

Module 2: Recommended Exercises

TMA4268 Statistical Learning V2023

Emma Skarstein, Daesoo Lee, Kenneth Aase, Stefanie Muff
Department of Mathematical Sciences, NTNU

January 19, 2023

Problem 1

- a) Describe a real-life application in which *classification* might be useful. Identify the response and the predictors. Is the goal inference or prediction?
- b) Describe a real-life application in which *regression* might be useful. Identify the response and the predictors. Is the goal inference or prediction?

Problem 2

Take a look at Figure 2.9 in the course book (p.31).

- a) Discuss whether a flexible or rigid method typically will have the highest test error.
- b) Does a small variance imply that the data has been under- or overfit?
- c) Relate the problem of over- and underfitting to the bias-variance trade-off.

Problem 3 – Exercise 2.4.9 from ISL textbook (modified)

This exercise involves the **Auto** dataset from the **ISLR** library. Load the data into your R session by running the following commands:

```
library(ISLR)
data(Auto)
```

PS: if the **ISLR** package is not installed (**library** function gives error) you can install it by running **install.packages("ISLR")** before you load the package the first time.

- a) View the data. What are the dimensions of the data? Which predictors are quantitative and which are qualitative?
- b) What is the range (min, max) of each quantitative predictor? Hint: use the **range()** function. For more advanced users, check out **sapply()**.
- c) What is the mean and standard deviation of each quantitative predictor?
- d) Now, make a new dataset called **ReducedAuto** where you remove the 10th through 85th observations. What is the range, mean and standard deviation of the quantitative predictors in this reduced set?
- e) Using the full dataset, investigate the quantitative predictors graphically using a scatterplot. Do you see any strong relationships between the predictors? Hint: try out the **ggpairs()** function from the **GGally** package.
- f) Suppose we wish to predict gas milage (**mpg**) on the basis of the other variables (both quantitative and qualitative). Make some plots showing the relationships between **mpg** and the qualitative predictors (hint: **geom_boxplot()**). Which predictors would you consider helpful when predicting **mpg**?

g) The correlation of two variables X and Y are defined as

$$\text{cor}(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}.$$

Both the correlation matrix and covariance matrix are easily assessed in R with the `cor()` and `cov()` functions. Use only the covariance matrix to find the correlation between `mpg` and `displacement`, `mpg` and `horsepower`, and `mpg` and `weight`. Do your results coincide with the correlation matrix you find using `cor(Auto[,quant])`?

```
quant = c(1,3,4,5,6,7)
covMat = cov(Auto[,quant])
```

Problem 4 – Multivariate normal distribution

The pdf of a multivariate normal distribution is on the form

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\Sigma|} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\},$$

where \mathbf{x} is a random vector of size $p \times 1$, $\boldsymbol{\mu}$ is the mean vector of size $p \times 1$ and Σ is the covariance matrix of size $p \times p$.

a) Use the `mvrnorm()` function from the `MASS` library to simulate 1000 values from multivariate normal distributions with

i)

$$\boldsymbol{\mu} = \begin{pmatrix} 2 \\ 3 \end{pmatrix} \quad \text{and} \quad \Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix},$$

ii)

$$\boldsymbol{\mu} = \begin{pmatrix} 2 \\ 3 \end{pmatrix} \quad \text{and} \quad \Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 5 \end{pmatrix},$$

iii)

$$\boldsymbol{\mu} = \begin{pmatrix} 2 \\ 3 \end{pmatrix} \quad \text{and} \quad \Sigma = \begin{pmatrix} 1 & 2 \\ 2 & 5 \end{pmatrix},$$

iv)

$$\boldsymbol{\mu} = \begin{pmatrix} 2 \\ 3 \end{pmatrix} \quad \text{and} \quad \Sigma = \begin{pmatrix} 1 & -2 \\ -2 & 5 \end{pmatrix}.$$

b) Make a scatterplot of the four sets of simulated datasets. Can you see which plot belongs to which distribution?

Problem 5 – Theory and practice: training and test MSE; bias-variance

We will now look closely into the simulations and calculations performed for the training error (`trainMSE`), test error (`testMSE`), and the bias-variance trade-off in lecture 1 of module 2.

Below, the code to run the simulation is included. The data is simulated according to the following specifications:

- True function $f(x) = x^2$ with normal noise $\varepsilon \sim N(0, 2^2)$.
- $x = -2.0, -1.9, \dots, 4.0$ (grid with 61 values).
- Parametric models are fitted (polynomials of degree 1 to degree 20).
- M=100 simulations.

a) Problem set-up

Look at the code below, copy it and run it yourself. Explain roughly what is done (you do not need not understand the code in detail), for example by commenting the code after copying it into your report.

