

# Computational Statistics - Lab 01

*Annalena Erhard (anner218) and Maximilian Pfundstein (maxpf364)*

*2019-02-10*

## Contents

1	Question 1: Be Careful When Comparing	1
2	Question 2: Derivative	2
3	Question 3: Variance	3
4	Question 4: Linear Algebra	6
5	Source Code	8

## 1 Question 1: Be Careful When Comparing

```
#####  
# Question 1 - Be Careful When Comparing  
#####  
  
x1 = 1/3  
x2 = 1/4  
  
if (x1-x2 == 1/12) {  
  print("Substraction is correct.")  
} else {  
  print("Substraction is wrong.")  
}  
  
## [1] "Substraction is wrong."  
  
x1 = 1  
x2 = 1/2  
  
if (x1-x2 == 1/2) {  
  print("Substraction is correct.")  
} else {  
  print("Substraction is wrong.")  
}
```

```
## [1] "Substraction is correct."
```

### Questions:

1. Check the results of the snippets. Comment what is going on.
2. If there are any problems, suggest improvements.

### Answers:

1. The first subtraction is not wrong - it is perfectly working as defined in IEEE\_754 which is a commonly used definition for floating point numbers. To make a long story short: You have an infinite amount of

real numbers for instance between 0.0 and 1.0 but just 32 or 64 bits for the representation (so  $2^{32}$  states or  $2^{64}$  states) so it's impossible to represent every number (t.ex. try writing down  $1/3$  in the decimal system, at one point you simply must stop). The second subtraction does not have a floating point error as you can represent multiples of the power of two in the binary system  $2^{-1} = 0.5$ .

2. You can either just use more bits for representing the numbers until you have the desired precision or use/write a class which can handle a specific amount of numbers behind the decimal point (which will be slower for sure).

Another way to handle this issue can be seen in the following code snippet. The used method `all.equal()` will try to handle computational errors allowing an error tolerance of  $1.5e-8$  by default. This value is a parameter of the function and thus can easily be changed. Applying this function both values are treated as equal and makes the expected result of the computation "correct".

```
x1 = 1/3
x2 = 1/4

if (all.equal(x1-x2, 1/12)) {
  print("Substraction is correct.")
} else {
  print("Substraction is wrong.")
}

## [1] "Substraction is correct."
```

## 2 Question 2: Derivative

**Question:** Write your own R function to calculate the derivative of  $f(x) = x$  in this way with  $\epsilon = 10^{-15}$ .

```
#####
# Question 2: Derivative
#####

epsilon = 10^(-15)

f_prime = function(x) {
  return( (f(x + epsilon) - f(x)) / epsilon)
}

f = function(x) {
  return(x)
}
```

