

# Computer Lab 1 Computational Statistics - Report

## Group 6

*Zijie Feng (zijfe244), Phillip Hoelscher (phiho267)*

*29 1 2019*

## Contents

<b>Question 1: Be careful when comparing</b>	<b>1</b>
1. Check the results of the snippets. Comment what is going on. . . . .	2
2. If there are any problems, suggest improvements. . . . .	2
<b>Question 2: Derivative</b>	<b>2</b>
1. Write your own R function . . . . .	2
2. Evaluate your derivative function at $x = 1$ and $x = 100000$ . . . . .	3
3. What values did you obtain? What are the true values? Explain the reasons behind the discovered differences. . . . .	3
<b>Question 3: Variance</b>	<b>3</b>
1. Write your own R function, myvar, to estimate the variance in this way. . . . .	3
2. Generate a vector $x = (x_1, \dots, x_{10000})$ with 10000 random numbers with mean $10^8$ and variance 1. . . . .	3
3. Compute the difference $Y_i = myvar(X_i) - var(X_i)$ . . . . .	4
4. Better implementation of variance estimator . . . . .	5
<b>Question 4: Linear Algebra</b>	<b>6</b>
1. Import the data set to R . . . . .	6
2. Compute optimal regression coefficients . . . . .	6
3. Solve with default solver <i>solve()</i> . . . . .	6
4. Check the condition number of the matrix A (function kappa()) and consider how it is related to your conclusion in step 3. . . . .	6
5. Scale the data set and repeat steps 2-4. How has the result changed and why? . . . . .	7

## Question 1: Be careful when comparing

Consider the following two R code snippets

```
# snippets 1
x1 = 1/3 ; x2 = 1/4
if(x1 - x2 == 1/12){
  print("Subtraction is correct")
} else{
  print("Subtraction is wrong")
}
```

```
## [1] "Subtraction is wrong"
```

```
# snippets 2
x1 = 1 ; x2 = 1/2
if(x1 - x2 == 1/2){
  print("Subtraction is correct")
}
```

```

} else{
print("Subtraction is wrong")
}

```

```
## [1] "Subtraction is correct"
```

## 1. Check the results of the snippets. Comment what is going on.

We can see that the output of the two snippets is different.

- 1- Snippet:  $x_1$  is stored with the value  $1/3$ , which is a consecutive number in decimal form (0.33333...). Especially floating point number are not stored with the “correct” value, but one that is very similar to the “real” value. This is due to the memory capacity of numbers in 32 or 64 bits. Thus information is lost due to regression, which leads to a rounding error. This is also called underflow.
- 2 - snippet: In this case the output shows that the calculation is correct. In this case, no rounding error occurred due to memory optimization.

## 2. If there are any problems, suggest improvements.

```

# snippets 1 - improved
x1 = 1/3 ; x2 = 1/4
if(isTRUE(all.equal(x1 - x2, 1/12))){
  print("Subtraction is correct")
} else{
  print("Subtraction is wrong")
}

```

```
## [1] "Subtraction is correct"
```

One way to fix the output error of snippet 1 is to use the function *all.equal*. This function tests if two objects are “near equality”. We can now see from the output that the calculation is correct.

## Question 2: Derivative

From the definition of a derivative a popular way of computing it at a point  $x$  is to use a small  $\epsilon$  and the formula:

$$f'(x) = \frac{f(x + \epsilon) - f(x)}{\epsilon}$$

### 1. Write your own R function

to calculate the derivative of  $f(x) = x$  in this way with  $\epsilon = 10^{-15}$

The function:

```

epsilon = 10**-15
derivate = function(x){
  deriv = ((x + epsilon) - x) / epsilon
  return(deriv)
}

```

2. Evaluate your derivative function at  $x = 1$  and  $x = 100000$ .

```
x1= 1
x2= 100000
derivate(x1)

## [1] 1.110223

derivate(x2)

## [1] 0
```

3. What values did you obtain? What are the true values? Explain the reasons behind the discovered differences.

Since  $f(x) = x$ , the true value should be 1. The reason why the results are wrong might be:

When  $x_1 = 1$ ,  $x_2 = 10^{-15}$  exceeds the expression ability of  $R$ . It becomes  $1 + 1.110^{-15} \neq 1 + 10^{-15}$ , which is also wrong and underflow. Then the numerator thereby becomes

$$1 + 1.1 * 10^{-15} - 1 \approx 1.110223^{-15} \neq 1.110^{-15}$$

(underflow again) with the result  $1.110223 \neq 1$ .

