a.      Plot `bodyfat` vs. `Height` (code, plot) Which is the dependent variable? Which is the independent variable?
The dependent variable is bodyfat and the independent variable is height.

b.  There is one obvious outlier in the Height column. Remove the corresponding row from the data. (Show: plot, code to remove the row). This will be the data used for the following questions. Confirm that the mean Height is now 70.31076.

c.  Create a linear model of `bodyfat` vs. `Height`. (code, output of summary(model))
    1.  What is the R2 value?
        **The R2 value is 0.008009**
    2.  Is this a "good" model? Why or why not?
        **This is not a good model as the best line of fit is not optimal.**
    3.  What is the linear equation relating bodyfat and Height according to this model?
        **Bodyfat = 33.4945 + Height * -0.2045**

d.  Create a linear model of `bodyfat` vs. **Weight**. (code, output of summary(model))
    1.  What is the R2 value?
        **The R2 value is 0.3751**
    2.  Is this a better model than that based on Height? Why or why not?
        **This model is better than the model based on Height as the best line of fit is optimal.**
    3.  What is the linear equation relating bodyfat and Weight according to this model?
        **Bodyfat = -12.05158 + Weight * 0.17439**
    4.  Plot bodyfat vs. Weight and overlay the best fit line. Use a different color for the line. (plot, code)
    5.  Plot the histogram of residuals (plot, code). Does this show an approximately normal distribution?
        **This does show an approximate normal distribution**
    6.  From the model, predict the bodyfat for two persons: Person A weighs 150 lbs, Person B weighs 300lbs. Include the 99% **confidence** intervals for the predictions. In which prediction (for Person A or B), are you more confident? Why?
        **I'm more confident in Person B as the fit value is 40.3, lower value is 35.7 and upper value is 44.9. Whereas for Person a the fit value is 14.1, lower value is 12.58 and upper value is 15.6.**

e.  Create a linear model of `bodyfat` vs. **Weight and Height**. (code, output of summary(model))
    1.  What is the R2 value?
        **The R2 value is 0.4606.**
    2.  Is this a better model than that based only on Weight or Height? Why or why not?
        **This model is better than the models only on weight or height because it is fitted much better**
    3.  What is the linear equation relating bodyfat, Weight, and Height according to this model?
        **Bodyfat = 32.404646 + 0.20138*Weight – 0.70260*Height**

4. From the model, predict the bodyfat for two persons: Person A weighs 150 lbs, Person B weighs 300lbs. Both persons have height=70". Include the 99% confidence intervals for the predictions. In which prediction (for Person A or B), are you more confident? Why?

**I'm more confident in Person B as the fit value is 43.6, lower value is 39.1 and upper value is 48.1. Whereas for Person a the fit value is 13.4, lower value is 12 and upper value is 14.9.**

f. Add a new transformed variable **BMI** = **Weight/Height$_2$** to the dataset. Create a linear model of `bodyfat` vs. **BMI**.
   1. Give R code, output of summary(model)
   2. Is this a better model than the previous models? Why or why not?
      **This model is not better than the previous model as there is a systematic behavior in the residuals.**
   3. What is the equation relating bodyfat, Weight, and Height according to this model? Is this a linear or nonlinear equation?
      **Bodyfat = 10.715 + BMI * 228.616. According to this model, this seems to be a linear equation**
   4. Plot `bodyfat` vs. `BMI` and overlay the best fit model as a straight line. (code, plot)
   5. From the model, predict the bodyfat for two persons: Person A weighs 150 lbs, Person B weighs 300lbs. Both persons have height=70". Include the 99% confidence intervals for the predictions.

**I'm more confident in Person B as the fit value is 43.6, lower value is 39.1 and upper value is 48.1. Whereas for Person a the fit value is 13.4, lower value is 12 and upper value is 14.9.**

   6. Body Mass Index (BMI) is actually defined as a person's weight in kilograms divided by the square of height in meters but your data has Weight in pounds and Height in inches. Thus, the correct BMI transformation should have been BMI = (Weight/2.20)/(Height*0.0254)$_2$. Would using this correct BMI transformation result in a different model from what was calculated? Why or why not?
      **No this would not result in a different model from what was calculated as it would still produce a linear equation and a relative best line of fit.**

g. Add a new categorical variable (factor) **AgeGroup** to the dataset. AgeGroup should have three values: "Young" for Age<40, "Middle" for Age between 40 and 60, and "Older" for Age>60.
