

STK4990 - Statistical Methods and Applications

How to write a mandatory assignment

Vinnie Ko

February 14, 2020

This document is an example of how one can write a mandatory assignment in a reasonable way. We solve *Exercise 6* from week 1 as if it is a mandatory assignment that has to be handed in.

For those of you who are not familiar with \LaTeX , but wish to use it, this document also functions as a template that you can follow and modify. Please note that the use of \LaTeX is not compulsory for the mandatory assignment. You are free to use any other word processor (e.g. Microsoft Word). Yet, regardless of which program you use, you have to make sure that your R code, its output and plots are clearly and properly displayed.

Dos and don'ts

- You are only allowed to hand in a single PDF file.
- The PDF file should contain your name, course and assignment number.
- All texts, tables and plots should be clearly readable.

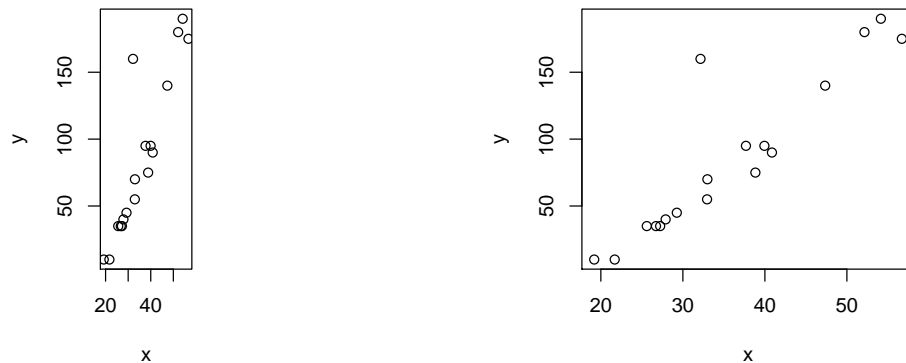


Figure 1: Left: A plot that is badly scaled. Right: Same plot as left, but now properly scaled.

- You should include all the codes that you used, including its output.
- Your code should be placed in the corresponding part of the problems. It should not be included as an appendix at the end of the document.
- Comment your code. So that other people can easily understand your code.
- “normal texts” (e.g. interpretation of the result) should not be written as a comment in R code.
- Answer all the questions that are asked.

Problem 1

a)

I used the following code to create a scatter plot.

```
1 > # Read the data into R from the web.
2 > HERS.data = read.table("http://www.uio.no/studier/emner/matnat/math/STK4900/v17/
   hers.sample.txt", header=T)
3 >
4 > # First observations in the data.
5 > head(HERS.data)
6   age sbp
7 1  75 168
8 2  73 116
9 3  69  96
10 4  64 119
11 5  54 137
12 6  72 163
13 >
14 > # Last observations in the data.
15 > tail(HERS.data)
16   age sbp
17 271  64 126
18 272  56 115
19 273  53 126
20 274  74 156
21 275  66 143
22 276  62 175
23 >
24 > # Create a scatter plot
25 > # Open a pdf device.
26 > pdf("./Exercise_6_a.pdf", width = 6, height = 6)
27 > # Create a scatter plot
28 > plot(
29 +   x=HERS.data[, "age"],
30 +   y=HERS.data[, "sbp"],
31 +   xlab="age",
32 +   ylab="sbp"
33 + )
34 > # Turn off the device.
35 > dev.off()
36 null device
37      1
```

Figure 2 shows the scatter plot between **sbp** (systolic blood pressure) and **age** generated by the code above. We can observe a weak positive correlation between **sbp** and **age**. However, at a first glance, the relationship looks quite weak. We need to run more analysis such as Pearson correlation coefficient and linear regression to determine whether there is a (positive) relationship between the two variables.

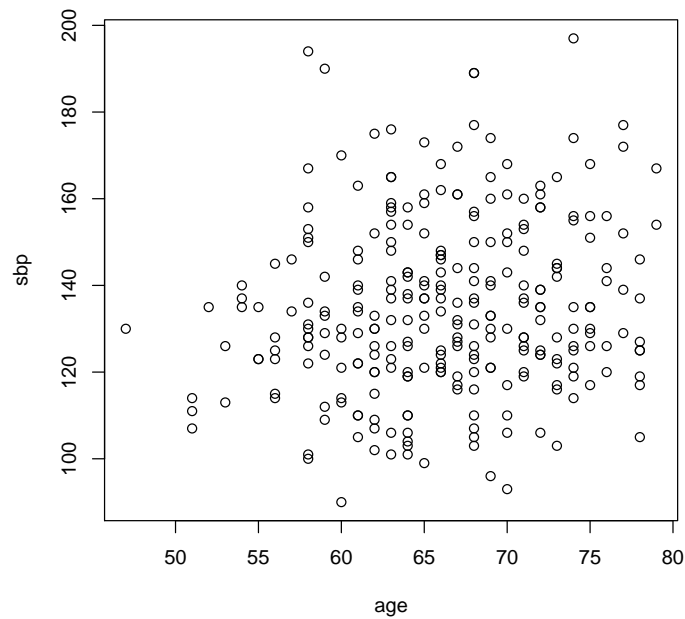


Figure 2: A scatter plot between `sbp` and `age` from a subset of HERS data.

