

FYS-STK4155 Project 1: Regression

Jan Ole Åkerholm
(Dated: October 7, 2019)

Abstract
Hi

Link to github repository: [Click here](#)

I. INTRODUCTION

Regression is widely used in science as a means to build models based on observed data. It allows for fitting a continuous function to discrete data sets, which can be used to predict and explain various phenomena in a range of sciences. The methods also have well studied tools for validating their accuracy.

This project explores three various methods for linear regression: ordinary least squares, Ridge and Lasso regression. First these methods will be applied on the Franke function, a continuous function which has some similarities with terrain. The methods will then be applied on real terrain data from Møsvatn Austfjell in Norway. For each case, the error in the models will be studied in some detail.

II. THEORY

A. Franke function

The Franke function is a two-dimensional weighted sum of exponentials which visibly shows some similarities to terrain, it is defined as:

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

A plot of the function is shown in FIG. 4

B. Ordinary least squares (OLS)

OLS and the other methods discussed here are also explained in detail by Hastie et. al. [1]

Suppose a data set of some observed phenomenon is given as:

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

where y_i is the observed data for some predictor \mathbf{x}_i and ϵ_i is random noise with mean value 0 and a variance of σ^2 . The ordinary least squares model is then used to model the underlying function $\mathbf{f}(\mathbf{x})$ as:

$$\tilde{\mathbf{y}} = \mathbf{X}\beta$$

where \mathbf{X} is the so-called "design matrix", and β are coefficients specific to the data. The design matrix can then be designed in some way that seems sensible for the problem at hand and depending on the predictors, and the model can then be created by finding the optimal parameters of β for some data set.

Note that the predictor \mathbf{x}_i may be multi-dimensional, for example in the case where the data point y_i is some terrain data (e.g. height) dependant on two coordinates.

Finding the optimal parameters of β is typically done by minimizing some cost function $C(\beta)$. In OLS, the *mean squared error* is typically used as the cost function, and is defined as:

$$C(\beta) = MSE(\beta) = \frac{1}{N} \sum_{i=0}^{N-1} (y_i - \tilde{y}_i)^2 \quad (2)$$

where N is the number of data points. This can be written in matrix notation as

$$C(\beta) = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \quad (3)$$

The cost function is minimized by taking its derivative and setting it equal to 0:

$$\frac{\partial C(\beta)}{\partial \beta} = 0 \quad (4)$$

$$\mathbf{X}^T (\mathbf{y} - \mathbf{X}\beta) = 0 \quad (5)$$

Solving for β gives the normal equation for OLS:

$$\beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (6)$$

The prediction can then be found by

$$\tilde{\mathbf{y}} = \mathbf{X}\beta$$

1. Confidence intervals of β

Having calculated an optimal set of β -values, the confidence interval of β_i can then be found from [1]:

$$(\beta_i - z^{1-\alpha} \sqrt{v_i} \hat{\sigma}, \beta_i + z^{1-\alpha} \sqrt{v_i} \hat{\sigma}) \quad (7)$$

where v_i is the i -th diagonal element of $(\mathbf{X}^T \mathbf{X})^{-1}$ and $\hat{\sigma}^2$ is the variance of $\tilde{\mathbf{y}}$. $z^{1-\alpha}$ depends on how wide the confidence interval is, and for a 95% confidence interval $z^{1-\alpha} = 1.96$.

2. Error measurement

In addition to using the MSE (equation 2) to measure the quality of the model, the R^2 score is also used. It is defined as:

$$R^2(\mathbf{y}, \tilde{\mathbf{y}}) = 1 - \frac{\sum_{i=0}^{N-1} (y_i - \tilde{y}_i)^2}{\sum_{i=0}^{N-1} (y_i - \bar{y})^2}$$

where \bar{y} is the mean value of \mathbf{y} . For a perfect fit, the R^2 score function gives a value of 1.

C. Resampling

When developing and testing the model the data is split into training and test data, so that the model can be tested on data which it has not been trained on. This is done to avoid training the model in such a way that it might fit the training data very well, but be completely useless as soon as it is used for data which is not part of the training (overfitting).

1. K-fold cross validation

In order to test the model on as much test data as possible, a method called K-fold cross validation is used. The method splits the data set randomly into k sections (folds), and trains on $k - 1$ of these sets while testing on the one which has not been part of the training. This is repeated k times and each time provides one estimate for the accuracy of the model and the β parameters. Usual values for k is 5 or 10, and I have chosen to stick with 5 for this project.

D. Ridge regression

Ridge regression is a method which adds *regularization* to the regression model. Regularization is a tool used to prevent the regression method failing if the data points are highly correlated, or to prevent overfitting of the data.

In ridge regression, the values of β are constrained so that $\|\beta\|_2^2 \leq t$ for some finite number $t > 0$.

