



**Deadline:** Thursday, November 17, 2022, 16:00

Before you start solving the exercise, please consult the instructions in the course website.

Importantly, note that the submission to each problem must be contained in an individual file.

- For each theoretical problem you must provide a single **pdf** file containing your answer to the respective problem. This file may be a scan of your (legible) handwriting.
- For each practical problem, you need to provide
  - the completed jupyter notebook (**.ipynb**) file, and
  - any necessary files required to reproduce your results. Additionally,
  - a **pdf** report generated from the jupyter notebook that shows all your results.

All the above files must be compressed into a **zip** archive and uploaded using the content management system of the course.

**Please remember, that each team member has to submit the signed Code of Conduct for each assignment. The Code of Conduct can be found in the material section on the course website.**

**Problem 1** (T, 5 Points). **Principles of statistical learning**

Explain in short and concise words the concepts of these pairs of terms:

1. unsupervised and supervised learning
2. prediction and inference
3. classification and regression
4. training and test data
5. parametric and non-parametric methods

You are allowed to use a maximum of **300 words** for this exercise, every **ten** more words will lead to losing **one** point.

*Solution.*

1. [1pts] **Supervised** learning describes a setting where we learn a mapping from features to outcomes, for which we are given pairs  $(\mathbf{x}_i, y_i)$  consisting of a feature vector  $\mathbf{x}_i$  and its associated label  $y_i$ .  
In contrast, **unsupervised** learning describes the setting where we are only given a set of unlabelled feature vectors, with the goal of learning useful structure or patterns in the data.
2. [1pts] In supervised learning there we can define two distinct tasks (goals) of learning models from data. We can either aim to **predict** outcomes for unseen inputs (new data points) or aim to **infer** the relationship between certain features and the response (e.g. understanding the way that  $Y$  is affected as features  $X$  change, i.e. the data generation process). *The learning task affects the model choice: inference requires more interpretable models than prediction.*
3. [1pts] Variables can be quantitative (i.e. numerical values - *continuous*, e.g. *income*, or *discrete*, e.g., *number of children*) or qualitative (i.e. values in one of  $K$  different categories - *nominal*, e.g., *sex*, or *ordinal*, e.g., *health diagnosis*). If the outcome of a supervised learning task is quantitative, the learning task is called **regression**, if the outcome is qualitative, it is called **classification**.
4. [1pts] *The combination of inputs (and outputs in the supervised case) that is used to learn the model is called **training set**. One usually evaluates the derived model using an independent set of inputs (and outputs in the supervised case), the **test set**.*



5. [1pts] **Parametric** methods make the strong assumption that the underlying model of the method lies in a well-defined class which is parameterised by a set of parameters whose size does not depend on the number of training points; here, picking the correct model amounts to simply specifying the correct set of parameters. As an example, linear regression makes a strong assumption that the outcome is a linear function of a (given number of) features, after which we only need to specify the linear coefficients of each feature.

In contrast, the underlying models of **Non-parametric** methods have a flexible number of parameters which depends on the complexity of the data; note that this does not refer to the dimensions of each observation, but rather to the arrangement of the observations in space. An extreme case of a non-parametric method is the k-NN classifier, which essentially stores the entire training set.

Importantly, parametric methods tends to outperform non-parametric ones given the same observations, whenever the assumptions of the former are accurate. Otherwise, and given enough data, the flexibility of the non-parameteric methods allows them to outperform parametric ones whose assumptions are not entirely correct.

*The relative performance of these two classes of methods can also be seen as yet another instance of the bias-variance tradeoff: parametric methods typically contain a smaller class of models than those of non-parametric ones, and therefore their variance is typically lower. However, if the true model lies far from the best parametric model—as is the case when its assumptions are wrong—the bias of the corresponding parametric method will eventually dominate the variance of a non-parameteric one, when given enough data.*

## Problem 2 (T, 6 Points). Bias-variance trade-off

1. [5pts] Prove that the expected test mean square error (MSE), for a given value  $x_0$ , can always be decomposed into the sum of three fundamental quantities: the variance of  $\hat{f}(x_0)$ , the squared bias of  $\hat{f}(x_0)$  and the variance of the error terms  $\epsilon$ . Please note that you should prove both qualities (both equals relation), enlist all assumptions and define all parameters and variables.

$$\begin{aligned}\mathbb{E}[(y_0 - \hat{f}(x_0))^2] &= \mathbb{E}[(\hat{f}(x_0) - \mathbb{E}(\hat{f}(x_0)))^2] + [\mathbb{E}(\hat{f}(x_0) - f(x_0))]^2 + \text{Var}(\epsilon) \\ &= \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\epsilon).\end{aligned}$$

Note that this is Equation (2.7) from ISLR (p. 34), which contains Equation (2.3) from ISLR (p. 19).

