# Assignment 5

## Jacob O

### 2023-11-09

## Question 1

1. The `infmort` data set from the package `faraway` gives the infant mortality rate for a variety of countries. The information is relatively out of date (from 1970s?), but will be fun to graph. Visualize the data using by creating scatter plots of mortality vs income while faceting using `region` and setting color by `oil` export status. Utilize a $\log_{10}$ transformation for both `mortality` and `income` axes. This can be done either by doing the transformation inside the `aes()` command or by utilizing the `scale_x_log10()` or `scale_y_log10()` layers. The critical difference is if the scales are on the original vs log transformed scale. Experiment with both and see which you prefer.
   a) The `rownames()` of the table gives the country names and you should create a new column that contains the country names. *rownames
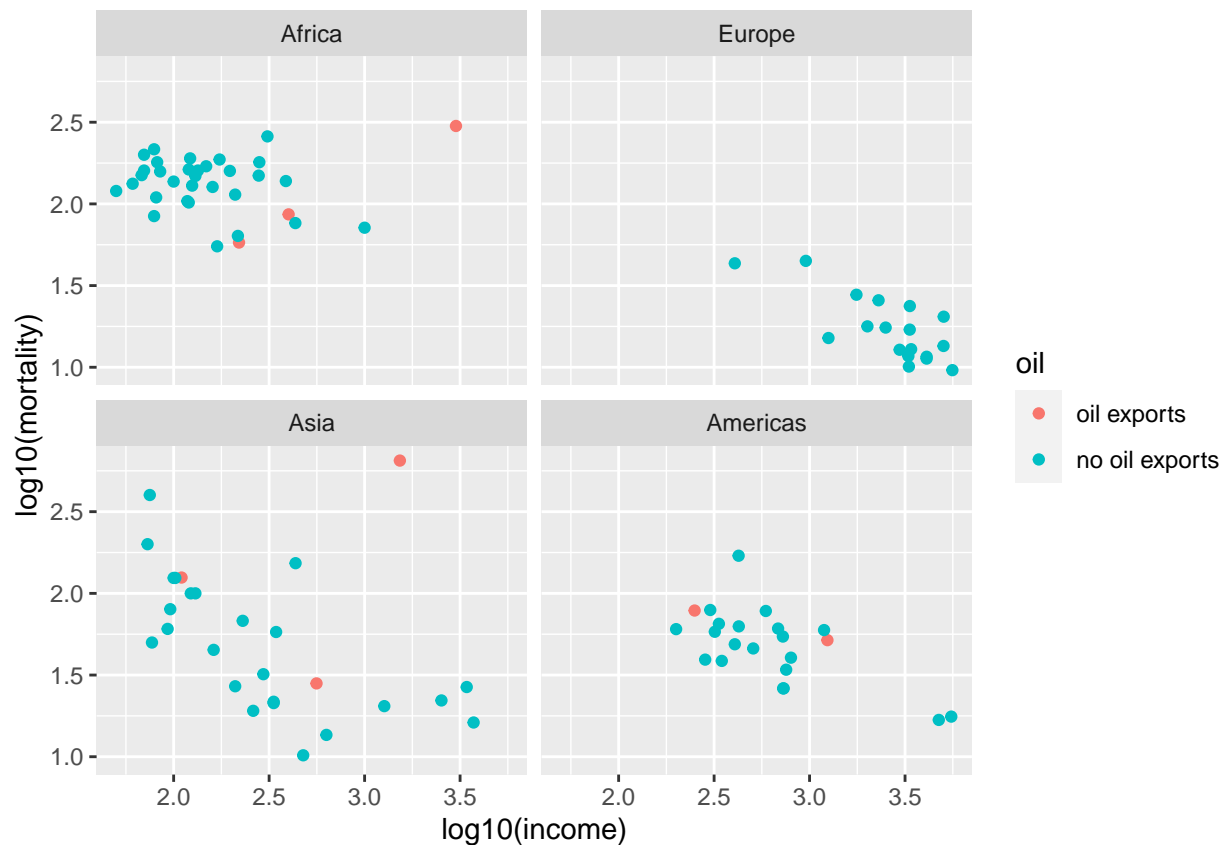
```
data(infmort)
yesssirrr <- infmort %>%
  mutate(rownames = rownames(infmort))
head(yesssirrr)
```

```
##                    region income mortality          oil
## Australia            Asia   3426      26.7 no oil exports
## Austria            Europe   3350      23.7 no oil exports
## Belgium            Europe   3346      17.0 no oil exports
## Canada            Americas   4751      16.8 no oil exports
## Denmark            Europe   5029      13.5 no oil exports
## Finland            Europe   3312      10.1 no oil exports
##                    rownames
## Australia        Australia
## Austria            Austria
## Belgium            Belgium
## Canada              Canada
## Denmark            Denmark
## Finland            Finland
```

