

FYS-STK4155 - Applied data analysis and machine learning

Project 1

Even M. Nordhagen

October 5, 2018

- Github repository containing programs and results:

<https://github.com/evenmn/FYS-STK4155>

Abstract

The aim of this project is to study the performance of linear regression in order to fit a two dimensional polynomial to terrain data. Both Ordinary Least Square (OLS), Ridge and Lasso regression methods were implemented, and for minimizing Lasso's cost function Gradient Descent (GD) was used. A fourth method was to minimize the cost function of Ridge using GD. The fitted polynomial was visualized and compared with the data, the Mean Square Error (MSE) and R^2 -score were analyzed, and finally the polynomial coefficients were studied applying visualization tools and Confidence Intervals (CI). To benchmark the results, we used Scikit Learn.

We found the self-implemented OLS and Ridge regression functions to reproduce the benchmarks, and Lasso was close to reproducing the benchmark as well. Anyway, the difference between results gotten from standard Ridge regression and when minimizing its cost function is large. The OLS regression method is considered as the most successful due to its small MSE and high R^2 -score.

Contents

1	Introduction	3
2	Theory	3
2.1	Linear regression	3
2.1.1	Ordinary Least Square (OLS)	4
2.1.2	Ridge regression	5
2.1.3	Lasso regression	5
2.1.4	General form	6
2.2	Multivariate linear regression	6
2.2.1	Terrain	6
2.2.2	Higher order	7
2.3	Some statistics	9
2.3.1	Confidence interval	10
2.4	Error analysis	10
2.4.1	Mean Square Error (MSE)	11
2.4.2	R ² score function	11
3	Methods	12
3.1	Resampling techniques	12
3.1.1	Bootstrap method	13
3.1.2	K-fold validation method	14
3.2	Minimization methods	14
3.2.1	Gradient Descent	15
3.3	Benchmark	16
4	Code	16
4.1	Code structure	17
4.2	Implementation and optimization	17
5	Results	18
5.1	Franke function	18
5.1.1	Visualization of graphes	18
5.1.2	Error	20
5.1.3	Coefficients β	21
5.2	Real data	24
5.2.1	Visualization of graphes	24
5.2.2	Error	26
5.2.3	Coefficients β	27

6	Discussion	29
7	Conclusion	30
8	References	31

1 Introduction

The linear regression methods were first introduced for more than two centuries ago, and have been used in a large number of fields throughout the years [**Gauss**][**Legende**]. In this project we will examine whether the methods are sufficient for fitting polynomials to real terrain data, or we need more complicated methods. To challenge the methods, we chose terrain data from the volcanic island of Lombok, Indonesia, where the contour lines are quite dense.

We developed our own software for ordinary least square (OLS), Ridge and Lasso linear regression, where the latter was based on minimization using gradient descent (GD). To verify the implementation, we tested it on data from the Franke function where we knew what the result should be. Further, the error was analyzed in order to decide which method that gave the best result, and all data was resampled using the bootstrap method to find the correct sampling variance.

For the results, see section *Results* (5), which again is discussed in section *Discussion* (6). The background theory can be found in section *Theory* (2), and all methods and techniques are presented in the section *Methods* (3). For code structure and implementation, see section *Code* (4), and finally, the conclusion is found in section (7) with the same name.

2 Theory

Regression analysis are widely used in different fields of natural sciences, data science and economics to mention some. In fact, as fast as we a data set that needs to be analyzed, regression is unavoidable. The simplest form is the linear regression, which is intuitive and easy to work with.

2.1 Linear regression

In linear regression, the dependent variable y_i is a linear combination of the parameters, and for a dependent variable this can be written as

$$y_i = \sum_j X_{ij}\beta_j. \tag{1}$$

In principle, $x_i j$ can be an arbitrary function of the arguments x_i , but in this project we will only fit the data with a polynomial function. [Hastie]

2.1.1 Ordinary Least Square (OLS)

Suppose we have a set of points $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, 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

$$\begin{aligned} y_1 &= \beta_0 x_1^0 + \beta_1 x_1^1 + \dots + \beta_p x_1^p + \varepsilon_1 \\ y_2 &= \beta_0 x_2^0 + \beta_1 x_2^1 + \dots + \beta_p x_2^p + \varepsilon_2 \\ &\vdots \\ y_n &= \beta_0 x_n^0 + \beta_1 x_n^1 + \dots + \beta_p x_n^p + \varepsilon_n, \end{aligned}$$

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

$$\vec{y} = \hat{X}^T \vec{\beta} + \vec{\varepsilon}, \quad (2)$$

with

$$\hat{X} = \begin{pmatrix} x_1^0 & x_1^1 & x_1^2 & \dots & x_1^p \\ x_2^0 & x_2^1 & x_2^2 & \dots & x_2^p \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_n^0 & x_n^1 & x_n^2 & \dots & x_n^p \end{pmatrix} \quad (3)$$

as the *design* matrix and

$$\vec{\beta} = (\beta_0, \beta_1, \dots, \beta_p) \quad (4)$$

as the coefficients.

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

$$\vec{\beta} = (\hat{X}^T \hat{X})^{-1} \hat{X}^T \vec{y}, \quad (5)$$

which again corresponds to minimizing the cost function,

$$\vec{\beta} = \operatorname{argmin}_{\vec{\beta}} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \right\}. \quad (6)$$

where the standard cost function is used, given by least squares,

$$Q(\vec{\beta}) = \sum_{i=1}^n (y_i - \tilde{y}_i)^2 = (\vec{y} - \hat{X} \vec{\beta})^T (\vec{y} - \hat{X} \vec{\beta}). \quad (7)$$

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,

$$Q^{\text{ridge}}(\vec{\beta}) = \sum_{i=1}^N (y_i - \tilde{y}_i)^2 + \lambda \|\vec{\beta}\|_2^2, \quad (8)$$

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

$$\|\vec{v}\|_2 = \sqrt{\vec{v}^T \vec{v}} = \left(\sum_{i=1}^N v_i^2 \right)^{1/2}. \quad (9)$$

This will eliminate the singularity problem.

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

$$\vec{\beta}^{\text{ridge}} = \text{argmin}, \vec{\beta} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}, \quad (10)$$

or we could simply solve the equation

$$\vec{\beta}^{\text{ridge}} = (\hat{X}^T \hat{X} + \lambda I)^{-1} \hat{X}^T \vec{\beta}. \quad (11)$$

In the latter equation we can easily see why this solves our problem. For this case, we expect the coefficients to be smaller compared with OLS, since the cost function is directly dependent on $\vec{\beta}^2$.

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

$$Q^{\text{lasso}}(\vec{\beta}) = \sum_{i=1}^N (y_i - \tilde{y}_i)^2 + \lambda |\vec{\beta}|, \quad (12)$$

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

$$\vec{\beta}^{\text{lasso}} = \operatorname{argmin}, \vec{\beta} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}. \quad (13)$$

We expect the coefficients to be even smaller for Lasso regression than for Ridge regression, since the cost function is directly dependent on $\vec{\beta}$.

2.1.4 General form

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

$$\vec{\beta}^q = \operatorname{argmin}, \vec{\beta} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j|^q \right\}, \quad (14)$$

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 Multivariate linear regression

The same approach as above can be used 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 an arbitrary order to points.