**Question:** Evaluate your derivative function at  $x = 1$  and  $x = 100000$ .

```
first = f_prime(1)
second = f_prime(100000)

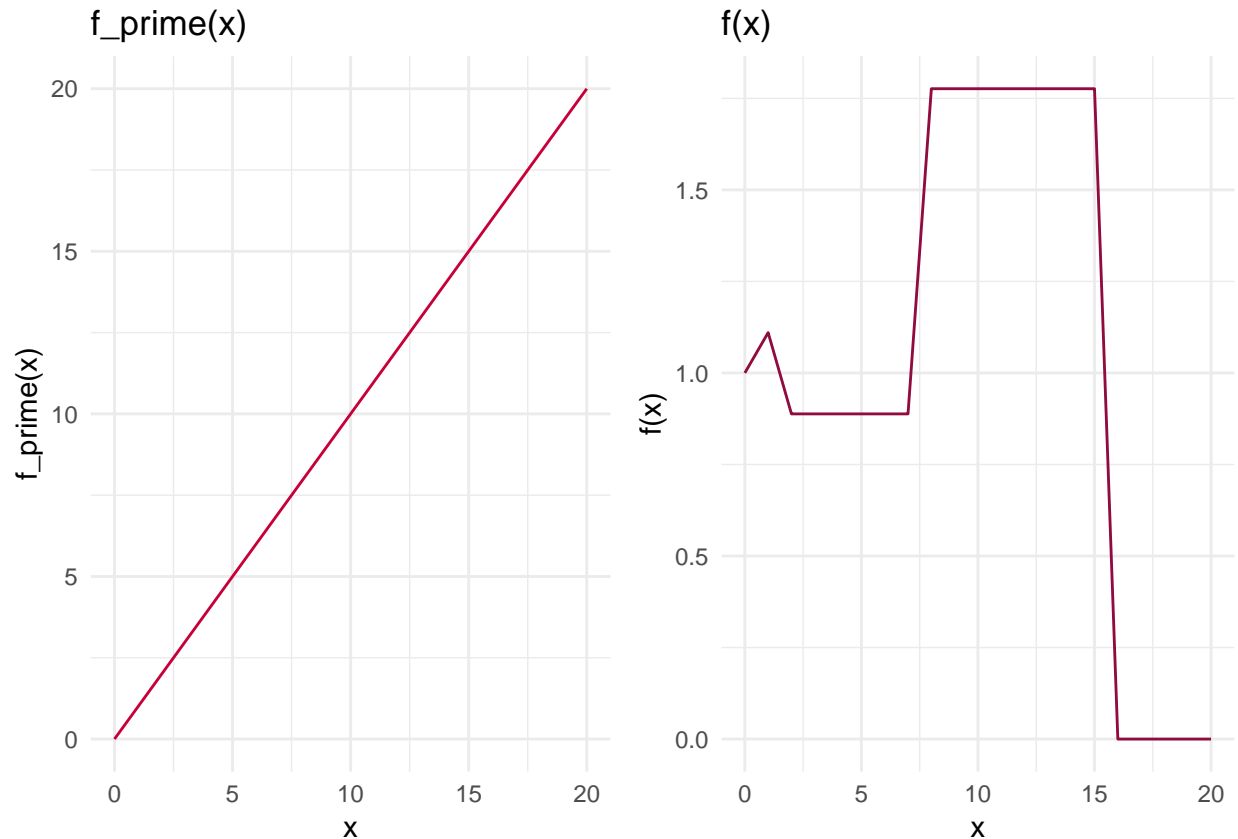
print(first)

## [1] 1.110223

print(second)

## [1] 0
```

The following plots show the function `f` and `f_prime` in the interval  $0 \dots 20$ . We observe that the derivative seems to take discrete values and is 0 around  $x = 16$ .



**Question:** What values did you obtain? What are the true values? Explain the reasons behind the discovered differences.

**Answer:** The values and plots can be seen above.

The derivate of  $f(x) = x$  is  $f'(x) = 1$  so the true slope is 1 at all spots.

As `epsilon` is a really small number and we do calculations with a rather big number (`x`) we run into precision problems which are more heavy if the magnitudes of the numbers is greatly different. Therefore we see that if we take  $x = 1$  the error is smaller, but still big enough to give an undesired result. As we take  $x = 100000$  the difference in magnitude increased further so we obtain the weird result 0 which is obviously totally wrong.

If you evaluate only the nominator you will see that it's 0 after  $x = 16$  so the result will always be 0 (assuming the denominator is unequal to 0).