b) Create scatter plots with the `log10()` transformation inside the `aes()` command.

```
ggplot(yesssirrr, aes(x=log10(income), y=log10(mortality), color = oil)) +
  geom_point() +
  facet_wrap(vars(region))
```

```
## Warning: Removed 4 rows containing missing values (`geom_point()`).
```
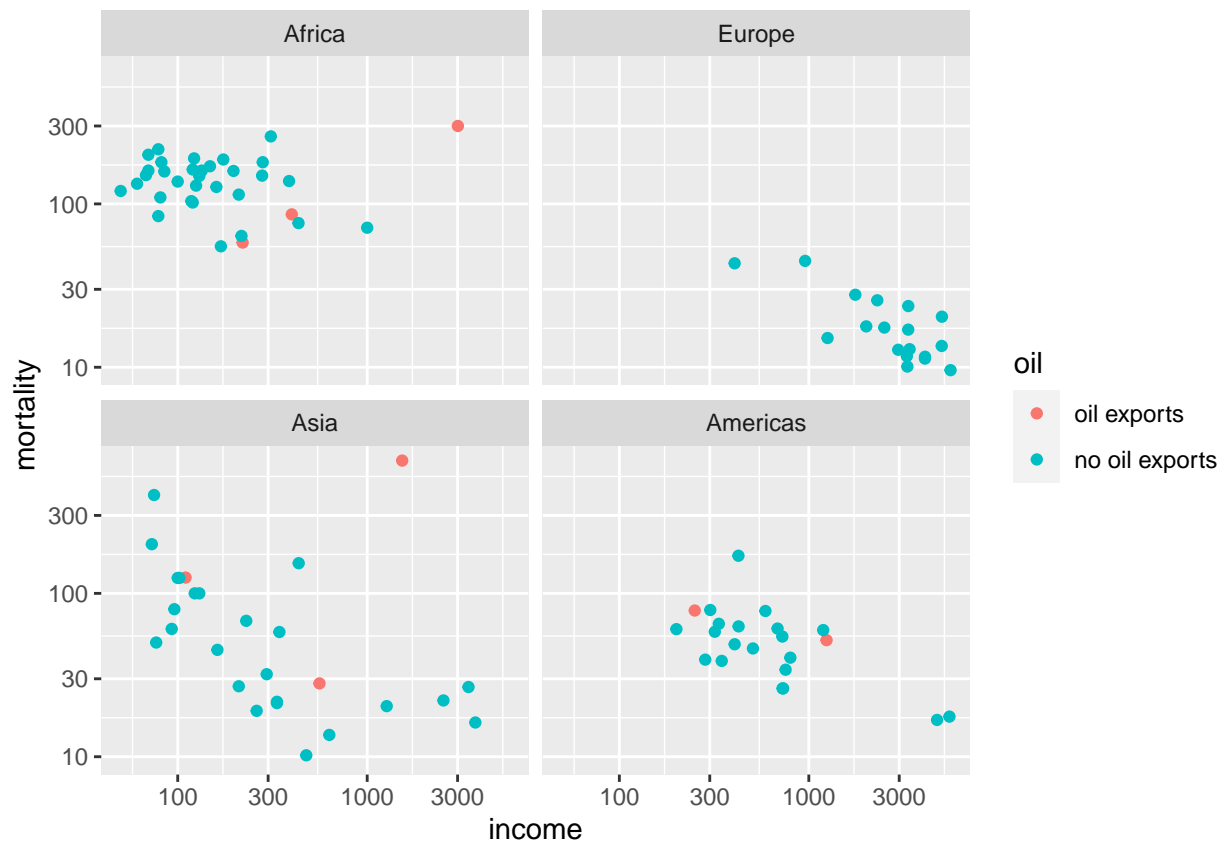
Create the scatter plots using the `scale_x_log10()` and `scale_y_log10()`. Set the major and minor breaks to be useful and aesthetically pleasing. Comment on which version you find easier to read.