   1. Show R code that adds the AgeGroup variable. This can be done with mutate and the cut() function like so: `cut (Age, breaks = c(-Inf,40,60,Inf), labels = c("Young", "Middle", "Older")`[Code]
   2. Create a linear model of `bodyfat` vs. **BMI and AgeGroup.**[Code, output of summary(model)]
   3. How many dummy (i.e., 0-1) variables were created in the model?
   4. Is this a better model than the previous models? Why or why not?
   5. What are the set of equations relating bodyfat, BMI, and AgeGroup according to this model?
   6. Plot `bodyfat` vs. `BMI` and overlay the model predictions (**multiple** lines: one for each value of the discrete variable). [Code, plot]

```
 1 Call:
 2 lm(formula = bodyfat ~ Height, data = Bodyfat)
 3
 4 Residuals:
 5     Min      1Q   Median      3Q      Max
 6 -19.5902  -6.7124   0.3966   6.0716  27.0919
 7
 8 Coefficients:
 9             Estimate Std. Error t value Pr(>|t|)
10 (Intercept)  33.4945    10.1096   3.313  0.00106 **
11 Height       -0.2045     0.1439  -1.421  0.15664
12 ---
13 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
14
15 Residual standard error: 8.352 on 250 degrees of freedom
16 Multiple R-squared:  0.008009,  Adjusted R-squared:  0.004041
17 F-statistic: 2.019 on 1 and 250 DF,  p-value: 0.1566
18
19 Saving 7 x 7 in image
20
21 Call:
22 lm(formula = bodyfat ~ Weight, data = Bodyfat)
23
24 Residuals:
25     Min      1Q   Median      3Q      Max
26 -17.8164  -4.7430   0.0746   4.9283  21.3605
27
28 Coefficients:
29             Estimate Std. Error t value Pr(>|t|)
30 (Intercept) -12.05158    2.58139  -4.669 4.95e-06 ***
31 Weight        0.17439    0.01424  12.249  < 2e-16 ***
32 ---
33 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
34
35 Residual standard error: 6.629 on 250 degrees of freedom
36 Multiple R-squared:  0.3751,    Adjusted R-squared:  0.3726
37 F-statistic:   150 on 1 and 250 DF,  p-value: < 2.2e-16
38
39 `geom_smooth()` using formula 'y ~ x'
40 Saving 7 x 7 in image
41 `geom_smooth()` using formula 'y ~ x'
42 `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
