

# Compulsory exercise 1: Group 5

TMA4268 Statistical Learning V2021

Hans Røhjel Odland and Aksel Haugen Madslien

2/9/2021

For some problems you will need to include some LaTeX code. Please install latex on your computer and then consult Compulsory1.Rmd for hints how to write formulas in LaTeX.

An example:

$$Y_i = f(x_i) + \varepsilon_i ,$$

Or the same formula  $Y_i = f(x_i) + \varepsilon_i$  in-line.

## Problem 1

a)

We consider  $Y = f(\mathbf{x}) + \varepsilon$ , where  $E(\varepsilon) = 0$  and  $\text{Var}(\varepsilon) = \sigma^2$ .

We find the expected value for  $\tilde{\beta}$  as

$$E(\tilde{\beta}) = E[(\mathbf{x}^T \mathbf{x} + \lambda \mathbf{I})^{-1} \mathbf{x}^T \mathbf{y}] = (\mathbf{x}^T \mathbf{x} + \lambda \mathbf{I})^{-1} \mathbf{x}^T E[\mathbf{y}] = (\mathbf{x}^T \mathbf{x} + \lambda \mathbf{I})^{-1} \mathbf{x}^T \mathbf{x} \beta + \lambda \mathbf{I} \beta - \lambda \mathbf{I} \beta = (\mathbf{x}^T \mathbf{x} + \lambda \mathbf{I})^{-1} (-\lambda \mathbf{I}) \beta + \mathbf{I} \beta = \beta - \lambda (\mathbf{x}^T \mathbf{x} + \lambda \mathbf{I})^{-1} \mathbf{x}^T \mathbf{x} \beta$$

b)

We let  $\tilde{f}(\mathbf{x}_0) = \mathbf{x}_0^T \tilde{\beta}$ . The variance for  $\tilde{f}(\mathbf{x}_0)$  then becomes

$$E[\tilde{f}(\mathbf{x}_0)] = E[\mathbf{x}_0^T \tilde{\beta}] = \mathbf{x}_0^T E[\tilde{\beta}] = \mathbf{x}_0^T (\beta - \lambda (\mathbf{x}^T \mathbf{x} + \lambda \mathbf{I})^{-1} \mathbf{x}^T \mathbf{x} \beta)$$

For the variation we get

$$\text{Var}[\tilde{f}(\mathbf{x}_0)] = \mathbf{x}_0^T \text{Var}[\tilde{\beta}] \mathbf{x}_0 = \mathbf{x}_0^T (\beta - \lambda (\mathbf{x}^T \mathbf{x} + \lambda \mathbf{I})^{-1} \mathbf{x}^T \mathbf{x} \beta) \mathbf{x}_0$$

## c)

```
id <- "1X_80KcoYbng1XvYFDirxjEW7LtpNr1m" # google file ID
values <- dget(sprintf("https://docs.google.com/uc?id=%s&export=download", id))
X = values$X
dim(X)
```

```
## [1] 100 81
```

```
x0 = values$x0
dim(x0)
```

```
## [1] 81 1
```

```
beta=values$beta
dim(beta)
```

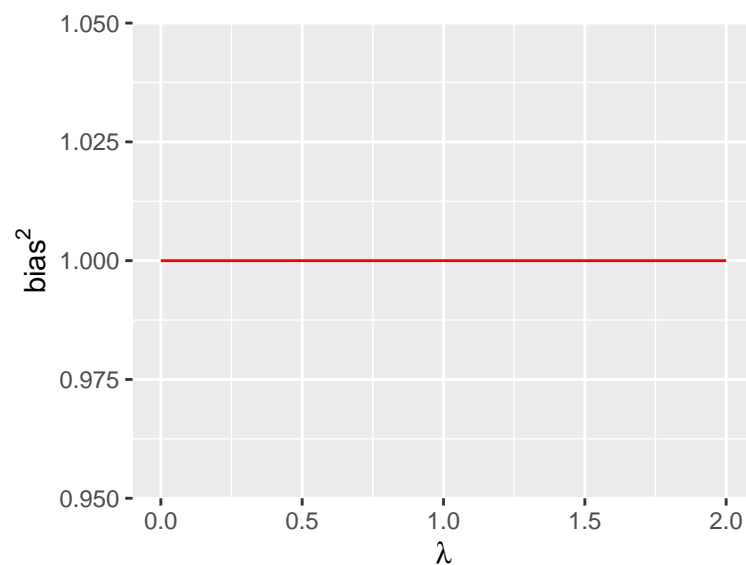
```
## [1] 81 1
```

```
sigma=values$sigma
sigma
```

```
## [1] 0.5
```

d)

```
library(ggplot2)
bias = function(lambda, X, x0, beta) {
  p = ncol(X)
  value = 1
  return(value)
}
lambdas = seq(0, 2, length.out = 500)
BIAS = rep(NA, length(lambdas))
for (i in 1:length(lambdas)) BIAS[i] = bias(lambdas[i], X, x0, beta)
dfBias = data.frame(lambdas = lambdas, bias = BIAS)
ggplot(dfBias, aes(x = lambdas, y = bias)) + geom_line(color = "red") + xlab(expression(lambda)) +
  ylab(expression(bias^2))
```