```
big.slagga <- ggplot(yesssirrr, aes(x=income, y=mortality, color=oil)) +
  geom_point() +
  scale_x_log10() +
  scale_y_log10() +
  facet_wrap(vars(region))
big.slagga
```

```
## Warning: Removed 4 rows containing missing values (`geom_point()`).
```
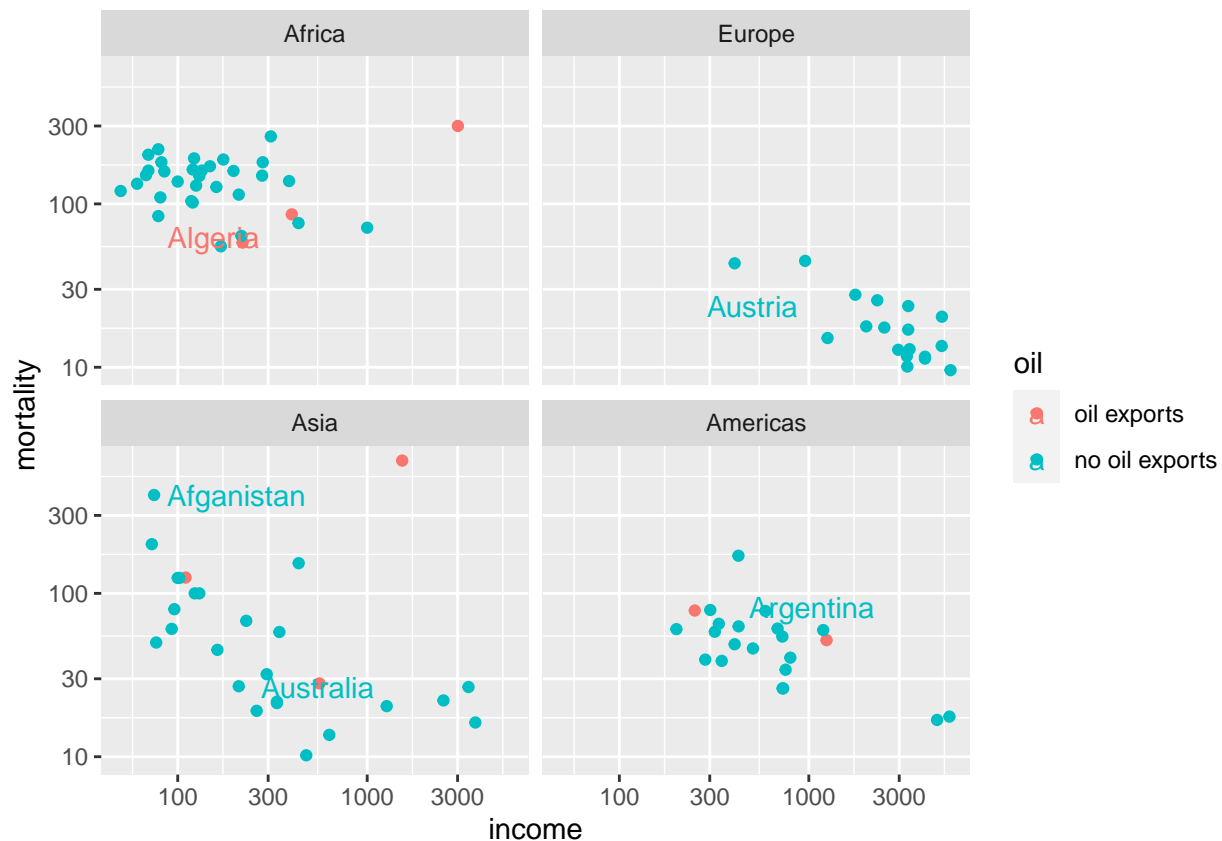
I didnt change the breaks cause I thought they were fine as is. I know how to do it though you just pass breaks = <vector of break #s> or do the same thing for minor_breaks.