43 Saving 7 x 7 in image
44 `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
45      fit      lwr      upr
46 1 14.10671 12.58441 15.62901
47      fit      lwr      upr
48 1 40.26499 35.66127 44.86872
49
50 Call:
51 lm(formula = bodyfat ~ Weight + Height, data = Bodyfat)
52
53 Residuals:
54     Min      1Q   Median      3Q      Max
55 -20.0617  -4.1356   0.0361   4.6262  15.9585
56
57 Coefficients:
```
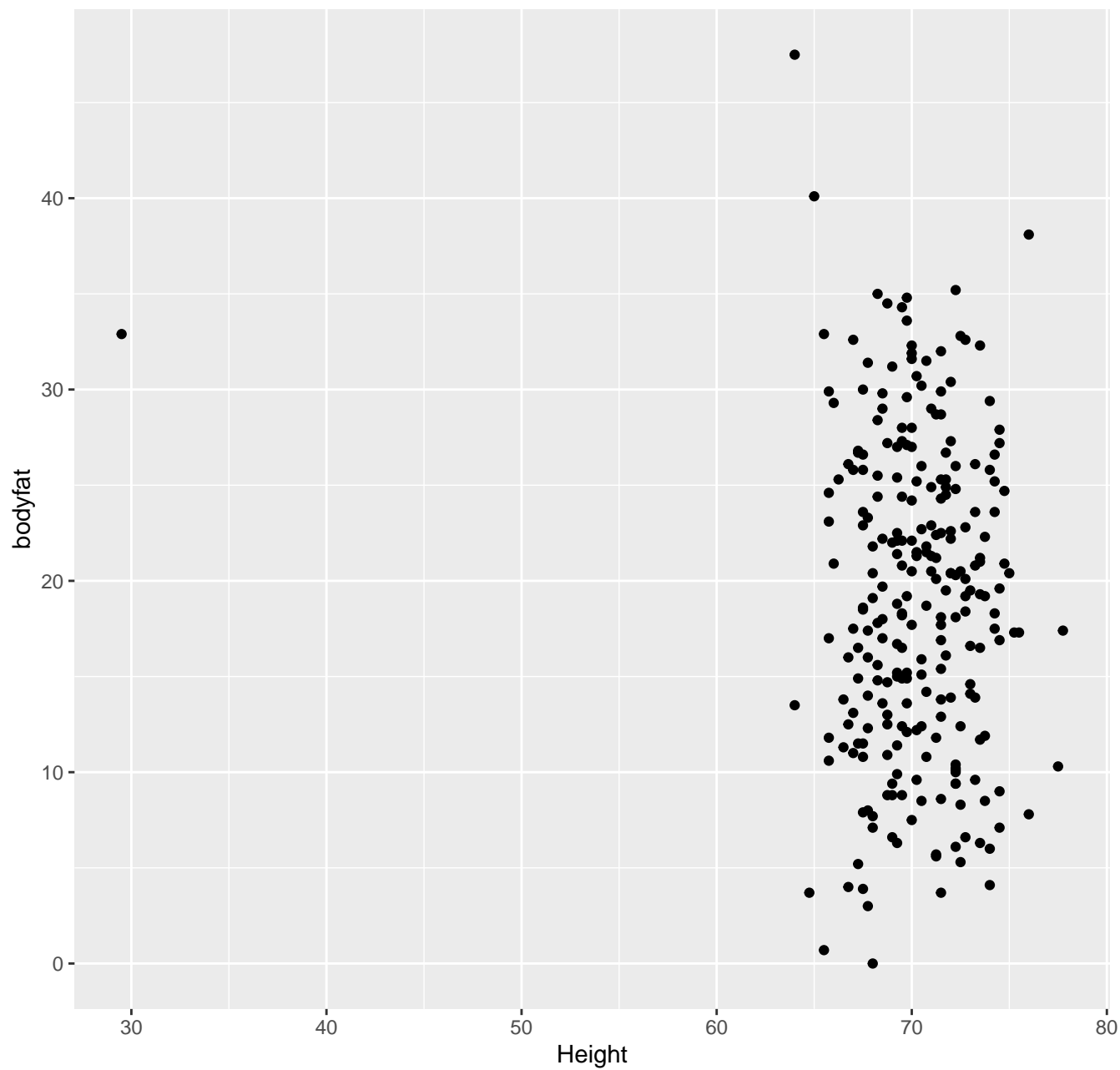
```
58              Estimate Std. Error t value Pr(>|t|)
59 (Intercept) 32.40464    7.46992   4.338 2.09e-05 ***
60 Weight       0.20138    0.01393  14.455  < 2e-16 ***
61 Height      -0.70260    0.11178  -6.285 1.45e-09 ***
62 ---
63 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
64
65 Residual standard error: 6.171 on 249 degrees of freedom
66 Multiple R-squared:  0.4606,    Adjusted R-squared:  0.4563
67 F-statistic: 106.3 on 2 and 249 DF,  p-value: < 2.2e-16
68
69       fit       lwr      upr
70 1 13.43045 11.98608 14.87482
71       fit       lwr      upr
