

# Analysis of Regression and Resampling Methods

Håkon Olav Torvik, Vette Vikenes and Sigurd Sørle Rustad



University of Oslo  
Norway  
October 3, 2021

## CONTENTS

I Introduction	1
II Ordinary Least Square	1
III Bias-variance trade-off and Bootstrapping	1
Bias-variance trade-off decomposition	1
IV Cross-validation	3
V Ridge Regression	3
VI Lasso regression	3
VII Analysis of real data	3
VIII Testing	3
References	4

## I. INTRODUCTION

Regression analysis is a statistical method for fitting a function to data. It is useful for building mathematical models to explain observations. There are several regression methods to achieve this, all with their strengths and weaknesses. We will in this paper study three different methods; ordinary least squares, Ridge and Lasso regression. All the code, results and instructions on running the code can be found in our GitHub repository here<sup>1</sup>

Larger datasets contain more information, giving more accurate models. However, real-world datasets usually have a fixed size, and getting more is practically impossible. For smaller datasets it is then useful to have tools mitigating the effects of little data. Resampling methods does this by running the same data through the regression, without over-fitting the model to the sample data. In addition to the regression methods, we will also study the effect of bootstrapping and cross-validating the data.

In order to study this, we need data to analyze. We are going to use two data sets to study the different methods. One data set we will generate ourselves using the Franke function, given by equation (1). To make the data more realistic we also add normally distributed noise. The other data is real terrain data from here<sup>2</sup>.

$$f(x, y) = \frac{3}{4} \exp\left(-\frac{(9x-2)^2}{4} - \frac{(9y-2)^2}{4}\right) + \frac{3}{4} \exp\left(-\frac{(9x+1)^2}{49} - \frac{(9x+1)^2}{49} - \frac{(9y+1)}{10}\right) \\ + \frac{1}{2} \exp\left(-\frac{(9x-7)^2}{4} - \frac{(9y-3)}{4}\right) - \frac{1}{5} \exp\left(-(9x-4)^2 - (9y-7)^2\right) \quad (1)$$

## II. ORDINARY LEAST SQUARE

## III. BIAS-VARIANCE TRADE-OFF AND BOOTSTRAPPING

### Bias-variance trade-off decomposition

We assume that our true data is generated from a noisy model with normally distributed noise  $\epsilon$  with a mean of zero and standard deviation  $\sigma^2$ , i.e.

$$\mathbf{y} = f(\mathbf{x}) + \epsilon$$

We have approximated this function with our design matrix  $\mathbf{X}$  and our parameters  $\beta$  such that our model becomes  $\tilde{\mathbf{y}} = \mathbf{X}\beta$ , where the values of  $\beta$  were obtained by optimizing the mean squared error via the cost function, given by

$$C(\mathbf{X}, \beta) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2 = \mathbb{E} \left[ (\mathbf{y} - \tilde{\mathbf{y}})^2 \right]$$

---

<sup>1</sup> <https://github.com/sigurdru/FYS-STK4155/tree/main/project1>

<sup>2</sup> <https://github.com/CompPhysics/MachineLearning/tree/master/doc/Projects/2021/Project1/DataFiles>

where  $\mathbb{E}$  is the expected value.

We want to show that the above expression can be written as

$$\mathbb{E} \left[ (\mathbf{y} - \tilde{\mathbf{y}})^2 \right] = \frac{1}{n} \sum_i (f_i - \mathbb{E}[\tilde{\mathbf{y}}])^2 + \frac{1}{n} \sum_i (\tilde{y}_i - \mathbb{E}[\tilde{\mathbf{y}}])^2 + \sigma^2$$

We begin by inserting our model expression for  $\mathbf{y}$  and adding and subtracting  $\mathbb{E}[\tilde{\mathbf{y}}]$  inside the expected value, before we square the expression.

