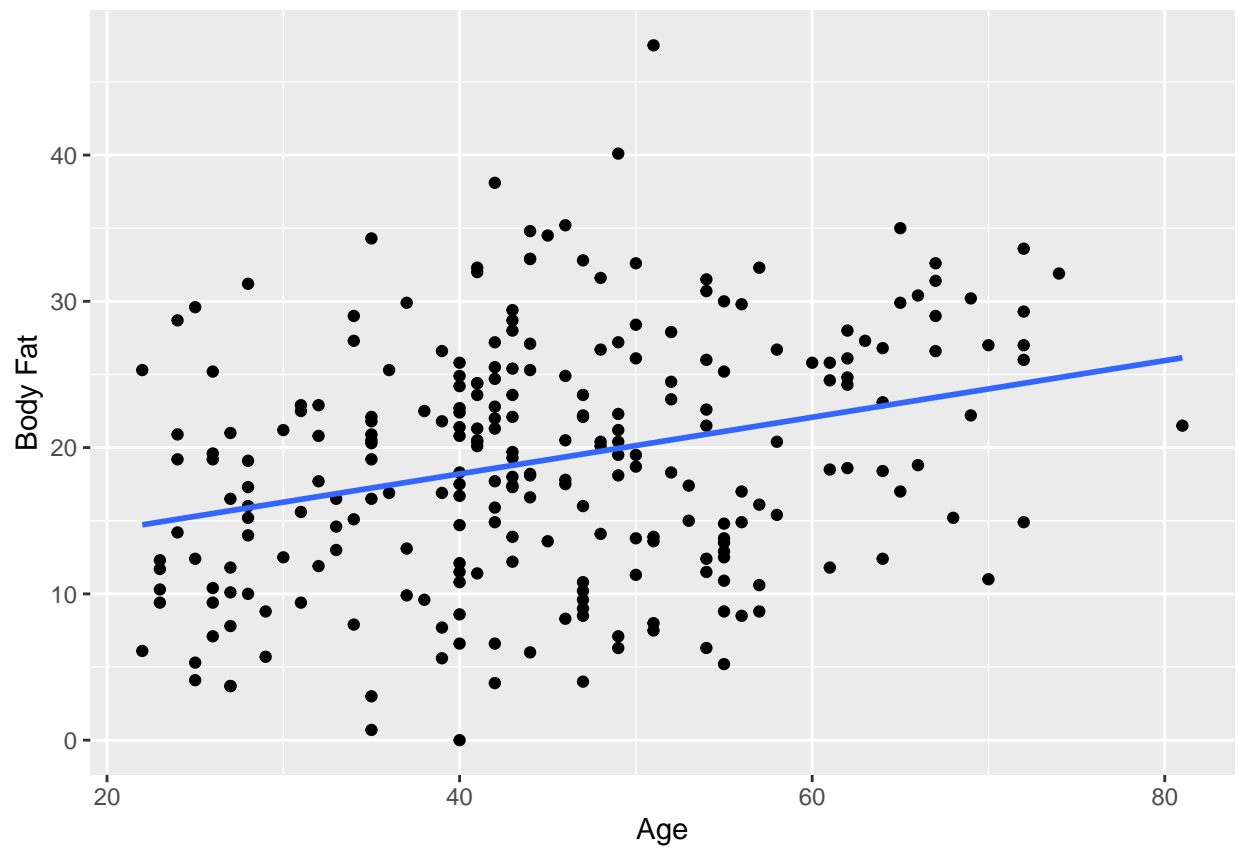
Figure 1: Correlation Plot

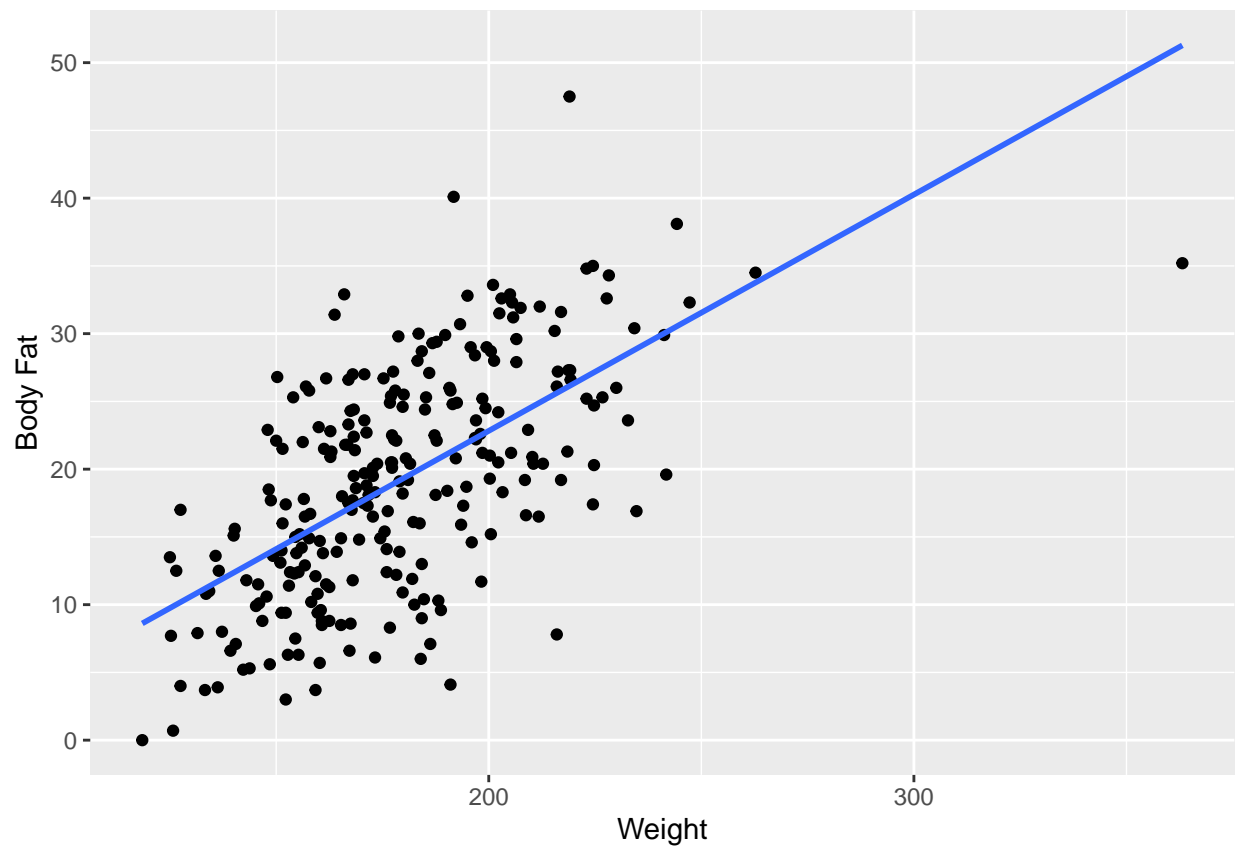
There is a strong positive correlation between all the variables, which implies that there is high multicollinearity. So, it will be better to use variable selection method to remove multicollinearity.

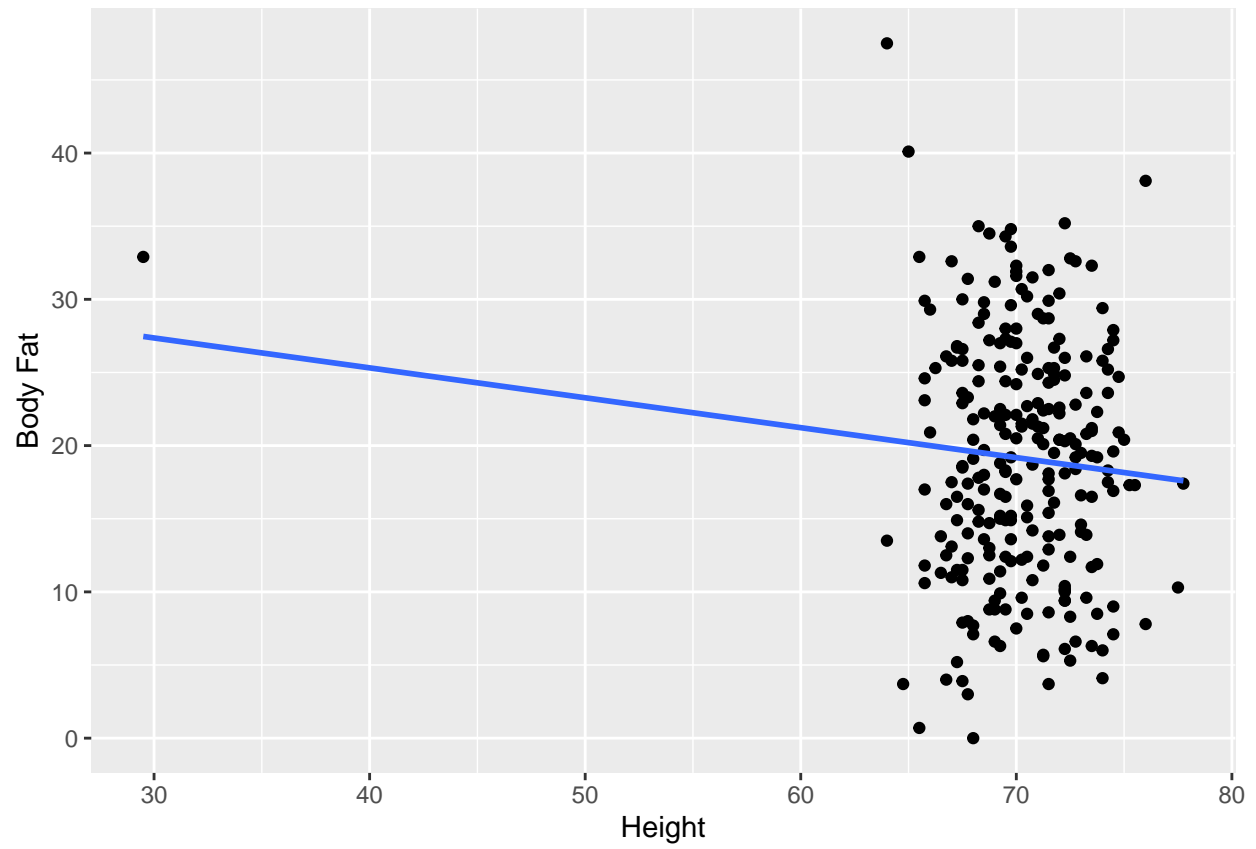
Let's take a look at the summary of our data.

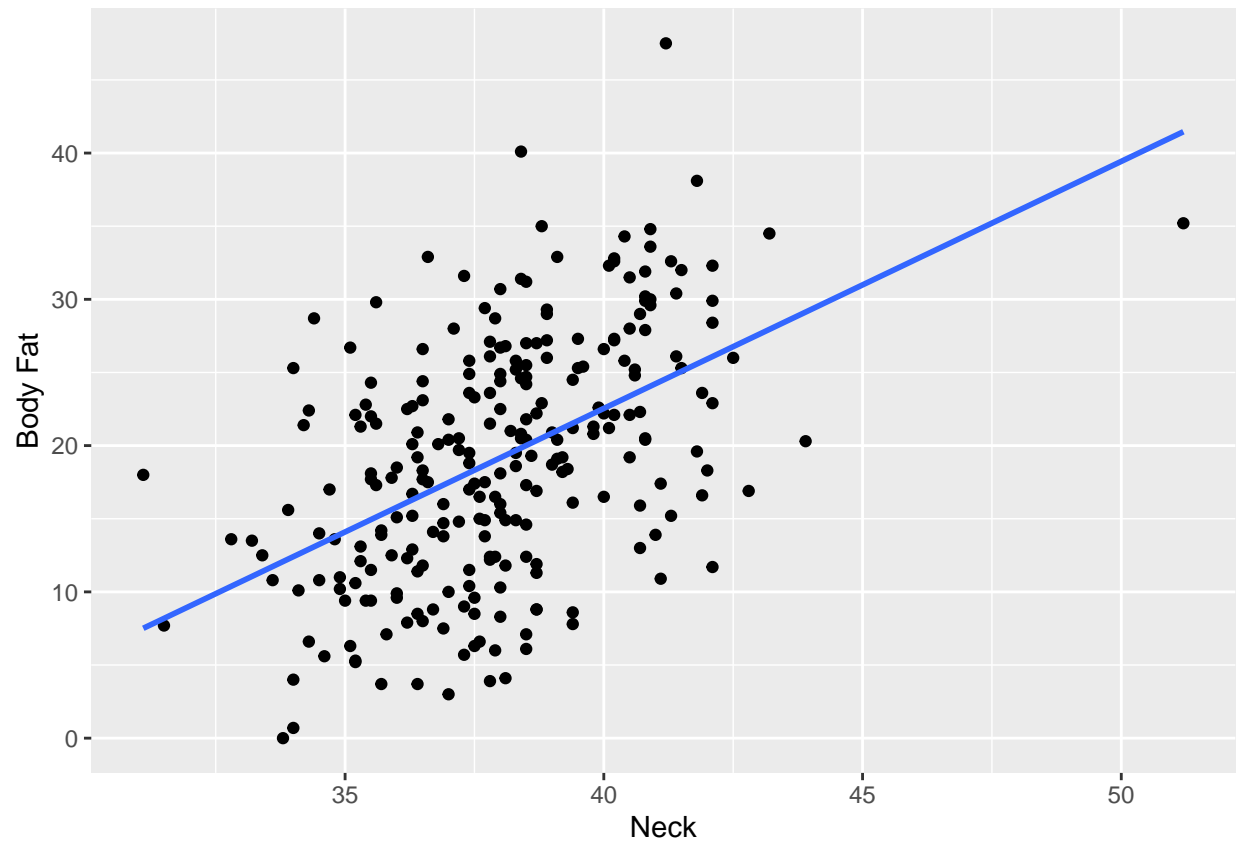
Variable	n	Mean	SD	Q1	Median	Q3
BodyFat	252	19.15	8.37	12.47	19.20	25.30
Age	252	44.88	12.60	35.75	43.00	54.00
Weight	252	178.92	29.39	159.00	176.50	197.00
Height	252	70.15	3.66	68.25	70.00	72.25
Neck	252	37.99	2.43	36.40	38.00	39.42
Chest	252	100.82	8.43	94.35	99.65	105.38
Abdomen	252	92.56	10.78	84.57	90.95	99.33
Hip	252	99.90	7.16	95.50	99.30	103.53
Thigh	252	59.41	5.25	56.00	59.00	62.35
Knee	252	38.59	2.41	36.98	38.50	39.92
Ankle	252	23.10	1.69	22.00	22.80	24.00
Biceps	252	32.27	3.02	30.20	32.05	34.32
Forearm	252	28.66	2.02	27.30	28.70	30.00
Wrist	252	18.23	0.93	17.60	18.30	18.80

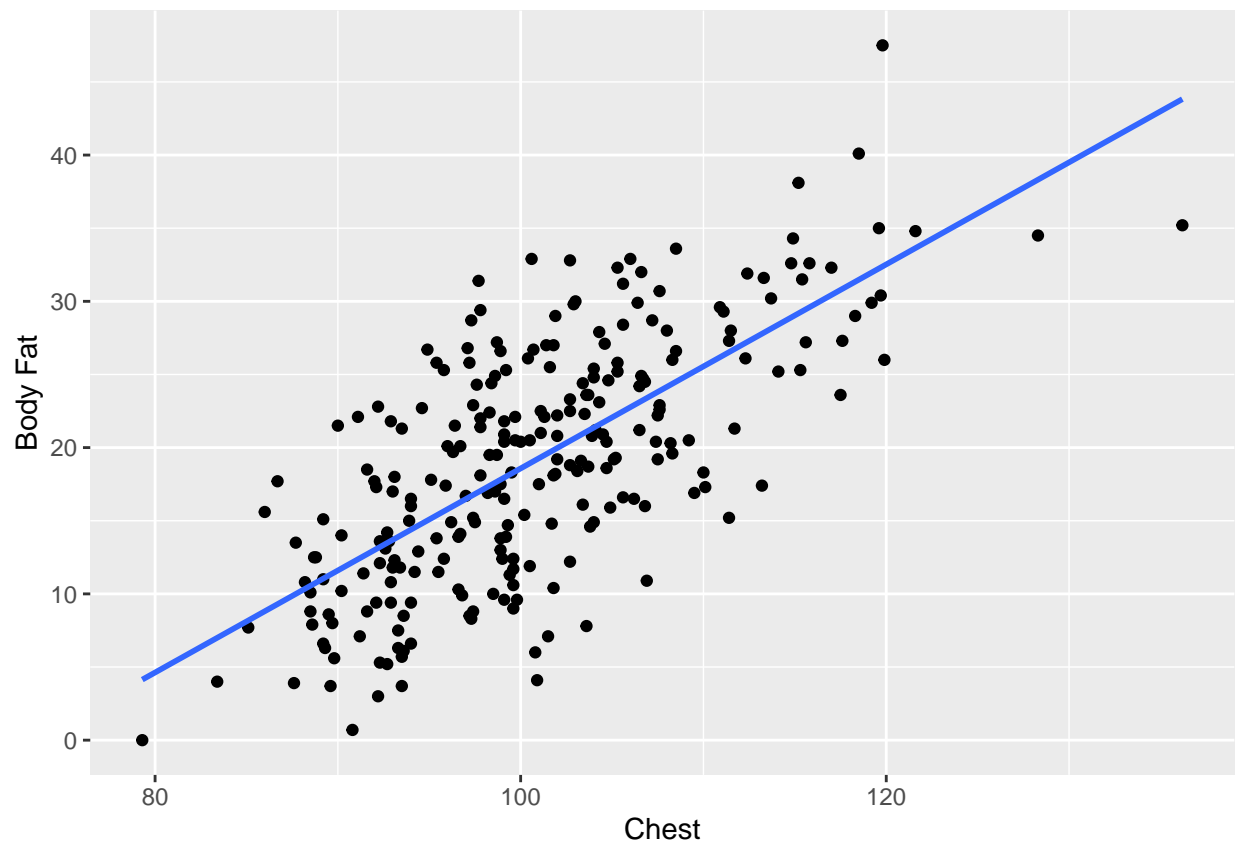
Since all the variables are in the same range, we shall look at the individual relationships between each explanatory variable and the response variable.