72 1 43.63797 39.13165 48.14429
73 Saving 7 x 7 in image
74
75 Call:
76 lm(formula = bodyfat ~ BMI, data = Bodyfat.copy)
77
78 Residuals:
79    Min      1Q  Median      3Q     Max
80 -31.669  -5.740  -0.169   5.968  24.561
81
82 Coefficients:
83             Estimate Std. Error t value Pr(>|t|)
84 (Intercept)   10.715      1.421   7.540 8.66e-13 ***
85 BMI          228.616     36.147   6.325 1.16e-09 ***
86 ---
87 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
88
89 Residual standard error: 7.786 on 250 degrees of freedom
90 Multiple R-squared:  0.1379,    Adjusted R-squared:  0.1345
91 F-statistic:    40 on 1 and 250 DF,  p-value: 1.162e-09
92
93 Saving 7 x 7 in image
94 `geom_smooth()` using formula 'y ~ x'
95 Saving 7 x 7 in image
96 `geom_smooth()` using formula 'y ~ x'
97       fit       lwr      upr
98 1 13.43045 11.98608 14.87482
99       fit       lwr      upr
100 1 43.63797 39.13165 48.14429
```
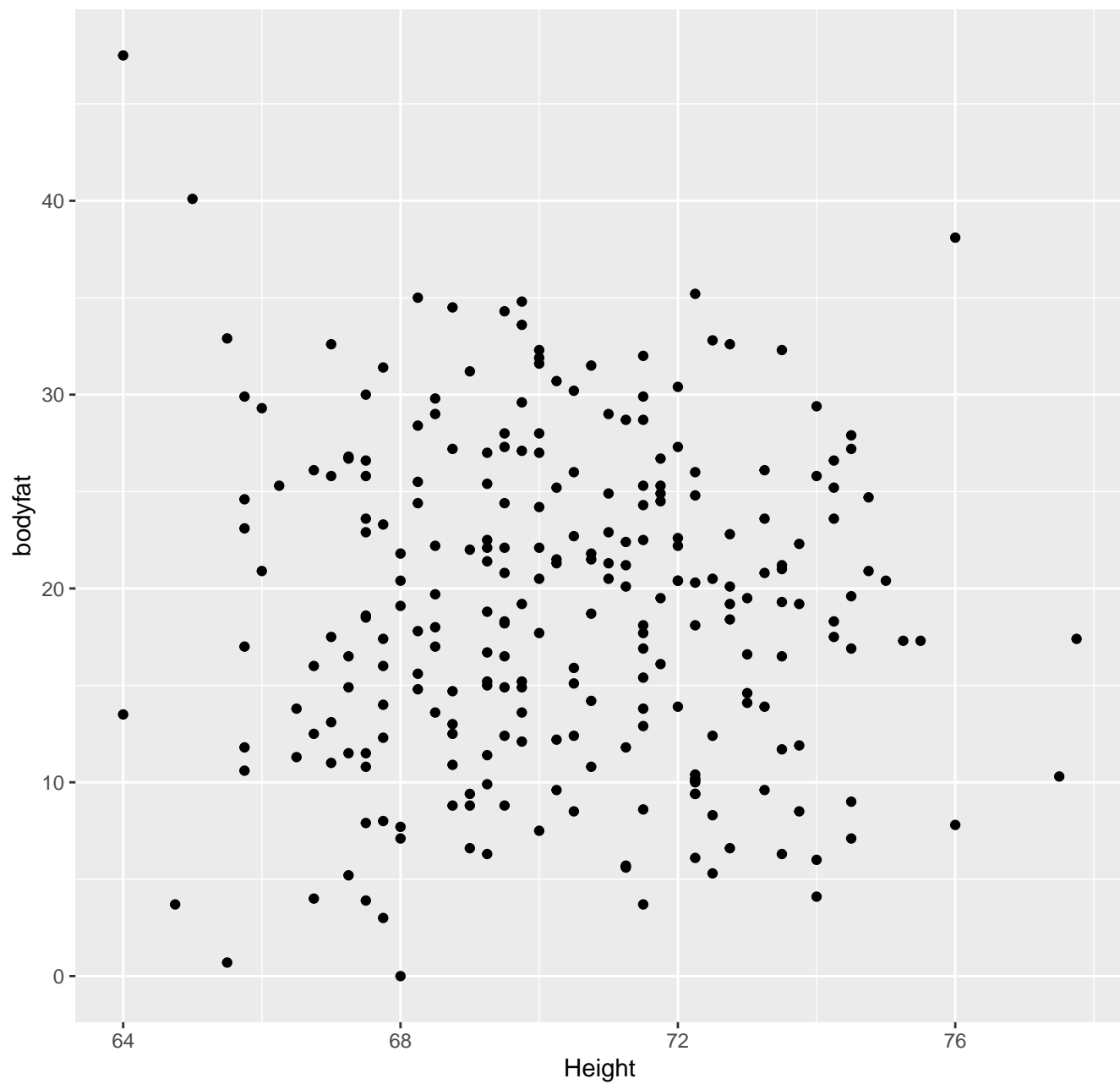
```r
 1 library(ggplot2)
 2 library(tidyverse)
 3 library(modelr)