We will learn more about the `lm` function in Module 3 - now just think of this as fitting a polynomial regression and then `predict` gives the fitted curve in our grid points. `predictions_list` is just a way to save M simulations of 61 gridpoints in x and 20 polynomial models.

```
set.seed(2) # to reproduce

M <- 100 # repeated samplings, x fixed
nord <- 20 # order of polynomials

#-----

x <- seq(from = -2, to = 4, by = 0.1)

truefunc <- function(x) {
  return(x ^ 2)
}
true_y <- truefunc(x)
error <- matrix(rnorm(length(x) * M, mean = 0, sd = 2),
               nrow = M,
               byrow = TRUE)
ymat <- matrix(rep(true_y, M), byrow = T, nrow = M) + error # Each row is a simulation

#-----

predictions_list <- lapply(1:nord, matrix, data = NA, nrow = M, ncol = ncol(ymat))
for(i in 1:nord){
  for(j in 1:M){
    predictions_list[[i]][j, ] <- predict(lm(ymat[j,] ~ poly(x, i, raw = TRUE)))
  }
}

# Plotting -----

library(tidyverse) # The tidyverse contains ggplot2, as well as tidyr and dplyr,
# which we can use for dataframe manipulation.

list_of_matrices_with_deg_id <-
  lapply(1:nord,
        function(poly_degree){cbind(predictions_list[[poly_degree]],
                                     simulation_num = 1:M, poly_degree)}
  )

# Now predictions_list is a list with 20 entries, where each entry is a matrix
# with 100 rows, where each row is the predicted polynomial of that degree.
# We also have a column for the simulation number, and a column for polynomial degree.

# Extract each matrix and bind them to one large matrix
stacked_matrices <- NULL
for (i in 1:nord) {
  stacked_matrices <-
    rbind(stacked_matrices, list_of_matrices_with_deg_id[[i]])
}
```

```

}
stacked_matrices_df <- as.data.frame(stacked_matrices)

# Convert from wide to long (because that is the best format for ggplot2)
long_predictions_df <- pivot_longer(stacked_matrices_df,
                                     !c(simulation_num, poly_degree),
                                     values_to = "y")

# Now we can use ggplot2!
# We just want to plot for degrees 1, 2, 10 and 20.

plotting_df <- cbind(long_predictions_df, x = x) %>% # adding the x-vector to the dataframe
  filter(poly_degree %in% c(1, 2, 10, 20)) # Select only the predictions using degree 1, 2, 10 or 20

ggplot(plotting_df, aes(x = x, y = y, group = simulation_num)) +
  geom_line(aes(color = simulation_num)) +
  geom_line(aes(x = x, y = x^2), size = 1.5) +
  facet_wrap(~ poly_degree) +
  theme_bw() +
  theme(legend.position = "none")

```

What do you observe in the produced plot? Which polynomial fits the best to the true curve?

b) Train and test MSE

First we produce predictions at each grid point based on our training data (x and y_{mat}). Then we draw new observations to calculate test MSE, see `testymat`.

Observe how `trainMSE` and `testMSE` are calculated, and then run the code.

```

set.seed(2) # to reproduce
M <- 100 # repeated samplings, x fixed but new errors
nord <- 20

x <- seq(from = -2, to = 4, by = 0.1)
truefunc <- function(x) {
  return(x^2)
}
true_y <- truefunc(x)
error <- matrix(rnorm(length(x) * M, mean = 0, sd = 2), nrow = M, byrow = TRUE)
testerror <- matrix(rnorm(length(x) * M, mean = 0, sd = 2), nrow = M,
                   byrow = TRUE)
ymat <- matrix(rep(true_y, M), byrow = T, nrow = M) + error
testymat <- matrix(rep(true_y, M), byrow = T, nrow = M) + testerror

predictions_list <- lapply(1:nord, matrix, data = NA, nrow = M, ncol = ncol(ymat))
for (i in 1:nord) {
  for (j in 1:M) {
    predictions_list[[i]][j, ] <- predict(lm(ymat[j, ] ~ poly(x,
      i, raw = TRUE)))
  }
}

```

```

trainMSE <- lapply(1:nord, function(poly_degree) {
  rowMeans((predictions_list[[poly_degree]] - ymat)^2)
})
testMSE <- lapply(1:nord, function(poly_degree) {
  rowMeans((predictions_list[[poly_degree]] - testymat)^2)
})