2.2.1 Terrain

A set of data points $\{(x_1, y_1, z_1), (x_2, y_2, z_2), \dots, (x_n, y_n, z_n)\}$ gives some coordinates in space, which for instance can describe the terrain. The system of linear equations to solve can then be lined up as

$$\begin{aligned} z_1 &= \beta_0 x_1^0 y_1^0 + \beta_1 x_1^1 y_1^0 + \beta_2 x_1^0 y_1^1 + \dots + \beta_{p^2} x_1^p y_1^p + \varepsilon_1 \\ z_2 &= \beta_0 x_2^0 y_2^0 + \beta_1 x_2^1 y_2^0 + \beta_2 x_2^0 y_2^1 + \dots + \beta_{p^2} x_2^p y_2^p + \varepsilon_2 \\ &\vdots \\ z_n &= \beta_0 x_n^0 y_n^0 + \beta_1 x_n^1 y_n^0 + \beta_2 x_n^0 y_n^1 + \dots + \beta_{p^2} x_n^p y_n^p + \varepsilon_n, \end{aligned}$$

when fitting a polynomial of p 'th order. In general the order associated with x -direction do not need to be the same as the order associated with y -direction.

We can set up a similar equation as for the first order case, presented in equation (2), but we will now (typically) have \vec{z} on the left hand side and \hat{X} will contain both x- and y-coordinates:

$$\vec{z} = \hat{X}^T \vec{\beta}. \quad (15)$$

The design matrix \hat{X} will now look like

$$\hat{X} = \begin{pmatrix} x_1^0 y_1^0 & x_1^1 y_1^0 & x_1^0 y_1^1 & x_1^1 y_1^1 & \dots & x_1^p y_1^p \\ x_2^0 y_2^0 & x_2^1 y_2^0 & x_2^0 y_2^1 & x_2^1 y_2^1 & \dots & x_2^p y_2^p \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ x_n^0 y_n^0 & x_n^1 y_n^0 & x_n^0 y_n^1 & x_n^1 y_n^1 & \dots & x_n^p y_n^p \end{pmatrix} \quad (16)$$

and we can again use the OLS method to find $\vec{\beta}$, i.e, calculating

$$\vec{\beta} = (\hat{X}^T \hat{X})^{-1} \hat{X}^T \vec{z}. \quad (17)$$

Similarly, we can use Ridge and Lasso regression in the same way as when fitting 1D polynomials.

In this project we will actually deal with terrain data. Firstly, we implement some regression tools which fit a 2D function to points in space. To verify the implementation, we pick points from a known function such that we know the result. The specific verification function used is the Franke Function, which looks like

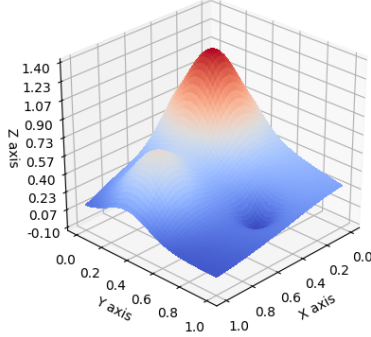
$$\begin{aligned} f(x, y) = & \frac{3}{4} \exp \left(-\frac{(9x-2)^2}{4} - \frac{(9y-2)^2}{4} \right) + \frac{3}{4} \exp \left(-\frac{(9x+1)^2}{49} - \frac{(9y+1)}{10} \right) \\ & + \frac{1}{2} \exp \left(-\frac{(9x-7)^2}{4} - \frac{(9y-3)^2}{4} \right) - \frac{1}{5} \exp \left(-(9x-4)^2 - (9y-7)^2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

and is a smooth 2D function as one can see in figure (1), part (a).

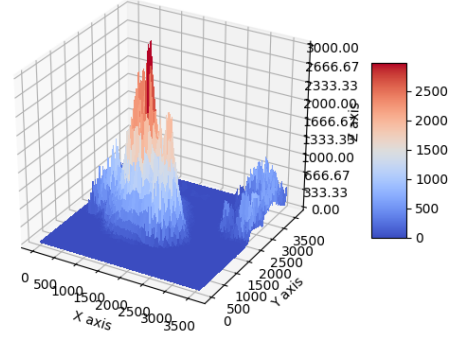
Secondly, we use the verified regression tools to fit a polynomial to real terrain data. The data is taken from United States Geological Survey (USGS) [USGS], and for investigation we picked the volcanic island of Lombok. Despite the island's small extent, it houses the second highest volcano in Indonesia which makes the terrain data interesting. See part (b) of figure (1) for terrain data.

2.2.2 Higher order

There should be no surprise that we can extend the theory above to even higher orders. Although we stick to 2D regression in this project, we add this section for completeness.



(a) Franke function



(b) Lombok terrain

Figure 1: The two data sets we are going to fit polynomials to.

Suppose now that we have a data set of n points in d dimensions, $\{(x_1^1, x_1^2, \dots, x_1^d), (x_2^1, x_2^2, \dots, x_2^d), \dots, (x_n^1, x_n^2, \dots, x_n^d)\}$, where the superscript indicates in which dimension the coordinate is, and the subscript is the coordinate number. We want to fit a polynomial of degree $p \in \mathbb{R}^d$ to the points, where we assume that the fitting polynomial has the same degree in all directions, which makes the notation neater, but is not essential. To be more specific, we want polynomial coefficients $\vec{\beta} = (\beta_0, \beta_1, \dots, \beta_{p^d})$ that minimizes the error in

$$\begin{aligned} x_1^d &= \beta_0(x_1^1)^0(x_1^2)^0 \dots (x_1^d)^0 + \beta_1(x_1^1)^1(x_1^2)^0 \dots (x_1^d)^0 + \dots + \beta_{p^d}(x_1^1)^p(x_1^2)^p \dots (x_1^d)^p + \varepsilon_1 \\ x_2^d &= \beta_0(x_2^1)^0(x_2^2)^0 \dots (x_2^d)^0 + \beta_1(x_2^1)^1(x_2^2)^0 \dots (x_2^d)^0 + \dots + \beta_{p^d}(x_2^1)^p(x_2^2)^p \dots (x_2^d)^p + \varepsilon_2 \\ &\vdots \\ x_n^d &= \beta_0(x_n^1)^0(x_n^2)^0 \dots (x_n^d)^0 + \beta_1(x_n^1)^1(x_n^2)^0 \dots (x_n^d)^0 + \dots + \beta_{p^d}(x_n^1)^p(x_n^2)^p \dots (x_n^d)^p + \varepsilon_n. \end{aligned}$$

This gives a design matrix

$$\hat{X} = \begin{pmatrix} (x_1^1)^0(x_1^2)^0 \dots (x_1^d)^0 & \dots & (x_1^1)^1(x_1^2)^0 \dots (x_1^d)^0 & \dots & (x_1^1)^p(x_1^2)^p \dots (x_1^d)^p \\ (x_2^1)^0(x_2^2)^0 \dots (x_2^d)^0 & \dots & (x_2^1)^1(x_2^2)^0 \dots (x_2^d)^0 & \dots & (x_2^1)^p(x_2^2)^p \dots (x_2^d)^p \\ \vdots & & \vdots & & \vdots \\ (x_n^1)^0(x_n^2)^0 \dots (x_n^d)^0 & \dots & (x_n^1)^1(x_n^2)^0 \dots (x_n^d)^0 & \dots & (x_n^1)^p(x_n^2)^p \dots (x_n^d)^p \end{pmatrix} \quad (18)$$

and we can find the optimal $\vec{\beta}$ in the same way as above.