 4
 5 Bodyfat <- read_csv("Bodyfat.csv")
 6
 7 ggplot(data=Bodyfat)+geom_point(mapping=aes(y=bodyfat,x=Height)) +
 8   ggtitle('a. Plot bodyfat vs. Height')
 9 ggsave("homework#6.1.pdf")
10
11 BodyfatRow <- Bodyfat[-c(42),]
12 ggplot(data=BodyfatRow)+geom_point(mapping=aes(y=bodyfat,x=Height)) +
13   ggtitle('b. Remove the corresponding row from the data')
14 ggsave("homework#6.2.pdf")
15
16 m <- lm(data=Bodyfat, formula=bodyfat~Height)
17 summary(m)
18
19 cf <- coef(m)
20
21 ggplot(data=Bodyfat)+geom_point(mapping=aes(y=bodyfat,x=Weight)) +
22   ggtitle('a. Plot bodyfat vs. Weight')
23 ggsave("homework#6.3.pdf")
24
25 m1 <- lm(data=Bodyfat, formula=bodyfat~Weight)
26 summary(m1)
27
28 ggplot(data = Bodyfat, aes(y = bodyfat,x = Weight)) +
29   geom_point(aes(color = "red")) +
30   geom_smooth(method = "lm") +
31   labs(title = "bodyfat vs Weight best line of it",)
32 ggsave("homework#6.4.pdf")
33
34 hello <- residuals(m1)
35 ggplot(data = Bodyfat) + geom_histogram(aes(x = hello))
36 ggsave("homework#6.5.pdf")
37
