

# TMA4268 V2024 Exam

## TMA4268 Statistical Learning V2024

Stefanie Muff and Sara Martino, Department of Mathematical Sciences, NTNU

May 11, 2024

Check in the end which libraries we need:

### Problem 1 (Fill-in-the-blank text, 5P)

0.5P per correct answer

Read the whole text and fill in the blanks such that the whole text makes sense (you might only understand which answer is correct after you continued reading):

In this course we learned about methods that can be used for statistical learning. We have mainly discussed supervised statistical learning methods, broadly divided into regression and classification problems. We were thereby discriminating between models that we use for prediction versus inference.

If the main goal is inference, we usually prefer (parametric, non-parametric, unsupervised, more advanced, more flexible, regression, classification) methods, because the goal is then to (predict, minimize the sum of bias and variance, understand, minimize prediction error). If the goal is pure prediction, we can use any method that performs well, whereas (regression, classification, clustering, supervised, non-parametric) methods are very flexible and thus more complex and (better interpretable, less interpretable, more appealing, better suited for inference) than, for example, linear models.

One important concept in the course was the bias-variance trade-off, and in that context we discussed variable selection and methods based on shrinkage. Lasso and ridge regression are two shrinkage methods we learned about. Lasso is an alternative to AIC-based model selection in linear models and should be preferred since it avoids (model selection bias, underestimated parameter estimates, large prediction error, irreducible error, over-fitting). In the context of trees, shrinkage-type regularization is performed via (bagging, random forests, tree pruning, boosting, classification trees, regression trees), whereas neural networks do this via (data augmentation, weight penalization, label smoothing, dropout, early stopping). More generally, any type of regularization has as a main goal to (reduce training error, reduce bias, increase the flexibility of the model, reduce test error, shrink parameters) by avoiding that the model over-fits the data.

In unsupervised learning the goal is to discover interesting aspects about the data when the (covariate, loss function, parametric model, reducible error, bias-variance trade-off, response variable) is not present. We have considered two types of methods: principal component analysis and clustering. Both methods have strong focus on finding (good predictions, groups, bias, variance) in the data, and on visualization.

### Problem 2 (Multiple choice and numeric answer, 7P)

#### a) (3P)

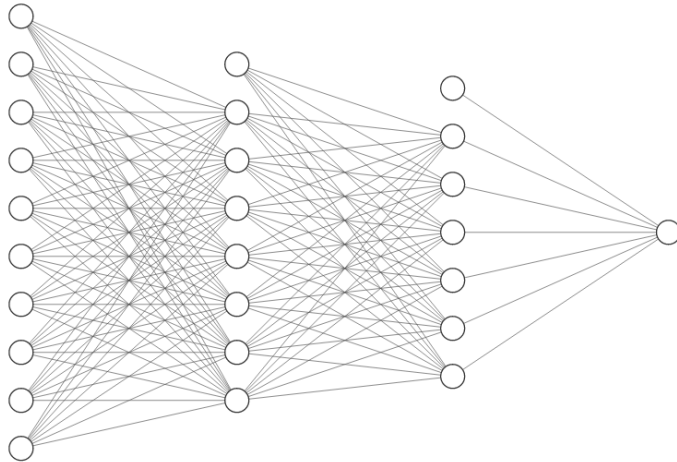
Which statements are correct? Here  $n$  are the number of data points and  $p$  the number of available variables.

- Forward selection requires that  $p < n$ .
- Backward selection requires that  $p < n$ .

- Lasso requires that  $p < n$ .
- Ridge regression is possible even if  $p > n$ .
- Neural networks must have  $p > n$ .
- Boosted regression trees allow for  $p > n$ .

**b) (2P)**

Look at the following architecture of a feed-forward neural network:



Which statements are correct?