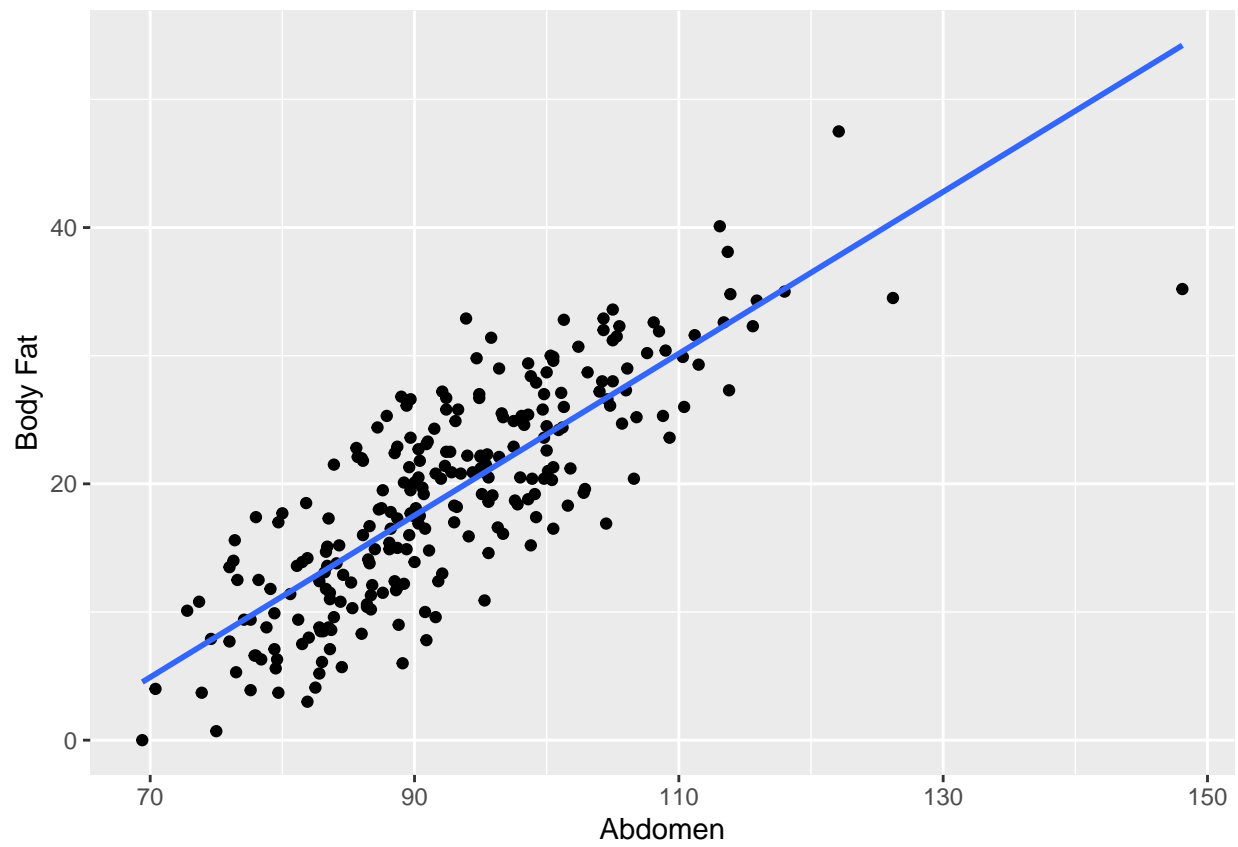


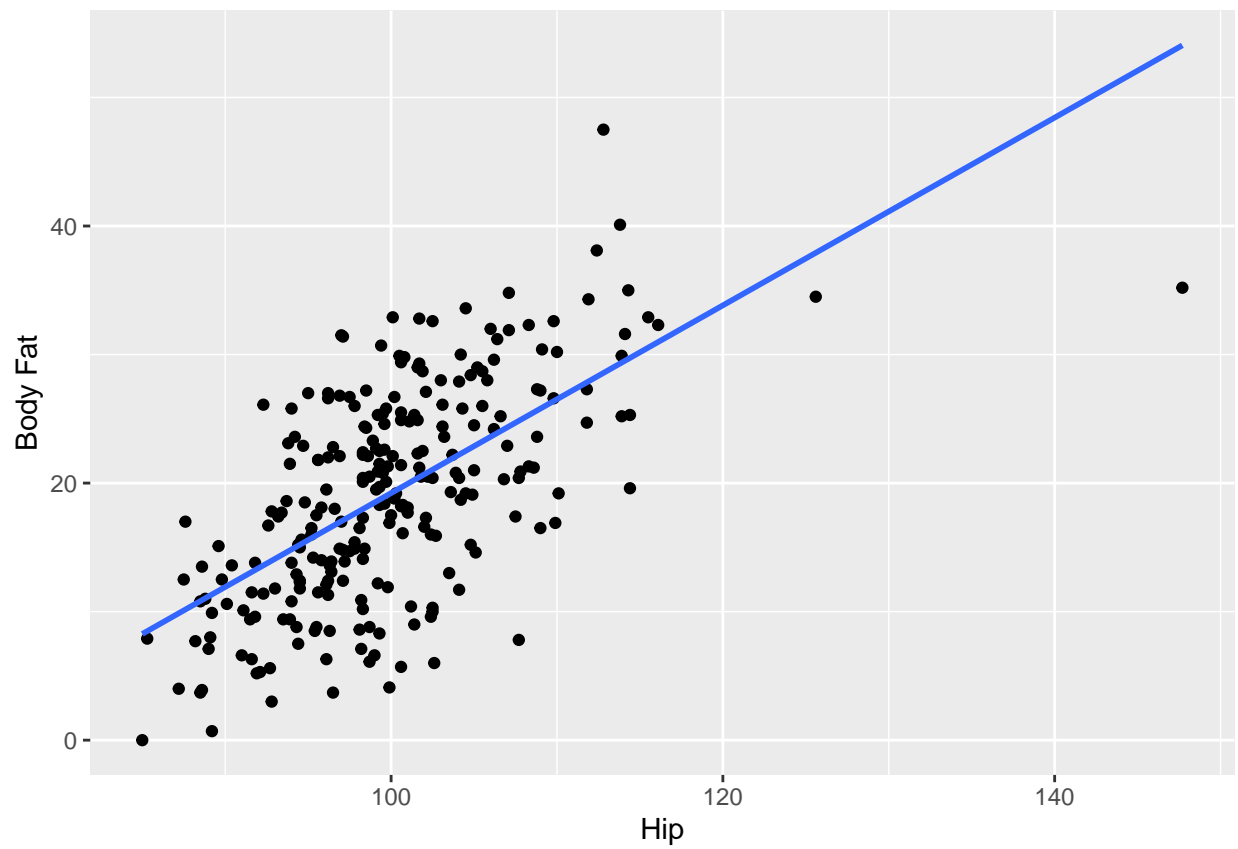


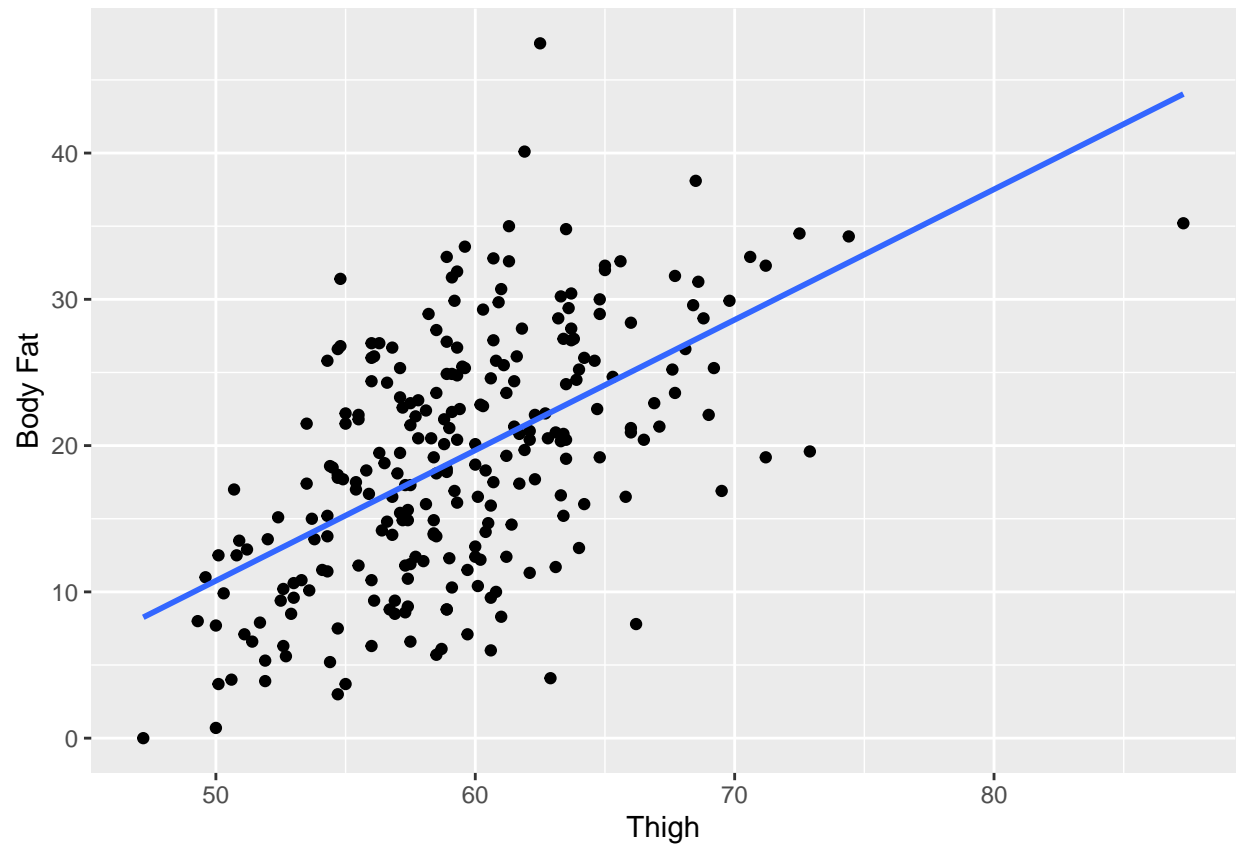


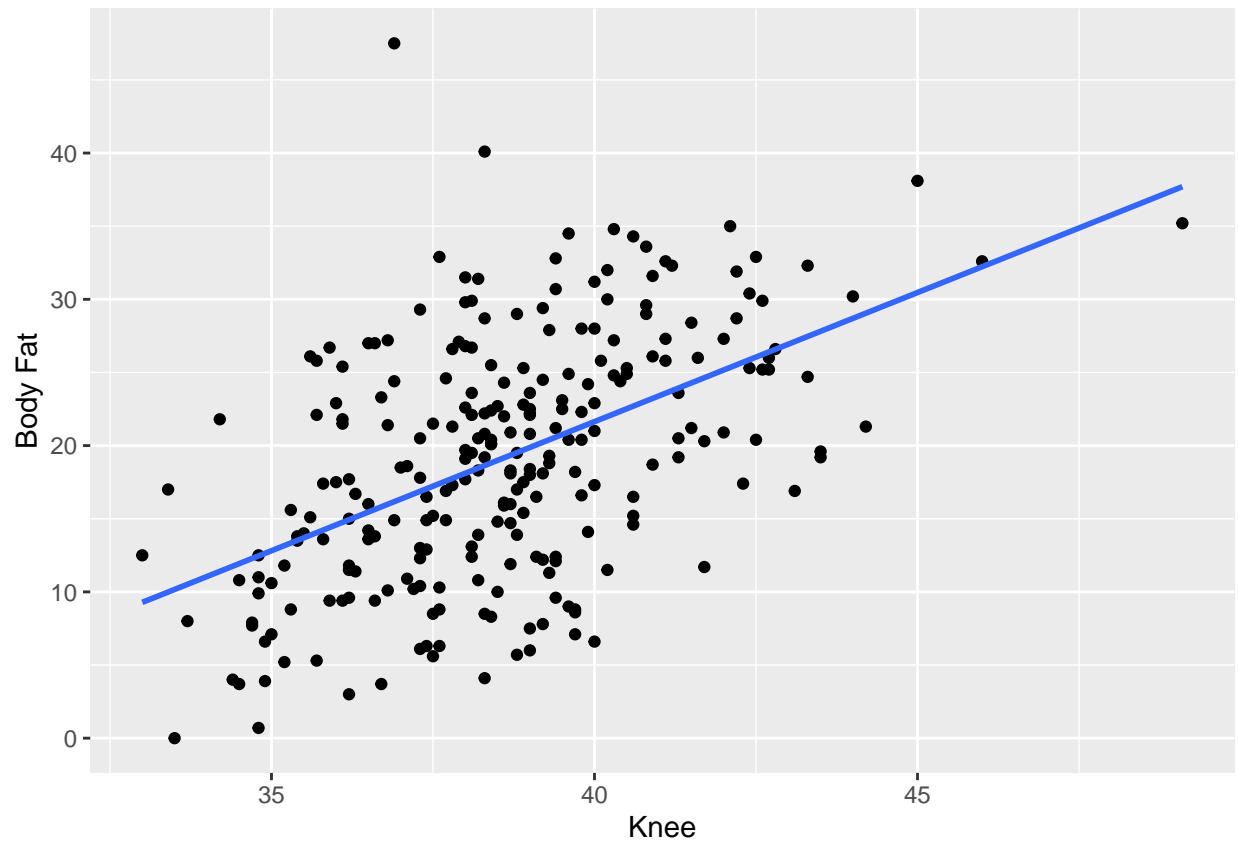


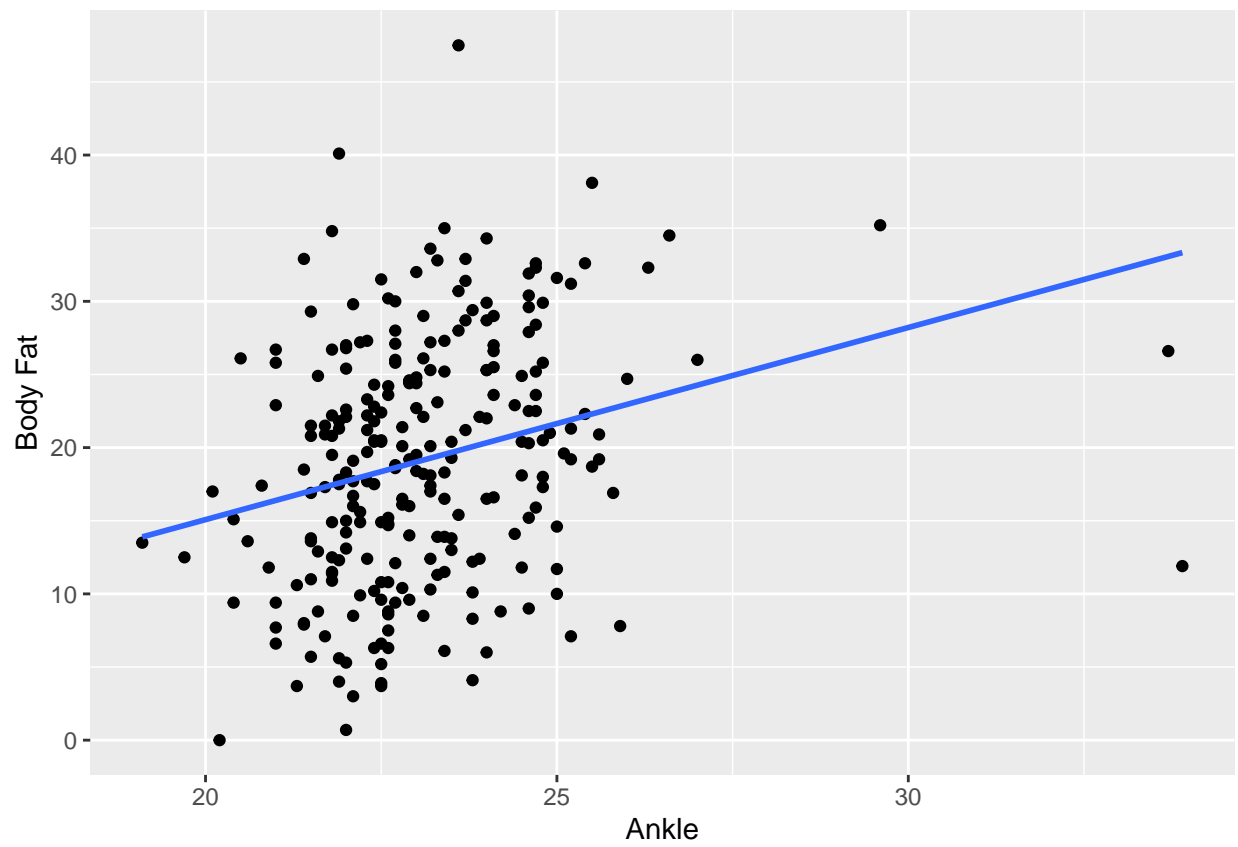


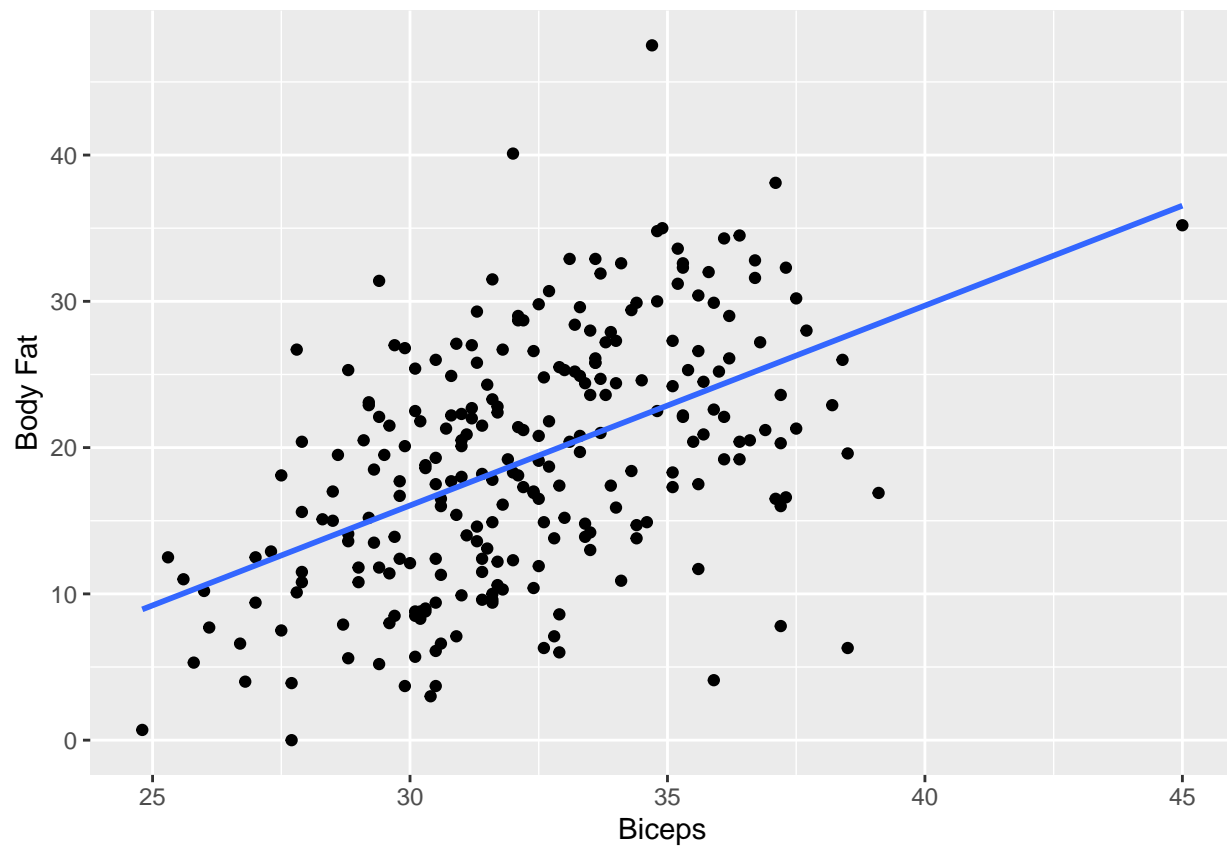


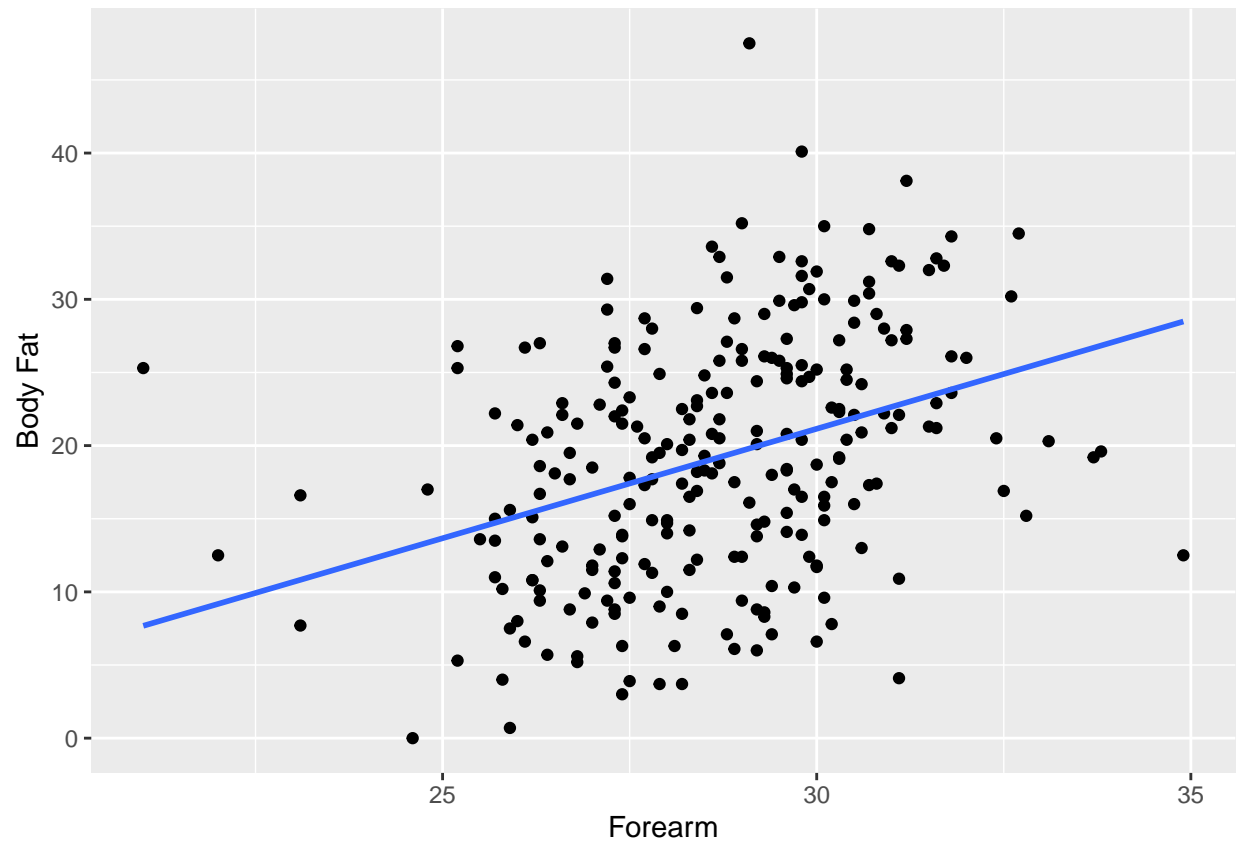


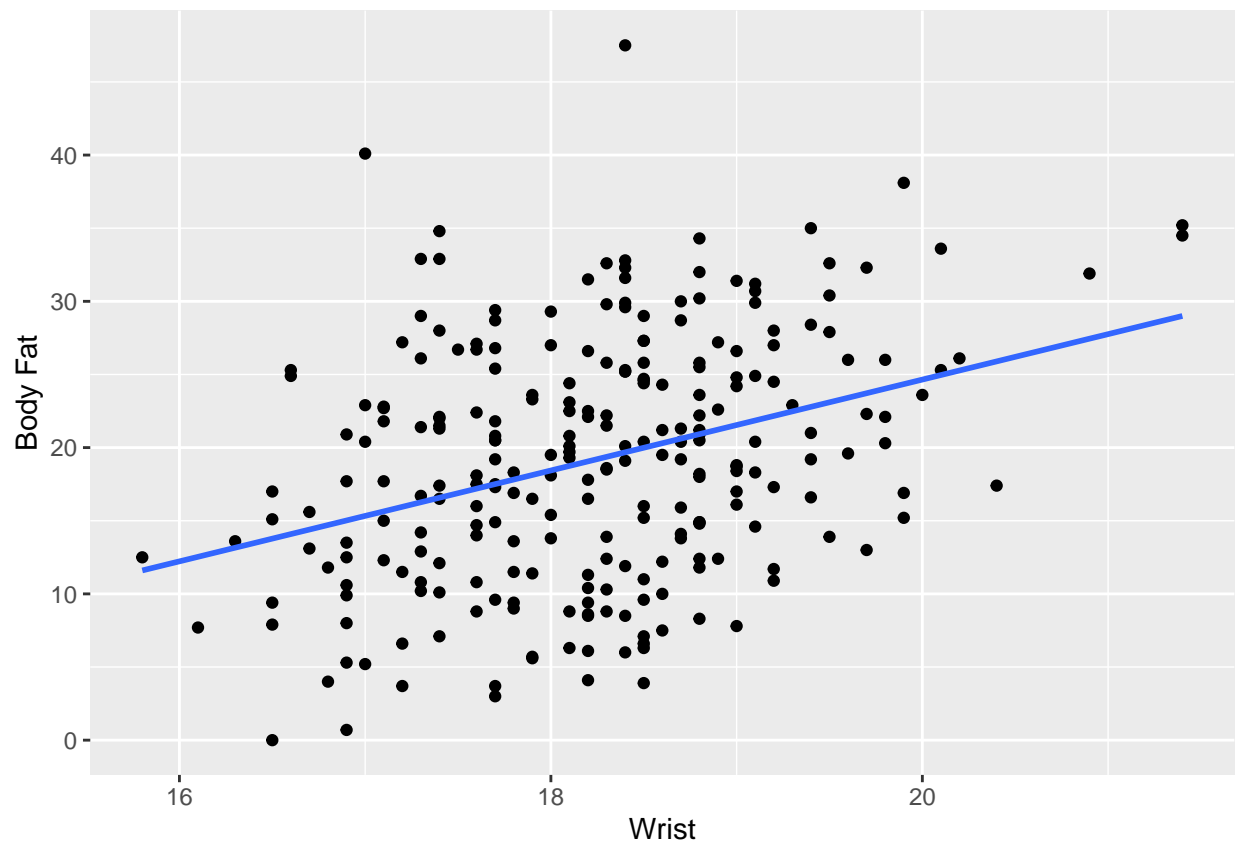


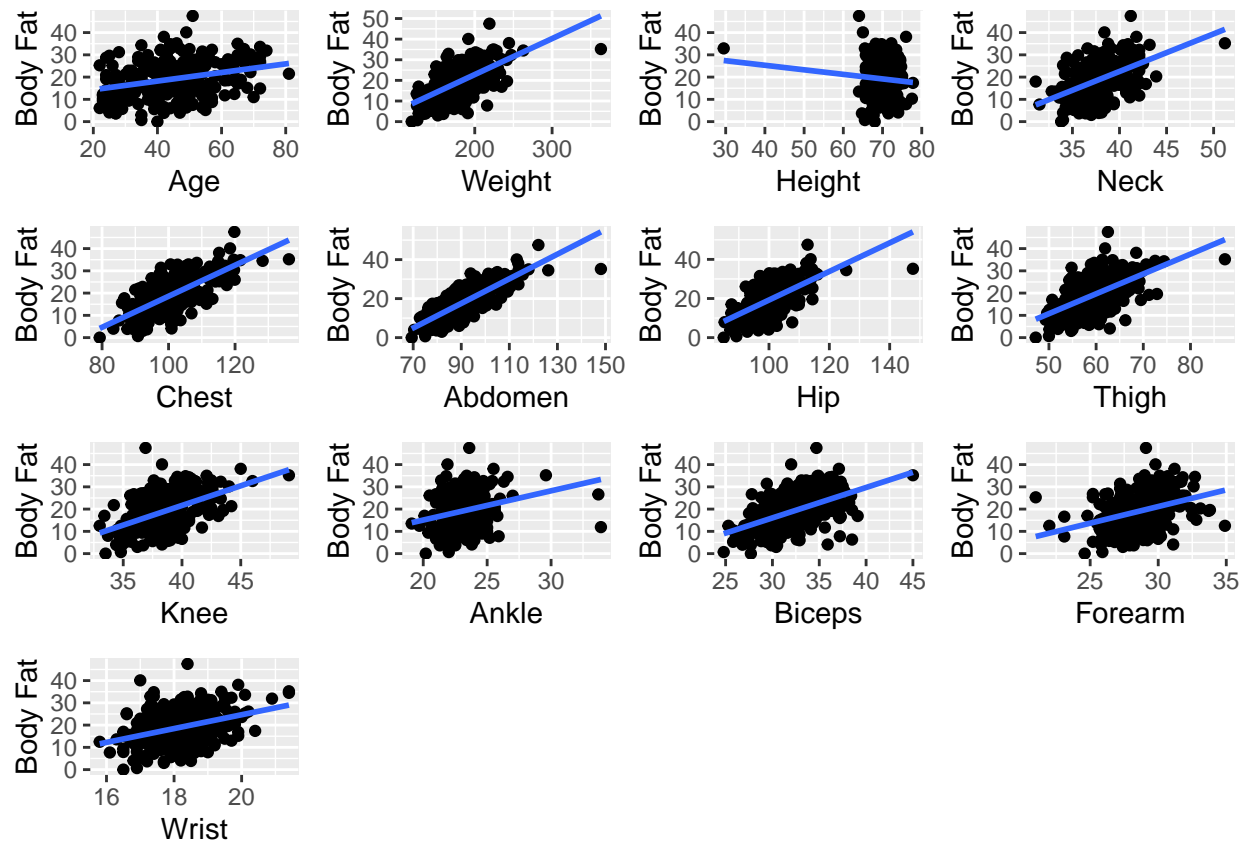








