

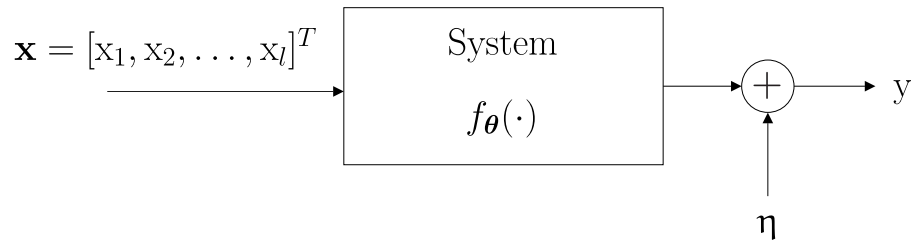
# 02471 Machine Learning for Signal Processing

## Week 2: Parameter Estimation

This exercise is based on S. Theodoridis: Machine Learning, A Bayesian and Optimization Perspective 2nd edition, section(s) 3.1–3.3, 3.5, 3.8–3.11.

The objective of the exercise is to get acquainted with parameter estimation, the bias–variance dilemma and implications with key results.

We consider the following system  $y = f_{\boldsymbol{\theta}}(\mathbf{x}) + \eta$  (see figure).



For this particular exercise, we ignore the choice of  $f_{\boldsymbol{\theta}}(\mathbf{x})$  and simply use linear regression models, and put focus on the estimation task, i.e. how to determine the parameters  $\boldsymbol{\theta}$ .

## Overview

The exercise have the following structure:

2.1 is an exercise for linear regression, where we get acquainted with matrix manipulations and the bias-variance tradeoff.

2.2 will derive central results from ML section 3.5 and 3.9. This exercise will train you to get more comfortable with manipulating matrices and expectations, as well as show how some of the key results are derived.

2.3 is an exercise for ridge regression which is used to explore the bias-variance trade-off further.

Exercises can be carried out in any order, however it is beneficial to complete 2.1 first.

## Notation

- Random variables are denoted with roman font, such as  $\mathbf{x}$ , or  $\boldsymbol{\theta}$ , and their corresponding values (i.e. an observation) are denoted in italic, such as  $x$  or  $\theta$ .
- Random vectors are denoted in bold roman font, such as  $\mathbf{x}$  or  $\boldsymbol{\theta}$ , and their corresponding values are denoted in bold italic, such as  $\mathbf{x}$  and  $\boldsymbol{\theta}$ .
- $J(\boldsymbol{\theta})$  is a cost function that we are seeking to minimize with respect to  $\boldsymbol{\theta}$ .
- $\theta$  denotes the parameters of a model.  $\hat{\theta}$  denotes a point estimate of theta and  $\theta_o$  denotes the optimal value for  $\theta$ .

## Code

The code can be found in the .m and .py files named in the same way as exercises, ie. the code for exercise 2.x.y is in the file 2\_x\_y.m (or .py).

For coding exercises that requires implementation we will usually write `complete this line` where the implementing should be done.

## Solutions

The solution is provided for all derivation exercises, and often hints are provided at the end of the document. If you get stuck, take a look at the hints, and if you are still stuck, take a look in the solution to see the approach being taken. Then try to do it on your own.

Solutions are also provided for some coding exercises. If you get stuck, take a look at the solution, and then try to implement it on your own.

### 2.1 Linear Models

To keep focus on the estimation task, we will use the simple linear regression model. The linear regression model, can for the  $n$ 'th data point be written as

$$y_n = \theta_0 + \theta_1 x_{n,1} + \dots + \theta_l x_{n,l} + \eta = \theta_0 + \boldsymbol{\theta}^T \mathbf{x} + \eta$$

#### Exercise 2.1.1

If you are unfamiliar with Linear Regression and Least Squares, work through example 3.1, page 74 in the ML book.