$$\begin{aligned} \mathbb{E} \left[ (\mathbf{y} - \tilde{\mathbf{y}})^2 \right] &= \mathbb{E} \left[ (f(\mathbf{x}) + \epsilon - \tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}] + \mathbb{E}[\tilde{\mathbf{y}}])^2 \right] = \mathbb{E} \left[ ((f(\mathbf{x}) - \mathbb{E}[\tilde{\mathbf{y}}]) + \epsilon + (\mathbb{E}[\tilde{\mathbf{y}}] - \tilde{\mathbf{y}}))^2 \right] \\ &= \mathbb{E} \left[ (f(\mathbf{x}) - \mathbb{E}[\tilde{\mathbf{y}}])^2 + \epsilon^2 + (\mathbb{E}[\tilde{\mathbf{y}}] - \tilde{\mathbf{y}})^2 \right] \\ &\quad + \mathbb{E} [2\epsilon(f(\mathbf{x}) - \mathbb{E}[\tilde{\mathbf{y}}]) + 2\epsilon(\mathbb{E}[\tilde{\mathbf{y}}] - \tilde{\mathbf{y}}) + 2(f(\mathbf{x}) - \mathbb{E}[\tilde{\mathbf{y}}])(\mathbb{E}[\tilde{\mathbf{y}}] - \tilde{\mathbf{y}})] \end{aligned}$$

where the cross terms have been written on a separate line since the expected value is linear. Next we will focus on the cross-terms. Since  $\epsilon$  is normally distributed, it's expected value is simply the mean, which is zero in our case. The two cross terms involving  $\epsilon$  is therefore zero, so we only need to consider

$$\mathbb{E} [(f(\mathbf{x}) - \mathbb{E}[\tilde{\mathbf{y}}])(\mathbb{E}[\tilde{\mathbf{y}}] - \tilde{\mathbf{y}})] = \mathbb{E} [f(\mathbf{x})\mathbb{E}[\tilde{\mathbf{y}}]] - \mathbb{E} [f(\mathbf{x})\tilde{\mathbf{y}}] - \mathbb{E} [\mathbb{E}[\tilde{\mathbf{y}}]\mathbb{E}[\tilde{\mathbf{y}}]] + \mathbb{E} [\tilde{\mathbf{y}}\mathbb{E}[\tilde{\mathbf{y}}]]$$

Since the expected value of an expected value is just the expected value itself the last two terms in the above equation both become  $\mathbb{E}[\tilde{\mathbf{y}}]^2$ , canceling each other out. Using that  $f(\mathbf{x})$  is a deterministic function, we have  $\mathbb{E}[f(\mathbf{x})] = f(\mathbf{x})$ . Expressing  $f(\mathbf{x})$  in terms of its expected value, we can write the first two terms in the above equation as

$$\begin{aligned} \mathbb{E} [f(\mathbf{x})\mathbb{E}[\tilde{\mathbf{y}}]] - \mathbb{E} [f(\mathbf{x})\tilde{\mathbf{y}}] &= \mathbb{E} [\mathbb{E} [f(\mathbf{x})] \mathbb{E}[\tilde{\mathbf{y}}]] - \mathbb{E} [\mathbb{E} [f(\mathbf{x})] \tilde{\mathbf{y}}] \\ &= \mathbb{E} [f(\mathbf{x})] \mathbb{E}[\tilde{\mathbf{y}}] - \mathbb{E} [f(\mathbf{x})] \mathbb{E}[\tilde{\mathbf{y}}] = 0 \end{aligned}$$

Hence, all the cross terms in the expected value cancel out, and we're left with

$$\mathbb{E} \left[ (\mathbf{y} - \tilde{\mathbf{y}})^2 \right] = \mathbb{E} \left[ (f(\mathbf{x}) - \mathbb{E}[\tilde{\mathbf{y}}])^2 \right] + \mathbb{E} \left[ (\mathbb{E}[\tilde{\mathbf{y}}] - \tilde{\mathbf{y}})^2 \right] + \mathbb{E} [\epsilon^2]$$