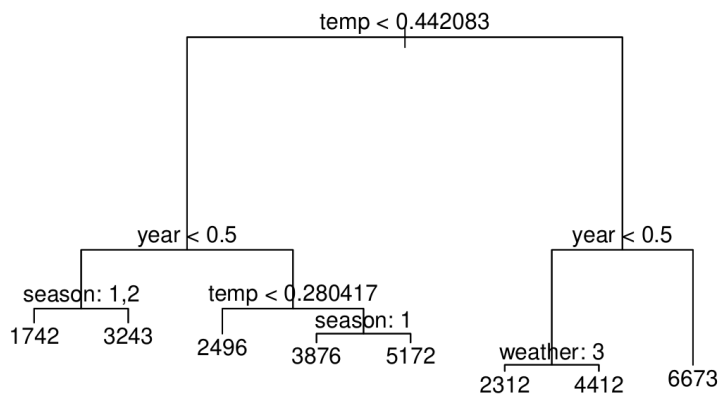
- This network has 9 input variables, two hidden layers and a continuous or binary outcome.
- This network has 10 input variables and two hidden layers with 8 and 7 nodes, respectively.
- This network can be used for binary classification.
- In this network we estimate 143 parameters.

**c) (1P) Numeric answer**

A covariate is included in a regression model as a natural cubic spline with three cut points. How many degrees of freedom does this spline term consume?

**d) (1P) Numeric answer**

We have fitted a tree to some training data to predict the number of bikes rented out in a given city on a specific day, and get the following result:



Using the fitted regression tree, what would you predict as the number of bikes rented on a day with the following covariates:

- year=0
- season=3
- holiday=0
- notworkday=0
- temp=0.3
- weather=3
- wind=0.3

### Problem 3 (theory, 8P)

#### a) (4P)

We learned about  $k$ -fold cross-validation (CV) as a way of doing model selection.

- (2P) Explain how  $k$ -fold CV is implemented and how the MSE is computed.
- (2P) State what the advantages and disadvantages of  $k$ -fold cross-validation (CV) are with respect to the Leave-one-out cross-validation (LOOCV) approach.

#### b) (4P)

Consider a regression setting and assume an additive error model:

$$Y_i = f_\theta(X_i) + \epsilon_i, \quad i = 1, \dots, N$$

where  $\theta$  denotes the model parameters defining the regression function  $f_\theta$ , and  $\epsilon$  is the error term.

- (2P) Describe the least squares method and the maximum likelihood method and how they are used to estimate  $\theta$ . It is enough to state the optimization problems, you do not need to solve them. State also the assumption that are made for each estimation method
- (2P) Show that, if you assume a Gaussian distribution for the error term, then the two methods are equivalent with respect to the estimate of  $\theta$ .

### Problem 4 – Data analysis 1 (16P)

Here we are looking at a regression problem, where we want to understand the factors that affect the sales of child carseats in 400 different stores. We use the `Carseats` data set from the ISLR package, which you can load using the code below

```
library(ISLR)
library(leaps)
library(glmnet)

data(Carseats)
```

It is useful to investigate the data for example using the code below and by typing `?Carseats` into the R console:

```
# Look at the data, for example using:
pairs(Carseats)
str(Carseats)
```

### a) (8P)

- (i) (1P) Fit a multiple regression model to predict **Sales** using **Price**, **ShelveLoc**, **US** and **Population** as predictors.
- (ii) (1P) Provide an interpretation of the estimated coefficient for **Price**. How much does the sales changes if price is increased by 10 dollars?
- (iii) (2P) Perform an F test to assess the significance of the **ShelveLoc** covariate and state your conclusion. If you move an item from a Medium to a Good shelf location, how does the expected number of sales changes?
- (iv) (2P) Compute the predicted sales for two observations both with a price of 100 dollars and population of 10 000 inhabitants but one with bad shelving location and not located in the US while the other with medium shelving and located in the US.
- (v) (2P) Look at the plots produced by

```
plot(mod)
```

where `mod` is your fitted model. Explain how each of the plots can be used to assess the fitted model. When relevant, state which assumptions are addressed and if those are met for the model under consideration.

### b) Model Selection (8P)

We now consider all predictors in the **Carseats** dataset except for **ShelveLoc**. Our goal is to build a model with as few predictor as possible but that still gives good predictions.