When  $x = 10^5$ ,  $10^{-15}$  is too small compared with  $x$ . So it is underflow and happens that  $10^5 + 10^{-15} = 10^5$ . The numerator becomes 0 and so does the derivative.

### Question 3: Variance

A known formula for estimating the variance based on a vector of  $n$  observations is

$$\text{Var}(\vec{x}) = \frac{1}{n-1} \left( \sum_{i=1}^n x_i^2 - \frac{1}{n} \left( \sum_{i=1}^n x_i \right)^2 \right)$$

1. Write your own R function, `myvar`, to estimate the variance in this way.

```
# write the preseted var function - input vector x
myvar = function(x){
  n = length(x)
  var_calculation = (1/(n-1)) * (sum(x^2) - (1/n)*(sum(x)^2))
  return(var_calculation)
}
```

2. Generate a vector  $x = (x_1, \dots, x_{10000})$  with 10000 random numbers with mean  $10^8$  and variance 1.

```
# initialize vector
n = 10000
x_vector = rnorm(n, mean = 10^8, sd = 1)
```

### 3. Compute the difference $Y_i = myvar(X_i) - var(X_i)$

For each subset  $X_i = \{x_1, \dots, x_i\}$ ,  $i=1, \dots, 10000$  compute the difference  $Y_i = myvar(X_i) - var(X_i)$ , where  $var(X_i)$  is the standard variance estimation function in R. Plot the dependence  $Y_i$  on  $i$ . Draw conclusions from this plot. How well does your function work? Can you explain the behaviour?

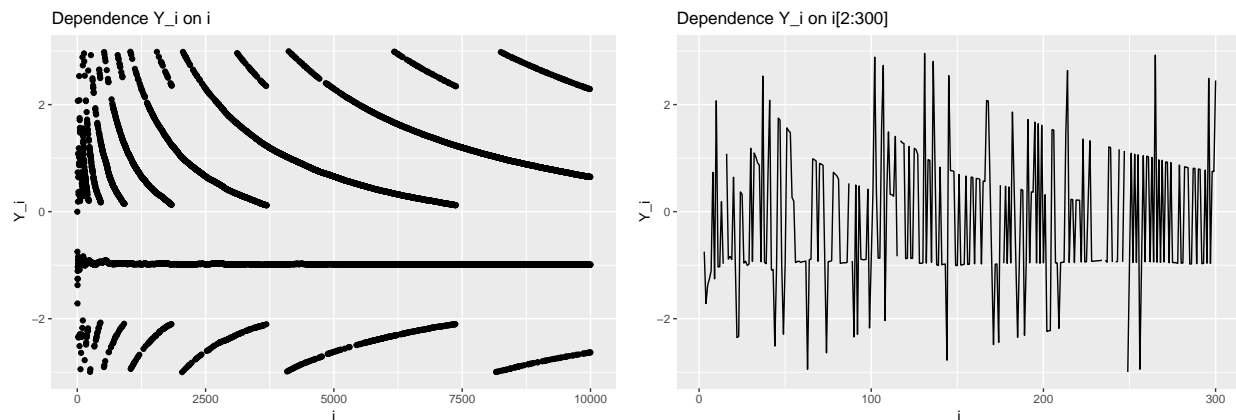
```
# try - 2
y = c()
for (i in 1:length(x_vector)) {
  y[i] = myvar(x_vector[1:i]) - var(x_vector[1:i])
}
# remove the first value - variance of one point can't be calculated - does not make sense
y = y[-1]

# create data framem for the plot
data_dependence_Y_i = data.frame(y_value = y,
                                  iteration = 1:length(y))

# create plot dependence Yi on i
plot_dependence_Y_i =
  ggplot(data = data_dependence_Y_i) +
  geom_point(aes(x = iteration, y = y_value)) +
  ggtitle("Dependence Y_i on i") + ylab("Y_i") + xlab("i") +
  ylim(-3,3)

plot_dependence_Y_i_2to300 =
  ggplot(data = data_dependence_Y_i[2:300,]) +
  geom_line(aes(x = iteration, y = y_value)) +
  ggtitle("Dependence Y_i on i[2:300]") + ylab("Y_i") + xlab("i") +
  ylim(-3,3)
```