38 new.dat <- data.frame(Weight=150)
39 predict(m1, newdata = new.dat, interval = 'confidence', level=0.99)
40
41 new.dat <- data.frame(Weight=300)
42 predict(m1, newdata = new.dat, interval = 'confidence', level=0.99)
43
44 m2 <- lm(data=Bodyfat, formula=bodyfat~Weight+Height)
45 summary(m2)
46
47 ggplot(data=BodyfatRow)+geom_point(mapping=aes(y=Weight,x=Height)) +
48   ggtitle('e.2 Bodyfat vs Weight and Height')
49
50 new.dat <- data.frame(Weight=150, Height=70)
51 predict(m2, newdata = new.dat, interval = 'confidence', level=0.99)
52
53 new.dat <- data.frame(Weight=300, Height=70)
54 predict(m2, newdata = new.dat, interval = 'confidence', level=0.99)
55 ggsave("homework#6.6.pdf")
56
57 Bodyfat.copy <- Bodyfat %>% mutate(BMI = Weight/Height^2)
```

```r
58
59  m3 <- lm(data=Bodyfat.copy, formula=bodyfat~BMI)
60  summary(m3)
61
62  ggplot(data=Bodyfat.copy)+geom_point(mapping=aes(y=bodyfat,x=BMI)) +
63    ggtitle('e.2 Bodyfat VS BMI')
64  ggsave("homework#6.7.pdf")
65
66  ggplot(data = Bodyfat.copy, aes(y = bodyfat,x = BMI)) +