In order to assess the robustness of our models, we split the data into a training and a test set as follows:

```
set.seed(1234)
# remove ShelveLoc covariate
car.all = Carseats[, -7]

# create train and test datasets
samples <- sample(1:400, 250, replace = F)
car.train <- car.all[samples, ]
car.test <- car.all[-samples, ]
```

- i) (1P) Fit a model on the train data with all **Sales** as response. Use all predictors except **ShelveLoc** and add also a quadratic effect of **Age**.
- ii) (1P) Look at the summary of your model. How do you interpret the F-statistics in the last row?
- iii) (2P) Carry out Ridge regression on the training set, choose the largest lambda within 1 standard error from the lambda with the minimal error in a 5-fold cross-validation. As above, include the quadratic term for **Age**. Report the MSE of test data.

*Requirement:* Use `set.seed(1100)` before running the cross-validation.

- iv) (2P) Carry out Lasso regression on the training set, choose the largest  $\lambda$  within 1 standard error from the  $\lambda$  with the minimal error in a 5-fold cross-validation. As above, include the quadratic term for **Age**. Report the MSE of test data.

*Requirement:* Use `set.seed(1100)` before running the cross-validation.

- v) (2P) Compare the estimated coefficients you obtained with Ordinary least square, Ridge and Lasso regression. What patterns do you notice? Are there some advantages/disadvantages of Lasso over ridge regression?

#### R-hints:

```
x.train <- model.matrix(Sales ~ ..., data = car.train)
set.seed(1100)
```

```
# ridge
cv_ridge <- cv.glmnet(x.train, car.train$Sales, alpha = ...)
plot(cv.ridge)
cv.ridge$...
fit_ridge = glmnet(..., ..., alpha = ..., lambda = ...)
```

## Problem 5 – Data analysis 2 (16P)

We are using a dataset on purchase of orange juice, where the costumers either purchased the brand “Citrus Hill” (CH) or “Minute Maid” (MM). The dataset is available in the ISLR package and can be loaded and modified as follows:

```
library(ISLR)
data(OJ)
# Select a subset of columns
d.OJ <- OJ[, c("Purchase", "WeekofPurchase", "PriceDiff", "PriceCH",
               "PriceMM", "SpecialMM", "LoyalCH", "PctDiscMM", "PctDiscCH", "Store7")]
d.OJ$Purchase <- ifelse(d.OJ$Purchase == "CH", 0, 1)
```

Note that we have coded the purchase of Citrus Hill as 0, and Minute Maid as 1.

You can look up what the different covariates in the dataset actually mean by typing `?OJ` in the R console.

Before starting, it is smart to investigate the data a little bit, for example by making **pairs** plots or looking at the structure using `str(d.OJ)` etc.

We are interested in understanding and predicting the purchase of Minute Maid vs Citrus Hill.

### a) (3P)

- (i) (1P) Fit a logistic regression model on the full data set with **Purchase** as response variable, using all the covariates.
- (ii) (2P) Since we have a model where the prices for both brands, as well as the price differences are included, we fit another logistic regression model where we use all covariates *except* **PriceCH** and **PriceMM**. Given the fitted model, quantify the effect of the price difference (**PriceDiff**) on purchase of MM when the price difference increases by 0.1 units.

### b) (5P)

Now split the dataset into a training and a test sample for prediction (assuming our aim is to predict purchase of MM). Split the dataset as in the code below, using the same seed. Then

- (i) (1P) Perform a quadratic discriminant analysis on the training data.
- (ii) (1P) Use the fitted model to predict purchase of MM in the test set using a probability cutoff of  $p = 0.5$ .
- (iii) (1.5P) Generate the confusion table and calculate the error rate, sensitivity and specificity for the prediction on the test set.
- (iv) (1.5P) Generate the confusion table and calculate the error rate, sensitivity and specificity also for the logistic regression model from a) i), when trained on the training data and then predicted on the test data. Use again  $p = 0.5$  as cutoff.