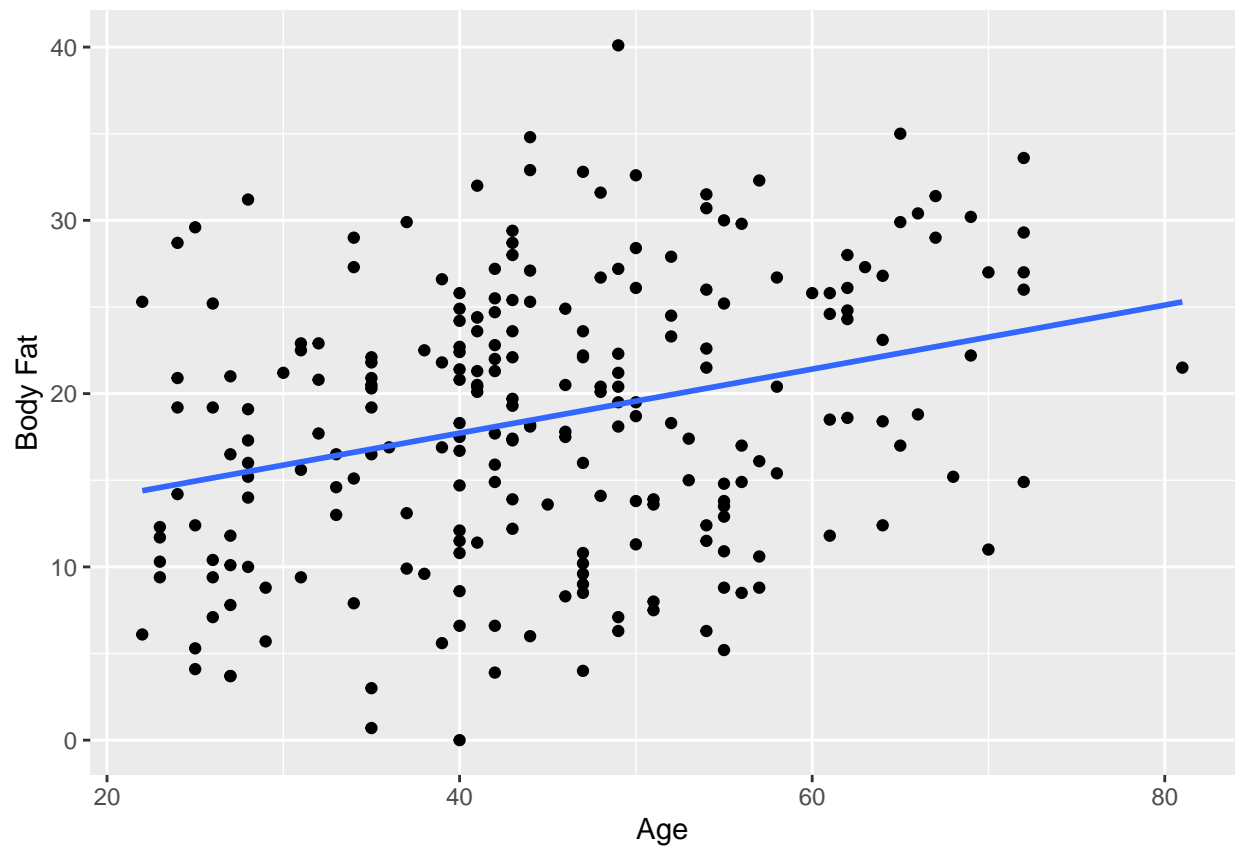


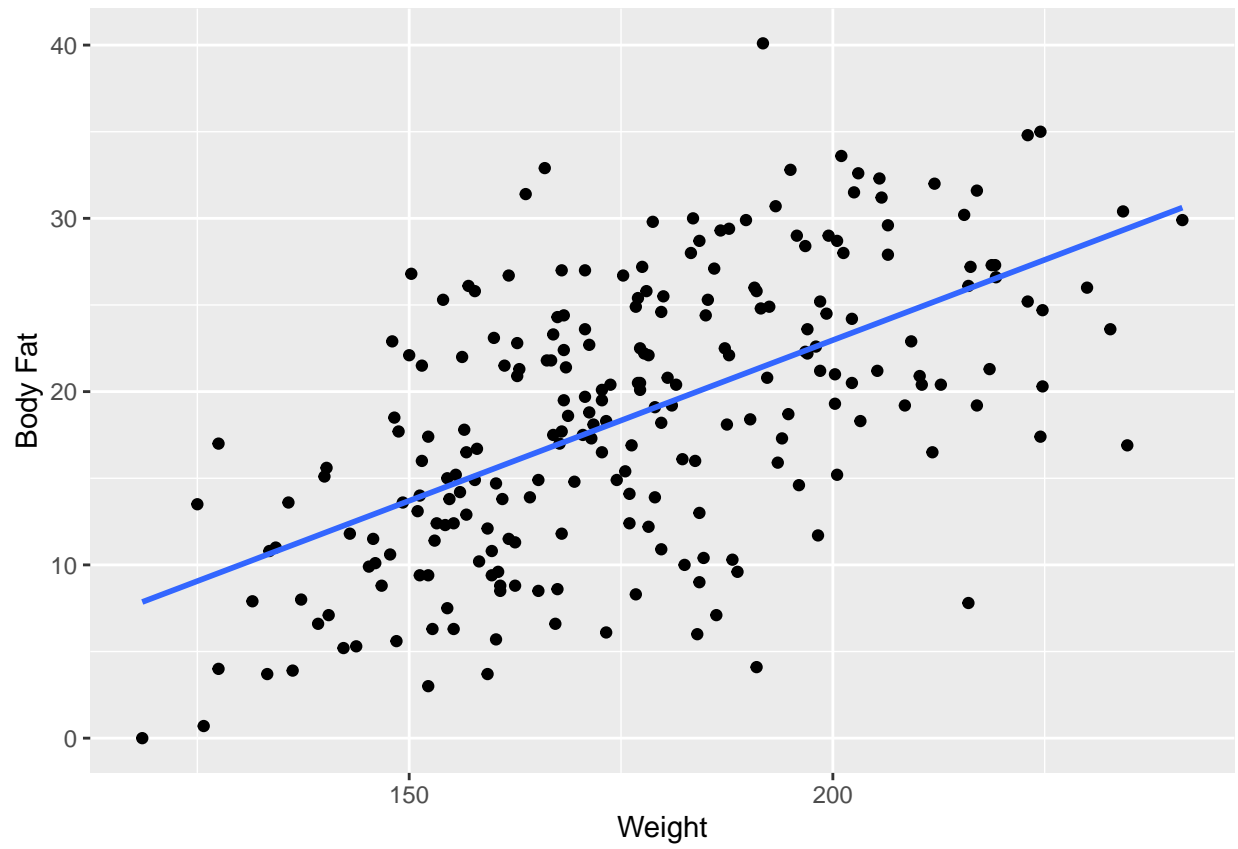


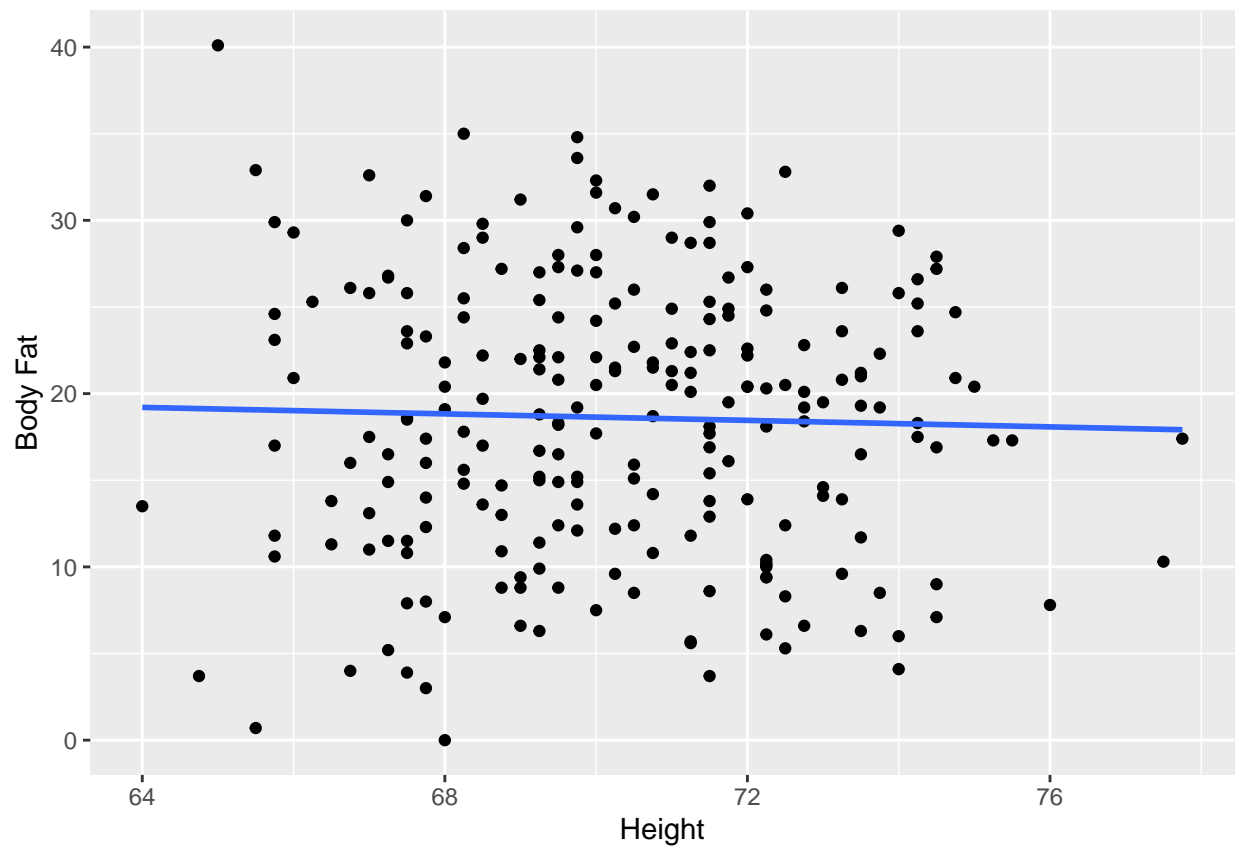
After plotting the relationship between the variables, we observe that there is indeed a good positive relationship between most of them except one of the explanatory variables. The plot between Height and BodyFat shows us that they have a very weak relationship between them. There