The plots:



At the y-axis we can recognize the  $Y_i$ , which can be calculated as follows:  $Y_i = myvar(X_i) - var(X_i)$ . On the x-axis there is the subset with the respective size  $i$ . If the functions would output the same output, then all points would be distributed on the 0 horizontal axis. However, we can't see this course. At the  $y_i$  value of about -1 there is a continuous line. In the upper part of the visualization it can be seen that the value runs recurrently towards 0. Similar approximate progression can also be seen in the lower part, but against the previously mentioned value of about -1. Since the initialize vector is a very small number, the calculation with the  $var()$  function may handle these values better than the  $myvar$  function.

It is obvious that the difference between estimated variance and real variance increases regularly, so the result is not good enough. The reason might be that  $\sum_{i=1}^n x^2 \approx 10^{20}$  and  $(\sum_{i=1}^n x)^2 \approx 10^{24}$  are two large numbers.

So it is overflow intermittently.

## 4. Better implementation of variance estimator

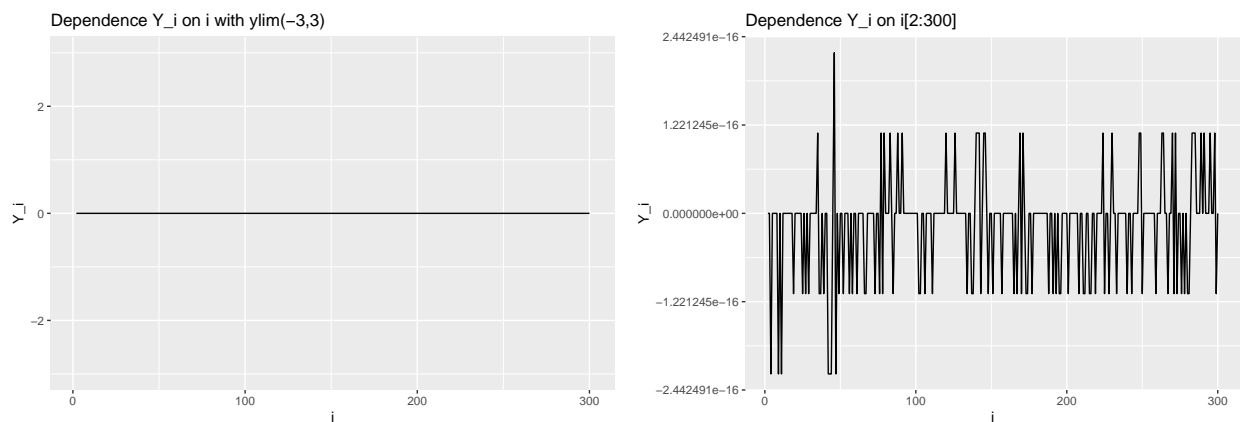
How can you better implement a variance estimator? Find and implement a formula that will give the same results as `var()`?

```
##3.4#####
myvar1 <- function(x){
  m <- mean(x)
  sum((x-m)**2)/(length(x)-1)
}
y1 <- sapply(1:n, function(i){
  myvar1(x_vector[1:i]) - var(x_vector[1:i]) # changed xbar to x_vector
})
y1 = y1[-1]
data_dependence_Y_i_improved = data.frame(y_value = y1,
                                             iteration = 1:length(y1))

plot_dependence_Y_i_2to300_improved_scale3 =
  ggplot(data = data_dependence_Y_i_improved[2:300,]) +
  geom_line(aes(x = iteration, y = y_value)) +
  ggtitle("Dependence Y_i on i with ylim(-3,3)") + ylab("Y_i") + xlab("i") +
  ylim(-3,3)

plot_dependence_Y_i_2to300_improved =
  ggplot(data = data_dependence_Y_i_improved[2:300,]) +
  geom_line(aes(x = iteration, y = y_value)) +
  ggtitle("Dependence Y_i on i[2:300]") + ylab("Y_i") + xlab("i")
```

Improved plot:



$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

Here we use the sample variance formula, which is an unbiased estimate of real variance for a finite sequence. Evidently, the result is better and the differences between estimated and real variances are much smaller than the previous formula.