d) The package `ggrepel` contains functions `geom_text_repel()` and `geom_label_repel()` that mimic the basic `geom_text()` and `geom_label()` functions in `ggplot2`, but work to make sure the labels don't overlap. Select 10-15 countries to label and do so using the `geom_text_repel()` function.

```
yesssirrr <- yesssirrr %>%
  mutate(Country = str_extract(rownames, pattern= '^[aA].*'))
big.slagga <- ggplot(yesssirrr, aes(x=income, y=mortality, color=oil)) +
  geom_point() +
  scale_x_log10() +
  scale_y_log10() +
  facet_wrap(vars(region)) +
  ggrepel::geom_text_repel(aes(label = Country))
big.slagga
```

```
## Warning: Removed 4 rows containing missing values (`geom_point()`).
```

```
## Warning: Removed 100 rows containing missing values (`geom_text_repel()`).
```

**Question 2**

3. Using the `datasets::trees` data, complete the following:

   a) Create a regression model for $y = $ `Volume` as a function of $x = $ `Height`.

```
data(trees)
str(trees) #FUCK your promise
```

```
## 'data.frame':    31 obs. of  3 variables:
##  $ Girth : num  8.3 8.6 8.8 10.5 10.7 10.8 11 11 11.1 11.2 ...
##  $ Height: num  70 65 63 72 81 83 66 75 80 75 ...
##  $ Volume: num  10.3 10.3 10.2 16.4 18.8 19.7 15.6 18.2 22.6 19.9 ...
```

```
model.guy <- lm(Volume ~ Height, data=trees)
trees <- mutate(trees, yhat=fitted(model.guy))
```