In practice the method is applied by changing the cost function so that the new problem to be minimized is:

$$C(\beta) = \frac{1}{N} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2 \quad (8)$$

where λ is some regularization parameter where multiple values should be tested in order to find the optimal result.

This leads to a new equation for the optimal values β^{Ridge} :

$$\beta^{\text{Ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (9)$$

Adding a small value to the diagonal of $\mathbf{X}^T \mathbf{X}$ avoids problems if $\mathbf{X}^T \mathbf{X}$ is singular, and also decreases the condition number of the matrix. If $\lambda = 0$ the method is identical to OLS as expected.

E. Lasso regression

Lasso regression is similar to Ridge regression, but instead of using the L_2 -norm for the λ term, the L_1 -norm is used instead, so that the cost function now is:

$$C(\beta) = \frac{1}{N} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1 \quad (10)$$

where the 1-norm is defined as:

$$\|\beta\|_1 = \sum_i |\beta_i|$$

There is no analytical solution for the optimal β values in this case, so when using the Lasso method an iterative method is used instead.

F. Bias-variance tradeoff

In order to gain more insight into the error of the model, the MSE can be decomposed into the bias, variance and irreducible error terms. Recall that $\mathbf{y} = \mathbf{f} + \epsilon$, and

$$MSE(\mathbf{y}, \tilde{\mathbf{y}}) = \frac{1}{N} \sum_{i=0}^{N-1} (y_i - \tilde{y}_i)^2 = \mathbb{E}[(\mathbf{y} - \tilde{\mathbf{y}})^2]$$

Inserting $\mathbf{y} = \mathbf{f} + \epsilon$, and adding and subtracting $\mathbb{E}[\tilde{\mathbf{y}}]$:

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

Where δ are terms on the form $\mathbb{E}[\epsilon]$ or $\mathbb{E}[\mathbb{E}[\tilde{\mathbf{y}}] - \tilde{y}]$ which tend to 0 for large enough samples. Using the definitions of the expectation values then we get:

$$\mathbb{E}[(\mathbf{y} - \tilde{\mathbf{y}})^2] = \frac{1}{N} \sum_{i=0}^{N-1} (f_i - \mathbb{E}[\tilde{y}])^2 + \frac{1}{N} \sum_{i=0}^{N-1} (\tilde{y}_i - \mathbb{E}[\tilde{y}])^2 + \sigma^2 \quad (11)$$

The first term is the bias squared, the second term is the variance of $\tilde{\mathbf{y}}$. σ^2 is the variance of ϵ , i.e. the noise of the measurement. It is also known as the irreducible error.

In general, one expects the bias to be higher when the complexity of the model is low, and the variance to gradually increase along with the complexity of the model. This gives a tradeoff between the bias and variance, where the lowest MSE is given by some complexity where the sum of the bias and the variance is low.

III. METHOD/IMPLEMENTATION

First the theory discussed will be applied to model a Franke function with noise. It will then be applied on real terrain data from Møsvatn Austfjell in Norway. For the rest of this report \mathbf{x} , and \mathbf{y} will be referring to the predictors (coordinates of the terrain), and \mathbf{z} will refer to the height of the terrain.

For both the Franke function and the terrain data, the x and y axes are scaled to $[0, 1]$. Additionally, for the terrain data the height of the terrain is normalized so that the highest point = 1.

A. Design matrix

For both the Franke function and the terrain data, the design matrix is set up with each row containing a polynomial on the form:

$$[x, y, x^2, xy, y^2, x^3, x^2y, xy^2, y^3, \dots]$$

This means a polynomial of degree m has

$$m \frac{m+1}{2}$$

terms in total.

B. Code structure

The core of the code is contained in the folder "resources". The regression methods are contained in the

file "regression.py" which contains a main class called "OLS()". This class contains methods for calculating the β parameters for a given design matrix and set of data. When an instance is initialized the fit is immediately calculated, and a prediction can be made.

There are also subclasses of "OLS()" for Ridge and Lasso regression, where the fitting methods are modified to use the Ridge and Lasso methods.

The "resources.py" file contains functions for plotting a 2D surface, setting up the design matrix, finding the confidence intervals for the β fit, finding the MSE, R^2 score, variance and bias. The "franke.py" file contains code for generating the Franke function.

C. Franke function

Before starting on the real terrain data, I first started by using the methods on the Franke function. The code for this can be found in "franke/ridge_lasso_fitting.py". The Franke function is generated and gaussian noise with $\mu = 0, \sigma = 0.1$ is generated and added to the data. The K-fold split is then set up and the code runs for all three models (OLS, Ridge and Lasso), for a given set of complexities m and lambda values λ (for OLS the loop ends after one λ value because the method does not depend on λ).

The file "testing/franke_fit.py" contains code for running the OLS regression on the Franke function and plotting the various fits. It also compares the results with Scikit-Learn's "LinearRegression" functionality, and calculates the confidence intervals of β .