### 3 Question 3: Variance

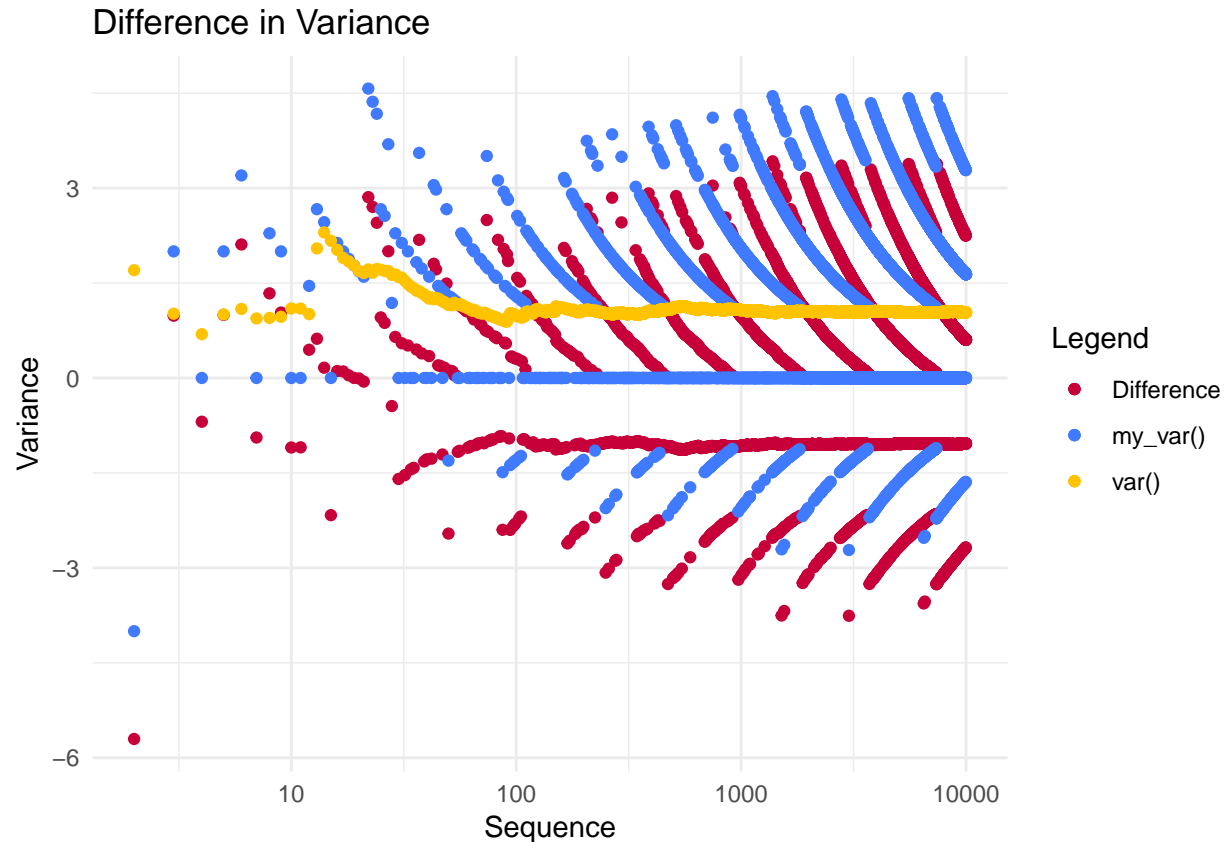
**Question:** Write your own R function, `myvar`, to estimate the variance in this way.

```
#####
# Question 3: Variance
#####
myvar = function(x) return(1/(length(x)-1) * (sum(x^2) - (sum(x)^2)/length(x)))
```

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

```
v = rnorm(10000, mean = 10^8, sd = 1)
```

**Question:** 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?



**Answer:** When interpreting the plot one has to keep in mind that each subset contains one value more than the previous one, so we should observe that we're getting better estimations for the variance which an increased index. Let's focus on the first values ( $index < 10$ ). They're highly scattered which is to be expected for such a small number of data points. With increasing  $index$  we see that some pattern is repeated. We have datapoints around -1 as well as trails towards the center. If we look at the trail that is on the upper half, ending at around an index of 100 we observe the following values:

```
kable(X[X$index > 70 & X$index < 100,])
```

	index	value	vec_myvar	vec_var
71	71	-1.0201689	0.000000	1.0201689
72	72	-1.0315800	0.000000	1.0315800
73	73	0.7537393	1.777778	1.0240385
74	74	2.4950624	3.506849	1.0117869
75	75	0.7297652	1.729730	0.9999645
76	76	-0.9926903	0.000000	0.9926903
77	77	-0.9810885	0.000000	0.9810885
78	78	-0.9724459	0.000000	0.9724459
79	79	-0.9711615	0.000000	0.9711615
80	80	0.6563458	1.620253	0.9639073

	index	value	vec_myvar	vec_var
81	81	0.6465441	1.600000	0.9534559
82	82	0.6283002	1.580247	0.9519468
83	83	2.1815748	3.121951	0.9403764
84	84	-0.9305039	0.000000	0.9305039
85	85	-0.9201338	0.000000	0.9201338
86	86	-0.9187333	0.000000	0.9187333
87	87	-2.3983854	-1.488372	0.9100133
88	88	2.0374444	2.942529	0.9050844
89	89	0.5571976	1.454546	0.8973479
90	90	0.5444620	1.438202	0.8937402
91	91	1.9538804	2.844444	0.8905640
92	92	1.8475873	2.813187	0.9655995
93	93	-0.9551112	0.000000	0.9551112
94	94	-2.3997272	-1.376344	1.0233831
95	95	0.3396715	1.361702	1.0220306
96	96	0.3346190	1.347368	1.0127494
97	97	-2.3408504	-1.333333	1.0075171
98	98	0.3133134	1.319588	1.0062742
99	99	-2.3022287	-1.306122	0.9961062