b)  Using the `summary` command, get the y-intercept and slope of the
    regression line.
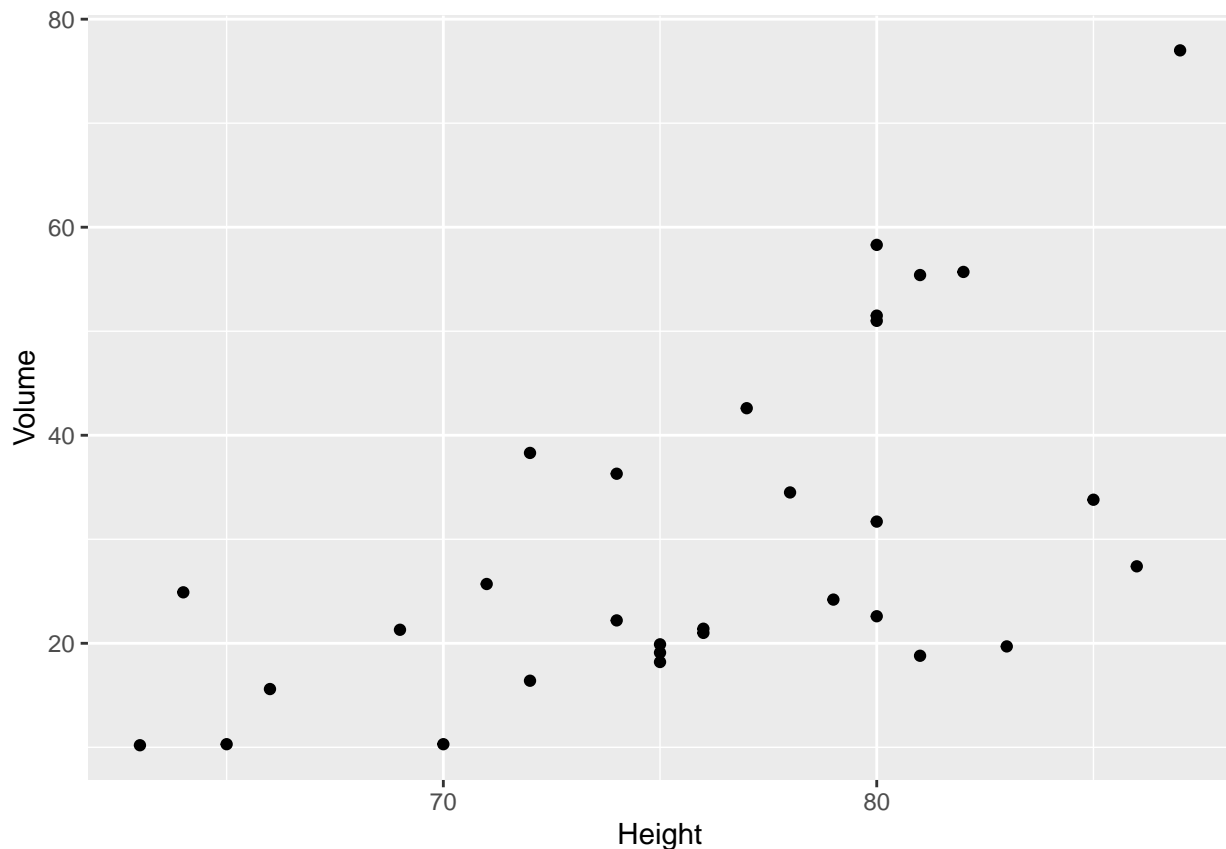
```
summary(model.guy)
```

```
##
## Call:
## lm(formula = Volume ~ Height, data = trees)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -21.274  -9.894  -2.894  12.068  29.852
##
```

4

```
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -87.1236    29.2731  -2.976 0.005835 **
## Height        1.5433     0.3839   4.021 0.000378 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 13.4 on 29 degrees of freedom
## Multiple R-squared:  0.3579, Adjusted R-squared:  0.3358
## F-statistic: 16.16 on 1 and 29 DF,  p-value: 0.0003784
```

