

# FYS-STK4155 - Applied data analysis and machine learning

## Project 1

Even M. Nordhagen

September 13, 2018

- Github repository containing programs and results:  
<https://github.com/evenmn/FYS-STK4155>

### **Abstract**

Do not forget to be specific

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Theory</b>	<b>3</b>
2.1	Regression . . . . .	3
2.1.1	Ordinary Least Square (OLS) . . . . .	3
2.1.2	Ridge regression . . . . .	4
2.1.3	Lasso regression . . . . .	5
2.1.4	General form . . . . .	5
2.2	Higher order regression . . . . .	5
2.2.1	Terrain . . . . .	5
2.2.2	Higher order . . . . .	5
2.3	Error analysis . . . . .	6
<b>3</b>	<b>Methods</b>	<b>6</b>
3.1	Resampling techniques . . . . .	6
3.1.1	Bootstrap method . . . . .	6
3.1.2	K-fold validation method . . . . .	6
3.2	Singular Value Decomposition (SVD) . . . . .	6
3.3	Minimization methods . . . . .	6
3.3.1	Gradient Descent . . . . .	7
<b>4</b>	<b>Code</b>	<b>8</b>
4.1	Code structure . . . . .	8
4.2	Implementation . . . . .	8
4.3	Optimalization . . . . .	8
<b>5</b>	<b>Results</b>	<b>8</b>
<b>6</b>	<b>Discussion</b>	<b>8</b>
<b>7</b>	<b>Conclusion</b>	<b>8</b>
<b>A</b>	<b>Appendix A</b>	<b>8</b>

# 1 Introduction

Write some motivating words about regression analysis

## 2 Theory

### 2.1 Regression

A few general words about regression

#### 2.1.1 Ordinary Least Square (OLS)

Suppose we have a set of points  $\{(x_0, y_0), (x_1, y_1), \dots, (x_{n-1}, y_{n-1})\}$ , and we want to fit a p'th order polynomial to them. The most intuitive way would be to find coefficients  $\vec{\beta}$  which minimize the error in

$$\begin{aligned}y_0 &= \beta_0 x_0^0 + \beta_1 x_0^1 + \dots + \beta_p x_0^p + \varepsilon_0 \\y_1 &= \beta_0 x_1^0 + \beta_1 x_1^1 + \dots + \beta_p x_1^p + \varepsilon_1 \\&\vdots \\y_{n-1} &= \beta_0 x_{n-1}^0 + \beta_1 x_{n-1}^1 + \dots + \beta_p x_{n-1}^p + \varepsilon_{n-1},\end{aligned}$$

which for OLS is defines as

$$\text{MSE} = \sum_i \varepsilon_i^2 \tag{1}$$

NEED TO REWRITE THIS + COST FUNCTION

Standard cost function

$$Q(\vec{\beta}) = \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2 = (\vec{y} - \hat{X}\vec{\beta})^T (\vec{y} - \hat{X}\vec{\beta}) \tag{2}$$

Instead of dealing with a set of equations, we can apply linear algebra. One can easily see that the equations above correspond to

$$\vec{y} = \hat{X}^T \vec{\beta} + \vec{\varepsilon}. \tag{3}$$

For a nonsingular matrix  $\hat{X}$  (but not necessary symmetric) we can find the optimal  $\vec{\beta}$  by solving

$$\vec{\beta} = (\hat{X}^T \hat{X})^{-1} \hat{X}^T \vec{y}, \tag{4}$$

which again corresponds to minimizing the cost function,

$$\vec{\beta} = \operatorname{argmin}, \vec{\beta} \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \right\}. \quad (5)$$

CONFIDENCE INTERVAL of  $\hat{\beta}$ :  $\operatorname{Var}(\hat{\beta})$ .

This works perfectly when all rows in  $\hat{X}$  are linearly independent, but this will generally not be the case for large data sets. If we are not able to diagonalize the matrix, we will not be able to calculate  $(\hat{X}^T \hat{X})^{-1}$ , so we need to do something smart.

Fortunately there is a simple trick we can do to make all matrices diagonalizable; we can add a diagonal matrix to the initial matrix.

### 2.1.2 Ridge regression

Ridge regression is a widely used method that can handle singularities in matrices. The idea is to modify the standard cost function by adding a small term,

$$Q^{\text{ridge}}(\vec{\beta}) = \sum_{i=1}^N (y_i - \tilde{y}_i)^2 + \lambda \|\vec{\beta}\|_2^2, \quad (6)$$

where  $\lambda$  is the so-called *penalty* and  $\|\vec{v}\|_2$  is defined as

$$\|\vec{v}\|_2 = \sqrt{\vec{v}^T \vec{v}} = \left( \sum_{i=1}^N v_i^2 \right)^{1/2}. \quad (7)$$

This will eliminate the singularity problem.

Further we find the optimal  $\vec{\beta}$  values by minimizing the function

$$\vec{\beta}^{\text{ridge}} = \operatorname{argmin}, \vec{\beta} \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}, \quad (8)$$

or we could simply solve the equation

$$\vec{\beta}^{\text{ridge}} = (\hat{X}^T \hat{X} + \lambda I)^{-1} \hat{X}^T \vec{y}. \quad (9)$$

In the latter equation we can easily see why this solves our problem.