are some influencer points in some of the plots, which we can try to eliminate to improve the model.

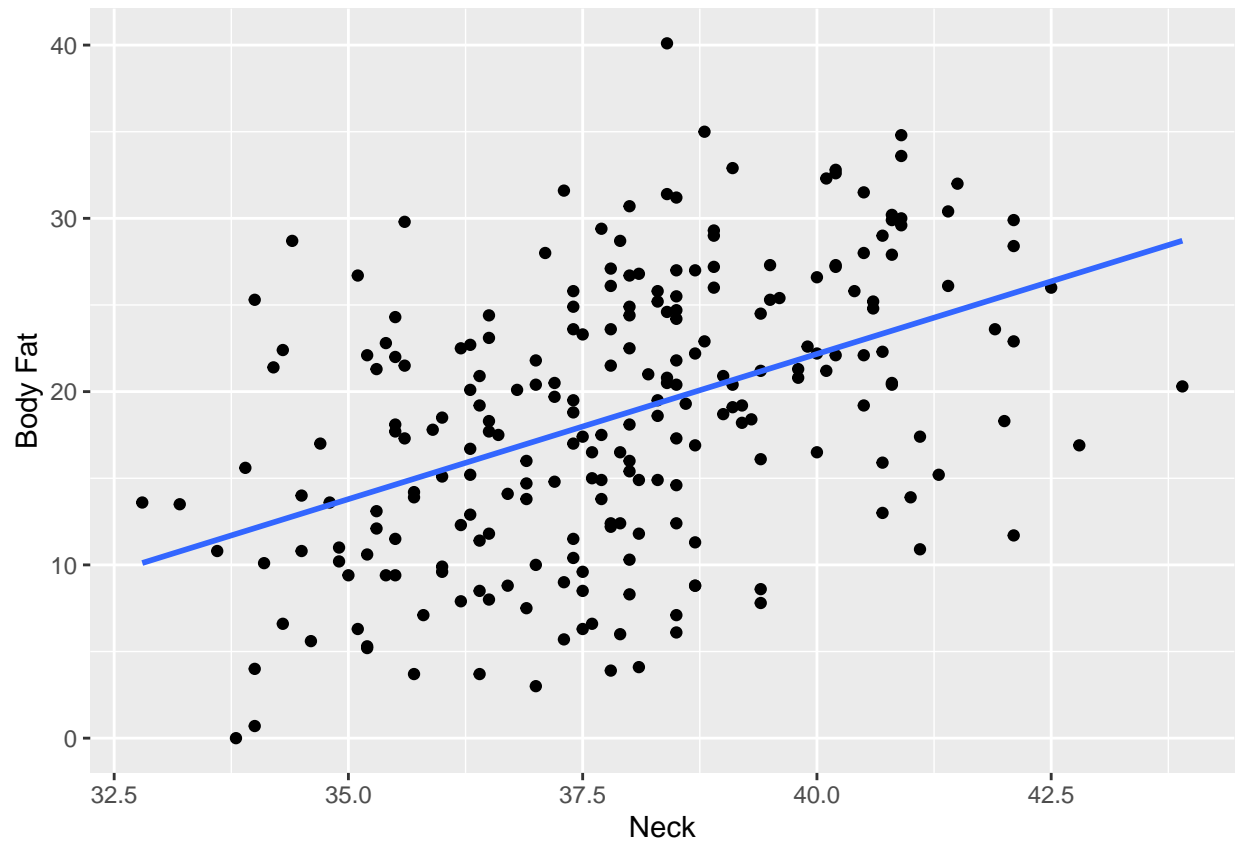
We are trying to find and remove the outliers to see if it makes the model any better.

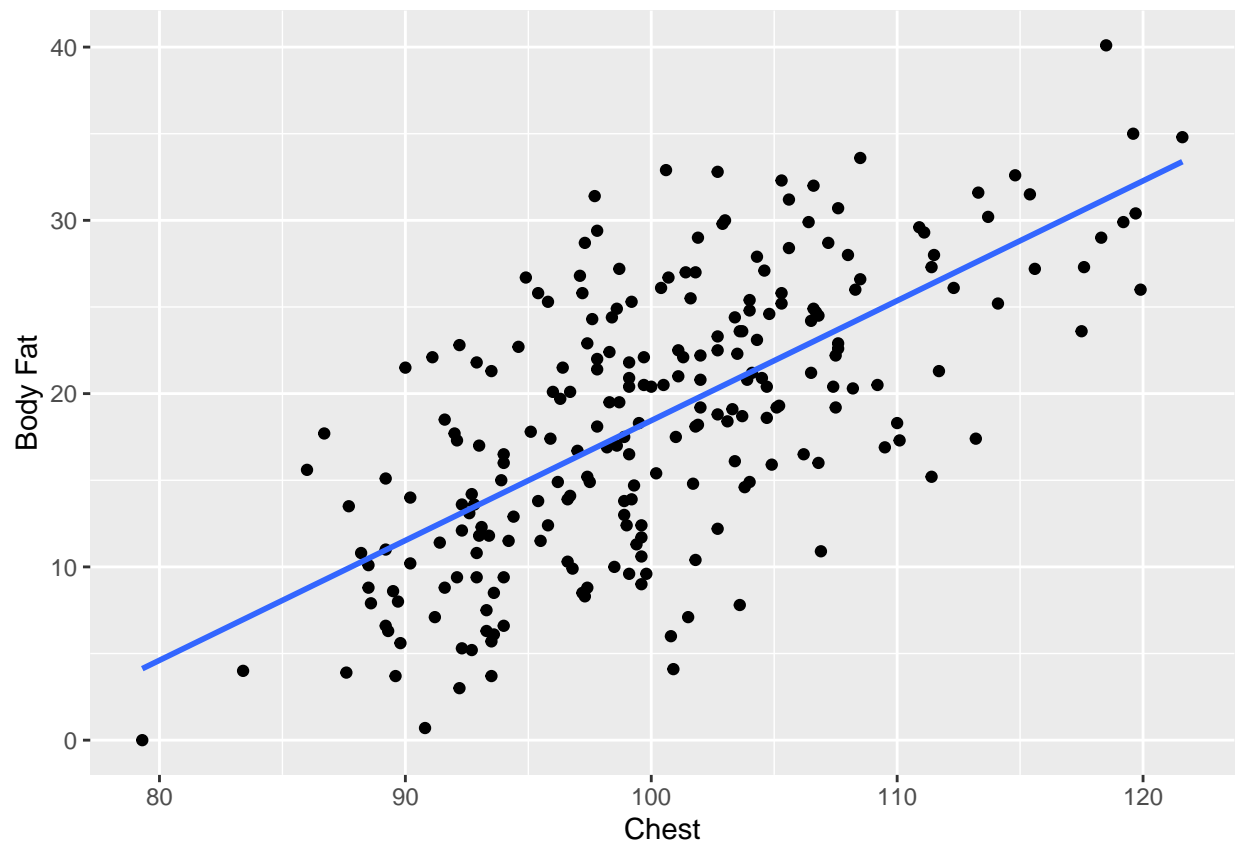
Let's look at the relationships between the explanatory variables and the response variable after eliminating the outliers