#### Exercise 2.1.2

The parameters of the linear regression model is estimated by minimizing the cost function

$$J(\boldsymbol{\theta}) = \sum_{n=1}^N (y_n - \boldsymbol{\theta}_b^T \mathbf{x}_n - \theta_0)^2$$

Let  $\boldsymbol{\theta} = [\boldsymbol{\theta}_b^T \ \theta_0]^T$  ( $\boldsymbol{\theta}$  is a column vector). Define a matrix  $X$  and vector  $\mathbf{y}$  such that the cost function can be rewritten to

$$J(\boldsymbol{\theta}) = (\mathbf{y} - X\boldsymbol{\theta})^T (\mathbf{y} - X\boldsymbol{\theta})$$

Validate the expressions are identical e.g. by using a 3-dimensional  $\boldsymbol{\theta}$  vector and 2 datapoints. Make sure to familiarize yourself with these types of rewritings.

#### Exercise 2.1.3

Show that

$$\frac{\partial}{\partial \boldsymbol{\theta}} J(\boldsymbol{\theta}) = -2X^T \mathbf{y} + 2X^T X \boldsymbol{\theta}$$

We can now find the optimal set of parameters  $\boldsymbol{\theta}$  by setting the derivative to zero, which then yields

$$(X^T X) \boldsymbol{\theta} = X^T \mathbf{y}$$

A few remarks to this solution.

- this is the preferred form for numerical solvers since these typically takes input on the form  $A\mathbf{x} = \mathbf{b}$  where we solve for  $\mathbf{x}$ . We will use this later in the coding exercises.
- the cost function  $J(\boldsymbol{\theta})$  is actually the squared error, and not the *mean* squared error (MSE), which is  $\text{MSE} = \frac{1}{N}(\mathbf{y} - X\boldsymbol{\theta})^T(\mathbf{y} - X\boldsymbol{\theta})$ . Note however, the factor  $\frac{1}{N}$  vanishes when the derivative is set to zero, so the two functions have the same optimal solution.

#### Exercise 2.1.4

Inspect the code associated with this exercise (`ex2_1_4`). Run the code to create a training-set with a 2-dimensional input variable and a 1-dimensional output variable.

Identify the different parameters for the linear regression, and identify the line where the model is applied.

#### Exercise 2.1.5

Inspect the code associated with this exercise (`ex2_1_5`) and complete the missing lines.

Use the script to evaluate the training and test errors on independent sets. The sets are generated by the same true weight vector and the same noise variance, for two models.

- Identify the number of parameters for the two models
- Discuss the reason why the models perform best for different training set sizes.

#### Exercise 2.1.6

Compare the training and test errors *per example* as function of the size of the training set for the two models.

#### Exercise 2.1.7

Compare the values of the training and test errors for large training sets with the value of the noise variance.

## 2.2 Estimation

This exercise is purely hand-derivation, where we will show three results. It is important to reflect on the results that you derive, don't let yourself be bogged down in pure algebraic manipulations.

#### Exercise 2.2.1

We will work with the Mean-squared error estimation. Make sure to read sec. 3.9.1 in the book before carrying out this exercise.

We assume the general form of regression

$$y = g(\mathbf{x}) + \eta$$

where  $\eta$  is zero-mean Gaussian noise with variance  $\sigma_\eta^2$ . If we have found an optimal estimate in the mean squared sense using data:

$$\hat{g}(\mathbf{x}) = \mathbb{E}[y|\mathbf{x}]$$

Then,

$$\text{MSE} = \mathbb{E}[(y - \mathbb{E}[y|\mathbf{x}])^2] = \sigma_\eta^2$$

Show the MSE relation. Consider what is the implications of the result?

### Exercise 2.2.2 (This is a difficult exercise)

As a first step towards understanding the parameter estimation task, we will analyze the variance of an unbiased estimator. This analysis does not impose any constraints on the model as such, only on the estimator.

From sec 3.5.1 we have the following result that we want to derive.