2.3 Some statistics

When dealing with data sets, we always want to know how much the data points vary from each other, in other words we want to calculate the variance of the data. The variance from one data set is given by

$$\sigma_y^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2 \quad (19)$$

where the sample mean is given by

$$\bar{y} = \sum_{i=1}^N y_i. \quad (20)$$

We will stress that sample mean is not necessarily the same as the distribution mean, but the sample mean will approximate the distribution mean for large data sets, known as the *central limit theorem*. [StatMek]

Often we study multiple data sets at the same time, and want to estimate the total variance. The sample mean of one data set α is still given by

$$\bar{y}_\alpha = \frac{1}{N} \sum_{i=1}^N y_{\alpha,i}, \quad (21)$$

which for M data sets gives a total sample mean of

$$\bar{y}_M = \frac{1}{M} \sum_{\alpha=1}^M \bar{y}_\alpha = \frac{1}{MN} \sum_{\alpha=1}^M \sum_{i=1}^N y_{\alpha,i}. \quad (22)$$

The total variance will then be a sum over the single set variances plus the cross terms

$$\sigma_M^2 = \frac{1}{M} \sum_{\alpha=1}^M \frac{\sigma_\alpha^2}{N} + \frac{2}{MN^2} \sum_{\alpha=1}^M \sum_{i < j}^N (y_{\alpha,i} - \bar{y}_M)(y_{\alpha,j} - \bar{y}_M). \quad (23)$$

The last term is called the covariance, and is a measure of how much the data sets are related. If they are totally uncorrelated, the covariance is zero. For large data sets this term is computational expensive to calculate, and we will rather use resampling techniques to estimate the variance. A few resampling techniques are presented in the *Method* section.

2.3.1 Confidence interval

The confidence interval is the interval that you are within with a probability given by the confidence level. The confidence interval (CI) means

$$\text{CI} = \bar{y} \pm z^* \frac{\sigma}{\sqrt{N}} \quad (24)$$

where \bar{y} is the sample mean, σ is the standard deviation of the PDF, N is the number of points and z^* is directly related to the confidence level and can be found in statistics tables.

As an example, suppose that you flip a coin N times and count number of heads. Before you start, you want to know how many heads you will count with a 95% confidence. The procedure for calculating the CI is then

1. Calculate the sample mean of your data set, which for this case is $N/2$.
2. Find z^* from <https://www.sumproduct.com/thought/simulation-stimulation> a table. For 95% confidence, the Z-score $z^* = 1.96$
3. Calculate the standard deviation of the PDF. The clue in this example, is that the PDF approximates a binomial distribution when N increases, with a known σ . **[Coin]**
4. Plug into the formula above

So, how can we make us of this in project 1? To predict which coefficients we could get from the regression, it could be useful to calculate their CI. This is possible, since we can find their variance from the analytical formula

$$\text{Var}(\vec{\beta}) = (\hat{X}^T \hat{X})^{-1} \sigma^2 \quad (25)$$

where the σ^2 is the sample variance of the points.

Another approach is to simply calculate N estimates of each beta, and sort them in increasing order. Leaving out the $p\%$ first elements of the data set and the $p\%$ last elements, we are left with the $q = 100\% - 2p\%$ most likely betas, and the confidence interval related to a confidence level of $q\%$ is within the smallest and largest elements of this new set.

2.4 Error analysis

To find out how good the fit is, we could potentially compare the polynomial plot to the data plot and decide whether the fit is good or not. However, that approach is risky in the sense that it is hard to determine how good a fit really is, and we therefore rely on error estimates. There are many ways to calculate the error, where some frequently used methods are

- the absolute error
- the relative error
- the mean square error (MSE)
- the R^2 score function

In this report we will study only the MSE and R^2 score function.

2.4.1 Mean Square Error (MSE)

The most popular error estimation method is the mean square error, also called least squares. [Hastie] We study the MSE in order to compute the reduction of the cost function, because it is basically the standard cost function, used in OLS.

$$\text{MSE}(\vec{\beta}) = \frac{1}{N} \sum_{i=1}^N (y_i - \tilde{y}(\vec{\beta}))^2 \quad (26)$$

Compared to least absolute value, the points far away from the fitted line are weighted stronger. We will also study a quantity called bias, which is just the sum over differences between the model and all points,

$$\text{bias}(\vec{\beta}) = \frac{1}{N} \sum_{i=1}^N (y_i - \tilde{y}(\vec{\beta})). \quad (27)$$

As one can see, the bias can both be positive and negative, where a large positive bias indicates underfitting and a large negative bias indicates overfitting. [Berkley] It can be shown that the MSE is a sum of the bias squared and the variance,

$$\text{MSE}(\vec{\beta}) = \text{bias}(\vec{\beta})^2 + \sigma^2, \quad (28)$$

and based on this we can show that the optimal complexity of our model (which is neither underfitted nor overfitted) is when the bias and variance are equal, see figure (2). From the figure, we observe that the bias decreases and the variance increases as the model gets more complex. The MSE is minimized when $\sigma^2 = \text{bias}$.

2.4.2 R^2 score function

The R^2 score function is a measure of how close the data are to the fitted regression line, and is a widely used quantity within statistics. [R2]

It is defined by how much of the variation that is not explained by the model, i.e.,

$$R^2 = \frac{\text{Explained variation}}{\text{Total variation}}$$

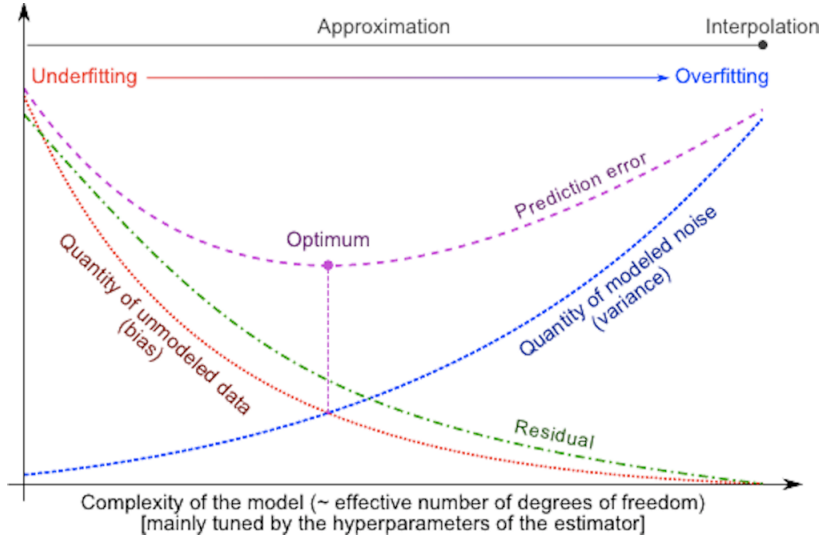


Figure 2: Bias-variance trade-off. Figure showing the bias (red graph), variance (blue graph) and MSE (purple graph) as functions of the complexity of the model. Figure taken from https://www.reddit.com/r/mlclass/comments/mmlfu/a_nice_alternative_explanation_of_bias_and/

where the explained variation is the difference between the variance and the MSE and the total variation is given by

$$\text{Total variation} = \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2.$$

The R^2 -score is usually between 0 and 1. If the model does not explain any of the variations it is 0 and if it explains all variations the score is 1. In that manner, we should fight for a high R^2 -score. It is in entirety given by

$$R^2(\vec{y}, \tilde{\vec{y}}) = 1 - \frac{\sum_{i=1}^N (y_i - \tilde{y})^2}{\sum_{i=1}^N (y_i - \bar{y})^2}. \quad (29)$$

3 Methods

3.1 Resampling techniques

A resampling technique is a way of estimating the variance of data sets without calculating the covariance. As we saw in section 2.4, the true covariance is given by a double sum which we will avoid calculating if possible. There are

different ways of doing this, and we have already went through several of them in this course:

- 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 only they will be covered here.