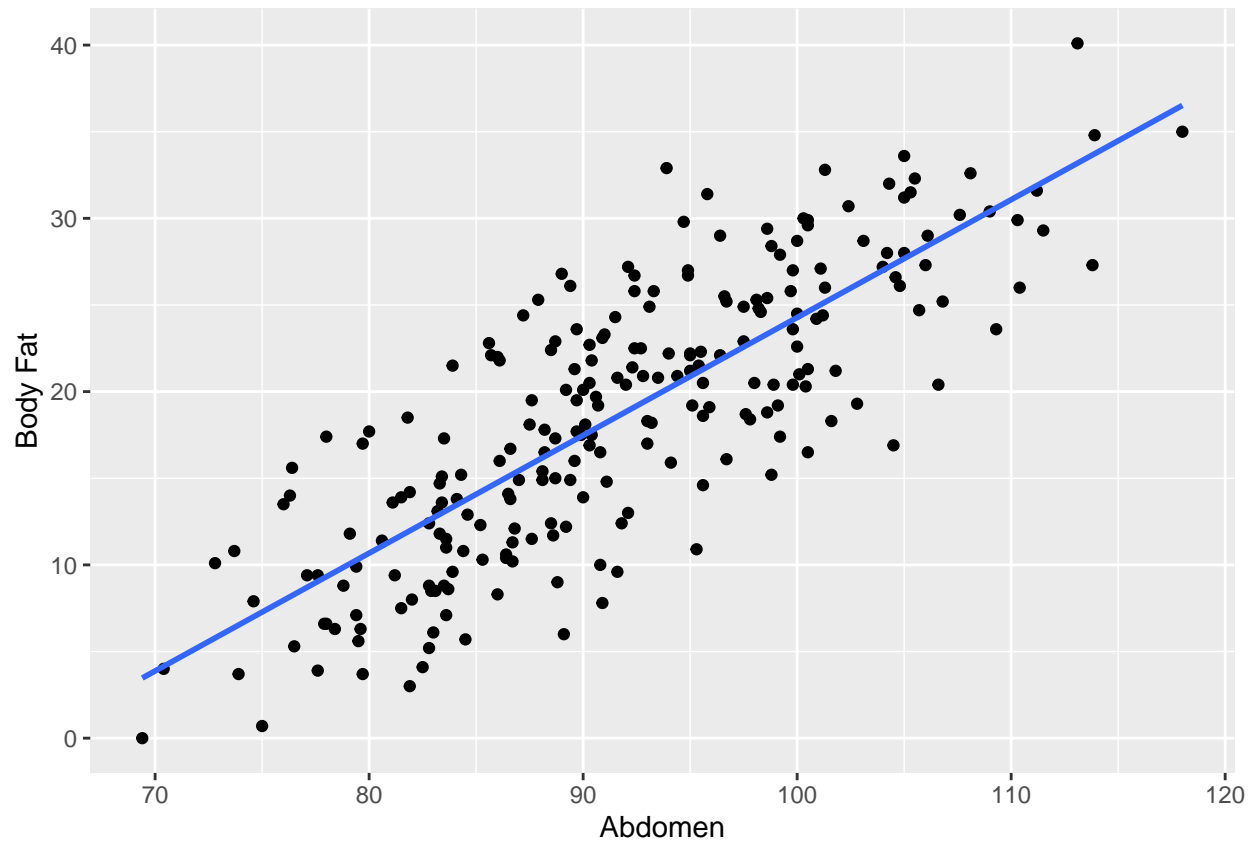


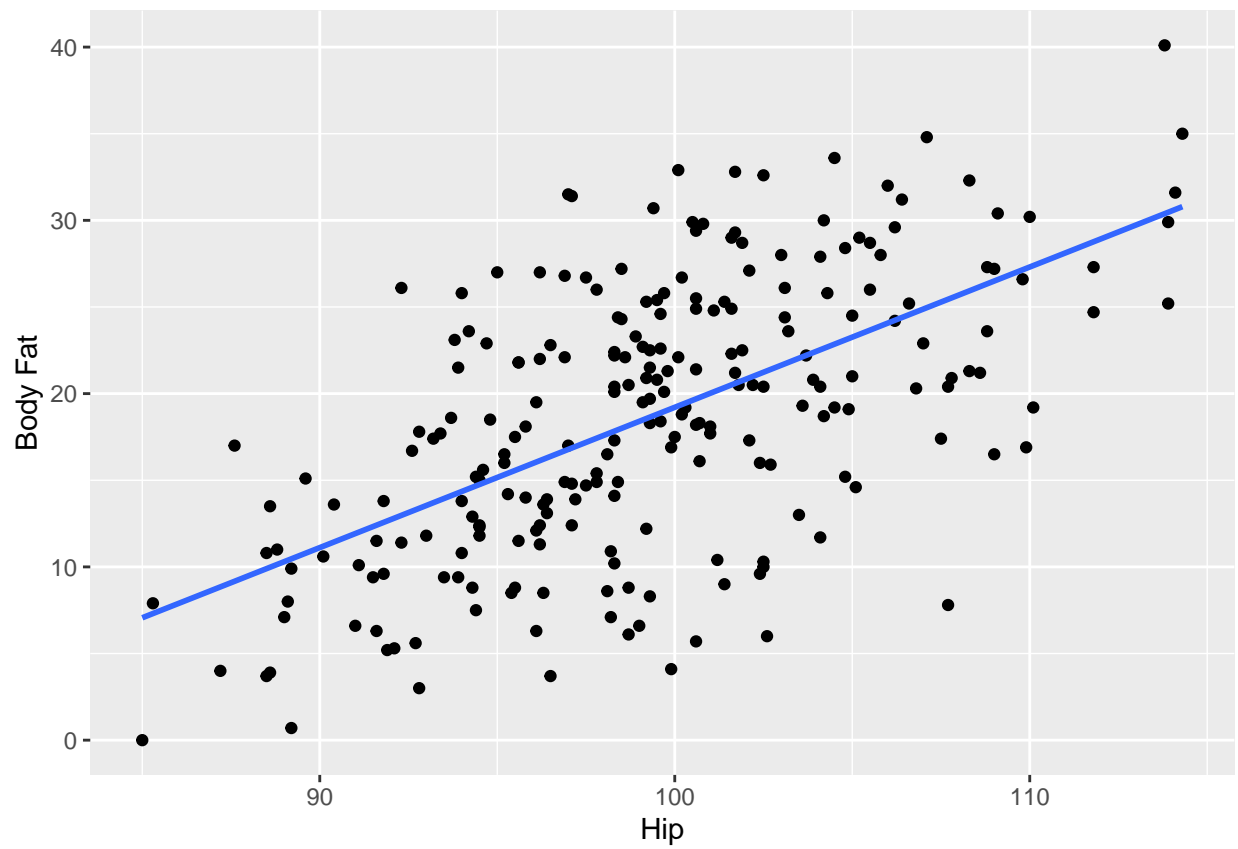


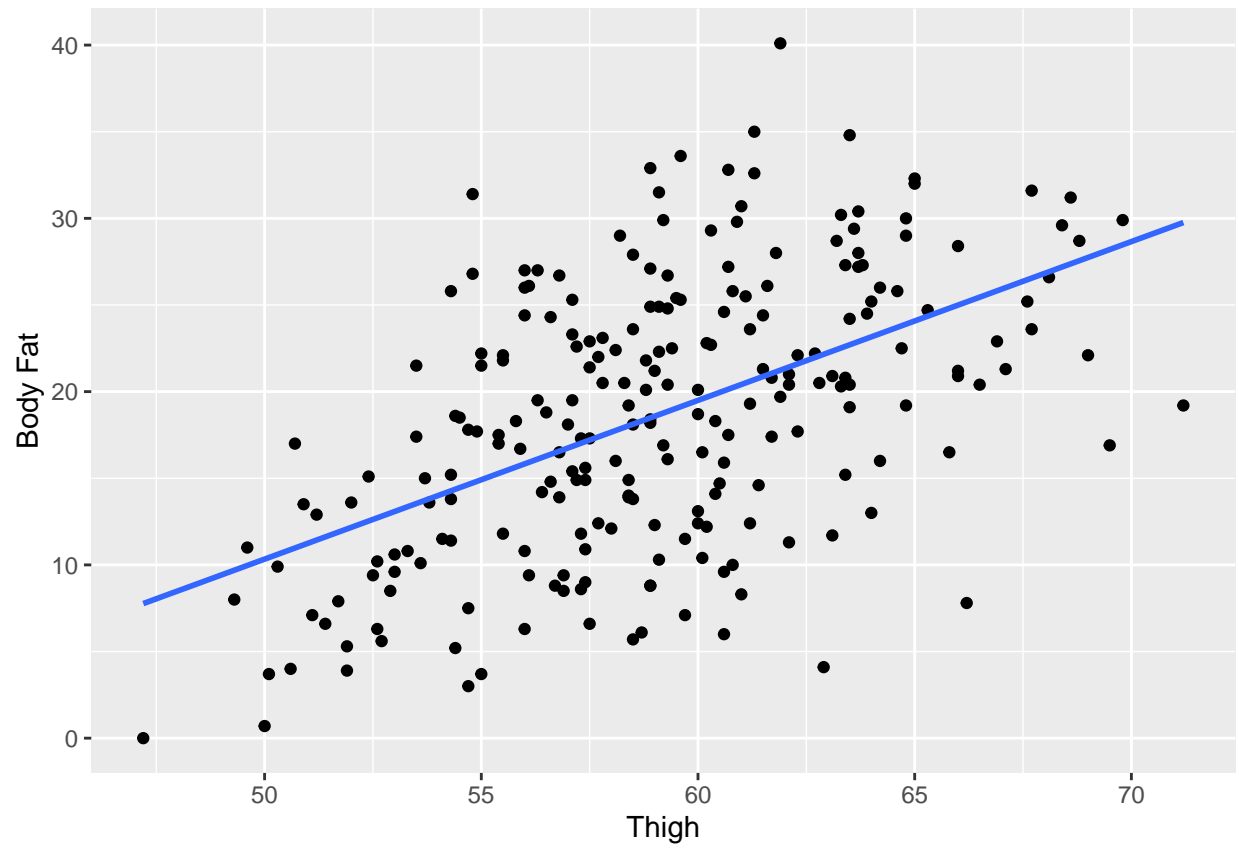


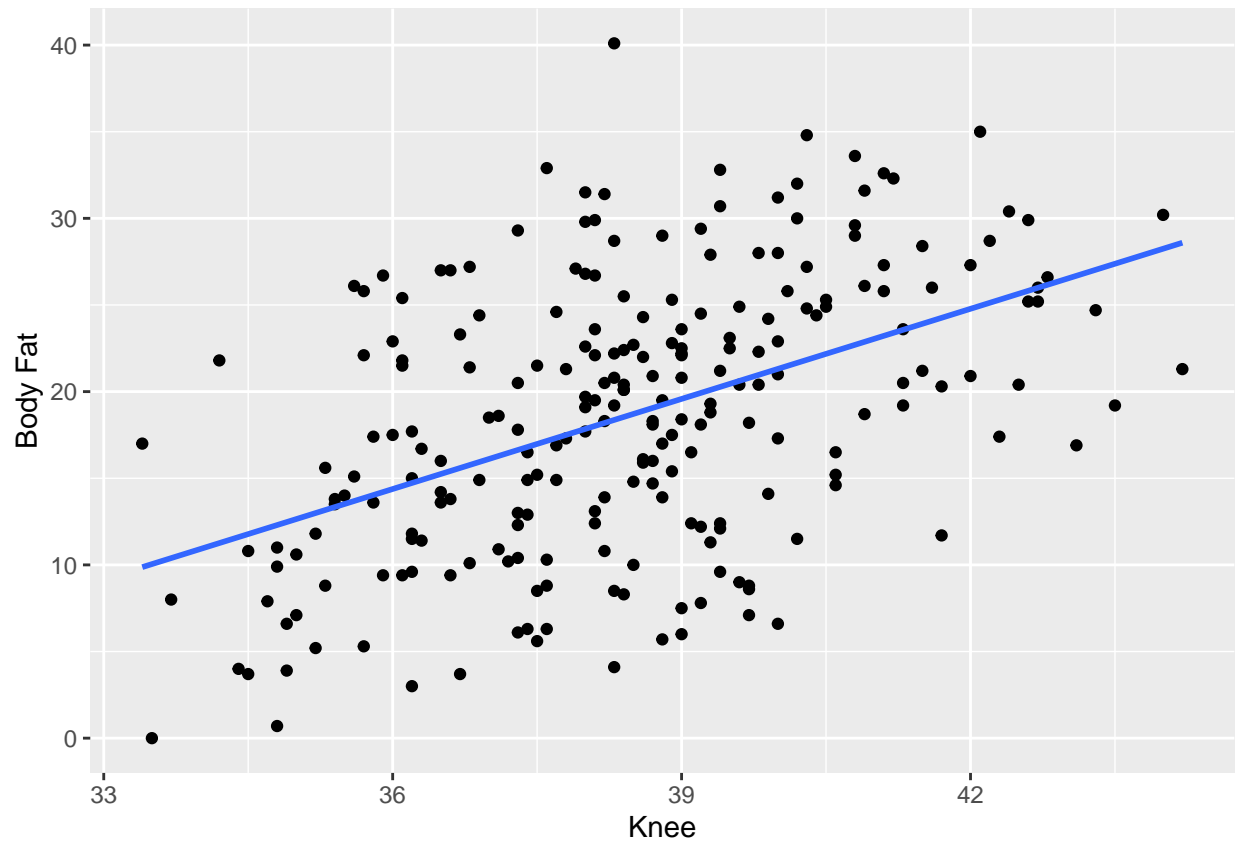


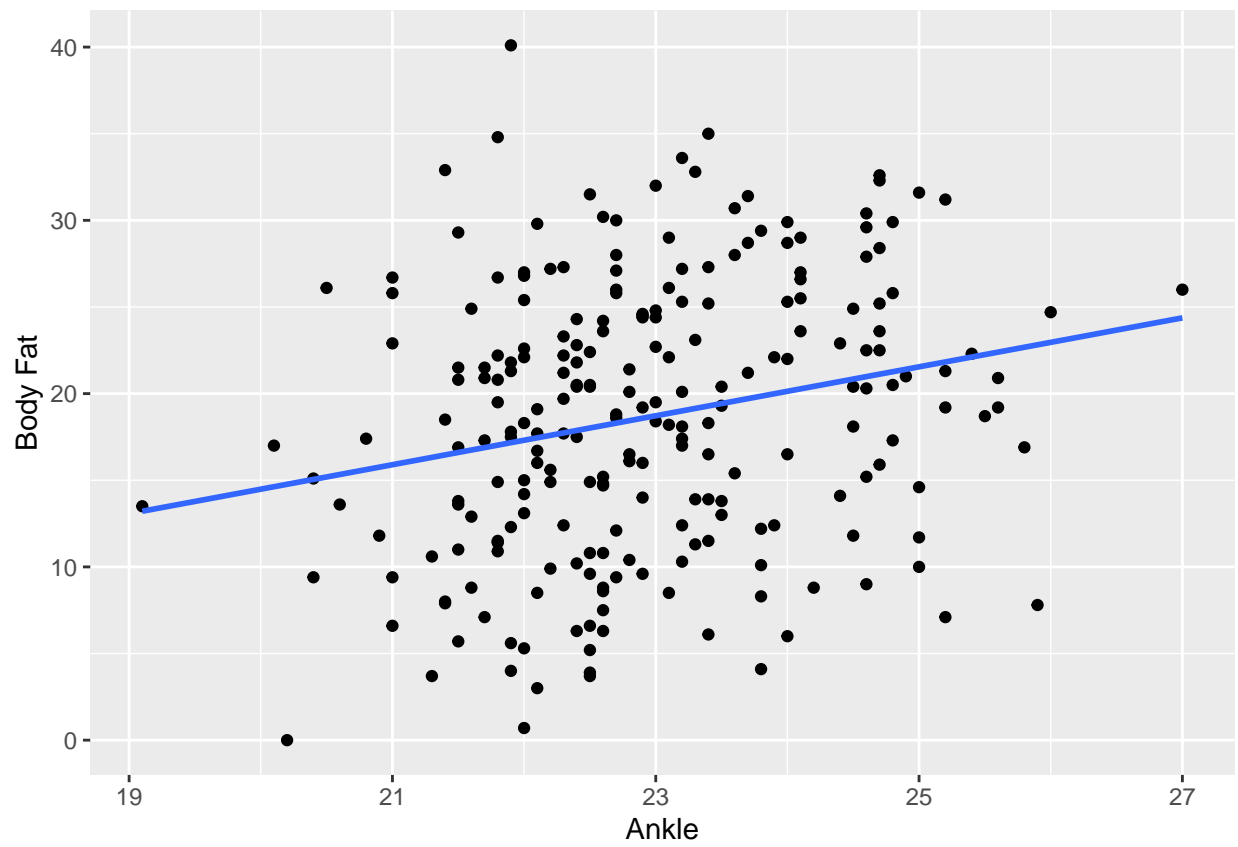


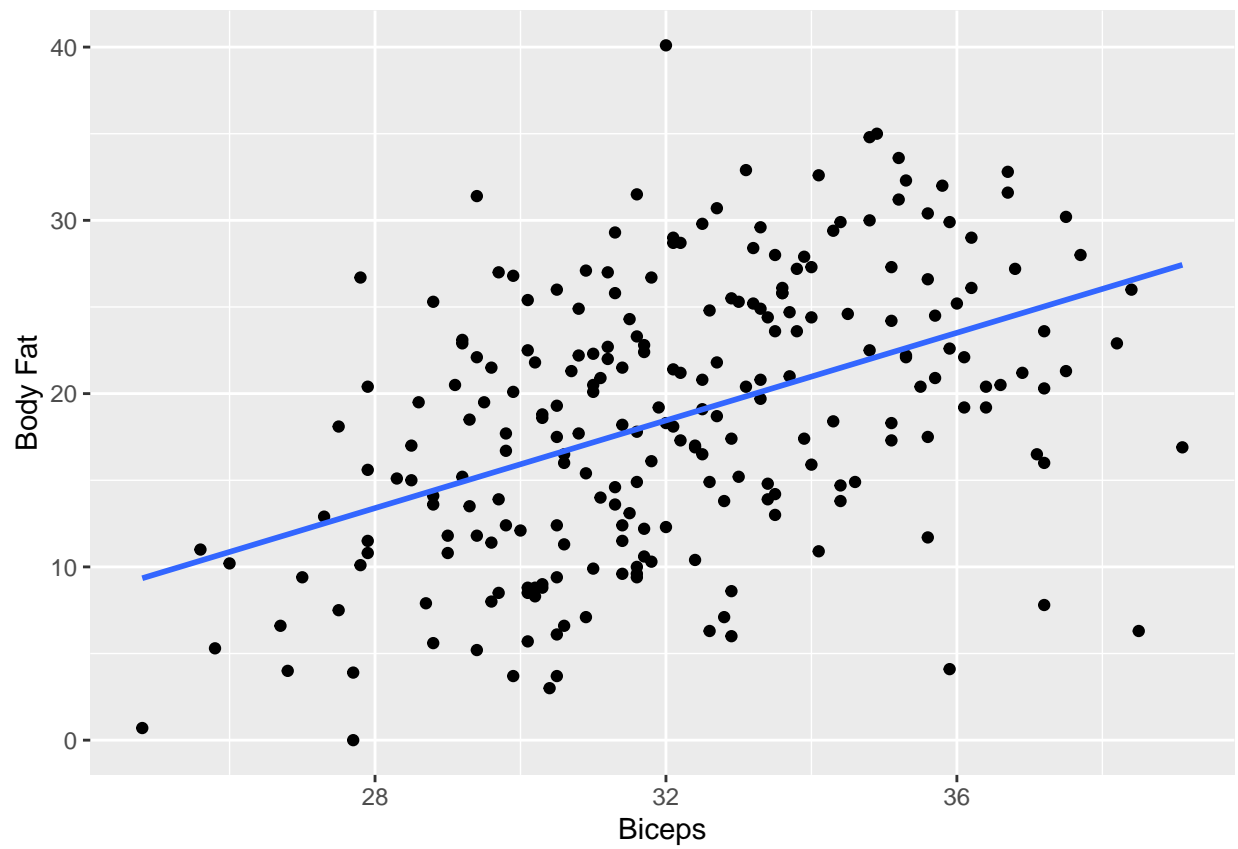


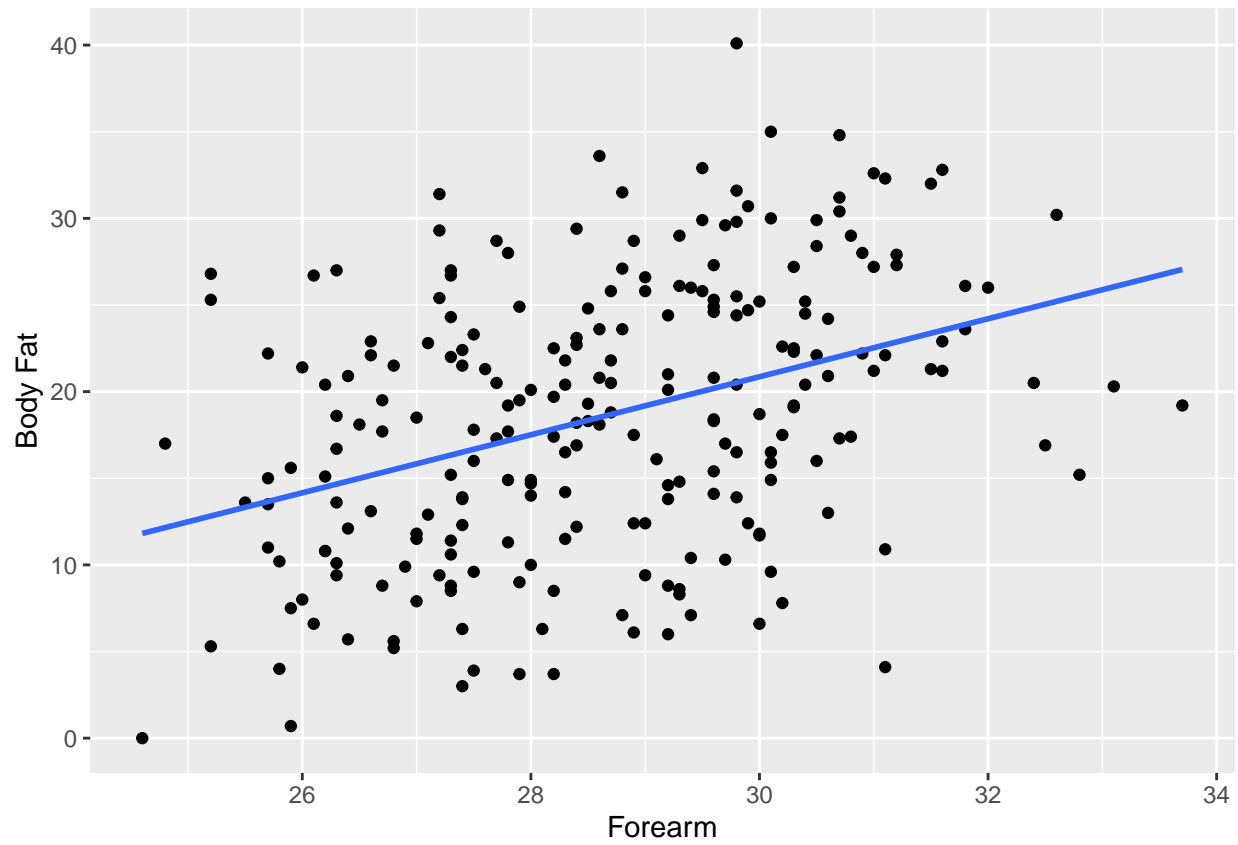


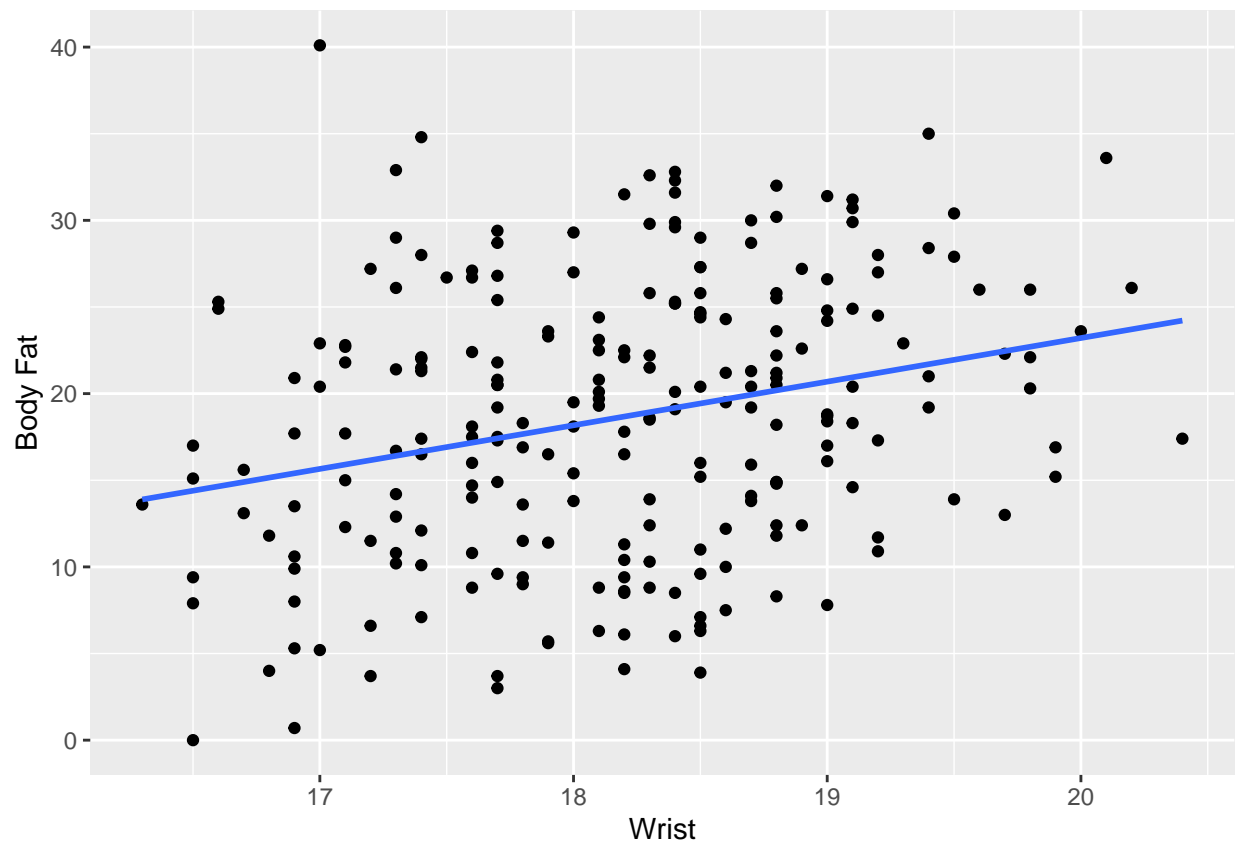


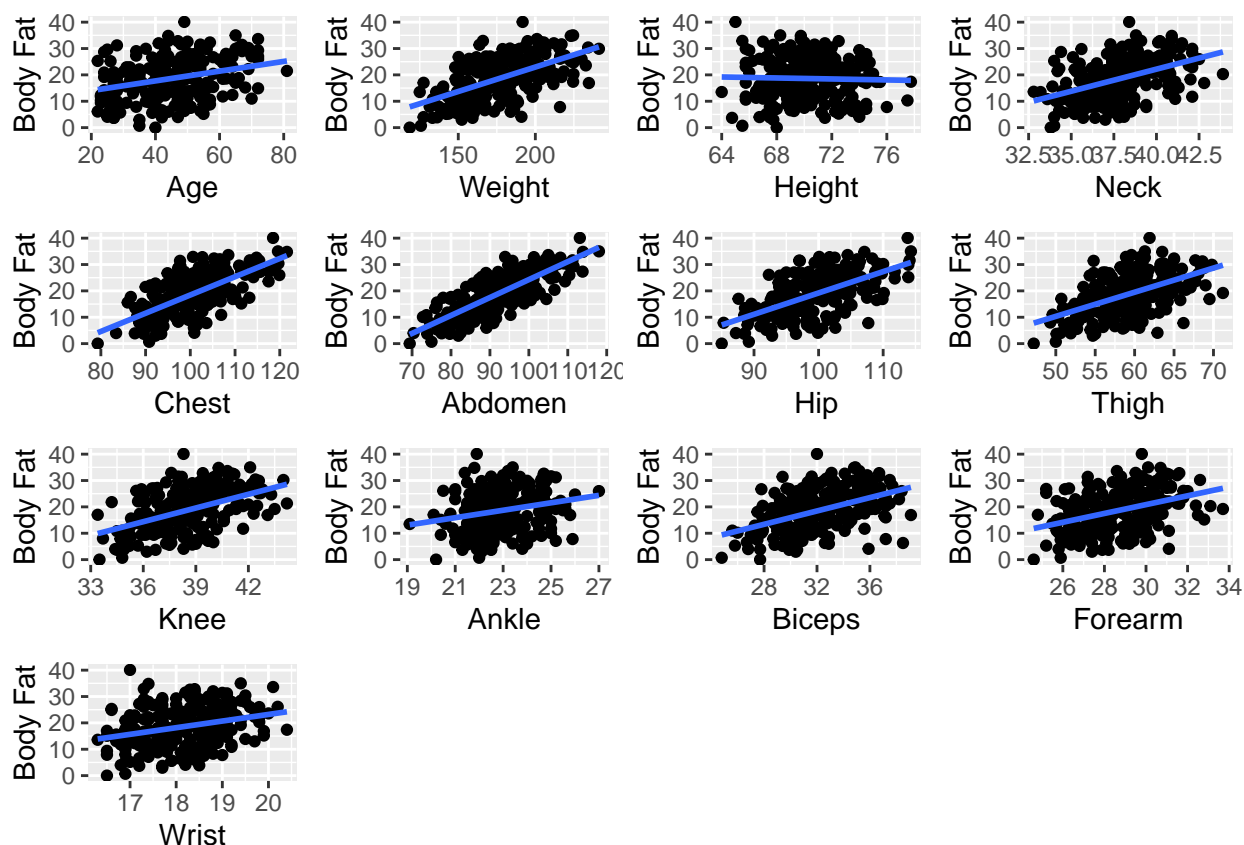












After removing the outlier points, we come to the conclusion, that there is no major significant change to the plots.

So firstly, we shall try to fit the model taking all the predictor variables and the response variable into consideration.

term	estimate	std_error	statistic	p_value	lower_ci	upper_ci
intercept	-18.188	17.349	-1.048	0.296	-52.365	15.988
Age	0.062	0.032	1.919	0.056	-0.002	0.126
Weight	-0.088	0.054	-1.652	0.100	-0.194	0.017
Height	-0.070	0.096	-0.725	0.469	-0.259	0.120
Neck	-0.471	0.232	-2.024	0.044	-0.929	-0.013
Chest	-0.024	0.099	-0.241	0.810	-0.219	0.171
Abdomen	0.955	0.086	11.044	0.000	0.784	1.125
Hip	-0.208	0.146	-1.422	0.156	-0.495	0.080
Thigh	0.236	0.144	1.636	0.103	-0.048	0.520
Knee	0.015	0.242	0.063	0.950	-0.461	0.492
Ankle	0.174	0.221	0.786	0.433	-0.262	0.610
Biceps	0.182	0.171	1.061	0.290	-0.156	0.519
Forearm	0.452	0.199	2.270	0.024	0.060	0.844
Wrist	-1.621	0.535	-3.030	0.003	-2.674	-0.567

As we can see from the above table, the p-values are quite high. This is due to high multicollinearity. To tackle this problem, we shall use AIC criterion to perform feature selection.