Let  $\hat{\theta}_i, i = 1, 2, \dots, m$ , be unbiased estimators of a parameter vector  $\theta$ , so that  $\mathbb{E}[\hat{\theta}_i] = \theta, i = 1, \dots, m$ . Moreover, assume that the respective estimators are uncorrelated to each other and that all have the same (total) variance,  $\sigma^2 = \mathbb{E}\left[(\hat{\theta}_i - \theta)^T (\hat{\theta}_i - \theta)\right]$ . Then by averaging the estimates,

$$\hat{\theta} = \frac{1}{m} \sum_{i=1}^m \hat{\theta}_i$$

the new estimator has total variance:

$$\sigma_c^2 = \mathbb{E}\left[(\hat{\theta} - \theta)^T (\hat{\theta} - \theta)\right] = \frac{1}{m} \sigma^2$$

Show that, by averaging the estimators  $\hat{\theta}_i$ , we obtain  $\sigma_c^2 = \frac{1}{m} \sigma^2$ .

### Exercise 2.2.3 (This is a difficult exercise)

In this exercise we will provide theoretic proof that a biased estimator can obtain a lower MSE compared to an unbiased estimator. To simplify the derivations, we assume that the estimator is a real scalar instead of a vector.

We define the biased estimator as a scaled estimate of the unbiased estimator,

$$\hat{\theta}_b = (1 + \alpha) \hat{\theta}_u$$

where  $\hat{\theta}_u$  is the unbiased estimator. We will show that  $\hat{\theta}_b$  has a lower MSE than the corresponding unbiased estimator for certain choices of  $\alpha$ , ie.  $\text{MSE}(\hat{\theta}_b) < \text{MSE}(\hat{\theta}_u)$ .

In addition, we also show that the norm of the biased estimator is lower than the corresponding unbiased estimator, ie.  $|\hat{\theta}_b| < |\hat{\theta}_u|$ .

The task is the following:

1. Compute  $\text{MSE}(\hat{\theta}_b)$  as a function of  $\alpha$ ,  $\text{MSE}(\hat{\theta}_u)$  and  $\theta_0$ , where  $\theta_0$  are the true optimal parameter (ignore the value for  $\alpha$  in this step).
2. Show that there exists an  $\alpha$  such that  $\text{MSE}(\hat{\theta}_b) \leq \text{MSE}(\hat{\theta}_u)$  by placing bounds on  $\alpha$  using the expression derived in step 1. What can you for instance say about the sign of  $\alpha$ ?
3. Conclude that  $\text{MSE}(\hat{\theta}_b) \leq \text{MSE}(\hat{\theta}_u)$  implies that  $|1 + \alpha| < 1$  and relate that to  $|\hat{\theta}_b|$ .

## 2.3 Bias-variance trade-off

Now we will consider a regularized version of the linear model (Ridge regression), where the parameters  $\boldsymbol{\theta}$  are learned by minimization for a fixed value of  $\lambda$ :

$$L(\boldsymbol{\theta}, \lambda) = \sum_{n=1}^N (y_n - \boldsymbol{\theta}^T \mathbf{x}_n)^2 + \lambda \|\boldsymbol{\theta}\|^2$$

where it is assumed that the bias  $\theta_0$  is encoded as part of the  $\mathbf{x}_n$  vector as in exercise 2.1.2.

### Exercise 2.3.1

Use the results from exercise 2.1 to write up the solution for Ridge regression on the form  $A\mathbf{x} = \mathbf{b}$ .

### Exercise 2.3.2

Inspect the code associated with this exercise (`ex2_3_2`) and complete the requested lines.

Consider why the function is not simultaneously minimized for  $(\boldsymbol{\theta}, \lambda)$ .