Using that  $\mathbb{E}[\epsilon^2] = \sigma^2$  and writing the expected values as sums we finally arrive at

$$\mathbb{E} \left[ (\mathbf{y} - \tilde{\mathbf{y}})^2 \right] = \frac{1}{n} \sum_i (f_i - \mathbb{E}[\tilde{\mathbf{y}}])^2 + \frac{1}{n} \sum_i (\tilde{y}_i - \mathbb{E}[\tilde{\mathbf{y}}])^2 + \sigma^2$$

which is what we wanted to show. git

#### IV. CROSS-VALIDATION

#### V. RIDGE REGRESSION

#### VI. LASSO REGRESSION

#### VII. ANALYSIS OF REAL DATA

#### VIII. TESTING

In order to make sure our algorithms are running correctly, it is necessary to test our algorithms. We did this by comparing our results to those produced by scikit-learn. First of all we generated some simpler data for testing, namely an exponential:

$$f_{\text{Test}}(x) = \exp(x) + \epsilon. \quad (2)$$

Here  $\epsilon$  denotes normally distributed noise, and  $x$  runs from  $x_{\min} = 0$  to  $x_{\max} = 1$  in  $N = 50$  randomly distributed steps. This generates the testing data visualized in figure 1.

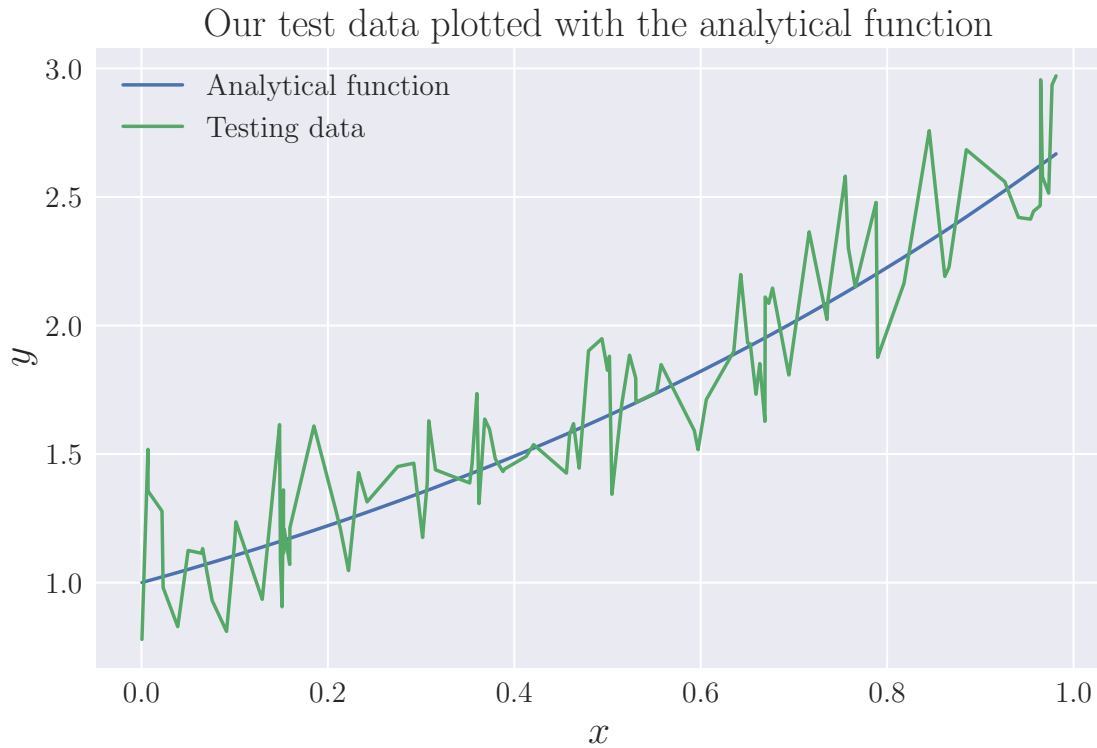


Figure 1. Here you we have plotted the testing data along with the analytical function.

First off we want to test the regression methods we have written.

---

[1] Ref.