b)

By using `lm()` we fitted a linear regression model to HERS data. As instructed in the problem text, we used `sbp` as outcome variable (y) and `age` as predictor (x).

```

1 > # Fit a linear model.
2 > lm.obj.1 = lm(sbp~age, data=HERS.data)
3 > # Show the result.
4 > summary(lm.obj.1)
5
6 Call:
7 lm(formula = sbp ~ age, data = HERS.data)
8
9 Residuals:
10      Min       1Q   Median       3Q      Max
11 -43.550  -13.520   -2.431   12.578   62.736
12
13 Coefficients:
14             Estimate Std. Error t value Pr(>|t|)
15 (Intercept)  105.7130    12.4024   8.524 1.05e-15 ***
16 age           0.4405     0.1865   2.363  0.0188 *
17 ---
18 Signif. codes:
19
20 Residual standard error: 19.76 on 274 degrees of freedom
21 Multiple R-squared:  0.01997, Adjusted R-squared:  0.01639
22 F-statistic: 5.582 on 1 and 274 DF, p-value: 0.01884
23
24 >
25 > # Extract estimated parameters
26 > beta.0.hat = lm.obj.1$coefficients[1]
27 > names(beta.0.hat) = NULL
28 > beta.1.hat = lm.obj.1$coefficients[2]
29 > names(beta.1.hat) = NULL
30 >

```

```

31 > # Create a scatter plot with fitted regression line
32 > # Open a pdf device.
33 > pdf("./Exercise_6_b.pdf", width = 6, height = 6)
34 > # Create a scatter plot
35 > plot(
36 +   x=HERS.data[, "age"],
37 +   HERS.data[, "sbp"],
38 +   xlab="age",
39 +   ylab="sbp"
40 + )
41 > # Plot the fitted linear model.
42 > abline(a=beta.0.hat, b=beta.1.hat, col="blue")
43 > # Turn off the device.
44 > dev.off()
45 null device
46      1
47 >

```

From the summary of the fitted linear model (`lm.obj.1`), we can read off the following estimated parameter values: $\hat{\beta}_0 = 105.7130$, $\hat{\beta}_1 = 0.4405$.

The fitted regression line is visualized with the original data points in Figure 3.

The slope ($\hat{\beta}_1$) can be interpreted as “expected change in *systolic blood pressure* (`sbp`) when the age of a woman increases by 1 year”. **We also have to mention what the unit of `sbp` is, but that information is not given in the exercise text.**

Further, we see, from the summary of the fitted linear model, that the p -value corresponding to $\hat{\beta}_1$ is 0.0188. This p -value comes from the t -test with following null and alternative hypothesis:

$$H_0 : \beta_1 = 0, \quad H_A : \beta_1 \neq 0.$$

If we would use the significance level of $\alpha = 0.05$, we can reject H_0 and conclude that age has a significant (linear) effect on systolic blood pressure.

However, if we use the significance level of $\alpha = 0.01$, we have $0.0188 > 0.01$ and we cannot reject H_0 . In this case, we conclude that age has no significant (linear) effect on systolic blood pressure.

The intercept ($\hat{\beta}_0$) can be interpreted as “expected value of *systolic blood pressure* (y) when age (x) is set to 0”.

However, this interpretation is not meaningful since age cannot be 0.

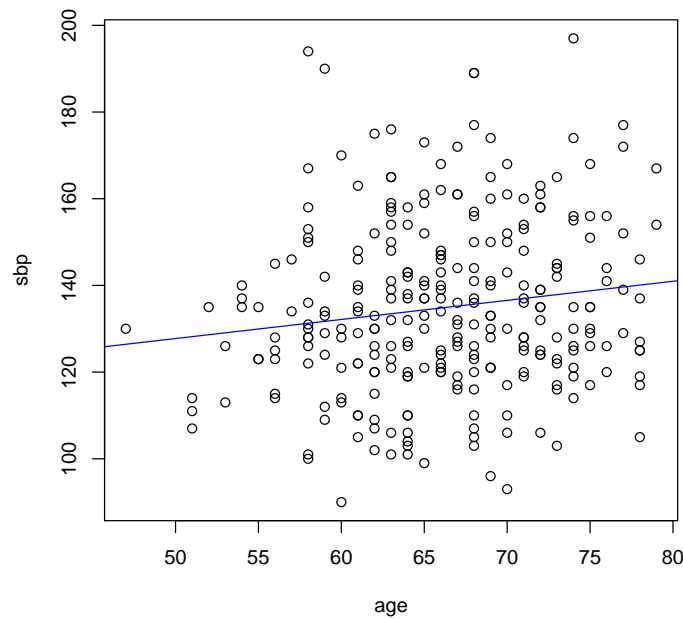


Figure 3: A scatter plot with fitted linear regression line from Problem 1 b).

c)