```

Next, we plot the training and test error for each of the 100 data sets we simulated, as well as two different plots that show the means across the simulations.

```

library(tidyverse) # The tidyverse contains ggplot2, as well as tidyr and dplyr,
# which we can use for dataframe manipulation.

# Convert each matrix in the list from wide to long (because that
# is the best format for ggplot2)
list_train_MSE <- lapply(1:nord, function(poly_degree) cbind(error = trainMSE[[poly_degree]],
  poly_degree, error_type = "train", simulation_num = 1:M))
list_test_MSE <- lapply(1:nord, function(poly_degree) cbind(error = testMSE[[poly_degree]],
  poly_degree, error_type = "test", simulation_num = 1:M))

# Now predictions_list is a list with 20 entries, where each entry
# is a matrix with 100 rows, where each row is the predicted
# polynomial of that degree.

stacked_train <- NULL
for (i in 1:nord) {
  stacked_train <- rbind(stacked_train, list_train_MSE[[i]])
}
stacked_test <- NULL
for (i in 1:nord) {
  stacked_test <- rbind(stacked_test, list_test_MSE[[i]])
}

stacked_errors_df <- as.data.frame(rbind(stacked_train, stacked_test))
# This is already on long format.
stacked_errors_df$error <- as.numeric(stacked_errors_df$error)
stacked_errors_df$simulation_num <- as.integer(stacked_errors_df$simulation_num)
stacked_errors_df$poly_degree <- as.integer(stacked_errors_df$poly_degree)

p.all_lines <- ggplot(data = stacked_errors_df, aes(x = poly_degree,
  y = error, group = simulation_num)) + geom_line(aes(color = simulation_num)) +
  facet_wrap(~error_type) + xlab("Polynomial degree") + ylab("MSE") +
  theme_bw() + theme(legend.position = "none")

p.bars <- ggplot(stacked_errors_df, aes(x = as.factor(poly_degree), y = error)) +
  geom_boxplot(aes(fill = error_type)) + scale_fill_discrete(name = "Error type") +
  xlab("Polynomial degree") + ylab("MSE") + theme_bw()

# Here we find the average test error and training error across the
# repeated simulations. The symbol '%>%' is called a pipe, and
# comes from the tidyverse packages, which provide convenient
# functions for working with data frames.
means_across_simulations <- stacked_errors_df %>%
  group_by(error_type, poly_degree) %>%

```

```

summarise(mean = mean(error))

p.means <- ggplot(means_across_simulations, aes(x = poly_degree, y = mean)) +
  geom_line(aes(color = error_type)) + scale_color_discrete(name = "Error type") +
  xlab("Polynomial degree") + ylab("MSE") + theme_bw()

library(patchwork) # The library patchwork is the best way of combining ggplot2 objects.
# You could also use the function ggarrange from the ggpubr
# package.

p.all_lines/(p.bars + p.means)

```

- Which value of the polynomial gives the smallest mean testMSE?
- Which gives the smallest mean trainMSE?
- Which would you use to predict a new value of y ?

c) Bias and variance - we use the truth!

Finally, we want to see how the expected quadratic loss can be decomposed into

- irreducible error: $\text{Var}(\varepsilon) = 4$
- squared bias: difference between mean of estimated parametric model chosen and the true underlying curve (truefunc)
- variance: variance of the estimated parametric model

Notice that the test data is not used – only predicted values in each x grid point.

Study and run the code. Explain the plots produced.

```

meanmat <- matrix(ncol = length(x), nrow = nord)
varmat <- matrix(ncol = length(x), nrow = nord)
for (j in 1:nord){
  meanmat[j,] <- apply(predictions_list[[j]], 2, mean) # we now take the mean over the M simulations -
  varmat[j,] <- apply(predictions_list[[j]], 2, var)
}

# nord times length(x)
bias2mat <- (meanmat - matrix(rep(true_y, nord), byrow = TRUE, nrow = nord))^2 #here the truth is final

```

Plotting the polys as a function of x :

```

df <- data.frame(x = rep(x, each = nord), poly_degree = rep(1:nord, length(x)),
  bias2 = c(bias2mat), variance = c(varmat),
  irreducible_error = rep(4, prod(dim(varmat)))) #irr is just 1

df$total <- df$bias2 + df$variance + df$irreducible_error

df_long <- pivot_longer(df, cols = !c(x, poly_degree), names_to = "type")

df_select_poly <- filter(df_long, poly_degree %in% c(1, 2, 10, 20))

ggplot(df_select_poly, aes(x = x, y = value, group = type)) +
  geom_line(aes(color = type)) +
  facet_wrap(~poly_degree, scales = "free", labeller = label_both) +
  theme_bw()

```

Now plotting effect of more complex model at 4 chosen values of x , compare to Figures in 2.12 on page 36 in

ISL (our textbook).

```
df_select_x <- filter(df_long, x %in% c(-1, 0.5, 2, 3.5))

ggplot(df_select_x, aes(x = poly_degree, y = value, group = type)) +
  geom_line(aes(color = type)) +
  facet_wrap(~x, scales = "free", labeller = label_both) +
  theme_bw()
```

Study the final plot you produced: when the flexibility increases (poly increase), what happens with i) the squared bias, ii) the variance, iii) the irreducible error?

d) Repeat a-c

Try to change the true function `truefunc` to something else - maybe order 3? What does this do the the plots produced? Maybe you then also want to plot `poly3`?

Also try to change the standard deviation of the noise added to the curve (now it is `sd=2`). What happens if you change this to `sd=1` or `sd=3`?

Or, change to the true function that is not a polynomial?

Acknowledgements

We thank Mette Langaas and her PhD students (in particular Julia Debik) from 2018 and 2019 for building up the original version of this exercise sheet.