2. [1pts] Explain in your own words irreducible and reducible error as well as their difference.

*Solution.*

### 1. Proof

**Definitions:**  $\mathbb{E}[(y_0 - \hat{f}(x_0))^2]$  refers to the **expected test mean squared error (MSE)**, i.e., the average test MSE that we would obtain, by first using a large number of different training sets and for each training set learning a model  $\hat{f}(x)$  estimating the true relationship  $y = f(x) + \epsilon$  with irreducible error  $\epsilon$ .  $\mathbb{E}[(y_0 - \hat{f}(x_0))^2]$  is then the average prediction error over all of these models for a fix data point  $(x_0, y_0)$ . We **assume**  $\mathbb{E}(\epsilon) = 0$ , i.e. the irreducible error to have a zero mean. The overall expected test MSE can be computed by averaging over all models and all possible values of  $x_0$  (i.e., all data points) in the test set. (1 Point)

**Step 0:** For the sake of brevity, use  $f = f(x_0)$ ,  $\hat{f} = \hat{f}(x_0)$ , thus  $y_0 = f + \epsilon$ .



**Step 1:** We first simplify  $\mathbb{E}[(y_0 - \hat{f})^2]$ :

$$\begin{aligned}\mathbb{E}[(y_0 - \hat{f})^2] &= \mathbb{E}[(f + \epsilon - \hat{f})^2] \\ &= \mathbb{E}\left[\left((f - \hat{f}) + \epsilon\right)^2\right] \\ &= \mathbb{E}\left[(f - \hat{f})^2 + 2(f - \hat{f})\epsilon + \epsilon^2\right] \\ &= \mathbb{E}[(f - \hat{f})^2] + 2\mathbb{E}[(f - \hat{f})\epsilon] + \mathbb{E}[\epsilon^2] \quad (1 \text{ Point})\end{aligned}$$

The last term is  $\text{Var}(\epsilon)$  since  $\text{Var}(\epsilon) = \mathbb{E}[(\epsilon - \mathbb{E}[\epsilon])^2] = \mathbb{E}[\epsilon^2] + [\mathbb{E}(\epsilon)]^2$  and we assume  $\mathbb{E}(\epsilon) = 0$ .

The second term equals to zero because  $\epsilon$  is independent of  $(f - \hat{f})$ , so

$$2\mathbb{E}[(f - \hat{f})\epsilon] = 2\mathbb{E}[(f - \hat{f})] \mathbb{E}[\epsilon] = 0 \text{ with again the assumption } \mathbb{E}(\epsilon) = 0.$$

Hence,

$$\mathbb{E}[(y_0 - \hat{f})^2] = \mathbb{E}[(f - \hat{f})^2] + \text{Var}(\epsilon) \quad (1 \text{ Point})$$

**Step 2:** We show now that

$$\text{Var}(\hat{f}) + [\text{Bias}(\hat{f})]^2 = \mathbb{E}[(f - \hat{f})^2]$$

with  $\text{Var}(\hat{f}) = \mathbb{E}(\hat{f}^2) - [\mathbb{E}(\hat{f})]^2$  and  $\text{Bias}(\hat{f}) = \mathbb{E}(\hat{f} - f)$

$$\begin{aligned}\text{Var}(\hat{f}) + [\text{Bias}(\hat{f})]^2 &= \mathbb{E}(\hat{f}^2) - [\mathbb{E}(\hat{f})]^2 + [\mathbb{E}(\hat{f} - f)]^2 \\ &= \mathbb{E}(\hat{f}^2) - [\mathbb{E}(\hat{f})]^2 + [\mathbb{E}(\hat{f}) - \mathbb{E}(f)]^2 \\ &= \mathbb{E}(\hat{f}^2) - [\mathbb{E}(\hat{f})]^2 + [\mathbb{E}(\hat{f})]^2 - 2\mathbb{E}(\hat{f})\mathbb{E}(f) + [\mathbb{E}(f)]^2 \quad \text{binomial formula} \\ &= \mathbb{E}(\hat{f}^2) - 2\mathbb{E}(\hat{f})\mathbb{E}(f) + \mathbb{E}(f)^2 \quad (\mathbb{E}(f) = f, \text{ since } f \text{ is deterministic}) \\ &= \mathbb{E}[(f - \hat{f})^2] \quad (1 \text{ Point})\end{aligned}$$