To get more meaningful interpretation of the intercept β_0 , we subtract 67 from the predictor (`age`) and fit linear regression with the following code.

```

1 > # Fit a linear regression with "age - 67" as predictor.
2 > lm.obj.2 = lm(sbp~I(age-67), data=HERS.data)
3 > # Show the result.
4 > summary(lm.obj.2)
5
6 Call:
7 lm(formula = sbp ~ I(age - 67), data = HERS.data)
8
9 Residuals:
10      Min       1Q   Median       3Q      Max
11 -43.550 -13.520  -2.431  12.578  62.736
12
13 Coefficients:
14             Estimate Std. Error t value Pr(>|t|)
15 (Intercept)  135.2284    1.1985  112.829  <2e-16 ***
16 I(age - 67)    0.4405    0.1865   2.363   0.0188 *
17 ---
18 Signif. codes:
19
20 Residual standard error: 19.76 on 274 degrees of freedom
21 Multiple R-squared:  0.01997, Adjusted R-squared:  0.01639
22 F-statistic: 5.582 on 1 and 274 DF, p-value: 0.01884
23
24 >
25 > # Summary of the result.
26 > result.table = data.frame(
27 +   lm.1 = lm.obj.1$coefficients,
28 +   lm.2 = lm.obj.2$coefficients
29 + )
30 > rownames(result.table) = c("beta.0.hat", "beta.1.hat")

```

```

31 > result.table
32           lm.1           lm.2
33 beta.0.hat 105.7129691 135.2283886
34 beta.1.hat  0.4405286  0.4405286

```

From the output, we have following least squares estimates.

	not centered	centered
$\hat{\beta}_0$	105.71	135.23
$\hat{\beta}_1$	0.44	0.44

Tip: <https://www.tablesgenerator.com/> is a handy website if you want to draw a table in L^AT_EX.

We can see that centering changed the intercept ($\hat{\beta}_0$), but the slope ($\hat{\beta}_1$) remains unchanged. This is logical since centering is nothing more than “shifting” all data points with the same amount.

The intercept ($\hat{\beta}_0$) can now be interpreted as “expected value of *systolic blood pressure* (y) when centered age ($x - 67$) is set to 0”.

In other words, the intercept can now be interpreted as “expected value of *systolic blood pressure* (y) when age (x) is set to 67”.

This is a meaningful interpretation, without any biological nonsense like in a).

d)

We fit a new linear regression with $\frac{\text{age}}{10}$ as predictor by using the following code.

```

1 > # Fit a linear regression with age/10 as predictor.
2 > lm.obj.3 = lm(sbp~I(age/10), data=HERS.data)
3 > # Show the result.
4 > summary(lm.obj.3)
5
6 Call:
7 lm(formula = sbp ~ I(age/10), data = HERS.data)
8
9 Residuals:
10      Min       1Q   Median       3Q      Max
11 -43.550 -13.520  -2.431  12.578  62.736
12
13 Coefficients:
14             Estimate Std. Error t value Pr(>|t|)
15 (Intercept)  105.713     12.402   8.524 1.05e-15 ***
16 I(age/10)     4.405       1.865   2.363  0.0188 *
17 ---
18 Signif. codes:
19
20 Residual standard error: 19.76 on 274 degrees of freedom
21 Multiple R-squared:  0.01997, Adjusted R-squared:  0.01639
22 F-statistic: 5.582 on 1 and 274 DF, p-value: 0.01884
23
24 >
25 > # Summary of the result.
26 > result.table = cbind(result.table, lm.3 = lm.obj.3$coefficients)
27 > result.table
28           lm.1           lm.2           lm.3
29 beta.0.hat 105.7129691 135.2283886 105.712969
30 beta.1.hat  0.4405286  0.4405286  4.405286
31 >

```

From the output, we have following least squares estimates.

	age	age/10
$\hat{\beta}_0$	105.71	105.71
$\hat{\beta}_1$	0.44	4.41

We can see that intercept remains unchanged. This is because we only “scaled” the predictor (x) and did not “shift” it.

For the slope ($\hat{\beta}_1$), we can observe that it is increased by the factor of 10.

On slide 43 from Lecture 2. $\hat{\beta}_1$ is defined as

$$\hat{\beta}_1 = \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{s_x^2}.$$

Dividing age by 10 is equal to: replacing x_i by $\frac{x_i}{10}$, replacing \bar{x} by $\frac{\bar{x}}{10}$ and replacing s_x by $\frac{s_x}{10}$.

So, the slope becomes

$$\begin{aligned} \frac{\frac{1}{n-1} \sum_{i=1}^n (\frac{1}{10}x_i - \frac{1}{10}\bar{x})(y_i - \bar{y})}{(\frac{1}{10}s_x)^2} &= \frac{\frac{1}{n-1} 10^2 \sum_{i=1}^n \frac{1}{10}(x_i - \bar{x})(y_i - \bar{y})}{s_x^2} \\ &= \frac{\frac{1}{n-1} 10^2 \cdot \frac{1}{10} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{s_x^2} \\ &= 10 \cdot \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{s_x^2} \\ &= 10 \hat{\beta}_1. \end{aligned}$$

Hence, we can see that dividing the predictor (x) by 10 results in that the slope increases by 10.

THE END