As we see the trails are **not** continuous as they heavily fluctuate. So, as it seems like they do not converge and are not continuous, this function (`my_var()`) seems not to be a good estimator for the variance. As we have rather big values with just a small variance from the created data and we have the sum of quite a few data, we might run into an underflow and precision problems.

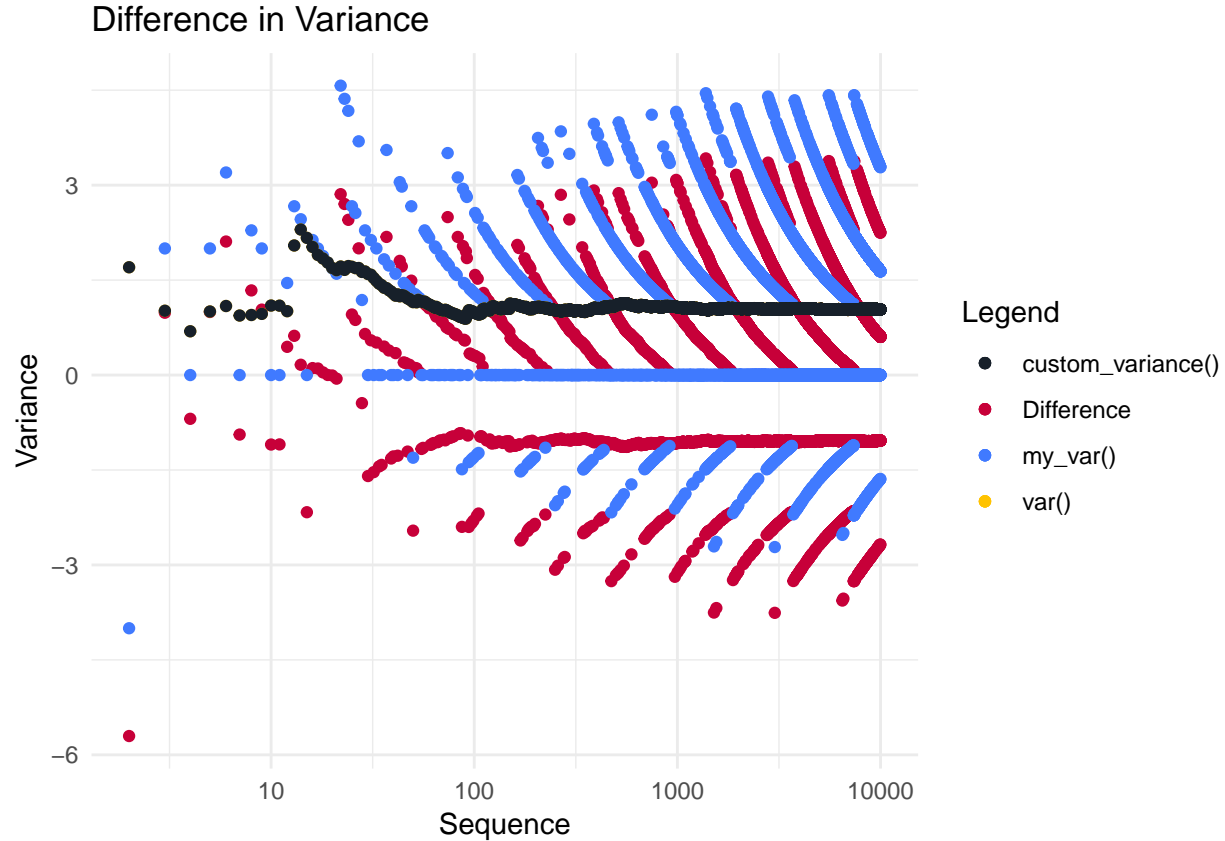
**Question:** How can you better implement a variance estimator? Find and implement a formula that will give the same results as `var()`?