## Question 4: Linear Algebra

The Excel file “tecator.xls” contains the results of a study aimed to investigate whether a near infrared absorbance spectrum and the levels of moisture and fat can be used to predict the protein content of samples of meat. For each meat sample the data consists of a 100 channel spectrum of absorbance records and the levels of moisture (water), fat and protein. The absorbance is  $-\log_{10}$  of the transmittance measured by the spectrometer. The moisture, fat and protein are determined by analytic chemistry. The worksheet you need to use is “data” (or file “tecator.csv”). It contains data from 215 samples of finely chopped meat. The aim is to fit a linear regression model that could predict protein content as function of all other variables.

### 1. Import the data set to R

```
##4.1#####
data <- as.matrix(read.csv("tecator.csv"))
data <- data[, -1]
X <- data[, -102]
Y <- data[, 102]
```

### 2. Compute optimal regression coefficients

Optimal regression coefficients can be found by solving a system of the type  $A\vec{\beta} = \vec{b}$  where  $A = X^T X$  and  $\vec{b} = X^T \vec{y}$  for the given data set. The matrix  $X$  are the observations of the absorbance records, levels of moisture and fat, while  $\vec{y}$  are the protein levels.

```
##4.2#####
Xe0 <- cbind(1,X)
A0 <- t(Xe0)%*%Xe0
b0 <- t(Xe0)%*%Y
```

### 3. Solve with default solver *solve()*

Try to solve  $A\vec{\beta} = \vec{b}$  with default solver `solve()`. What kind of result did you get? How can this result be explained?

```
##4.3#####
tryCatch(solve(A0,b0),
  error=function(e){cat("ERROR :\n",conditionMessage(e),"\n")})
```

“Error in `solve.default(A, b_vector)` : system is computationally singular: reciprocal condition number =  $7.13971\text{e-}17$ ”

An error occurs when the function `solve` is used. The reciprocal of condition number of matrix  $A$  is too small, and such value represents the tolerance *tol* for detecting linear dependencies in the columns of  $A$ . So some columns of  $A$  is linear dependent within default tolerance, which means the matrix  $A$  is not invertible here.

### 4. Check the condition number of the matrix $A$ (function `kappa()`) and consider how it is related to your conclusion in step 3.

```
##4.4#####
beta0 <- solve(A0,b0,tol = 7.78804e-17)
kappa(A0)
```

```
## [1] 8.523517e+14
```

We still make the model by such  $A$  and  $b$  with smaller tolerance. The function *kappa* can provide us the condition number of matrix  $A$ , which follows the following formula:

$$\kappa(A) = \|A\| * \|A^{-1}\|.$$

It is to measure for a matrix  $A$  of solving a linear-equation system and show the significance of scaling data. The larger the condition number is, the worse matrix  $A$  does. Because the tolerance would be much smaller, when the condition number increases.

## 5. Scale the data set and repeat steps 2-4. How has the result changed and why?

```
##4.5#####  
X <- scale(data[, -102])  
Xe1 <- cbind(1, X)  
A1 <- t(Xe1) %*% Xe1  
b1 <- t(Xe1) %*% Y  
beta1 <- solve(A1, b1)  
# beta1 <- solve(t(Xe1) %*% Xe1) %*% t(Xe1) %*% Y  
kappa(A1)
```

```
## [1] 490471520661
```

After scaling the data, the function *solve* works well for new matrix  $A$ . Its corresponding condition number ( $\approx 4.9 * 10^{11}$ ) is much smaller than the previous one ( $\approx 8.5 * 10^{14}$ ). Theoretically, the model trained by scaled data might be better than the one by un-scaled data.

The mse unscaled:

```
## [1] 0.0762971
```

The mse scaled:

```
## [1] 0.07629706
```

According to the MSE, the model with scaled data has smaller MSE and better performance confirmly. One of reason could be that when the range of matrix  $A$  is very large, it is easy to occur underflow when calculating the determinant and inverse of such matrix. This would have a negative effect to prediction. Additionally, the condition number might be large based on un-scaled data since its value is in accordance with the determinant of original matrix  $A$  and its inverse.