**Step 3:** Inserting the result of step 2 in the result of step 1 gives:

$$\mathbb{E}[(y_0 - \hat{f})^2] = \text{Var}(\hat{f}) + [\text{Bias}(\hat{f})]^2 + \text{Var}(\epsilon) \quad (1 \text{ Point})$$

We have proven that that the expected test mean square error (MSE), for a given value  $x_0$ , can always be decomposed into the sum of three fundamental quantities: the variance of  $\hat{f}(x_0)$ , the squared bias of  $\hat{f}(x_0)$  and the variance of the error terms  $\epsilon$ .

- The accuracy of a prediction  $\hat{Y}$  for ground truth  $Y$  depends on two quantities, the reducible and the irreducible error. The **reducible error** refers to the error resulting from the fact that a learned model is not be a perfect estimate for the true relationship. This error can be reduced by a better fit of the algorithm. The **irreducible error** refers to noise that cannot be reduced by a better fit of the algorithm (*even if it were possible to find the perfect true model*). This is, because  $Y$  is also a function of  $\epsilon$ , which, by definition, cannot be predicted using  $X$ . The *irreducible* noise may come from unmeasured variables. There might be useful features in predicting  $Y$ , but if we don't measure them, the model cannot use them for its prediction. For example, the risk of an adverse reaction of a drug may depend on the patient's general feeling of well-being on the day given. Note: the irreducible error provides an upper bound on the accuracy of our prediction for  $Y$ .



**Problem 3** (T, 6 Points). **Gauss-Markov theorem**

1. [2pts] Explain in your own words what the Gauss-Markov theorem is and why it is important.
2. [3pts] Explain in your own words and write down the formal math formulation of the three *Gauss-Markov error term assumptions*.
3. [1pts] Is it also the *best* (in terms of test error) linear unbiased estimate (argue with the bias-variance trade-off)?

*Solution.*

1.
  - [1pts] The Gauss-Markov theorem asserts that the **least squares estimates** of the parameters  $\beta$  have the **smallest variance among all linear unbiased estimates under certain conditions**. *However, there may well exist a biased estimator with smaller mean squared error. Such an estimator would trade a little bias for a larger reduction in variance. Biased estimates are commonly used.*
  - [1pts] **Importance:** When comparing different unbiased estimators for **model selection**, it is interesting to know which one has the highest precision, i.e. lowest variance: we want to make sure that the likelihood of obtaining an estimate very close to the true value is as high as possible. This means we want to use the estimator with the lowest variance of all unbiased estimators. *Note: When estimating regression models, we know that the results of the estimation procedure are random. When we use unbiased estimators, we know that on average, we estimate the true parameter.*
2. The **Gauss-Markov error term assumptions**:
  - [1pts] The error terms have zero mean:  $E[\epsilon_i] = 0$ .
  - [1pts] Homoscedasticity: The error terms have constant finite variance:  $\text{Var}(\epsilon_i) = \sigma^2 < \infty \forall i$
  - [1pts] Independence: The error terms are uncorrelated:  $\text{Cov}(\epsilon_i, \epsilon_j) = 0, \forall i \neq j$
3. [1pts] It is also the *best* linear unbiased estimate in terms of test error, since any unbiased estimate has bias 0. Thus, due to the bias-variance trade-off, variance is the only reducible term influencing test error, thus the least squares estimate is the best.

**Problem 4** (T, 8 Points). **Least Squares Error**

1. [2pts] Let  $Y$  be a random variable. Show that

$$\mathbb{E}(Y) = \underset{c}{\operatorname{argmin}} \mathbb{E}[(Y - c)^2].$$

Explain why this might be an important result for Linear Regression.

2. [6pts] The  $R^2$  statistic is a common measure of model fit corresponding to the fraction of variance in the data that is explained by the model. In general,  $R^2$  is given by the formula

$$R^2 = \frac{\text{TSS} - \text{RSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}.$$

- Show that for univariate regression,  $R^2 = \text{Cor}(X, Y)^2$  holds.
- Show that in the univariate case,  $R^2 = \text{Cor}(Y, \hat{Y})^2$  holds.