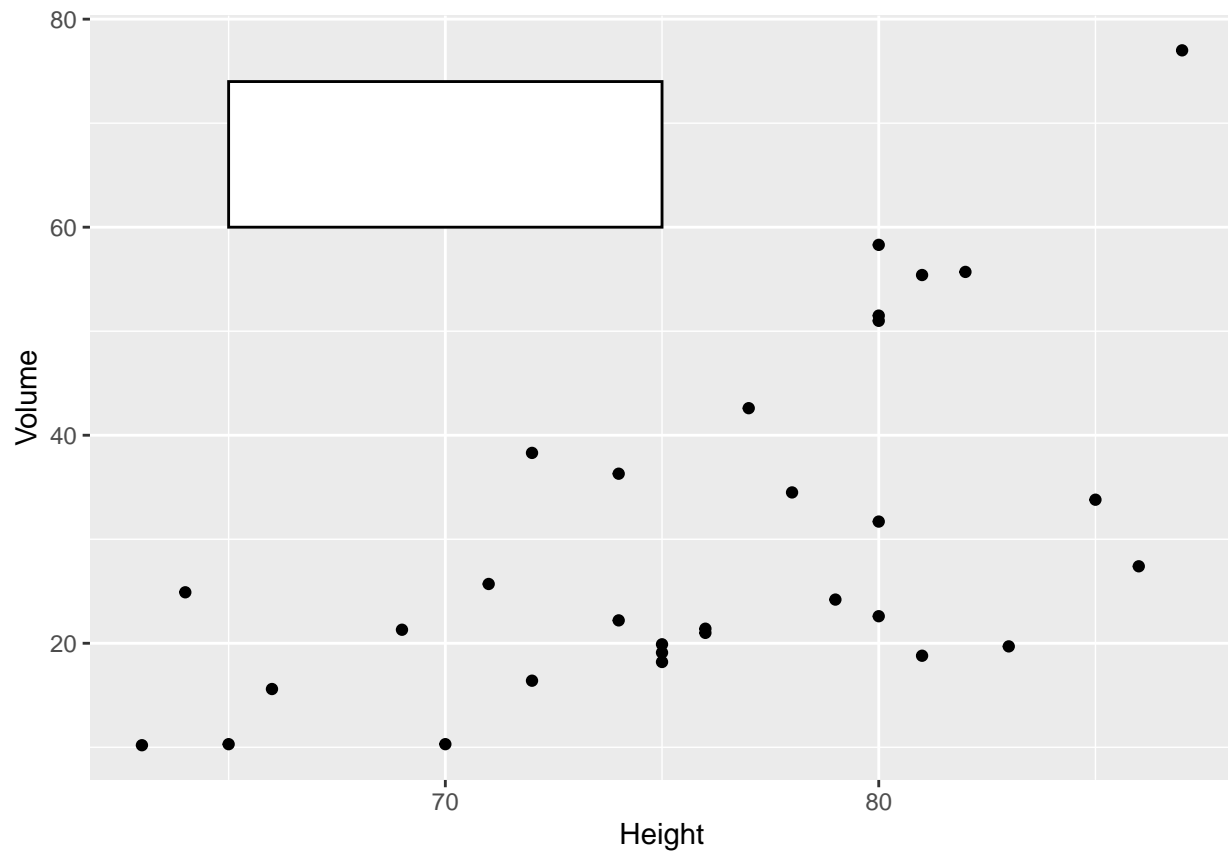c)  Using `ggplot2`, create a scatter plot of Volume vs Height.

```
ggplot(trees, aes(x=Height, y=Volume)) +
  geom_point()
```



d) Create a nice white filled rectangle to add text information to using by adding the following annotation layer.

```r