67    geom_point(aes(color = "red")) +
68    geom_smooth(method = "lm") +
69    labs(title = "Bodyfat VS BMI best line of it",)
70  ggsave("homework#6.8.pdf")
71
72  new.dat <- data.frame(Weight=150, Height=70)
73  predict(m2, newdata = new.dat, interval = 'confidence', level=0.99)
74
75  new.dat <- data.frame(Weight=300, Height=70)
76  predict(m2, newdata = new.dat, interval = 'confidence', level=0.99)
77
78  Bodyfat.copy <- Bodyfat %>% mutate(cut(Age, breaks = c(-Inf,40,60,Inf))) %>% mutate(BMI =
    Weight/Height^2)
79
80  m4 <- lm(data=Bodyfat.copy, formula=BMI~breaks)
81  summary(m4)
```
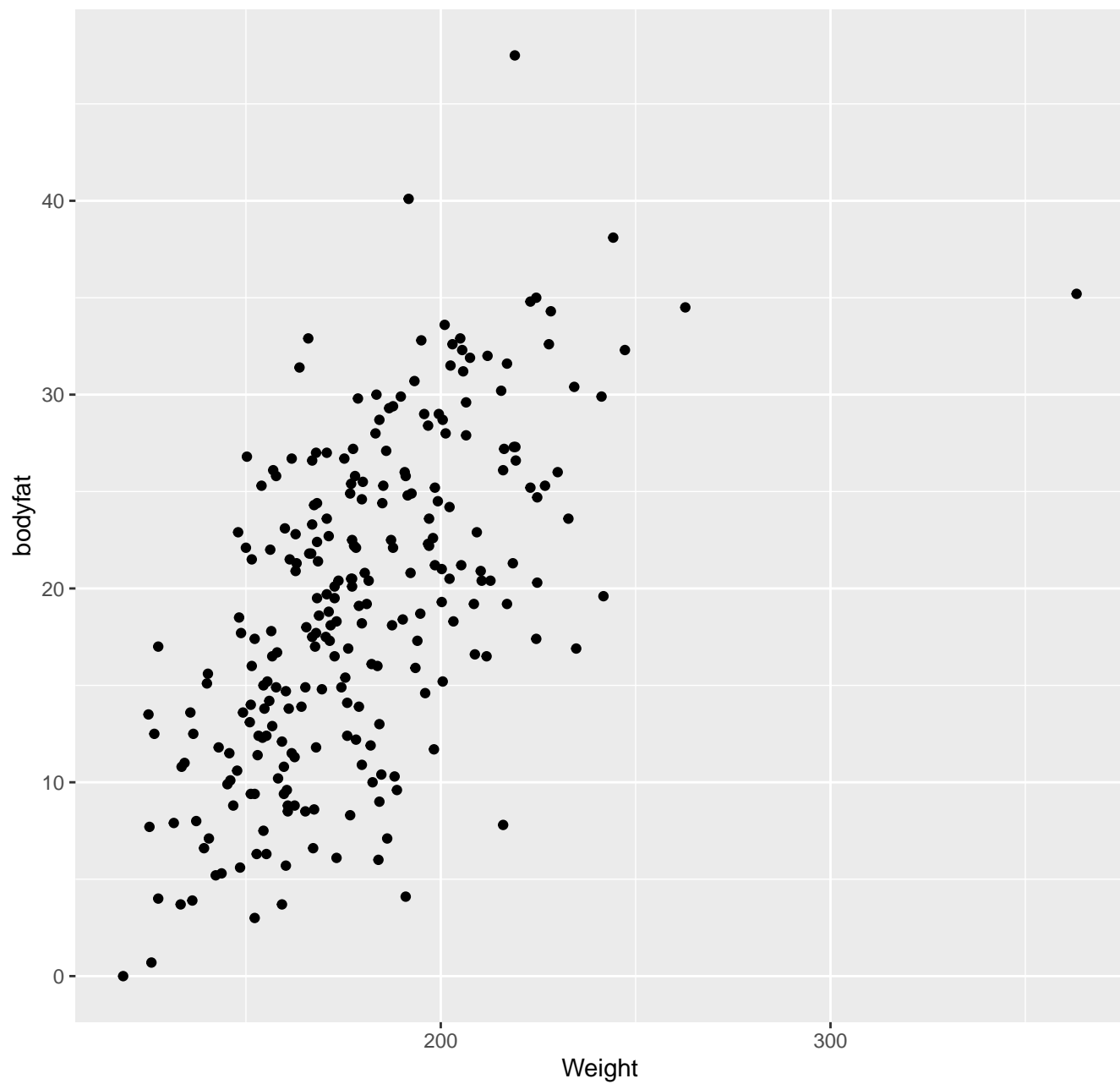
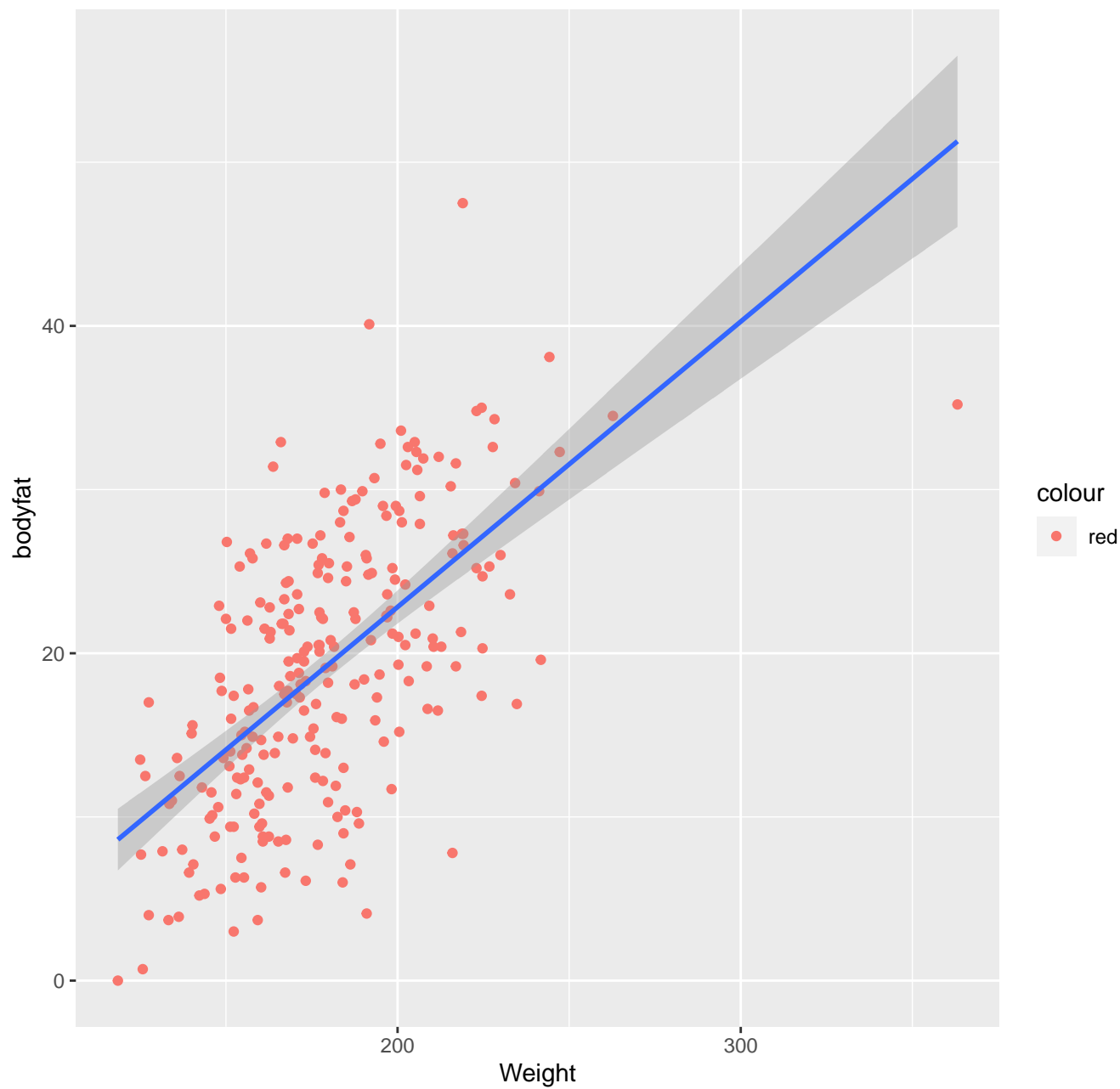a. Plot bodyfat vs. Height

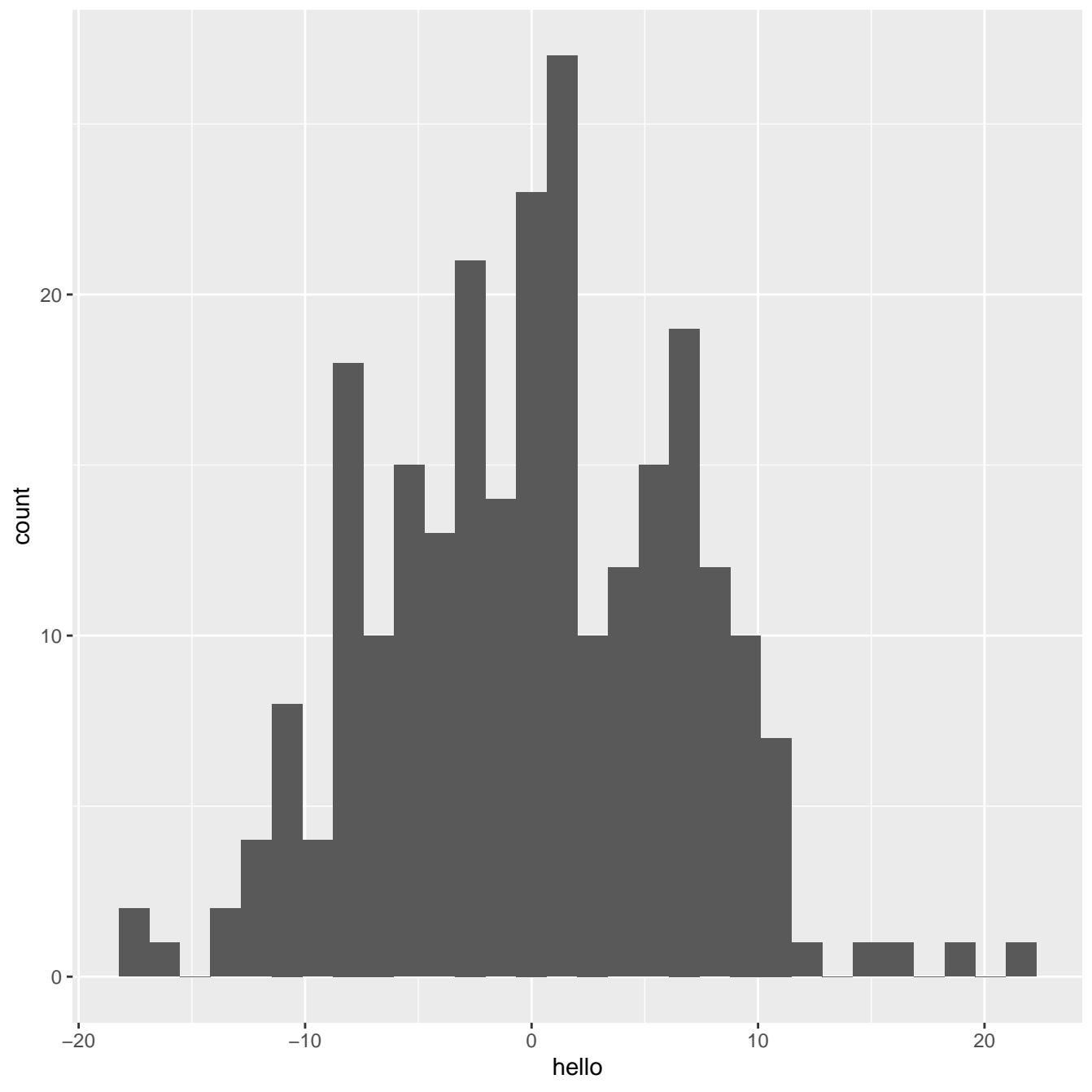# b. Remove the corresponding row from the data
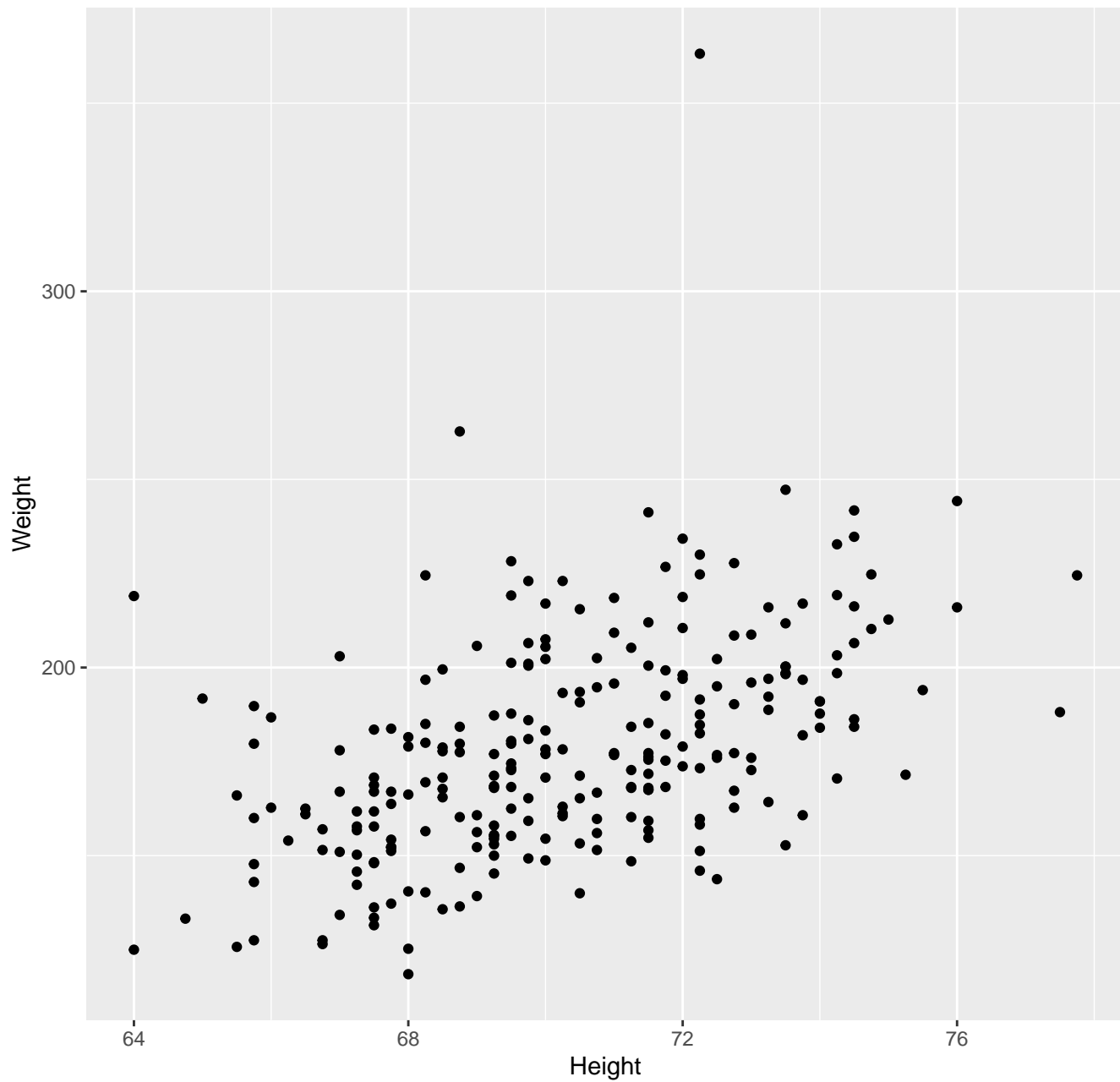
a. Plot bodyfat vs. Weight
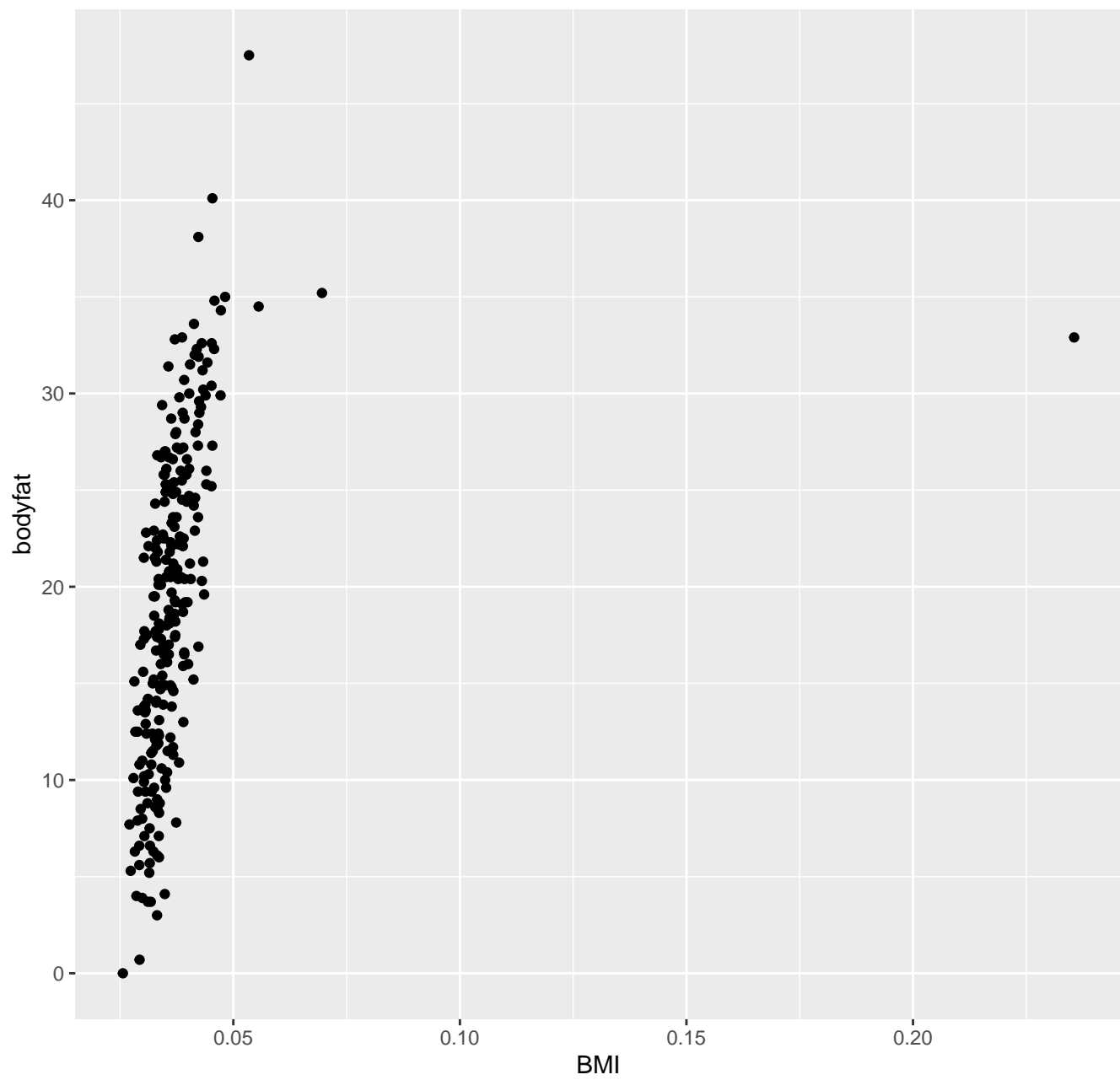
bodyfat vs Weight best line of it

e.2 Bodyfat vs Weight and Height

e.2 Bodyfat VS BMI

Bodyfat VS BMI best line of it