The training set averages generalization error in the point  $\mathbf{x}$  can be written as

$$\mathbb{E}_{\mathcal{D}}[(f(\mathbf{x}; \mathcal{D}) - \mathbb{E}[y|\mathbf{x}])^2] = \mathbb{E}_{\mathcal{D}}[(f(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[f(\mathbf{x}; \mathcal{D})])^2] + (\mathbb{E}_{\mathcal{D}}[f(\mathbf{x}; \mathcal{D})] - \mathbb{E}[y|\mathbf{x}])^2$$

where  $\mathbb{E}_{\mathcal{D}}[\cdot]$  is the expectation with respect to training sets.

Hence the average error is split into a variance part, quantifying the variation among solutions for different training sets and a bias part quantifying the performance of the average model with respect to best possible model.

The corresponding code snippet plots the relative amount of variance and bias for a linear model.

### Exercise 2.3.3

Plot the average generalization error, the bias error, and the variance error for a large range of weight decay values. Comment on the two regimes where the generalization error stems from variance and bias respectively.

### Exercise 2.3.4

Consider the following questions;

- What is the role of the weight decay in these two regimes?
- Which weight decay value would you recommend?

### Exercise 2.3.5 (Optional)

In this exercise we will analyze bias-variance trade-off on the example of a simple physical experiment which is tossing the ball up in the air. We consider an experiment in which a ball has been tossed into the air (from the height of an averaged human, 168cm) and its position is

measured as a function of time. The data obtained from this hypothetical experiment is stored in 'h' array.

Inspect the code associated with this exercise (`ex2_3_5`), and answer a following questions:

1. From mechanics we know that the height should vary as the square of the time, so instead of a linear fit to the data we should use a quadratic one:

$$h(t) = h_0 + v_0 t + \frac{1}{2} a t^2$$

Which degree of plotted polynomials would you choose? Explain why.

2. How does the bias and variance change with the degree of polynomials?
3. In which case do we see underfitting and in which overfitting. Explain why.
4. How does the model complexity change with the degree of polynomials?

## HINTS

For exercise 2.2.1

Relate the result to your observations from 2.1.7

For exercise 2.2.2

To solve this problem, first carry out substitution of  $\hat{\theta}$  in the expression  $\sigma_c^2 = \mathbb{E}[(\hat{\theta} - \theta_o)^T(\hat{\theta} - \theta_o)]$ , and then use the result that the estimators are uncorrelated, formally  $\mathbb{E}[(\hat{\theta}_i - \theta_o)^T(\hat{\theta}_j - \theta_o)] = \sigma^2 \delta_{ij}$ , where  $\delta_{ij} = 1$  when  $i=j$  and zero otherwise.

For exercise 2.2.3

To show this, we first need an expression for the MSE of the biased estimator

$$\text{MSE}(\hat{\theta}_b) = \mathbb{E}[(\hat{\theta}_b - \theta_o)^2]$$

Insert by substitution the expression for  $\hat{\theta}_b$ , and then add  $\alpha\theta_o - \alpha\theta_o$  to obtain

$$\text{MSE}(\hat{\theta}_b) = (1 + \alpha)^2 \text{MSE}(\hat{\theta}_u) + \alpha^2 \theta_o^2$$

Next we seek the solution for  $\alpha$  so that

$$\text{MSE}(\hat{\theta}_b) < \text{MSE}(\hat{\theta}_u)$$

We solve this inequality for  $\alpha$  first by substituting the expression for  $\text{MSE}(\hat{\theta}_b)$  and then divide both sides with  $\theta_o^2 + \text{MSE}(\hat{\theta}_u)$  and rearranging. Note you can assume  $\alpha \neq 0$ , since  $\hat{\theta}_b = \hat{\theta}_u$  for  $\alpha = 0$ . This will lead to the solution

$$-\frac{2\text{MSE}(\hat{\theta}_u)}{\text{MSE}(\hat{\theta}_u) + \theta_o^2} < \alpha < 0$$

The final step is to verify that the norm of  $\hat{\theta}_b < \hat{\theta}_u$ . This result is obtained by bounding the lower bound to get the result  $|1 + \alpha| < 1$ .

For exercise 2.3.1

The solution is written using sums in eq. (3.42) in the book.