### 2.1.3 Lasso regression

The idea behind Lasso regression is similar to the idea behind Ridge regression, and they differ only by the exponent factor in the last term. The modified cost function now writes

$$Q^{\text{lasso}}(\vec{\beta}) = \sum_{i=1}^N (y_i - \tilde{y}_i)^2 + \lambda \|\vec{\beta}\|_2, \quad (10)$$

and to find the optimal coefficients  $\vec{\beta}$ , we need to minimize

$$\vec{\beta}^{\text{lasso}} = \operatorname{argmin}, \vec{\beta} \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j \right\}. \quad (11)$$

### 2.1.4 General form

We can generalize the models above to a minimization problem where we have a  $q$  in the last exponent,

$$\vec{\beta}^q = \operatorname{argmin}, \vec{\beta} \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^q \right\}, \quad (12)$$

such that  $q = 2$  corresponds to the Ridge method and  $q = 1$  is associated with Lasso regression. It can also be interesting to try other  $q$ -values.

## 2.2 Higher order regression

We can use the same approach as above when dealing with regression of higher order, since the problem is to fit a function to points, no matter how many components they have. We will first take a look at how we can fit a 2D polynomial to some terrain data, before we briefly describe how to fit a function of arbitrary order to points.

### 2.2.1 Terrain

A set of data points  $\{(x_0, y_0, z_0), (x_1, y_1, z_1), \dots, (x_{n-1}, y_{n-1}, z_{n-1})\}$  gives some coordinates in space, which for instance can describe the terrain. Mention the Franke Function

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

### 2.2.2 Higher order

Although we stick to 2D regression in this project, I add this section for completeness.

## 2.3 Error analysis

Cost function (loss function)

Different methods to estimate error:

- Absolute error
- Relative error
- Mean square error (MSE)
- $R^2$  score function

## 3 Methods

### 3.1 Resampling techniques

A resampling technique is.. There are plenty of resampling techniques, and we have already went through several of them in this course:

- Validation set approaches
- Leave one out validation
- Jackknife resampling
- K-fold validation
- Bootstrap method
- Blocking method.

For this particular project we have been focusing on the bootstrap and the k-fold validation methods, so here I will cover them only

### 3.1.1 Bootstrap method

### 3.1.2 K-fold validation method

## 3.2 Singular Value Decomposition (SVD)

## 3.3 Minimization methods

When the interaction term is excluded, we know which  $\alpha$  that corresponds to the energy minimum, and it is in principle no need to try different  $\alpha$ 's. However, sometimes we have no idea where to search for the minimum point, and we need to try various  $\alpha$  values to determine the lowest energy. If we do not know where to start searching, this can be a time consuming activity. Would it not be nice if the program could do this for us?

In fact there are multiple techniques for doing this, where the most complicated ones obviously also are the best. Anyway, in this project we will have good initial guesses, and are therefore not in need for the most fancy algorithms.

### 3.3.1 Gradient Descent

Perhaps the simplest and most intuitive method for finding the minimum is the gradient descent method (GD), which reads

$$\alpha^+ = \alpha - \eta \cdot \frac{d\langle E(\alpha) \rangle}{d\alpha}. \quad (13)$$

where  $\alpha^+$  is the updated  $\alpha$  and  $\eta$  is a step size. The idea is that one finds the gradient of the energy with respect to a certain  $\alpha$ , and moves in the direction which minimizes the energy. This is repeated until one has found an energy minimum, where the energy minimum is defined as either where  $\frac{d\langle E(\alpha) \rangle}{d\alpha}$  is smaller than a given tolerance, or the energy fluctuates around a value are smaller than a tolerance, and thus changes minimally.

To implement equation 13, we need an expression for the derivative of  $E$  with respect to alpha:

$$\bar{E}_\alpha = \frac{d\langle E(\alpha) \rangle}{d\alpha}. \quad (14)$$

By using the expression for the expectation value for the energy  $\langle E(\alpha) \rangle$  in equation 15

$$\langle E(\alpha) \rangle = \frac{\langle \psi_T(\alpha) | H | \psi_T(\alpha) \rangle}{\langle \psi_T(\alpha) | \psi_T(\alpha) \rangle} \quad (15)$$

and applying the chain rule of differentiation, it can be shown that equation 14 is equal to equation 16

$$\bar{E}_\alpha = 2 \left[ \langle E_L(\alpha) \frac{\bar{\psi}_\alpha}{\psi_\alpha} \rangle - \langle E_L(\alpha) \rangle \langle \frac{\bar{\psi}_\alpha}{\psi_\alpha} \rangle \right] \quad (16)$$

where

$$\bar{\psi}_\alpha = \frac{d\psi(\alpha)}{d\alpha}. \quad (17)$$

The algorithm of this minimization method is thus as follows:

```

for (max number of iterations with minimizing)

    do M Monte Carlo cycles
    calculate E and dE/dalpha

    Check if dE/dalpha < eps or alpha fluctuation over the last 5 steps
        ↪ is < eps

        if yes, print optimal alpha and break loop
        if no, continue to next iteration

```

## 4 Code

### 4.1 Code structure

### 4.2 Implementation

### 4.3 Optimalization

## 5 Results

## 6 Discussion

## 7 Conclusion

## A Appendix A