3.1.1 Bootstrap method

When we construct a data set, we usually draw samples from a probability density function (PDF) and get a set of samples \vec{x} . If we draw a large number of samples, the sample variance will approach the variance of the PDF. The bootstrap method turns this upside down, and tries to estimate the PDF given a set data set, because if we know the PDF, we know in principle everything about the data set.

The assumption we need to make, is that the relative frequency of x_i equals $p(x_i)$, which is reasonable (for instance, think about how the histogram looks like when we draw samples from a normal distribution). In this project the vector that we want to find, $\vec{\beta}(x, y)$, is a function of two set of variables. Fortunately, they are independent of eachother, so we can safely apply the independent bootstrap on them separately. The independent bootstraps goes as

1. Draw n samplings from the data set \vec{x} with replacement and denote the new data set as $\vec{x}^* = \{x_1, x_2, \dots, x_n\}$
2. Compute $\vec{\beta}(\vec{x}^*) \equiv \vec{\beta}^*$
3. Repeat the procedure above K times
4. The average value of all K $\vec{\beta}^*$'s are stored in a new vector $\vec{\bar{\beta}}^*$
5. Finally, the variance of $\vec{\bar{\beta}}^*$ should correspond to the sample variance

[Bootstrap]

The implementation could look something like this

```
def bootstrap(data, K=1000):
    dataVec = np.zeros(K)
    for k in range(K):
        dataVec[k] = np.average(np.random.choice(data, len(data)))
    Avg = np.average(dataVec)
    Var = np.var(dataVec)
    Std = np.std(dataVec)

    return Avg, Var, Std
```

which is strongly inspired by the lecture notes [Lecture].

It is also worth to mention that we should hold back a small part of the original data set for testing, since we need to test the model on data that it has not seen before.

3.1.2 K-fold validation method

K-fold validation is a method that has more describing name than many of its fellow methods. The idea is to make the most use of our data by splitting it into k folds and training k times on it. Every time we train, we leave out one of the folds, which gonna be our validation data. This validation data needs to be different every time, and we are therefore restricted to k unique training sessions. A typical overview looks like this:

- Split data set into k equally sized folds
- Use the $k - 1$ first folds as training data, and leave the k 'th fold for validation
- Use the $k - 2$ first folds plus the k -th fold as training data, and leave the $(k - 1)$ 'th fold for validation
- Continue until all folds are used as training data

[Berkeley]

3.2 Minimization methods

Suppose we have a very simple model trying to fit a straight line to data points. In that case, we could manually vary the coefficients and find a line that fits the points quite good. However, when the model gets more complicated, this can be a time consuming activity. Would it not be good 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. We will stick to gradient descent in this project, which might be the simplest method for our purposes.

3.2.1 Gradient Descent

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

$$\beta_i^{\text{new}} = \beta_i - \eta \cdot \frac{\partial Q(\beta_i)}{\partial \beta_i} \quad (30)$$

where β_i^{new} is the updated β and η is a step size, in machine learning often referred to as the learning rate. The idea is to find the gradient of the cost function $Q(\vec{\beta})$ with respect to a certain β_i , and move in the direction which minimizes the cost function. This is repeated until a minimum is found, defined by either

$$Q(\beta_i) < \varepsilon \quad (31)$$

or that the change in β_i for the past x steps is small.

Before we can implement equation (30), we need an expression for the derivative of Q with respect to β_i . The general form of the cost function, as discussed in section 2.1.4, reads

$$Q(\vec{\beta}, \lambda, q) = \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j|^q, \quad (32)$$

and its derivative with respect to β_k is

$$\frac{\partial Q(\vec{\beta}, \lambda, q)}{\partial \beta_k} = -2 \sum_i \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right) x_{ik} + \frac{\beta_k}{|\beta_k|} q \lambda |\beta_k|^{q-1}. \quad (33)$$

The vectorized version looks like

$$\frac{\partial Q(\vec{\beta}, \lambda, q)}{\partial \vec{\beta}} = -2 \hat{X}^T (\vec{y} - \hat{X} \vec{\beta}) + \frac{\vec{\beta}}{|\vec{\beta}|} q \lambda |\vec{\beta}|^{q-1} \quad (34)$$

where all operations are element wise. The algorithm of this minimization method is thus as follows:

```
while dbeta > epsilon:
    e = z - X.dot(beta)
    debeta = 2*X.T.dot(e) - np.sign(beta)*q*lambda*np.power(abs(beta), q-1)
    beta += eta*dbeta
```

3.3 Benchmark

To benchmark the results in this project, we use Scikit Learn, which is a package for machine learning in Python. [Scikit] The built-in regression tools are made to give high performance and be easy to implement, and is therefore the ideal benchmark in this project. The procedure of calling the OLS, Ridge and Lasso functions looks like this:

```
from sklearn.linear_model import LinearRegression, Ridge, Lasso

reg = LinearRegression()
reg.fit(X, z)
beta_OLS = reg.coef_

reg = Ridge(alpha=lambda)
reg.fit(X, z)
beta_Ridge = reg.coef_

reg = Lasso(alpha=lambda)
reg.fit(X, z)
beta_Lasso = reg.coef_
```

where X is the design matrix and z is the dependent variables. See *regression_scikit.py* for actual implementation.

4 Code

This project is largely based on numerical calculation, so a few words about the code is necessary. The code is written in Python, which is easy to work with and considered fast enough. The most expensive part is without doubt the minimization using gradient descent, and for really large systems a low-level language might be preferred. Although no performance benchmarks will be provided in this article, we can reveal that the code was more than fast enough for our data sets. To run the code, the following packages are required:

- NumPy - Fundamental package for scientific works in Python
- Matplotlib - For plotting
- Scikit Learn - Benchmark for self-produced code
- Scipy - For reading terrain data
- tqdm - For progression bar

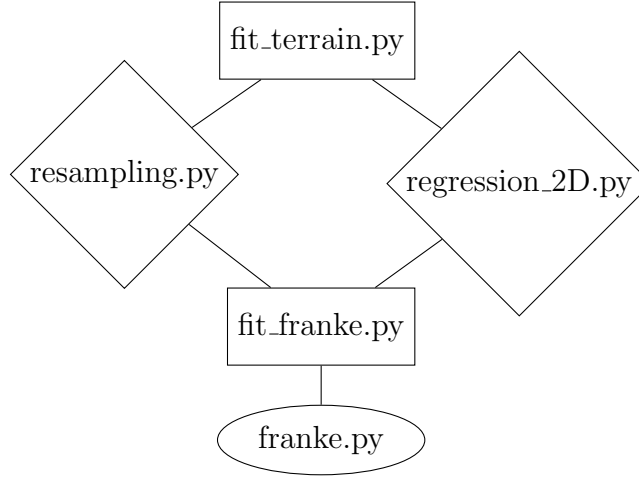


Figure 3: Code structure

4.1 Code structure

To keep the code neat and clean we decided to write object oriented code, although we do not have that many functions. The code of this project consists of two main functions `fit_franke.py` and `fit_terrain.py`, where the former is used to test the tools on a known function and the latter is used to fit a polynomial to terrain data. Both of them are calling the classes for 2D regression and resampling techniques, named `regression_2D.py` and `resampling.py` respectively, and `fit_franke.py` is also communicating with `franke.py`. For overview of the code structure, see figure (3).