annotate('rect', xmin=65, xmax=75, ymin=60, ymax=74,
         fill='white', color='black') +
```
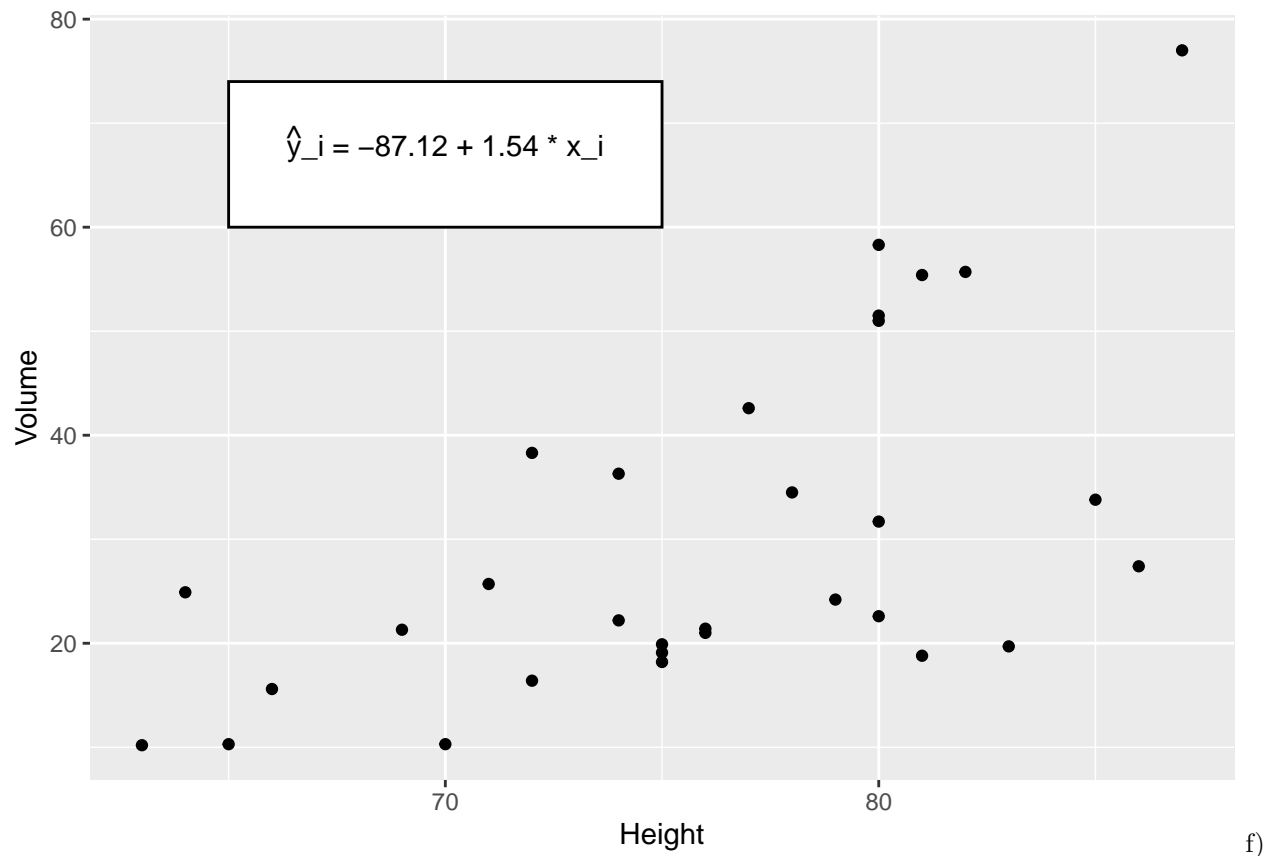
```
ggplot(trees, aes(x=Height, y=Volume)) +
  geom_point() +
  annotate('rect', xmin=65, xmax=75, ymin=60, ymax=74,
           fill='white', color='black')
```