**Answer:** We use this formula:

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

This is the sample standard deviance and thus is biased (if we assume our  $\vec{v}$  as a sample from a population).

```
custom_variance = function(x) {
  diff_mean = x - mean(x)
  return(sum(diff_mean^2 / (length(x) - 1)))
}
```



It can be seen that this function does a way better job at converging and returns almost (or exactly) the same values as the built in function and thus covering the `var()` plot almost perfectly.

## 4 Question 4: Linear Algebra

**Question:** Import the data set to R.

```
#####
# Question 4: Linear Algebra
#####

data = read.csv2("tecator.csv", sep=",", dec=".")
kable(head(data[, c(1, 101, 102, 103, 104)]))
```

Sample	Channel100	Fat	Protein	Moisture
1	2.81920	22.5	16.7	60.5
2	3.17942	40.1	13.5	46.0
3	2.54816	8.4	20.5	71.0
4	2.79622	5.9	20.7	72.8
5	3.13753	25.5	15.5	58.3
6	3.45307	42.7	13.7	44.0

**Question:** Optimal regression coefficients can be found by solving a system of the type  $A\vec{\beta} = \vec{b}$  where  $A = X^T X$  and  $\vec{b} = X^T \vec{y}$ . Compute  $A$  and  $\vec{b}$  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.

```
X = as.matrix(data[, c(1:102, 104)])
Y = as.matrix(data[, c(103)])
A = t(X) %*% X
b = t(X) %*% Y
```

**Question:** 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?

```
tryCatch(
  expr = {
    beta = solve(A) %*% b
  },
  error = function(e){
    paste("That escalated rather quickly: ", e)
  }
)
```

```
## [1] "That escalated rather quickly: Error in solve.default(A): system is computationally singular: "
```

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

```
kappa(A)
```

```
## [1] 4.286972e+15
```

As we mostly work with *rational* numbers we are used to the fact that almost every number has a inverse. An inverse  $a^{-1}$  is defined as that element that, multiplied with  $a$  results in the *neutral element*, at least for multiplications in arithmetics. In arithmetics we only have one number that does not have an multiplicative inverse which is 0 as  $1/0$  is undefined.

With matirces there are way more such matrices that do not have an inverse, exactly then when:

- The matrix is not a square.
- The determinant and thus the span is 0.

If we imagine a matrix as a linear transformation the determinant is the factor by which the space is streched or compressed. Thus a determinant of 0 tells us if the given linear transformation is squishing the amount of dimensions.

So it's reasonable that we have matrices that do not have an inverse. If we wanted to use the inverse for calculating or solving an equation, we have to find a different way (t.ex. QR-decomposition).

The `kappa()` function computes an estimate of the condition number of a matrix. Given a linear equation  $Ax = b$  the number gives us an estimation of how inaccurate the approximation of  $x$  is going to be. One can also say that it says how much  $x$  is going to change in respect to  $b$ . So if we have a large condition number it means that a small error in  $b$  is likely to cause a large error in  $x$ . As our number here is rather large, we can conclude that our features are linearly dependant.

**Question:** Scale the data set and repeat steps 2-4. How has the result changed and why?

```
data_scaled = scale(data)

X = as.matrix(data_scaled[, c(1:102, 104)])
Y = as.matrix(data_scaled[, c(103)])
A = t(X) %*% X
b = t(X) %*% Y

beta = solve(A) %*% b
```

The result has changed as we scaled the data. Before we ran into computational issues as the scale for each feature was on a different scope/scale which can lead to those errors. The `sclae()` made them of equal size and thus solved the problem.

Last but not least, let's look at the new conditional number, which should be smaller:

```
kappa(A)
```

```
## [1] 664322330621
```

It is smaller as we'd have expected. So everyone is happy and we can conclude this lab!