Stepwise Selection Method

```

## -----
##
## Candidate Terms:
##
## 1 . Age
## 2 . Weight
## 3 . Height
## 4 . Neck
## 5 . Chest
## 6 . Abdomen
## 7 . Hip
## 8 . Thigh
## 9 . Knee
## 10 . Ankle
## 11 . Biceps
## 12 . Forearm
## 13 . Wrist
##
## Step 0: AIC = 1788.893
## BodyFat ~ 1
##
##
## Variables Entered/Removed:
##
##                               Enter New Variables
## -----
## Variable      DF      AIC      Sum Sq      RSS      R-Sq      Adj. R-Sq
## -----
## Abdomen        1    1517.790    11631.527    5947.463    0.662      0.660
## Chest          1    1619.387     8678.314    8900.676    0.494      0.492
## Hip            1    1665.968     6871.209    10707.781    0.391      0.388
## Weight         1    1672.431     6593.016    10985.974    0.375      0.373
## Thigh          1    1696.228     5505.047    12073.943    0.313      0.310
## Knee           1    1715.443     4548.394    13030.596    0.259      0.256
## Biceps         1    1720.633     4277.257    13301.733    0.243      0.240
## Neck           1    1721.509     4230.918    13348.072    0.241      0.238
## Forearm        1    1755.625     2295.825    15283.165    0.131      0.127
## Wrist          1    1758.646     2111.485    15467.505    0.120      0.117
## Age            1    1768.522     1493.300    16085.689    0.085      0.081
## Ankle          1    1772.404     1243.536    16335.454    0.071      0.067
## Height         1    1788.866      140.798    17438.192    0.008      0.004
## -----
##
## - Abdomen added
##
##
## Step 1 : AIC = 1517.79
## BodyFat ~ Abdomen
##
##                               Enter New Variables
## -----
## Variable      DF      AIC      Sum Sq      RSS      R-Sq      Adj. R-Sq
## -----
## Weight         1    1473.185    12635.745    4943.245    0.719      0.717

```



```

## Wrist      1    1487.794    12340.703    5238.287    0.702    0.700
## Neck       1    1492.307    12246.047    5332.942    0.697    0.694
## Hip        1    1495.420    12179.763    5399.227    0.693    0.690
## Height     1    1499.561    12090.306    5488.683    0.688    0.685
## Knee       1    1505.912    11950.225    5628.765    0.680    0.677
## Ankle      1    1509.705    11864.851    5714.139    0.675    0.672
## Age        1    1511.129    11832.460    5746.530    0.673    0.670
## Chest      1    1511.370    11826.979    5752.011    0.673    0.670
## Thigh      1    1512.282    11806.111    5772.879    0.672    0.669
## Biceps     1    1513.992    11766.822    5812.168    0.669    0.667
## Forearm    1    1517.479    11685.830    5893.159    0.665    0.662
## -----
##
## - Weight added
##
##
## Step 2 : AIC = 1473.185
## BodyFat ~ Abdomen + Weight
##
## Remove Existing Variables
## -----
## Variable    DF      AIC      Sum Sq      RSS      R-Sq      Adj. R-Sq
## -----
## Weight      1    1517.790    11631.527    5947.463    0.662    0.660
## Abdomen     1    1672.431    6593.016    10985.974    0.375    0.373
## -----
##
## Enter New Variables
## -----
## Variable    DF      AIC      Sum Sq      RSS      R-Sq      Adj. R-Sq
## -----
## Wrist       1    1467.041    12792.936    4786.054    0.728    0.724
## Neck        1    1470.714    12722.674    4856.316    0.724    0.720
## Thigh       1    1471.003    12717.100    4861.889    0.723    0.720
## Forearm     1    1471.753    12702.597    4876.393    0.723    0.719
## Biceps      1    1471.911    12699.554    4879.436    0.722    0.719
## Height      1    1473.122    12676.037    4902.953    0.721    0.718
## Knee        1    1474.689    12645.463    4933.527    0.719    0.716
## Age         1    1475.086    12637.683    4941.307    0.719    0.716
## Ankle       1    1475.108    12637.258    4941.731    0.719    0.715
## Chest       1    1475.184    12635.755    4943.235    0.719    0.715
## Hip         1    1475.185    12635.750    4943.240    0.719    0.715
## -----
##
## - Wrist added
##
##
## Step 3 : AIC = 1467.041
## BodyFat ~ Abdomen + Weight + Wrist
##
## Remove Existing Variables
## -----
## Variable    DF      AIC      Sum Sq      RSS      R-Sq      Adj. R-Sq
## -----

```

```

## Wrist          1    1473.185    12635.745    4943.245    0.719    0.717
## Weight         1    1487.794    12340.703    5238.287    0.702    0.700
## Abdomen        1    1665.591    6971.728    10607.262    0.397    0.392
## -----
##
##                                     Enter New Variables
## -----
## Variable      DF      AIC      Sum Sq      RSS      R-Sq      Adj. R-Sq
## -----
## Forearm       1    1462.220    12920.754    4658.236    0.735    0.731
## Biceps        1    1464.325    12881.667    4697.323    0.733    0.728
## Thigh         1    1466.902    12833.398    4745.592    0.730    0.726
## Neck          1    1467.712    12818.118    4760.871    0.729    0.725
## Height        1    1467.805    12816.349    4762.641    0.729    0.725
## Age           1    1467.925    12814.088    4764.902    0.729    0.725
## Knee          1    1467.957    12813.481    4765.509    0.729    0.725
## Ankle         1    1468.252    12807.906    4771.084    0.729    0.724
## Hip           1    1468.555    12802.165    4776.825    0.728    0.724
## Chest         1    1468.975    12794.192    4784.798    0.728    0.723
## -----
##
## - Forearm added
##
##
## Step 4 : AIC = 1462.22
## BodyFat ~ Abdomen + Weight + Wrist + Forearm
##
##                                     Remove Existing Variables
## -----
## Variable      DF      AIC      Sum Sq      RSS      R-Sq      Adj. R-Sq
## -----
## Forearm       1    1467.041    12792.936    4786.054    0.728    0.724
## Wrist         1    1471.753    12702.597    4876.393    0.723    0.719
## Weight        1    1489.138    12354.321    5224.669    0.703    0.699
## Abdomen       1    1667.587    6971.898    10607.091    0.397    0.389
## -----
##
##                                     Enter New Variables
## -----
## Variable      DF      AIC      Sum Sq      RSS      R-Sq      Adj. R-Sq
## -----
## Neck          1    1461.442    12971.820    4607.169    0.738    0.733
## Age           1    1462.136    12959.116    4619.874    0.737    0.732
## Biceps        1    1462.380    12954.637    4624.353    0.737    0.732
## Thigh         1    1462.743    12947.970    4631.020    0.737    0.731
## Knee          1    1463.145    12940.583    4638.407    0.736    0.731
## Ankle         1    1463.235    12938.912    4640.078    0.736    0.731
## Height        1    1463.241    12938.803    4640.187    0.736    0.731
## Hip           1    1464.029    12924.286    4654.704    0.735    0.730
## Chest         1    1464.193    12921.242    4657.747    0.735    0.730
## -----
##
## - Neck added
##

```

```

##
## Step 5 : AIC = 1461.442
## BodyFat ~ Abdomen + Weight + Wrist + Forearm + Neck
##
## Remove Existing Variables
## -----
## Variable      DF      AIC      Sum Sq      RSS      R-Sq      Adj. R-Sq
## -----
## Neck          1      1462.220    12920.754    4658.236    0.735      0.731
## Wrist         1      1466.596    12839.142    4739.848    0.730      0.726
## Forearm       1      1467.712    12818.118    4760.871    0.729      0.725
## Weight        1      1481.547    12549.427    5029.563    0.714      0.709
## Abdomen       1      1669.382     6980.536    10598.454    0.397      0.387
## -----
##
## Enter New Variables
## -----
## Variable      DF      AIC      Sum Sq      RSS      R-Sq      Adj. R-Sq
## -----
## Age           1      1460.806    13019.755    4559.235    0.741      0.734
## Biceps        1      1460.917    13017.748    4561.242    0.741      0.734
## Thigh         1      1462.065    12996.921    4582.069    0.739      0.733
## Height        1      1462.408    12990.691    4588.299    0.739      0.733
## Hip           1      1462.840    12982.815    4596.175    0.739      0.732
## Ankle         1      1462.858    12982.485    4596.505    0.739      0.732
## Knee          1      1462.873    12982.217    4596.773    0.739      0.732
## Chest         1      1463.441    12971.830    4607.160    0.738      0.731
## -----
##
## - Age added
##
## Step 6 : AIC = 1460.806
## BodyFat ~ Abdomen + Weight + Wrist + Forearm + Neck + Age
##
## Remove Existing Variables
## -----
## Variable      DF      AIC      Sum Sq      RSS      R-Sq      Adj. R-Sq
## -----
## Age           1      1461.442    12971.820    4607.169    0.738      0.733
## Neck          1      1462.136    12959.116    4619.874    0.737      0.732
## Weight        1      1467.800    12854.091    4724.899    0.731      0.726
## Forearm       1      1468.402    12842.791    4736.199    0.731      0.725
## Wrist         1      1468.468    12841.553    4737.436    0.731      0.725
## Abdomen       1      1604.115     9463.485    8115.505    0.538      0.529
## -----
##
## Enter New Variables
## -----
## Variable      DF      AIC      Sum Sq      RSS      R-Sq      Adj. R-Sq
## -----
## Thigh         1      1459.054    13087.141    4491.849    0.744      0.737
## Biceps        1      1460.131    13067.893    4511.097    0.743      0.736
## Height        1      1461.755    13038.740    4540.250    0.742      0.734

```

```

## Ankle          1    1461.988    13034.531    4544.459    0.741    0.734
## Knee           1    1462.444    13026.307    4552.683    0.741    0.734
## Hip            1    1462.628    13022.984    4556.006    0.741    0.733
## Chest          1    1462.760    13020.598    4558.392    0.741    0.733
## -----
##
## - Thigh added
##
##
## Step 7 : AIC = 1459.054
## BodyFat ~ Abdomen + Weight + Wrist + Forearm + Neck + Age + Thigh
##
##                      Remove Existing Variables
## -----
## Variable      DF      AIC      Sum Sq      RSS      R-Sq      Adj. R-Sq
## -----
## Neck           1    1460.490    13025.470    4553.520    0.741    0.735
## Thigh          1    1460.806    13019.755    4559.235    0.741    0.734
## Age            1    1462.065    12996.921    4582.069    0.739    0.733
## Forearm        1    1466.142    12922.190    4656.800    0.735    0.729
## Wrist          1    1466.195    12921.208    4657.782    0.735    0.729
## Weight         1    1469.469    12860.293    4718.697    0.732    0.725
## Abdomen        1    1592.969     9875.999    7702.991    0.562    0.551
## -----
##
##                      Enter New Variables
## -----
## Variable      DF      AIC      Sum Sq      RSS      R-Sq      Adj. R-Sq
## -----
## Hip            1    1458.996    13123.666    4455.324    0.747    0.738
## Biceps         1    1459.620    13112.629    4466.361    0.746    0.738
## Ankle          1    1460.337    13099.909    4479.081    0.745    0.737
## Height         1    1460.811    13091.466    4487.523    0.745    0.736
## Chest          1    1461.011    13087.897    4491.093    0.745    0.736
## Knee           1    1461.054    13087.144    4491.846    0.744    0.736
## -----
##
## - Hip added
##
##
## Step 8 : AIC = 1458.996
## BodyFat ~ Abdomen + Weight + Wrist + Forearm + Neck + Age + Thigh + Hip
##
##                      Remove Existing Variables
## -----
## Variable      DF      AIC      Sum Sq      RSS      R-Sq      Adj. R-Sq
## -----
## Hip            1    1459.054    13087.141    4491.849    0.744    0.737
## Neck           1    1461.431    13044.562    4534.428    0.742    0.735
## Age            1    1461.690    13039.905    4539.085    0.742    0.734
## Weight         1    1462.200    13030.713    4548.277    0.741    0.734
## Thigh          1    1462.628    13022.984    4556.006    0.741    0.733
## Forearm        1    1464.820    12983.180    4595.810    0.739    0.731
## Wrist          1    1466.261    12956.816    4622.173    0.737    0.730

```

```
## Abdomen      1      1592.182      9960.704      7618.285      0.567      0.554
```

```
## -----
```

```
##
```

```
##
```

Enter New Variables

```
## -----
```

```
## Variable      DF      AIC      Sum Sq      RSS      R-Sq      Adj. R-Sq
```

```
## -----
```

```
## Biceps        1      1459.822      13144.377      4434.613      0.748      0.738
```

```
## Height        1      1460.331      13135.415      4443.575      0.747      0.738
```

```
## Ankle         1      1460.338      13135.285      4443.704      0.747      0.738
```

```
## Knee          1      1460.994      13123.702      4455.288      0.747      0.737
```

```
## Chest         1      1460.996      13123.666      4455.324      0.747      0.737
```

```
## -----
```

```
##
```

```
##
```

```
## No more variables to be added or removed.
```

```
##
```

Final Model Output

```
## -----
```

```
##
```

```
##
```

Model Summary

```
## -----
```

```
## R              0.864      RMSE              4.282
```

```
## R-Squared      0.747      Coef. Var      22.359
```

```
## Adj. R-Squared 0.738      MSE              18.335
```

```
## Pred R-Squared 0.725      MAE              3.439
```

```
## -----
```

```
## RMSE: Root Mean Square Error
```

```
## MSE: Mean Square Error
```

```
## MAE: Mean Absolute Error
```

```
##
```

```
##
```

ANOVA

```
## -----
```

```
##              Sum of  
##              Squares      DF      Mean Square      F      Sig.
```

```
## -----
```

```
## Regression    13123.666      8      1640.458      89.473      0.0000
```

```
## Residual      4455.324      243      18.335
```

```
## Total        17578.990      251
```

```
## -----
```

```
##
```

```
##
```

Parameter Estimates

```
## -----
```

```
##      model      Beta      Std. Error      Std. Beta      t      Sig      lower      upper
```

```
## -----
```

```
## (Intercept)  -22.656      11.714      -1.934      0.054      -45.730      0.417
```

```
## Abdomen      0.945      0.072      1.217      0.000      0.803      1.087
```

```
## Weight      -0.090      0.040      -0.316      0.025      -0.168      -0.011
```

```
## Wrist       -1.537      0.509      -0.171      0.003      -2.540      -0.533
```

```
## Forearm      0.516      0.186      0.125      0.006      0.149      0.883
```

```
## Neck        -0.467      0.225      -0.136      0.039      -0.909      -0.024
```

```
## Age          0.066      0.031      0.099      0.034      0.005      0.126
```

```
## Thigh        0.302      0.129      0.190      0.020      0.048      0.557
```

```
## Hip         -0.195      0.138      -0.167      0.159      -0.468      0.077
```

```
## -----
```

After performing AIC, we observe that most of the p-values except that of Hip are significant, which is good. So, we shall remove the variable, Hip to improve our model.