4.2 Implementation and optimization

To get a code which provides good performance, we need to be thoughtful when implementing it. The loops are often the bottlenecks, so by replacing loops with vector operations we can save a lot of time. An example on this, is when we calculate the MSE

$$\text{MSE}(\vec{\beta}) = \sum_i \left(y_i - \beta_0 - \sum_j X_{ij} \beta_j \right)^2, \quad (35)$$

where \vec{y} and $\vec{\beta}$ are vectors and \hat{X} is a matrix. For implementing this directly, we need a double loop, which will be slow for large systems. A better solution would be to exploit the linear algebra properties of vectors:

$$\text{MSE}(\vec{\beta}) = (\vec{y} - \hat{X}^T \vec{\beta})^T \cdot (\vec{y} - \hat{X}^T \vec{\beta}) \quad (36)$$

which is usually much faster.

5 Results

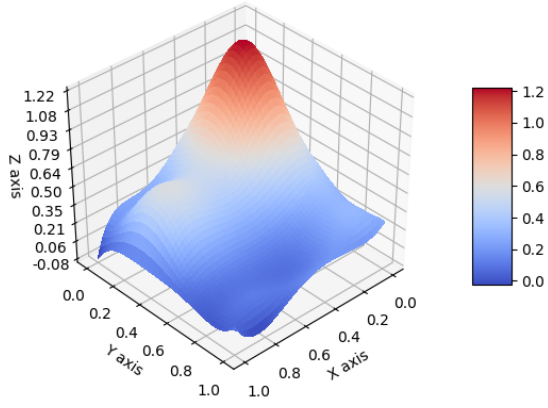
Here the results will be presented. We start looking at how good our linear regression is, by fitting a polynomial to points withdrawn from a known PDF. Thereafter we will see how good the regression is when using real terrain data. To give an intuition on how good the fit is, we will in both cases show 3D plots of the polynomials given by OLS, Ridge and Lasso methods, and then present a proper error analysis.

5.1 Franke function

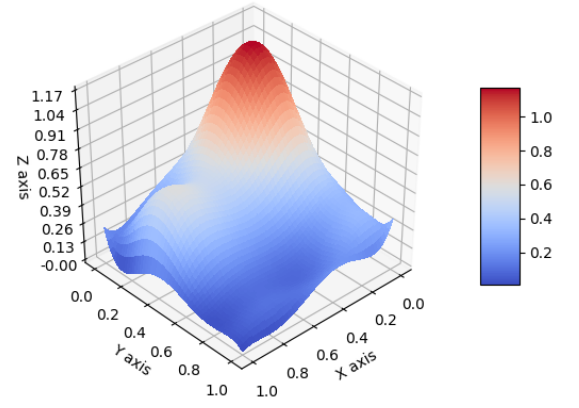
Below one can find the results when fitting the Franke function. All results in this section was obtained from the same data set, consisting of 1000 points, withdrawn from a uniform distribution ($x, y \in [0, 1]$). Further, the polynomial degrees in x- and y-directions were set to $P_x = P_y = 5$. For the methods that require minimization, we used gradient descent with the parameters $\lambda = 1e - 5$ (penalty), $\eta = 1e - 4$ (learning rate) and $niter = 1e5$ (number of iterations) when nothing else is specified. For benchmarking the results (mainly to verify the code), Scikit learn was used. In the first place, the studies were done without noise, and then everything was repeated with adding a noise of $\mathcal{N}(0, \sigma^2 = 0.01)$

5.1.1 Visualization of graphes

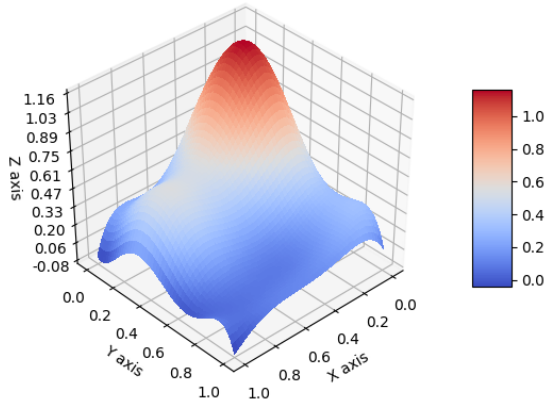
In figure (4) the polynomial gotten from regression on noisy data is compared with the polynomial gotten from regression on smooth data.



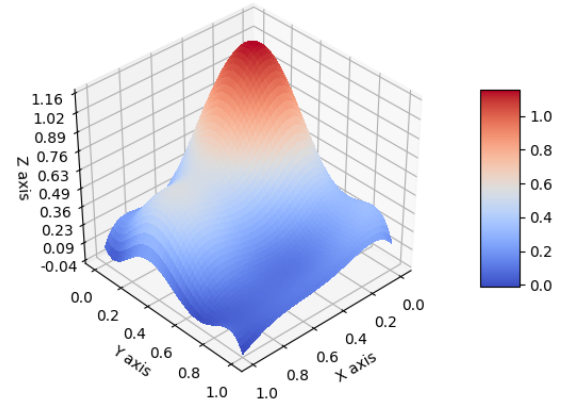
(a) OLS without noise



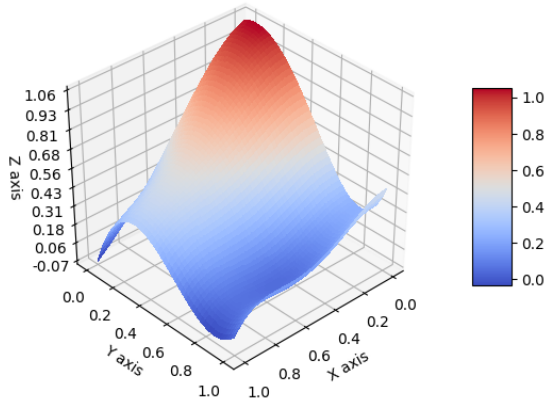
(b) OLS with noise



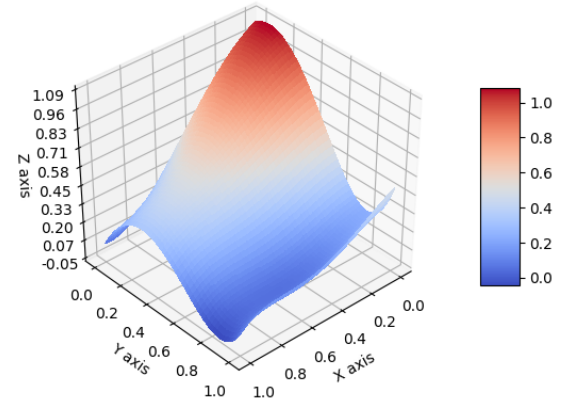
(c) Ridge without noise



(d) Ridge with noise



(e) Lasso without noise



(f) Lasso with noise

Figure 4: Fitted polynomial by OLS, Ridge and Lasso with and without noise. For Ridge and Lasso, we used a low penalty of $\lambda = 1e - 5$. Lasso was performed with $\eta = 1e - 4$ and $1e5$ iterations. The noise was sampled from a normal distribution with $\sigma = 0.1$.