e) Add some annotation text to write the equation of the line
$\hat{y}_i = -87.12 + 1.54 * x_i$ in the text area.

```
ggplot(trees, aes(x=Height, y=Volume)) +
  geom_point() +
  annotate('rect', xmin=65, xmax=75, ymin=60, ymax=74,
              fill='white', color='black') +
  annotate('text', x=70, y= 68, label=latex2exp::TeX('\\hat{y}_i = -87.12 + 1.54 * x_i'))
```

```
## Warning in is.na(x): is.na() applied to non-(list or vector) of type
## 'expression'
```

$\hat{y}\_i = -87.12 + 1.54 * x\_i$
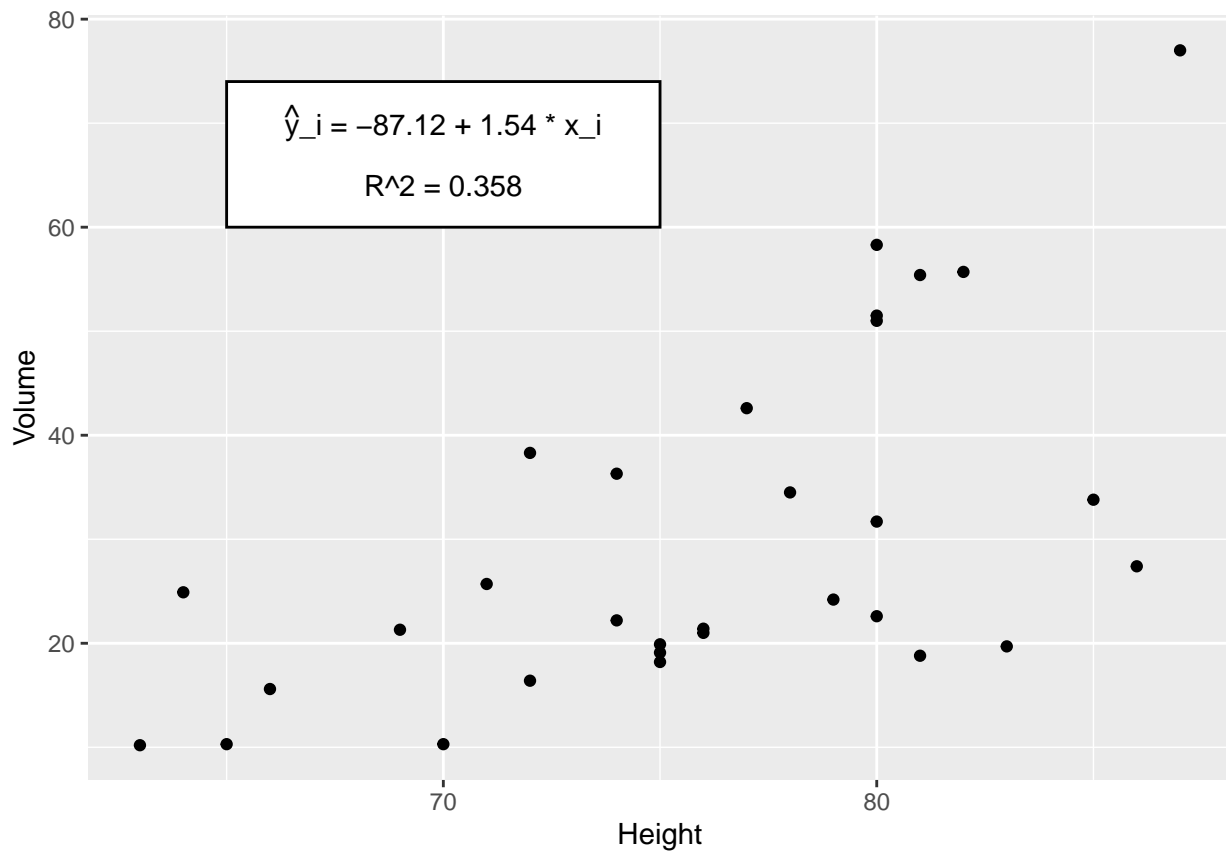
f)

Add annotation to add $R^2 = 0.358$

```
ggplot(trees, aes(x=Height, y=Volume)) +
  geom_point() +
  annotate('rect', xmin=65, xmax=75, ymin=60, ymax=74,
               fill='white', color='black') +
  annotate('text', x=70, y= 70, label=latex2exp::TeX('\\hat{y}_i = -87.12 + 1.54 * x_i')) +
  annotate('text', x= 70, y= 64, label=latex2exp::TeX('R^2 = 0.358'))
```

```
## Warning in is.na(x): is.na() applied to non-(list or vector) of type
## 'expression'
```

```
## Warning in is.na(x): is.na() applied to non-(list or vector) of type
## 'expression'
```

The plot shows a scatter plot with Height on the x-axis and Volume on the y-axis, with an annotation box containing:

$$\hat{y}\_i = -87.12 + 1.54 * x\_i$$

$$R^2 = 0.358$$

g)  Add the regression line in red. The most convenient layer function to uses is `geom_abline()`. It appears that the `annotate` doesn't work with `geom_abline()` so you'll have to call it directly.

```
ggplot(trees, aes(x=Height, y=Volume)) +
  geom_point() +
  annotate('rect', xmin=65, xmax=75, ymin=60, ymax=74,
              fill='white', color='black') +
  annotate('text', x=70, y= 70, label=latex2exp::TeX('\\hat{y}_i = -87.12 + 1.54 * x_i')) +
  annotate('text', x= 70, y= 64, label=latex2exp::TeX('R^2 = 0.358')) +
  geom_line(aes(y=yhat), color='red')
```

```
## Warning in is.na(x): is.na() applied to non-(list or vector) of type
## 'expression'
```

```
## Warning in is.na(x): is.na() applied to non-(list or vector) of type
## 'expression'
```

$$\hat{y}\_i = -87.12 + 1.54 * x\_i$$

$$R\^2 = 0.358$$