*Solution.*

1. Let

$$f(c) = \mathbb{E}[(Y - c)^2] = \mathbb{E}(Y^2) - 2c\mathbb{E}(Y) + \mathbb{E}(c^2) = c^2 - 2\mathbb{E}(Y)c + \mathbb{E}(Y^2)$$

We find the global minimum showing that the first derivative is 0 if  $c = \mathbb{E}(Y)$  and the second derivative is positive at  $c = \mathbb{E}(Y)$ :

$$f'(c) = 2c - 2\mathbb{E}(Y)$$

$$f''(c) = 2$$

Clearly  $f'(\mathbb{E}(Y)) = 0$  and  $f''(\mathbb{E}(Y)) > 0$ .

Thus the best achievable prediction when using squared error loss is the expected value. In practice the expected value is estimated by the mean of all  $Y$ s for a given  $X$ .

This shows, that if we do regression with squared error loss and the distribution of the data is known, the mean gives the best prediction. Note, that the analogue to the mean in the classification setting is the Bayes classifier.

2. (a) First, we list what is given:

$$\begin{aligned} \text{Cor}(X, Y) &= \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \\ \hat{y}_i &= \hat{\beta}_0 + \hat{\beta}_1 x_i \\ \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 &= \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{aligned}$$

We then have a closer look to the RSS:

$$\begin{aligned} \text{RSS} &= \sum_i (y_i - \hat{y})^2 \\ &= \sum_i (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2 \\ &= \sum_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \\ &= \sum_i (y_i - (\bar{y} - \hat{\beta}_1 \bar{x}) - \hat{\beta}_1 x_i)^2 \\ &= \sum_i ((y_i - \bar{y}) + (\hat{\beta}_1 \bar{x} - \hat{\beta}_1 x_i))^2 \\ &= \sum_i ((y_i - \bar{y}) - \hat{\beta}_1 (x_i - \bar{x}))^2 \\ &= \underbrace{\sum_i (y_i - \bar{y})^2}_{=TSS} - 2\hat{\beta}_1 \underbrace{\sum_i (y_i - \bar{y})(x_i - \bar{x})}_* + \hat{\beta}_1^2 \sum_i (x_i - \bar{x})^2 \end{aligned}$$



Since  $R^2 = 1 - \frac{RSS}{TSS}$ , we now investigate the fraction.

$$\begin{aligned}
 \frac{RSS}{TSS} &= \frac{TSS + *}{TSS} \\
 &= 1 + \frac{*}{TSS} \\
 &= 1 - \frac{2\hat{\beta}_1 \sum_i (y_i - \bar{y})(x_i - \bar{x}) - \hat{\beta}_1^2 \sum_i (x_i - \bar{x})^2}{\sum_i (y_i - \bar{y})^2} \\
 &= 1 - \frac{\frac{(\sum_i (y_i - \bar{y})(x_i - \bar{x}))^2}{\sum_i (x_i - \bar{x})^2}}{\sum_i (y_i - \bar{y})^2} \\
 &= 1 - Cor(X, Y)^2 \\
 Cor(X, Y)^2 &= 1 - \frac{RSS}{TSS} = R^2
 \end{aligned}$$

(b)

$$\begin{aligned}
 Cor(Y, \hat{Y})^2 &= \left( \frac{Cov(Y, \hat{Y})}{\sqrt{Var(Y)}\sqrt{Var(\hat{Y})}} \right)^2 \\
 &= \left( \frac{Cov(Y, \hat{\beta}_0 + \hat{\beta}_1 X)}{\sqrt{Var(Y)}\sqrt{Var(\hat{\beta}_0 + \hat{\beta}_1 X)}} \right)^2 \\
 &= \left( \frac{\beta_1 Cov(Y, X)}{\sqrt{Var(Y)}\hat{\beta}_1\sqrt{Var(X)}} \right)^2 \\
 &= Cor(Y, X)^2
 \end{aligned}$$



**Problem 5** (P, 14 Points). Solving this problem will make it easier to solve the following programming one. This exercise uses the *Auto* data set which is provided in the *zip* file. For this problem you need to write your solutions in the Jupyter notebook provided with the assignment `Practical_Problem_1.ipynb`. Reminder:

- Please rename submission files to include your name and matriculation number in your pdf and .pynb file (example: `john_123456_practical_problem_1.ipynb`).
  - Please answer each question in a few sentences as precise and short as possible and focus on the most outstanding facts. Explicitly describe and name the information you retain from plots and tables.
  - Make sure your code runs by running each cell in the Jupyter notebook and save all results.
1. [2pts] Using the `plot()` function in `matplotlib` and create scatterplots between all the variables. Is the relationship between those variables linear? Describe exactly four different connections between the variables. (Exclude the *name* variable, which is qualitative.)
  2. [2pts] Detect the two variable pairs **in the scatterplots** that appear to be **the most highly** correlated and anti-correlated, respectively. Justify your choice using the `np.corrcoef()` function. Do the results from `np.corrcoef()` differ from what you see in the plot?
  3. [4pts] Perform simple linear regression with *mpg* as the response using the variables *cylinders*, *displacement*, *horsepower* and *year*, respectively, as features using the `LinearRegression()` function provided in the `sklearn` package. Which predictors appear to have a statistically significant relationship to the outcome (use an appropriate measurement)? How good are the resulting models (provide all four  $R^2$  values)?
  4. [4pts] Use the `LinearRegression()` function to perform one multiple linear regression with *mpg* as the response and all other variables except *name* as the predictors. Use the `score()` and `get_params()` functions to print the results. Compare the full model to those generated in 5.3: (a) How is the model fit (using  $R^2$ )? (b) What can you observe in the different models concerning the significance of the relationship between response and individual predictors? What does the sign of the coefficient (i.e. of the estimate) tell you about the relationship between the predictor and the response? Provide an example from your fitted model.
  5. [2pts] Use the `plot()` function to produce the residual versus fitted plot of the [multiple](#) linear regression fit. Identify the residual plot. Does the residual plot suggest any non-linearity in the data (provide an explanation)?



**Problem 6 (Bonus). Chebychev Polynomials** So far in this assignment we have used least squares to find linear coefficients of a set of provided features. Here, we explore the case where we can enhance the available features with new ones, that are functions of the given ones, for the task of fitting a smooth function on a set of noisy scalars. For this problem consult notebook `Bonus_Problem_1.ipynb`.

From a regression point of view, we are given pairs  $(x_i, y_i)$ , with  $x_i \in [-1, 1]$ ,  $y_i \in \mathbb{R}$ ,  $i = 1, \dots, n$  and we seek to predict the scalar values  $y_i$  in `y_target.csv`. Here, we assume these points to lie on a polynomial curve on the domain  $[-1, 1]$ . Hence, we assume that the observation follows the functional form

$$y = w_0 + w_1 f_1(x) + \dots + w_D f_D(x) + \epsilon, \quad x \in [-1, 1], \quad \epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2), \quad (6.1)$$

where we use as features  $f_1(x_i), \dots, f_D(x_i)$  the Chebychev polynomials, with  $f_d(x_i)$  of degree  $d = 1, \dots, D$  computed on each input point  $x_i$ .

1. You first need to create a function that generates the feature vector for each input point. For this, you have to compute the Chebychev polynomials on the given input  $x_i$  for each degree from  $1, \dots, D$ . Remember to also add the constant value  $f_0(x_i) = 1$  that corresponds to the intercept  $w_0$  in your feature vectors, too.

*Hint: You might find the commands `np.linspace` and `np.polynomial.chebyshev` useful, check the documentation for details.*

2. Using the value of each Chebychev Polynomial as a feature, perform least squares estimation for the dependent variable  $y$ . For degree  $D \geq 1$ , this will be a Multiple Linear Regression problem. For this create

- a function `fit_cheb`, that is given a vector of scalar inputs, their labels, and the maximum degree of the polynomials to use and computes the best linear coefficients  $w_d$ .
- a function `predict_cheb` that is given the computed weights  $w_i$ , the input vector and the maximum degree of the polynomials to use, and computes the prediction  $\hat{y}$  for each input point.

3. We now analyse the error of our predictions. First, create a function `evaluate_model_on_dataset` that computes the mean squared error of your method, given a training and a testing set.

Use the function `split_data` to split the dataset into a training and a test set. On the same split compute the prediction error for varying degree  $D$  and plot the results.

What is the behaviour of the prediction error with respect to the degree? Based on your observation, which maximum polynomial degree would you use and why?

*Hint: Increasing the degree of the polynomial makes the model less rigid as discussed in the lecture.*

4. Now compare the above splitting scheme with that of the code `split_data_around_point`. Then plot, once again, the prediction error for varying degree  $D$ . What is the behaviour of the prediction error with respect to the two different splitting schemes? Based on your observation, what is a good requirement for the statistics of the training and the testing set?