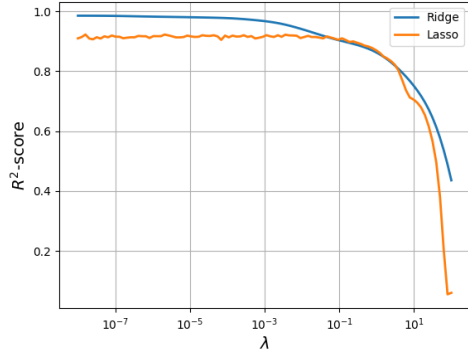
5.1.2 Error

To be able to determine which method is best, we need to analyze the error. In table (1), the MSE and R^2 -score are given for OLS, Ridge, Lasso and Ridge with gradient descent (RidgeGD).

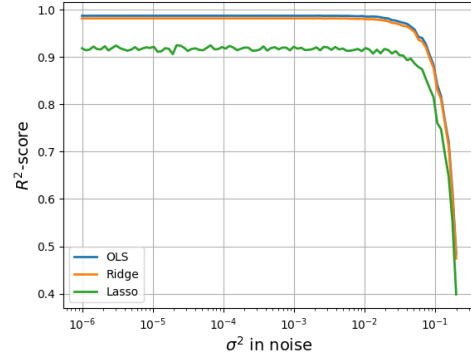
Table 1: Mean Square Error and R^2 -score presented for OLS, Ridge, Lasso and Ridge + gradient descent (RidgeGD), where noise was added to the data. The parameters used were $\lambda = 1e-5$ (penalty), $\eta = 1e-4$ (learning rate), $niter = 1e5$ (number of iterations) and $\mathcal{N}(0, \sigma^2 = 0.01)$ (noise).

	MSE		R2	
	self	scikit	self	scikit
OLS	0.001215	0.001215	0.9854	0.9853
Ridge	0.001628	0.001628	0.9804	0.9804
Lasso	0.009476	0.006651	0.8861	0.9201
RidgeGD	0.009067	0.001628	0.8911	0.9804

We have also studied how the R^2 score is dependent on various parameters. In figure (5), the R^2 -score is plotted as a function of the penalty for Ridge and Lasso (a), and as a function of the noise for OLS, Ridge and Lasso (b).



(a) Lambda vs. R2-score



(b) Variance vs. R2-score

Figure 5: R^2 -score plotted as a function of the penalty λ (a) and as a function of the noise (b). The penalty was varied in the interval $\lambda \in [10^{-8}, 10^2]$ and the variance of the normal distribution the noise was withdrawn from was varied in the interval $\sigma^2 \in [10^{-6}, 10^{-0.7}]$. The other parameters used were $\lambda = 1e - 5$ (penalty, was held constant for (b) only), $\eta = 1e - 4$ (learning rate), $niter = 1e5$ (number of iterations) and $\mathcal{N}(0, \sigma^2 = 0.01)$ (noise, was held constant for (a) only).

5.1.3 Coefficients β

It can also be interesting to see how the coefficients actually differ between the self-built functions and the Scikit Learn function. In figure (6), the beta values are visualized such that large numbers got strong color, negative numbers are blue and positive numbers are red. The methods OLS, Ridge, Lasso and Ridge produced using gradient descent (RidgeGD) are presented.

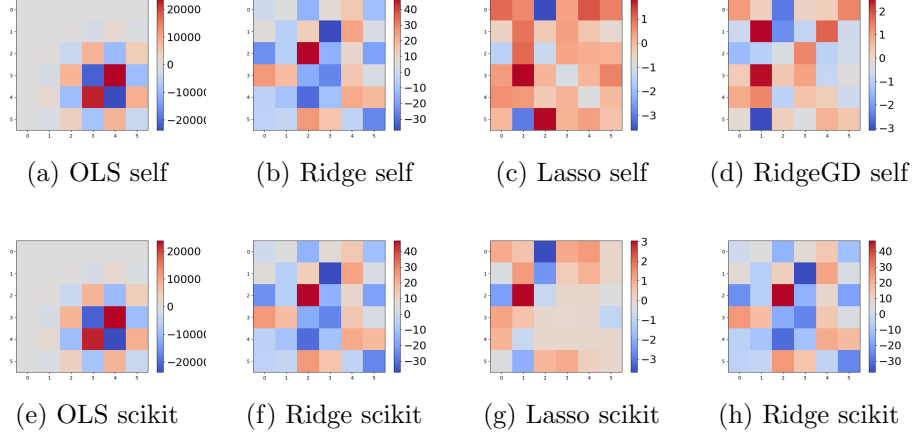


Figure 6: Beta values visualized for different methods. The upper images are gotten from the self-built functions, where the last one (Ridge2) is Ridge found by minimizing beta using gradient descent. The lower images are the benchmarks which are produced by Scikit learn. The parameters used were $\lambda = 1e-5$ (penalty), $\eta = 1e-4$ (learning rate), $niter = 1e5$ (number of iterations) and $\mathcal{N}(0, \sigma^2 = 0.01)$ (noise).

Furthermore, the confidence intervals are given in table (2) for all betas and all methods, when noise is included.

Table 2: Confidence intervals. The parameters used were $\lambda = 1e - 5$ (penalty), $\eta = 1e - 4$ (learning rate), $niter = 1e5$ (number of iterations) and $\mathcal{N}(0, \sigma^2 = 0.01)$ (noise).

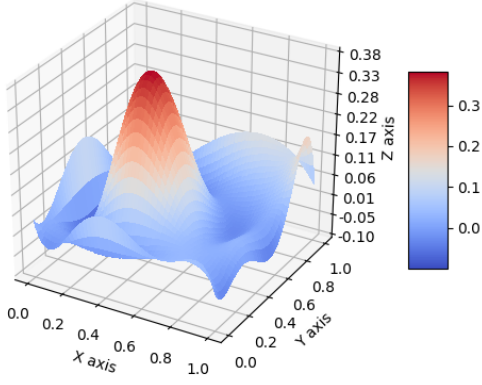
	OLS	Ridge	Lasso	RidgeGD
β_1	0.07853	0.86104	0.05386	1.7483
β_2	0.07853	0.86104	0.05386	1.7483
β_3	0.07853	0.86104	0.05386	1.7483
β_4	0.07853	0.86104	0.05386	1.7483
β_5	0.07853	0.86104	0.05386	1.7483
β_6	0.07853	0.86104	0.05386	1.7483
β_7	0.07853	0.86104	0.05386	1.7483
β_8	0.07853	0.86104	0.05386	1.7483
β_9	0.07853	0.86104	0.05386	1.7483
β_{10}	0.07853	0.86104	0.05386	1.7483
β_{11}	0.07853	0.86104	0.05386	1.7483
β_{12}	0.07853	0.86104	0.05386	1.7483
β_{13}	0.07853	0.86104	0.05386	1.7483
β_{14}	0.07853	0.86104	0.05386	1.7483
β_{15}	0.07853	0.86104	0.05386	1.7483
β_{16}	0.07853	0.86104	0.05386	1.7483
β_{17}	0.07853	0.86104	0.05386	1.7483
β_{18}	0.07853	0.86104	0.05386	1.7483
β_{19}	0.07853	0.86104	0.05386	1.7483
β_{20}	0.07853	0.86104	0.05386	1.7483
β_{21}	0.07853	0.86104	0.05386	1.7483
β_{22}	0.07853	0.86104	0.05386	1.7483
β_{23}	0.07853	0.86104	0.05386	1.7483
β_{24}	0.07853	0.86104	0.05386	1.7483
β_{25}	0.07853	0.86104	0.05386	1.7483
β_{26}	0.07853	0.86104	0.05386	1.7483
β_{27}	0.07853	0.86104	0.05386	1.7483
β_{28}	0.07853	0.86104	0.05386	1.7483
β_{29}	0.07853	0.86104	0.05386	1.7483
β_{30}	0.07853	0.86104	0.05386	1.7483
β_{31}	0.07853	0.86104	0.05386	1.7483
β_{32}	0.07853	0.86104	0.05386	1.7483
β_{33}	0.07853	0.86104	0.05386	1.7483
β_{34}	0.07853	0.86104	0.05386	1.7483
β_{35}	0.07853	0.86104	0.05386	1.7483
β_{36}	0.07853	0.86104	0.05386	1.7483

