

# REGRESSION AND RESAMPLING TECHNIQUES

## FYS-STK4155

### ASSIGNMENT ONE

Ivar Haugaløkken Stangeby

September 10, 2019

## 1 Introduction

Regression analysis is the act of estimating relationships among variables. In this project, we study various regression methods in more detail. In particular, we compare the *ordinary least squares* (OLS) method with the *Ridge regression* and *Lasso regression* techniques. As we shall see, these three methods are all variations over the same theme. We start by testing the methods on noisy data sampled from a function known as *Franke's* bivariate test function, which has been in widespread use in the evaluation of interpolation and data fitting techniques. Finally, we run regression on real terrain data, comparing the aforementioned methods.

## 2 Regression Techniques

In general, the goal of regression analysis is to fit a *model function*  $f(\mathbf{x}, \beta)$  to a set of  $n$  data points  $\Omega = (\mathbf{x}_i, y_i)_{i=1}^n$ . A simple example is that of a linear polynomial with two parameters:

$$f(x, \beta) = \beta_0 + \beta_1 x. \quad (1)$$

The *model parameters*  $\beta$  are determined in order to minimize a suitable *cost function*  $C(\Omega, \beta)$  which measures to which extent the model function manages to capture trends in the data  $\Omega$ . It is the choice of cost function  $C(\Omega, \beta)$  which distinguishes the three regression techniques, OLS, Lasso regression, and Ridge regression.

*Remark 1.* It is often assumed a priori that the data is infact generated from a noisy model, such that each  $y_i$  can be described as

$$y_i = f(\mathbf{x}_i) + \varepsilon_i \quad (2)$$

where each  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$  is normally distributed with zero mean and variance  $\sigma^2$ . This assumption on the error gives rise to what is known as *general linear models*.

**The design matrix.** We are often interested in finding the best model function in a specific function space. Assuming this space has a basis  $\mathcal{B} = \{\varphi_i\}_{i=1}^M$ , we may write our model function in terms of the basis functions and the model parameters as

$$f(\mathbf{x}, \beta) := \sum_{i=1}^M \beta_i \varphi_i(\mathbf{x}). \quad (3)$$

As an example, if we were to use the space  $\Pi_2^2$  of bi-quadratic polynomials, our basis functions would be

$$\mathcal{B} = \{1, x, x^2, y, y^2, xy\}. \quad (4)$$

With this formulation, we can represent the approximation  $\hat{\mathbf{y}}$  to  $\mathbf{y}$  as a matrix product

$$\hat{\mathbf{y}} = \mathbf{X}\beta \quad (5)$$

where  $\mathbf{X}$  is the *design matrix* defined by components  $X_{ij} = \varphi_j(\mathbf{x}_i)$ , and  $\beta = [\beta_1, \dots, \beta_M]$  is the model parameters.

**Performance metrics.** In order to evaluate the how accurately the model function  $f(\Omega, \beta)$  captures trends in the target data  $\mathbf{y}$ , a few standard performance metrics are used. Firstly, the *mean squared error*:

$$\text{MSE}(\mathbf{y}, \hat{\mathbf{y}}) := \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (6)$$

which simply averages the squared error over all samples and estimates. Secondly, we have the *coefficient of determination* or  $R^2$ -score:

$$R^2(\mathbf{y}, \hat{\mathbf{y}}) := 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}. \quad (7)$$

The  $R^2$ -score measures much of the variation in  $\mathbf{y}$  which can be attributed to a simple linear relation between  $\mathbf{x}$  and  $\mathbf{y}$ . It is a ratio, thus a value of one tells us that all the variation in the data can be attributed to an approximate linear relationship between the data  $\mathbf{x}$  and the response variable  $\mathbf{y}$ . A value of zero indicates that a non-linear model may be preferable.

## 2.1 Ordinary Least Squares

One common cost function is the one which involves the sum of squared residuals (or squared errors):

$$C(\Omega, \beta) = \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 := \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (8)$$

where  $\hat{y}_i := f(\mathbf{x}_i, \beta)$ . The method involving the minimization of this specific cost function is known as *least squares*. With the matrix notation from above in mind, we can also write the cost function as

$$C(\Omega, \beta) = \frac{1}{n} (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}). \quad (9)$$

**Optimizing the parameters in OLS.** Assume now that our model is of the form given in Equation (3) and that our cost function is the mean squared error defined in Equation (8). We are interested in finding the parameters  $\beta$  that minimize the cost function  $C(\Omega, \beta)$ . Since  $C(\Omega, \beta)$  is convex, it suffices to differentiate with respect to  $\beta$  and equating to zero.

We have that

$$\frac{\partial C(\Omega, \beta)}{\partial \beta} = \mathbf{X}^T (\mathbf{y} - \mathbf{X}\beta), \quad (10)$$

and equating this to zero yields the familiar *normal equations*

$$\mathbf{X}^T \mathbf{y} = \mathbf{X}^T \mathbf{X} \beta. \quad (11)$$

If the matrix  $X^T X$  is invertible, we may obtain the solution by direct numerical matrix inversion. In this case, the optimal model parameters are found directly by

$$\beta = (X^T X)^{-1} X^T y. \quad (12)$$

However, these matrices may be ill-conditioned when the number of equations are very large, and it is therefore common to apply approximate solvers for the inverse, for instance using low-rank approximation based on SVD-decomposition, which we will briefly turn to in the following.

**Singular value decomposition.** Recall that any matrix  $A \in \mathbb{C}^{n,m}$  can be decomposed as

$$A = U \Sigma V^T, \quad (13)$$

where  $U$  and  $V$  are comprised of the eigenvectors of  $AA^T$  and  $A^T A$  respectively. As these eigenvectors are orthonormal, it follows that both  $U$  and  $V$  are unitary matrices. Furthermore,  $\Sigma$  is a matrix

$$\Sigma = \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{bmatrix}, \quad (14)$$

where  $\Sigma_1$  is a square diagonal matrix of size  $r \times r$  with the non-zero singular values of  $A$ . The integer  $r$  is the *rank* of  $A$ . As the matrix  $\Sigma$  is mostly containing zeros, the information stored in  $A$  is attributed to only some parts of  $U$  and  $V$ . We can remove the redundant parts, and more compactly express  $A$  as  $A = U_1 \Sigma_1 V_1^T$  without loss of accuracy in the decomposition. Here,  $U_1$  is  $m \times r$  and  $V_1$  is  $n \times r$ .

Applying the singular value decomposition to the matrix  $X = U \Sigma V^T$  we can analyze the expression for the prediction  $\hat{y}$  in terms of the matrix  $X^T X$ . First of all, we have that

$$\begin{aligned} X^T X &= V \Sigma^T U^T U \Sigma V^T \\ &= V D V^T, \end{aligned} \quad (15)$$

where we have used that  $U U^T = I$  and defined  $D = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$ . Thus, plugging this into Equations (5) and (12), we obtain the expression

$$\hat{y} = X \beta = X (V D V^T)^{-1} X^T y. \quad (16)$$

Finally, substituting  $X = U\Sigma V^T$  and noting that diagonal matrices always commute, we end up with

$$\hat{y} = UU^T y. \tag{17}$$