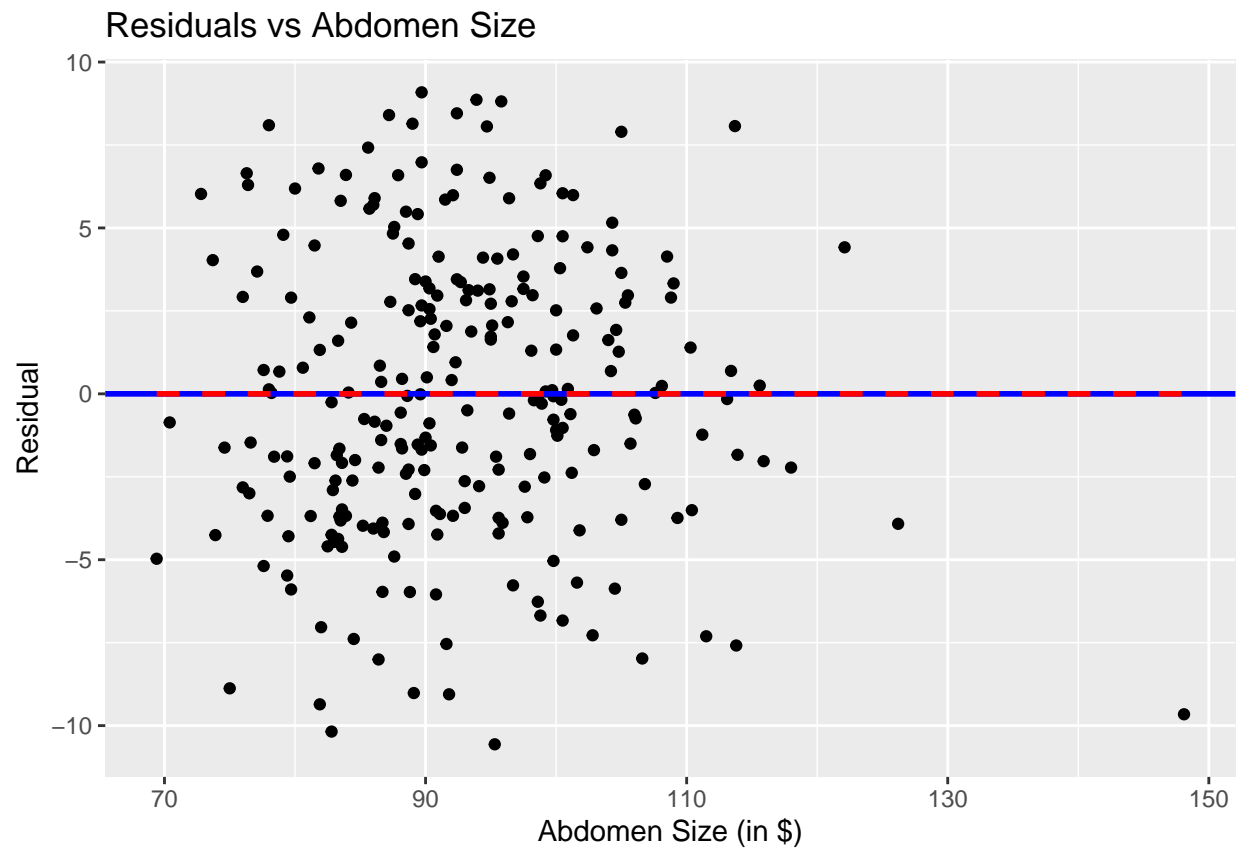
```
##
## Call:
## lm(formula = BodyFat ~ Abdomen + Weight + Wrist + Forearm, data = body_fat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.5626  -3.1235  -0.1461   3.1313   9.0867
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -34.85407     7.24500  -4.811 2.62e-06 ***
## Abdomen         0.99575     0.05607  17.760 < 2e-16 ***
## Weight        -0.13563     0.02475  -5.480 1.05e-07 ***
## Wrist         -1.50556     0.44267  -3.401 0.000783 ***
## Forearm        0.47293     0.18166   2.603 0.009790 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.343 on 247 degrees of freedom
## Multiple R-squared:  0.735, Adjusted R-squared:  0.7307
## F-statistic: 171.3 on 4 and 247 DF, p-value: < 2.2e-16
```

Table 1: Estimates of parameters from fitted model

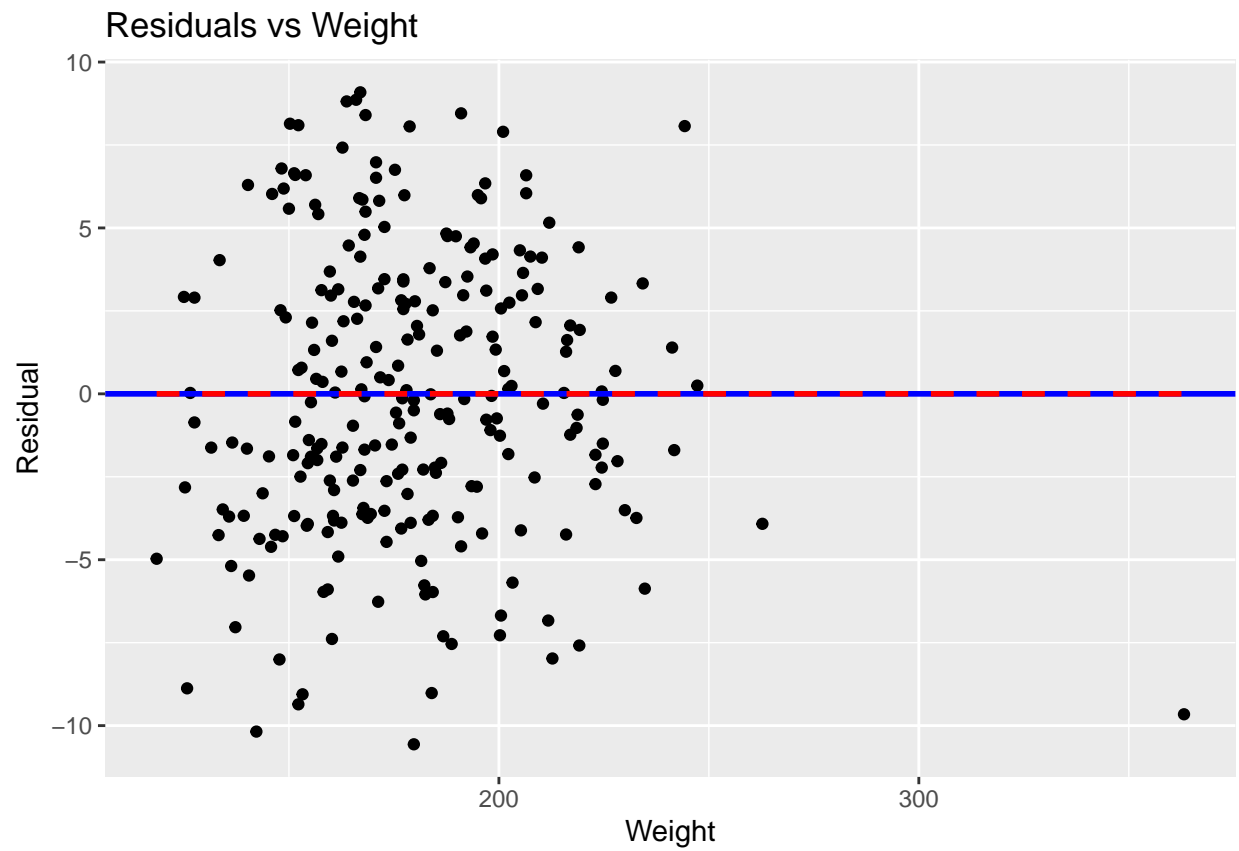
term	estimate	p_value
intercept	-34.854	0.000
Abdomen	0.996	0.000
Weight	-0.136	0.000
Wrist	-1.506	0.001
Forearm	0.473	0.010

The model seems pretty good. Now, let us assess the model fit.

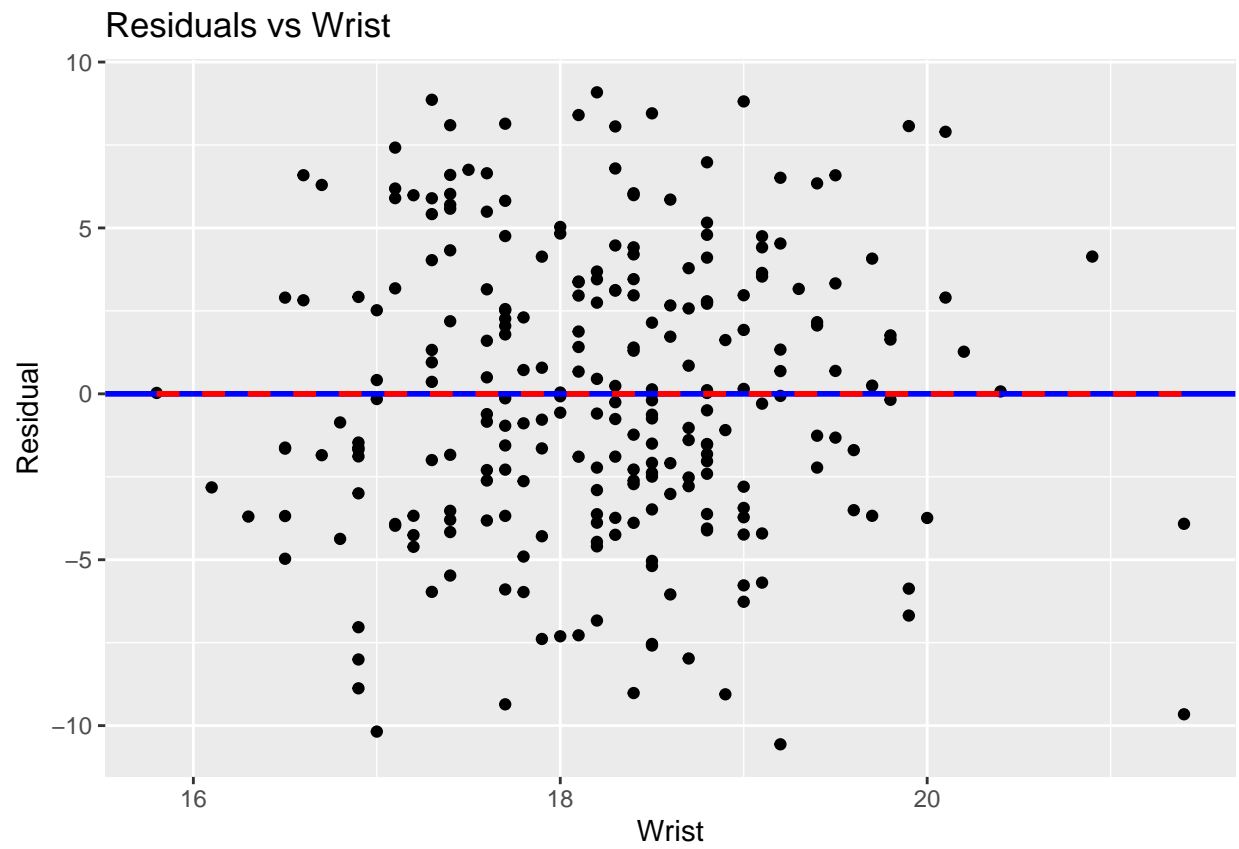
We can assess our first two model assumptions by producing scatterplots of our residuals against each of our explanatory variables. First, let's begin with the scatterplot of the residuals against Abdomen



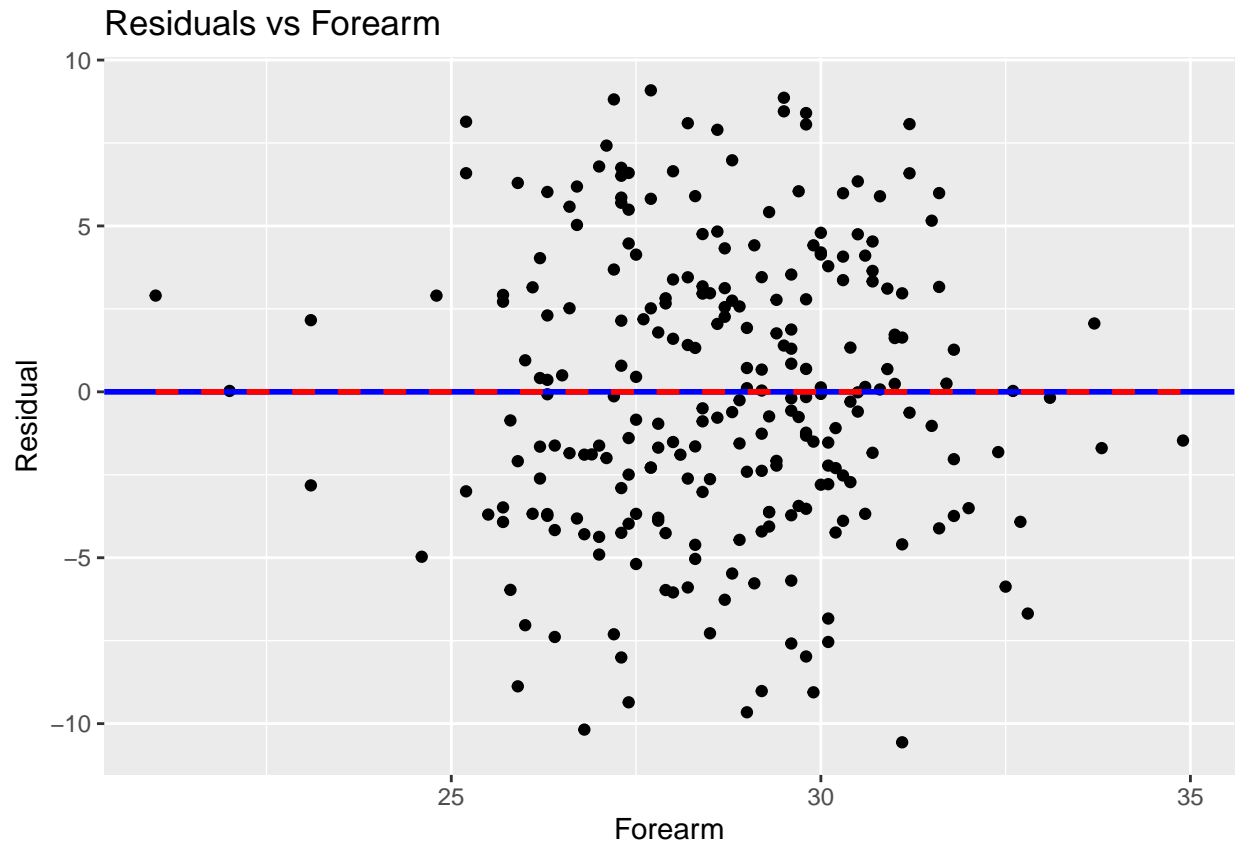
Now, let's plot a scatterplot of the residuals against Weight:



Next, let's plot a scatterplot of the residuals against Wrist:

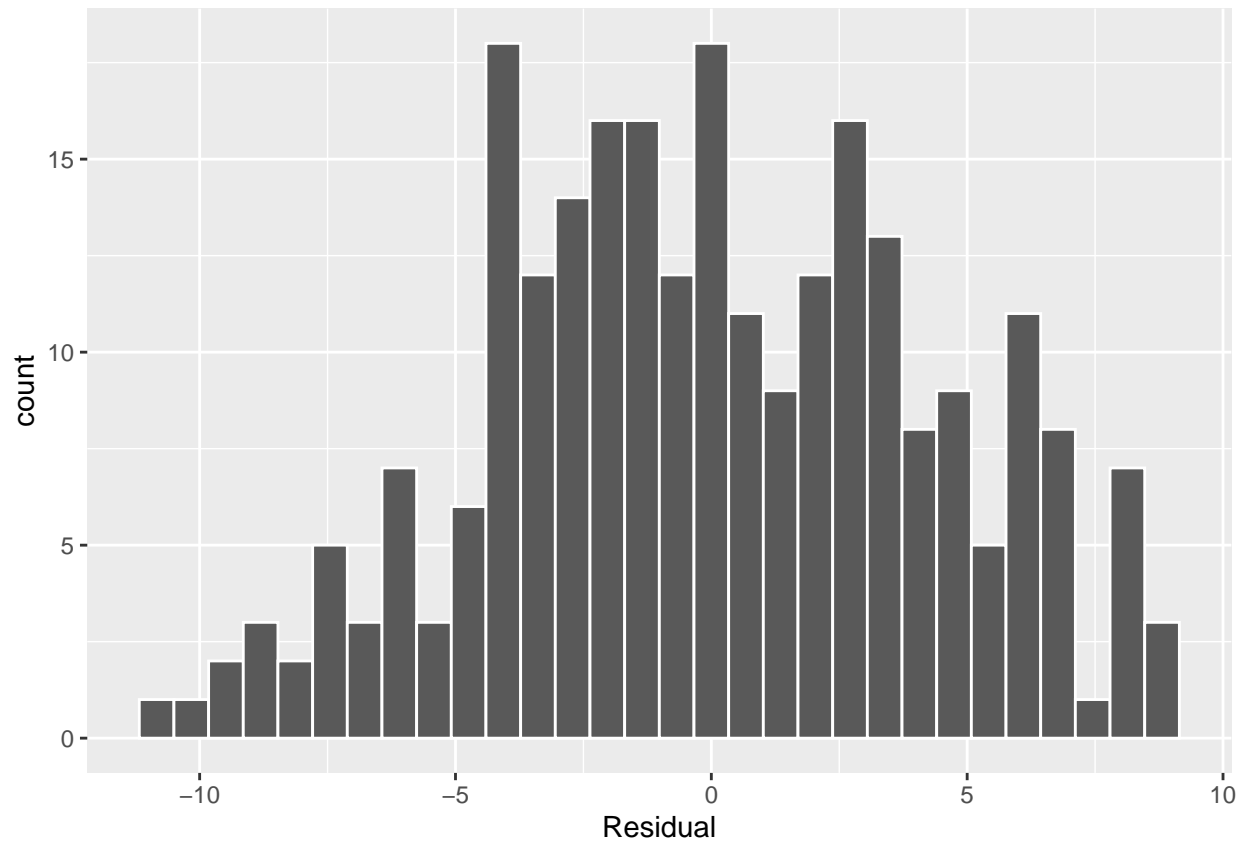


Let's take a look at the scatterplot of the residuals against Forearm:



From the above residuals vs fitted values graphs, we observe that the residuals are randomly scattered around the zero line. This suggests that the residuals have constant variance and mean zero.

Finally, we can check if the residuals are normally distributed by producing a histogram:



We can see from the above plot that the residuals are approximately consistent with a normal distribution and our data roughly fits the bell shaped curve.