D. Terrain data

The terrain data is gathered from the Github project repository. I chose to use the data from Møsvatn Austfjell, as it seemed fairly interesting, but not too noisy and complex. The code in "terrain_fitting/terrain_testing.py" is structured in much the same way as the code for the Franke function fitting, fitting the data with various values for m and λ and plotting the results.

Additionally the file "terrain_fitting/final_fit.py" uses OLS only to fit the data with various complexities and plotting the resulting model. In this case I chose not to use K-fold because I suspect there might be something wrong with my cross validation code, which will be explained further in the results section.

IV. RESULTS AND DISCUSSION

A. Franke function

1. Ordinary Least Squares

The result of the training of the OLS model on the Franke function with noise, up to a complexity of $m = 15$ can be seen in FIG. 1. The results are somewhat surprising, as the bias and variance are both above the MSE. As the MSE is expected to be a sum of the bias,

variance and the irreducible error, this is not realistic, and I suspect it is due to some error in my code. I got the same result both with and without the cross validation, so the error must be somewhere else, but I have not had time to find it.

A clue might be that the MSE and variance appear to almost exactly balance each other out, and the bias appears to be almost completely constant. Another clue is that the MSE drops to 0.01, or $0.1^2 = \sigma^2$, indicating that it is displaying the variance of the noise that was added to the data set.

OLS Franke function error

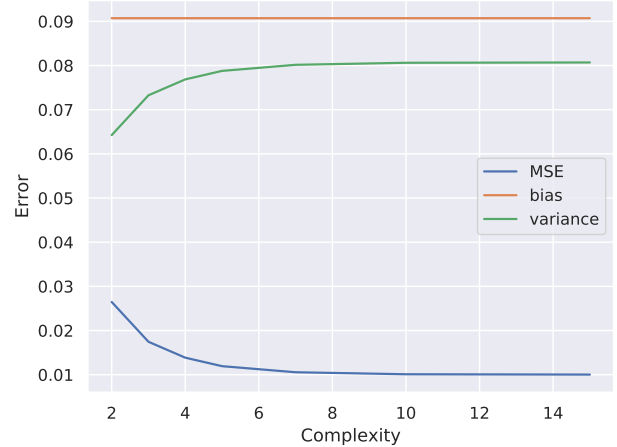


FIG. 1: MSE, bias and variance of the OLS model on the Franke function for complexities ranging from 2 to 15

Scikit Learn comparison

When comparing with the OLS method in Scikit Learn, the results for the fit and R2 score are the same. For $m = 2$ complexity:

	β_0	β_1	β_2	β_3	β_4	β_5
My method	1.17	-1.031	-0.794	0.100	0.857	-0.301
Scikit	1.17	-1.031	-0.794	0.100	0.857	-0.301

β with confidence intervals

The β values and their confidence intervals were calculated as follows (for $m = 2$):

$$\begin{aligned}
 \beta_0 &= 1.17141889 \pm 0.00255787 \\
 \beta_1 &= -1.03123099 \pm 0.00744324 \\
 \beta_2 &= -0.79432508 \pm 0.00752144 \\
 \beta_3 &= 0.10020426 \pm 0.00660813 \\
 \beta_4 &= 0.85686316 \pm 0.005875 \\
 \beta_5 &= -0.30056054 \pm 0.00661286
 \end{aligned}$$

2. Ridge regression

When making a models using Ridge regression, multiple values of λ were used, the full result for the MSE can be seen in the heatmap in FIG. 2. A similar heatmap for the R^2 score can be found in the github repository under "franke_fitting/fig".

From the heatmap, it appears that the lower the λ -value, the better the MSE. The R^2 heatmap shows the

same. Because the MSE and R^2 scores are calculated using cross validation, this appears to imply that the OLS is a better model than Ridge in this case, and that I have not experienced any overfitting yet up to a complexity of $m = 15$, but keeping in mind that there is some error in the code, making any conclusion on this is probably impossible with the current results. The Ridge regression

V. CONCLUSION

Hello

REFERENCES

[1] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition (Springer Series in Statistics)*. Springer, 2016. ISBN: 0387848576.