## Problem 2

a)

```
id <- "1yYlE15gYY3BEtJ4d7KWaFGIOEweJIn_" # google file ID
d.corona <- read.csv(sprintf("https://docs.google.com/uc?id=%s&export=download", id), header=T)

table(d.corona$deceased)
```

```
##
##      0      1
## 1905  105
```

```
table(d.corona$country, d.corona$sex)
```

```
##
##           female male
##   France         60   54
## indonesia        30   39
##   japan         120  174
##   Korea         879  654
```

```
table(d.corona$sex, d.corona$deceased)
```

```
##
##           0      1
##   female 1046   43
##   male   859   62
```

```
francedf <- subset(d.corona, country == "France")
table(francedf$sex, francedf$deceased)
```

```
##
##           0      1
##   female  55     5
##   male   43    11
```

b)

i)

```
# Multiple Linear Regression Example
fit <- lm(as.numeric(deceased) ~ sex + country + age, data=d.corona)
summary(fit)$coef # show results
```

```
##               Estimate Std. Error  t value    Pr(>|t|)
## (Intercept)    0.043861643 0.0252284684  1.738577 8.226271e-02
## sexmale        0.030814710 0.0099018024  3.112030 1.884264e-03
## countryindonesia -0.053478088 0.0335835996 -1.592387 1.114555e-01
## countryjapan    -0.097524595 0.0242695290 -4.018397 6.076062e-05
## countryKorea   -0.071966403 0.0215415300 -3.340821 8.506420e-04
## age            0.001304608 0.0002180126  5.984098 2.571005e-09
```

```
#predict(fit)
```

```
deceasedmale <- fit$coefficients[6]*75 + fit$coefficients[2] + fit$coefficients[5]
```

```
deceasedmale
```

```
##          age
## 0.05669394
```

The probability of dying of Covid-19 for a male age in Korea is found to be 5.7%.

ii)

Does males have higher probability to die than females?

```
fit <- glm(deceased ~ sex, data = d.corona, family = binomial)
summary(fit)$coef
```

```
##               Estimate Std. Error  z value    Pr(>|z|)
## (Intercept) -3.191529  0.1556015 -20.510908 1.720507e-93
## sexmale      0.562894  0.2037278  2.762971 5.727782e-03
```

The estimate readings for men dying of corona is positive, which means that we can conclude that men have a higher probability of dying of corona than women.

###iii)

```
fit <- glm(deceased ~ country, data = d.corona, family = 'binomial')
summary(fit)$coef
```

```
##               Estimate Std. Error  z value    Pr(>|z|)
## (Intercept)   -1.8123788  0.2696369 -6.721552 1.797990e-11
## countryindonesia -0.7370664  0.5369628 -1.372658 1.698586e-01
## countryjapan    -1.4351729  0.4088358 -3.510389 4.474509e-04
## countryKorea    -1.1833535  0.2951062 -4.009925 6.073807e-05
```

From these readings we can conclude that there is not enough evidence to say that there is a higher risk of dying of corona in Indonesia than in France, since the p-value is not significant. For Japan and Korea the p-value is much more significant, and also has a negative estimate, which means that there is a higher risk of dying of corona in Japan and Korea than in France.

iv)

A person is 10 years older than another person. The probability of dying is linear in terms of age because of the logistic regression, so we can see the probability of a person dying at an age of 65 and an age of 75 from task i and see that there is an increase of risk to die in case of higher age.

```
deceasedmale75 <- fit$coefficients[6]*75 + fit$coefficients[2] + fit$coefficients[5]
deceasedmale65 <- fit$coefficients[6]*65 + fit$coefficients[2] + fit$coefficients[5]

diff <- (deceasedmale75 - deceasedmale65)*100
```

This gives an age difference in 1.3% and we see that the model is linear.

c)

i)

```
fit <- glm(deceased ~ age*sex, data = d.corona, family = 'binomial')
summary(fit)$coef
```

	Estimate	Std. Error	z value	Pr(> z )
## (Intercept)	-4.7759607615	0.484626626	-9.85492852	6.526311e-23
## age	0.0278572418	0.007278416	3.82737708	1.295160e-04
## sexmale	0.5588617879	0.628776551	0.88880825	3.741061e-01
## age:sexmale	0.0005070244	0.009476422	0.05350378	9.573305e-01

Here we see that the `age:sexmale` coefficients has a positive estimate, but doesn't have a significant p-value. Age has a slightly lower p-value, and we can see that age is not a greater risk factor for males than for females.

ii)

```
fit <- glm(deceased ~ age*country, data = d.corona, family = 'binomial')
summary(fit)$coef
```

	Estimate	Std. Error	z value	Pr(> z )
## (Intercept)	-6.65232422	1.68621990	-3.945111	7.976311e-05
## age	0.06637008	0.02081162	3.189088	1.427225e-03
## countryindonesia	4.27185119	2.11130583	2.023322	4.303998e-02
## countryjapan	2.05835682	2.00125741	1.028532	3.036998e-01
## countryKorea	2.33159626	1.72031415	1.355332	1.753119e-01
## age:countryindonesia	-0.06981659	0.03206143	-2.177588	2.943672e-02
## age:countryjapan	-0.04403013	0.02627239	-1.675909	9.375603e-02
## age:countryKorea	-0.04173146	0.02148983	-1.941917	5.214712e-02

No, since the coefficient for the `age:countryindonesia` interaction is negative, which means that Indonesia is lower than it is for France. The p-value is slightly significant, which gives a low but greater risk factor for the Indonesian population than for the French.

### Problem 3

a)