**R-hints:**

```
set.seed(4268)
samples <- sample(1:1070, 1070 * 0.7, replace = F)
d.OJ.train <- d.OJ[samples, ]
d.OJ.test <- d.OJ[-samples, ]
```

```
library(MASS)
qdaMod <- qda()
postQDA = predict(qdaMod, newdata = ..)$class
table(...)
```

### c) (4P)

- (i) (3P) We continue analyzing the same that, focusing on the task to obtain good predictions of what the customers purchase. To this end, use a generalized additive model (still for the binary outcome **Purchase**) that only contains smoothed versions of **PriceDiff** and **LoyalCH** (no other covariates). Fit it on the training data and explain the details of your choice, for example how many degrees of freedom the smoothed terms consume and what functional form they have. Use maximum 5 sentences.
- (ii) (1P) Calculate the misclassification error rate on the test data. Is it lower than for logistic regression and QDA?

### d) (4P)

- (i) (2P) Finally, you should use a gradient tree boosting method. Take a large enough number of trees for a given learning rate and then choose the tree number that gives the lowest error rate for a 10-fold CV. Explain your choices. Use the R-hints below and fit the model on the training data.
- (ii) (2P) Calculate the misclassification error on the test data. Compare to the findings from b) and c) and interpret in 1-2 sentences.

#### R-hints:

```
library(gbm)
# Check the help file:
`?`(gbm())
gbm(..., n.trees = ..., shrinkage = ..., interaction.depth = ..., cv.folds = 10)
```

## Multiple and single choice and numerical answer questions

### Problem 6 (7.5P)

#### a)(1P) (numerical answer )

We have a dataset that we want to use to predict the starting salary after graduation, based on three predictors: the GPA (Grade Point Average, a number that indicates how high you scored in your courses on average), the IQ and the gender.

Assume we use R to fit the following linear model

```
formula = salary ~ GPA + IQ + GENDER + GPA:IQ + GPA:GENDER
mod = lm(formula, data = dataset)
```

Suppose we use least squares to fit the model, and get the following estimates

mod\$coefficients					
(Intercept)	GPA	IQ	GENDERFemale	GPA:IQ	GPA:GENDERFemale
50	20	0.07	35	0.01	-10

- i) What is the predicted salary for a female with IQ of 110 and GPA of 4.0? Give your answer with a precision of one decimal after the period.

## b) Single Choice (1P)

In the same problem as **a)** , which sentence is correct:

- For a fixed value of IQ and GPA, males earn more on average than females.
- For a fixed value of IQ and GPA, females earn more on average than males.
- For a fixed value of IQ and GPA, males earn more on average than females provided that the GPA is high enough.
- For a fixed value of IQ and GPA, females earn more on average than males provided that the GPA is high enough.

## c) (1P)

For which of the following techniques it is important to standardize the predictors:

- Multiple Linear Regression
- Ridge Regression
- Principal Component Analysis
- K-nearest neighbour classification

## d) (2P)

For each optimization criterion choose the correct method from the drop down menu:

- $\operatorname{argmin}_{\beta} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \beta)^2$ :
- $\operatorname{argmin}_{\beta} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \beta)^2 + \lambda \sum_{j=1}^p |\beta_j|$ :
- $\operatorname{minimize}_{C_1, \dots, C_K} \left\{ \sum_{k=1}^K \frac{1}{|C_k|} \sum_{i, i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2 \right\}$ :
- $\operatorname{argmin}_{R_1(j,s), R_2(j,s)} \left[ \sum_{i: x_i \in R_1(j,s)} (y_i - \hat{y}_{R_1})^2 + \sum_{i: x_i \in R_2(j,s)} (y_i - \hat{y}_{R_2})^2 \right]$ :

## e) (2.5P)

We are looking at the `state.x77` dataset given in R. This dataset consists of data related to the 50 states of the United States of America. The dataset contains 8 variables regarding different aspect of the state. You can check the data in R typing:

```
data(state)
`?`(state.x77)
```

to see what the different variables mean.

Here we carried out a principal component analysis and give the biplot and the scree plot below. In the biplot we also color the states according to their region (Northeast, South, North, Central, West). Which of the following statements are correct?

- (i) Population, income and area contribute most to the second PC.
- (ii) The first component explains 45% of the variability of the response variable.
- (iii) California (CA) has a very high loading for the second principal component.
- (iv) The first three PCs explain about 80% of the variability in the data.
- (v) Illiteracy has a low loading on the second PC.