## Appendix

```
knitr::opts_chunk$set(echo = TRUE, eval = TRUE)
# the libraries used in this lab
library(ggplot2)
# snippets 1
x1 = 1/3 ; x2 = 1/4
if(x1 - x2 == 1/12){
  print("Subtraction is correct")
} else{
  print("Subtraction is wrong")
}
# snippets 2
x1 = 1 ; x2 = 1/2
if(x1 - x2 == 1/2){
  print("Subtraction is correct")
} else{
  print("Subtraction is wrong")
}
# snippets 1 - improved
x1 = 1/3 ; x2 = 1/4
if(isTRUE(all.equal(x1 - x2, 1/12))){
  print("Subtraction is correct")
} else{
  print("Subtraction is wrong")
}
epsilon = 10**-15
derivate = function(x){
  deriv = ((x + epsilon) - x) / epsilon
  return(deriv)
}
x1= 1
x2= 100000
derivate(x1)
derivate(x2)
# write the preseted var function - input vector x
myvar = function(x){
  n = length(x)
  var_calculation = (1/(n-1)) * (sum(x^2)-(1/n)*(sum(x)^2))
  return(var_calculation)
}
# initialize vector
n = 10000
x_vector = rnorm(n, mean = 10^8, sd = 1)
# try - 2
y = c()
for (i in 1:length(x_vector)) {
  y[i] = myvar(x_vector[1:i]) - var(x_vector[1:i])
}
# remove the first value - variance of one pont can't be calculated - does not make sense
y = y[-1]

# create data framem for the plot
data_dependence_Y_i = data.frame(y_value = y,
```



```

                                iteration = 1:length(y))
# create plot dependence Yi on i
plot_dependence_Y_i =
  ggplot(data = data_dependence_Y_i) +
  geom_point(aes(x = iteration, y = y_value)) +
  ggtitle("Dependence Y_i on i") + ylab("Y_i") + xlab("i") +
  ylim(-3,3)

plot_dependence_Y_i_2to300 =
  ggplot(data = data_dependence_Y_i[2:300,]) +
  geom_line(aes(x = iteration, y = y_value)) +
  ggtitle("Dependence Y_i on i[2:300]") + ylab("Y_i") + xlab("i") +
  ylim(-3,3)
plot_dependence_Y_i

plot_dependence_Y_i_2to300
##3.4#####
myvar1 <- function(x){
  m <- mean(x)
  sum((x-m)**2)/(length(x)-1)
}
y1 <- sapply(1:n, function(i){
  myvar1(x_vector[1:i])-var(x_vector[1:i]) # changed xbar to x_vector
})
y1 = y1[-1]
data_dependence_Y_i_improved = data.frame(y_value = y1,
                                iteration = 1:length(y1))

plot_dependence_Y_i_2to300_improved_scale3 =
  ggplot(data = data_dependence_Y_i_improved[2:300,]) +
  geom_line(aes(x = iteration, y = y_value)) +
  ggtitle("Dependence Y_i on i with ylim(-3,3)") + ylab("Y_i") + xlab("i")+
  ylim(-3,3)

plot_dependence_Y_i_2to300_improved =
  ggplot(data = data_dependence_Y_i_improved[2:300,]) +
  geom_line(aes(x = iteration, y = y_value)) +
  ggtitle("Dependence Y_i on i[2:300]") + ylab("Y_i") + xlab("i")
plot_dependence_Y_i_2to300_improved_scale3

plot_dependence_Y_i_2to300_improved
##4.1#####
data <- as.matrix(read.csv("tecator.csv"))
data <- data[,-1]
X <- data[,-102]
Y <- data[,102]
##4.2#####
Xe0 <- cbind(1,X)
A0 <- t(Xe0)%*%Xe0
b0 <- t(Xe0)%*%Y
##4.3#####
tryCatch(solve(A0,b0),
          error=function(e){cat("ERROR :\n",conditionMessage(e),"\n")})

```

```

##4.4#####
beta0 <- solve(A0,b0,tol = 7.78804e-17)
kappa(A0)
##4.5#####
X <- scale(data[, -102])
Xe1 <- cbind(1,X)
A1 <- t(Xe1)%*%Xe1
b1 <- t(Xe1)%*%Y
beta1 <- solve(A1,b1)
# beta1 <- solve(t(Xe1)%*%Xe1)%*%t(Xe1)%*%Y
kappa(A1)
mse_unscaled <- mean((Xe0%*%beta0-Y)**2)
mse_scaled <- mean((Xe1%*%beta1-Y)**2)
mse_unscaled
mse_scaled

```