5.2 Real data

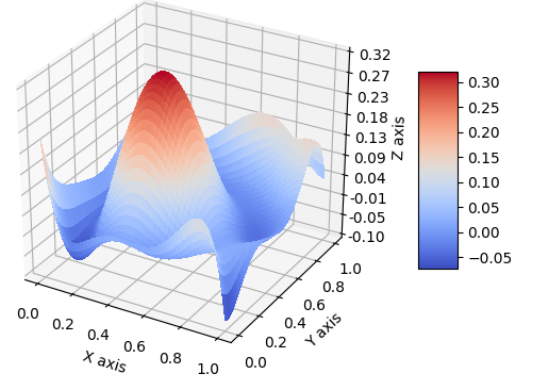
Below one can find the results when fitting the terrain data with a polynomial of a maximum degree of 5 in x- and y-directions. For the methods that require minimization, we used gradient descent with the parameters $\lambda = 1e - 5$ (penalty), $\eta = 1e - 4$ (learning rate) and $niter = 1e5$ (number of iterations) when nothing else is specified. For benchmarking the results (mainly to verify the code), Scikit learn was used. In the first place, the studies were done without noise, and then everything was repeated with adding a noise of $\mathcal{N}(0, \sigma^2 = 0.01)$

5.2.1 Visualization of graphes

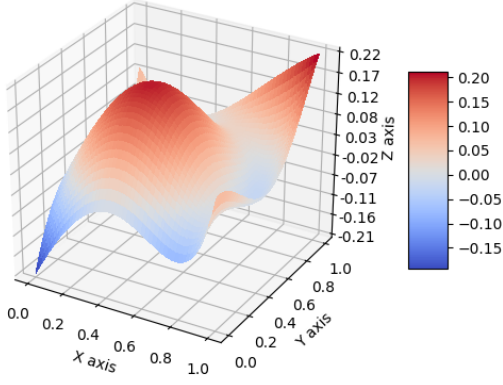
In figure (7) the polynomial gotten from regression on noisy data is compared with the polynomial gotten from regression on smooth data. In figure (7), we can see how good we were able to fit the polynomials using various methods



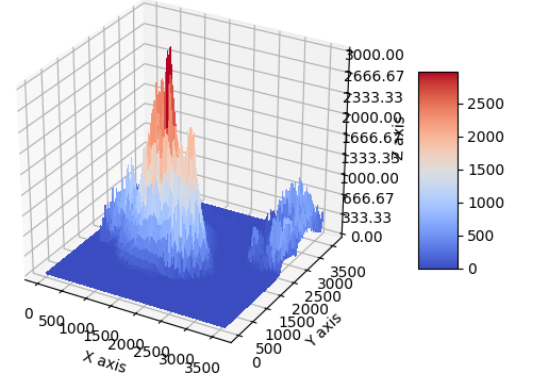
(a) OLS



(b) Ridge



(c) Lasso



(d) Terrain data

Figure 7: Fitted polynomial by OLS, Ridge and Lasso. For Ridge and Lasso, we used a low penalty of $\lambda = 1e - 5$. Lasso was performed with $\eta = 1e - 4$ and 1e5 iterations. The terrain data used in regression was normalized, but the data used in the terrain plot was not.

5.2.2 Error

To be able to determine which method is best, we need to analyze the error. In table (3), the MSE and R^2 -score are given for OLS, Ridge, Lasso and Ridge with gradient descent (RidgeGD).

Table 3: Mean Square Error and R^2 -score presented for OLS, Ridge, Lasso and Ridge + gradient descent (RidgeGD), where noise was added to the data. The parameters used were $\lambda = 1e-5$ (penalty), $\eta = 1e-4$ (learning rate), $niter = 1e5$ (number of iterations) and $\mathcal{N}(0, \sigma^2 = 0.01)$ (noise).

	MSE		R2	
	self	scikit	self	scikit
OLS	0.006282	0.006282	0.9267	0.9267
Ridge	0.007886	0.007886	0.9080	0.9080
Lasso	0.01294	0.01340	0.8491	0.8437
RidgeGD	0.01299	0.007885	0.8485	0.9080

We have also studied how the R^2 score is dependent on various parameters. In figure (8), the R^2 -score is plotted as a function of the penalty for Ridge and Lasso (a), and as a function of the noise for OLS, Ridge and Lasso (b).

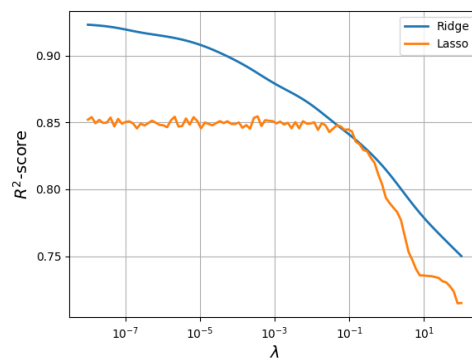


Figure 8: R^2 -score plotted as a function of the penalty λ . The penalty was varied in the interval $\lambda \in [10^{-8}, 10^2]$, and the other parameters used were $\eta = 1e-4$ (learning rate) and $niter = 1e5$ (number of iterations).

5.2.3 Coefficients β

It can also be interesting to see how the coefficients actually differ between the self-built functions and the Scikit Learn function. In figure (6), the beta values are visualized such that large numbers got strong color, negative numbers are blue and positive numbers are red. The methods OLS, Ridge, Lasso and Ridge produced using gradient descent (RidgeGD) are presented.

Furthermore, the confidence intervals are given in table (4) for all betas and all methods, when noise is included.

Table 4: Confidence intervals. The parameters used were $\lambda = 1e - 5$ (penalty), $\eta = 1e - 4$ (learning rate), $niter = 1e5$ (number of iterations) and $\mathcal{N}(0, \sigma^2 = 0.01)$ (noise).