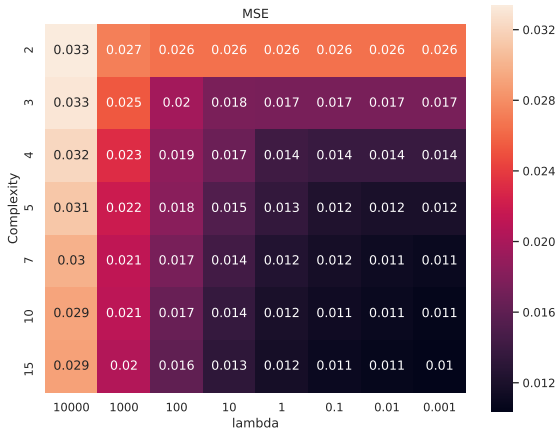


FIG. 2: MSE values for Ridge regression on the Franke function with noise. The lowest value is for $m = 15$ and $\lambda = 0.001$. Full size plots can be found in the Github repository.

shows the same problem with the bias-variance tradeoff as the OLS, as can be seen from FIG. 3

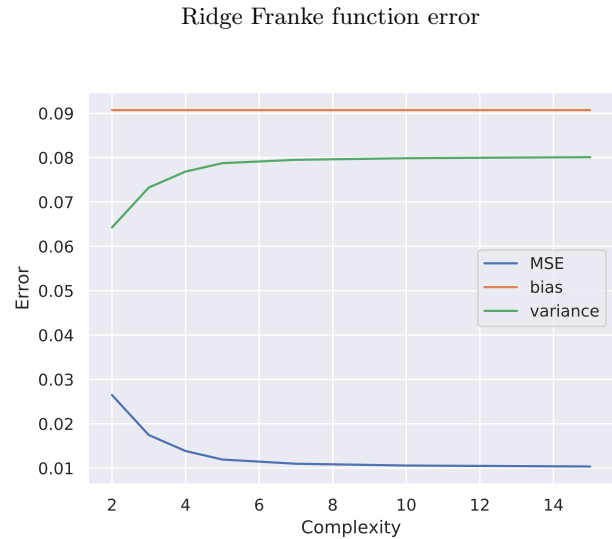


FIG. 3: MSE, bias and variance of the Ridge model on the Franke function for complexities ranging from 2 to 15

Scikit Learn comparison

When comparing with the Ridge method in Scikit Learn, the results for the fit and R2 score are the same, like in the case for OLS. For $m = 2$ complexity and $\lambda = 10$:

	β_0	β_1	β_2	β_3	β_4	β_5
My method	1.17	-1.027	-0.791	0.098	0.854	-0.302
Scikit	1.17	-1.027	-0.791	0.098	0.854	-0.302

VI. APPENDIX

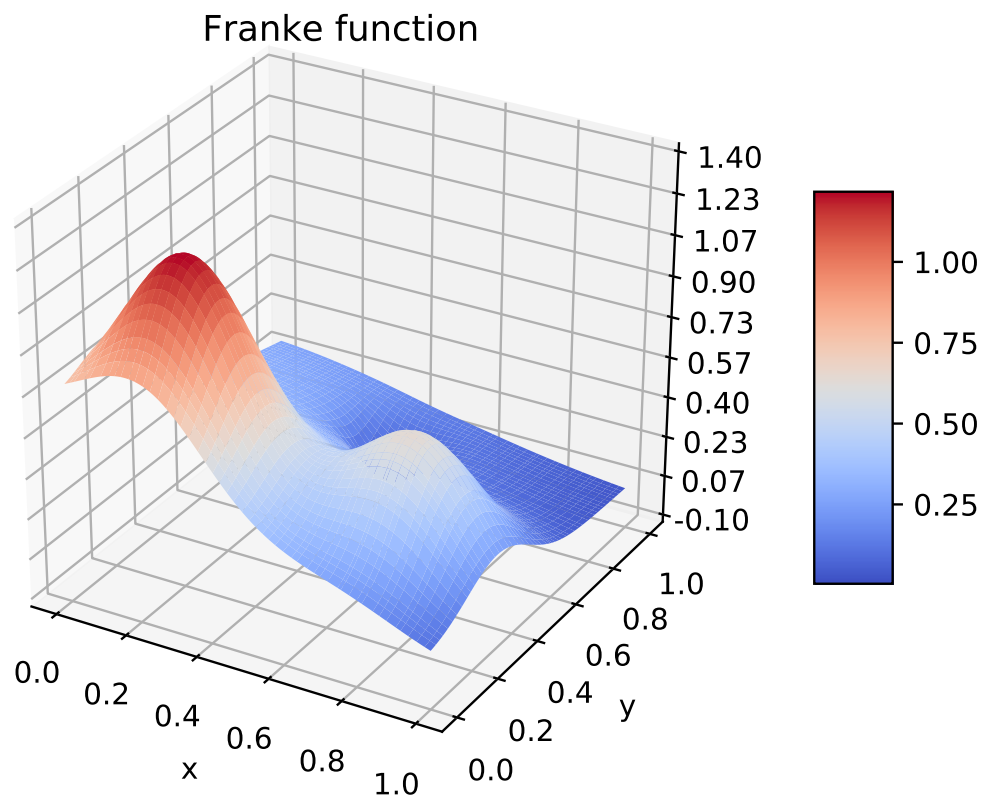


FIG. 4: Franke function on the domain $x, y \in [0, 1]$