## 5 Source Code

```
knitr::opts_chunk$set(echo = TRUE, cache = FALSE, include = TRUE, eval = TRUE)
library(ggplot2)
library(knitr)
library(gridExtra)

#####
# Question 1 - Be Careful When Comparing
#####

x1 = 1/3
x2 = 1/4

if (x1-x2 == 1/12) {
```



```

    print("Substraction is correct.")
  } else {
    print("Substraction is wrong.")
  }

x1 = 1
x2 = 1/2

if (x1-x2 == 1/2) {
  print("Substraction is correct.")
} else {
  print("Substraction is wrong.")
}

x1 = 1/3
x2 = 1/4

if (all.equal(x1-x2, 1/12)) {
  print("Substraction is correct.")
} else {
  print("Substraction is wrong.")
}

#####
# Question 2: Derivative
#####

epsilon = 10-15

f_prime = function(x) {
  return( (f(x + epsilon) - f(x)) / epsilon)
}

f = function(x) {
  return(x)
}

first = f_prime(1)
second = f_prime(100000)

print(first)
print(second)

sequence = seq(from = 0, to = 20, by = 1)
func = f(sequence)
deri = f_prime(sequence)
df = data.frame(sequence, deri)

p1 = ggplot(df) +

```

```

geom_line(aes(x = sequence, y = func), color = "#C70039") +
labs(title = "f_prime(x)", y = "f_prime(x)", x = "x") +
theme_minimal()

p2 = ggplot(df) +
geom_line(aes(x = sequence, y = deri), color = "#900C3F") +
labs(title = "f(x)", y = "f(x)", x = "x") +
theme_minimal()

grid.arrange(p1, p2, nrow = 1)

#####
# Question 3: Variance
#####

myvar = function(x) return(1/(length(x)-1) * (sum(x^2) - (sum(x)^2)/length(x)))

v = rnorm(10000, mean = 10^-8, sd = 1)

X = data.frame()

for (i in 1:length(v)) {
  Xi = v[1:i]
  vec_myvar = myvar(as.vector(Xi))
  vec_var = var(Xi)
  Yi = vec_myvar - vec_var
  Yi_index = list(index = i, value = Yi, vec_myvar = vec_myvar, vec_var = vec_var)
  X = rbind(X, Yi_index)
}

ggplot(X[2:nrow(X),]) +
geom_point(aes(x = index, y = value, colour = "Difference")) +
geom_point(aes(x = index, y = vec_myvar, colour = "my_var()")) +
geom_point(aes(x = index, y = vec_var, colour = "var()")) +
labs(title = "Difference in Variance", y = "Variance",
x = "Sequence", color = "Legend") +
scale_color_manual(values = c("#C70039", "#407AFF", "#FFC300")) +
scale_x_log10() +
theme_minimal()

kable(X[X$index > 70 & X$index < 100,])

custom_variance = function(x) {
  diff_mean = x - mean(x)
  return(sum(diff_mean^2 / (length(x) - 1)))
}

```

```

for (i in 1:length(v)) {
  Xi = v[1:i]
  X$vec_customvar[[i]] = custom_variance(as.vector(Xi))
}

ggplot(X[2:nrow(X),]) +
  geom_point(aes(x = index, y = value, colour = "Difference")) +
  geom_point(aes(x = index, y = vec_myvar, colour = "my_var()")) +
  geom_point(aes(x = index, y = vec_var, colour = "var()")) +
  geom_point(aes(x = index, y = vec_customvar, colour = "custom_variance()")) +
  labs(title = "Difference in Variance", y = "Variance",
       x = "Sequence", color = "Legend") +
  scale_color_manual(values = c("#17202A", "#C70039", "#407AFF", "#FFC300")) +
  scale_x_log10() +
  theme_minimal()

#####
# Question 4: Linear Algebra
#####

data = read.csv2("tecator.csv", sep=",", dec=".")
kable(head(data[, c(1, 101, 102, 103, 104)]))

X = as.matrix(data[, c(1:102, 104)])
Y = as.matrix(data[, c(103)])
A = t(X) %*% X
b = t(X) %*% Y

tryCatch(
  expr = {
    beta = solve(A) %*% b
  },
  error = function(e){
    paste("That escalated rather quickly: ", e)
  }
)

kappa(A)

data_scaled = scale(data)

X = as.matrix(data_scaled[, c(1:102, 104)])
Y = as.matrix(data_scaled[, c(103)])
A = t(X) %*% X
b = t(X) %*% Y

beta = solve(A) %*% b

```

kappa(A)