	OLS	Ridge	Lasso	RidgeGD
β_1	0.07853	0.86104	0.05386	1.7483
β_2	0.07853	0.86104	0.05386	1.7483
β_3	0.07853	0.86104	0.05386	1.7483
β_4	0.07853	0.86104	0.05386	1.7483
β_5	0.07853	0.86104	0.05386	1.7483
β_6	0.07853	0.86104	0.05386	1.7483
β_7	0.07853	0.86104	0.05386	1.7483
β_8	0.07853	0.86104	0.05386	1.7483
β_9	0.07853	0.86104	0.05386	1.7483
β_{10}	0.07853	0.86104	0.05386	1.7483
β_{11}	0.07853	0.86104	0.05386	1.7483
β_{12}	0.07853	0.86104	0.05386	1.7483
β_{13}	0.07853	0.86104	0.05386	1.7483
β_{14}	0.07853	0.86104	0.05386	1.7483
β_{15}	0.07853	0.86104	0.05386	1.7483
β_{16}	0.07853	0.86104	0.05386	1.7483
β_{17}	0.07853	0.86104	0.05386	1.7483
β_{18}	0.07853	0.86104	0.05386	1.7483
β_{19}	0.07853	0.86104	0.05386	1.7483
β_{20}	0.07853	0.86104	0.05386	1.7483
β_{21}	0.07853	0.86104	0.05386	1.7483
β_{22}	0.07853	0.86104	0.05386	1.7483
β_{23}	0.07853	0.86104	0.05386	1.7483
β_{24}	0.07853	0.86104	0.05386	1.7483
β_{25}	0.07853	0.86104	0.05386	1.7483
β_{26}	0.07853	0.86104	0.05386	1.7483
β_{27}	0.07853	0.86104	0.05386	1.7483
β_{28}	0.07853	0.86104	0.05386	1.7483
β_{29}	0.07853	0.86104	0.05386	1.7483
β_{30}	0.07853	0.86104	0.05386	1.7483
β_{31}	0.07853	0.86104	0.05386	1.7483
β_{32}	0.07853	0.86104	0.05386	1.7483
β_{33}	0.07853	0.86104	0.05386	1.7483
β_{34}	0.07853	0.86104	0.05386	1.7483
β_{35}	0.07853	0.86104	0.05386	1.7483
β_{36}	0.07853	0.86104	0.05386	1.7483

6 Discussion

The first thing we can see from the results, is that the fitted functions in figure (4) look very similar with and without noise. This is probably caused by too small noise. However, we observe that both the OLS regression and Ridge are really good fits compared to the actual function in figure (1). Lasso is also similar, but with the naked eye we can see that it is obviously not as good as OLS and Ridge.

This is also reflected in the error analysis, where the MSE is way larger for Lasso compared to OLS and Ridge. OLS seems to have a slightly lower MSE than Ridge again, but for Ridge with gradient descent the MSE is almost as high as for Lasso. From this, it may seem like the gradient descent method does not work properly, because the cost function is obviously not fully minimized. When it comes to the R^2 -score, we observe that also here OLS got the best score barely in front of Ridge, with RidgeGD and Lasso far behind.

Furthermore, from the R^2 vs. λ plot in figure (5), we can see that Ridge mostly got a better R^2 -score than Lasso, only beaten around $\lambda = 1$. What surprises me, is that the methods do not get more similar when the penalty decreases, since they should both approach OLS when the penalty gets smaller. In the second plot in the same figure, we again find Lasso to give the lowest R^2 -score.

In figure (6), we can see that the beta values produced by Scikit's OLS function and the self-built OLS-function are more or less identical. The same applies to Ridge, so we can conclude that those two methods are implemented correctly (also based on the error analysis). The situation is unlike for the gradient descent methods, where the coefficients differ. Especially for RidgeGD, the difference is significantly both in sign and magnitude. For Lasso, we can reveal some similarities, and at least the magnitudes are similar. Again it is tempting to conclude that the minimization is wrong. Something I find strange, is that every second coefficient is negative and every second is positive for OLS, and obtain a cool pattern which I cannot explain.

We have earlier stated that we expect the coefficients in Ridge regression to be smaller than the coefficients in OLS regression, and the coefficients in Lasso regression should be even smaller than Ridge again. This is exactly what we see in figure (6).

Now over to the real terrain data. In figure (7), we can see that OLS reproduces the sharp peak of the volcano best among the methods, followed by Ridge and finally Lasso. This is expected from regression on the Franke function. Another thing we observe, is that none of the models is able to reproduce the flat ocean around the volcano. This is because we have limited the models to 5 degrees in

each direction, and that is not enough degrees to fit those areas.

Again we observe that the self-built OLS function and Scikit learn's OLS function give identical MSE and R^2 -score, and exactly the same applies for Ridge. When it comes to Lasso, the self-built function now gives a lower MSE and higher R^2 -score than the function of Scikit learn, which is a bit surprising. Ridge produced from gradient descent is still weak compared to all the other methods.

When it comes to how the R^2 -score is varying as a function of the penalty, we can see that also here the R^2 -score from Lasso crosses the R^2 -score from Ridge. This happens around $\lambda = 0.1$, and that is the only point where Lasso can be considered as better than Ridge. For small penalties, Ridge is way better than Lasso. In the second plot in figure (8), we can see that Lasso gives the definitely lowest R^2 -score among the methods.

7 Conclusion

After what we have seen, we can conclude that OLS gives the best results for our particular data set, but Ridge is not a bad choice either. Lasso gives both higher mean square error (MSE) and lower R^2 -score than the other methods, and since the implemented Lasso gives similar results to the Lasso function of Scikit Learn, the Lasso is considered as the least good method among the methods that we have tested. We have also seen that the lower penalty the higher R^2 -score, which is expected since OLS is the best method.

As an afterthought, the choice of terrain data was maybe too ambitious, and I do not find the results sufficient. To obtain a better fit, the complexity could be increased by increasing the polynomial degree, or we could turn to other regression methods such as logistic regression. Apropos the polynomial degree, I misunderstood what polynomial of degree 5 was, and used maximum degree of 5 in each spatial directions (should of course have fixed this). Despite this misunderstanding, the mathematical operations I have used are still as valid as for a polynomial of maximum total degree of 5.

8 References

- [1] A.M.Legendre. Nouvelles méthodes pour la détermination des orbites des comètes. Libraire pour les Mathématiques. (1805).
- [2] C.F. Gauss. Theoria Motus Corporum Coelestium in Sectionibus Conicis Solem Ambientum. Hamburg: Frid. Perthes and I.H. Besser. (1809).
- [3] T. Hastie, R. Tibshirani, J. Friedman. The Elements of Statistical Learning. Springer-Verlag, New York. (2009).
- [4] United States Geological Survey (USGS). <https://earthexplorer.usgs.gov/>
- [5] Lecture notes in Statistical Physics: Statistical Physics - a second course. F. Ravndal, E. G. Flekkøy. (2014).
- [6] Simulation stimulation. L. Bastick. <https://www.sumproduct.com/thought/simulation-stimulation>
- [7] Machine Learning Crash Course: Part 4 - The Bias-Variance Dilemma. D. Geng, S. Shih. <https://ml.berkeley.edu/blog/2017/07/13/tutorial-4/> (2017).
- [8] Regression Analysis: How Do I Interpret R-squared and Assess the Goodness-of-Fit? <http://blog.minitab.com/blog/adventures-in-statistics-2/regression-analysis-how-do-i-interpret-r-squared-and-assess-the-goodness-of-fit> (2013).
- [9] Bootstrap Methods: Another Look at the Jackknife. B. Efron. Ann. Statist., Volume 7, Number 1, 1-26. (1979).
- [10] Lecture notes FYS-STK4155: Regression Methods. <https://compphysics.github.io/MachineLearning/doc/pub/Regression/pdf/Regression-minted.pdf> M. Hjorth-Jensen. (2018).
- [11] Scikit-learn: Machine Learning in Python. F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thiron, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, E. Duchesnay. Journal of Machine Learning Research, Volume 12. (2011).