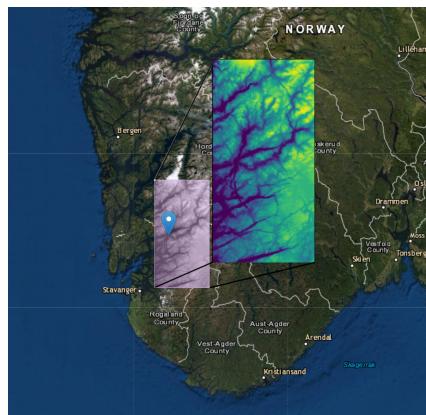


# Regression analysis and resampling methods

Jens Due, Mohamed Ismail & Oliver Hebnes

October 9, 2019



## Abstract

In this project, three regression methods are developed and tested with the aim of approximating two different data sets and generating models that fit these data. The first data set is generated from the Franke function, and the second is terrain data from an area in the south-western part of Norway. The regression methods used are the ordinary least squares method (OLS), the Ridge-regression and the LASSO-regression. For analyzing the quality of the models, several values are calculated and plotted as a function of model complexity. These include the mean squared error (MSE),  $R^2$ -score, confidence interval, bias and variance. Out of these, the MSE, and even more so, the  $R^2$ -score (with ideal values 0 and 1, respectively) are used as the determining factors for deciding upon which model is the best performer. It is shown that the models trained on the data from the Franke function has close to ideal values for the mentioned quality-parameters ( $MSE = 0.07$ ,  $R^2 = 0.92$ ), while the models do not perform as well on the terrain data ( $MSE = 0.325$ ,  $R^2 = 0.68$ ). This is concluded to be because of the high complexity of the terrain data, and inability to create a model with corresponding complexity. In the end, it is determined that the OLS is the best model for the Franke-function, and the Ridge-regression is the best for the terrain data. This is, however, not a straightforward answer, as all models seem to perform similarly in both data sets, showing that there is no 'one model to rule them all'.

# Contents

<b>I</b>	<b>Introduction</b>	<b>3</b>
<b>II</b>	<b>Theoretical Methods and Technicalities</b>	<b>3</b>
i	Ordinary Least Squares . . . . .	3
ii	Ridge regression . . . . .	4
iii	Lasso regression . . . . .	4
iv	Confidence interval . . . . .	5
v	Bias-Variance Tradeoff . . . . .	5
vi	MSE and R2-score . . . . .	7
vii	Resampling . . . . .	8
i	K-Fold Cross Validation . . . . .	8
<b>III</b>	<b>Data Sets</b>	<b>8</b>
i	Franke's Function . . . . .	8
ii	Terrain Data . . . . .	9
<b>IV</b>	<b>Implementation</b>	<b>10</b>
<b>V</b>	<b>Results</b>	<b>10</b>
i	Data from Franke's function . . . . .	10
i	Choosing hyperparameter $\lambda$ . . . . .	10
ii	Train and test model . . . . .	11
iii	Bias and variance . . . . .	11
iv	Confidence Interval . . . . .	12
v	Regression analysis . . . . .	12
ii	Terrain data . . . . .	13
i	Choosing hyperparameter $\lambda$ . . . . .	13
ii	Train and test model . . . . .	14
iii	Bias and variance . . . . .	14
iv	Confidence Interval . . . . .	15
v	Regression analysis . . . . .	15
<b>VI</b>	<b>Discussion</b>	<b>16</b>
i	Franke's Data . . . . .	16
ii	Terrain Data . . . . .	17
<b>VII</b>	<b>Conclusion</b>	<b>18</b>
<b>Appendices</b>		<b>20</b>
<b>A</b>	<b>Derivation of The Bias -Variance Tradeoff</b>	<b>20</b>
<b>B</b>	<b>Ordinary least square method fit on Franke's function.</b>	<b>21</b>
<b>C</b>	<b>Model fit as a function of data points</b>	<b>22</b>

## I Introduction

The ability to learn from our mistakes is of huge importance to human beings. This is how the world has been formed, by trial and error. Humans do not possess the ability to program each other, but since newer, modern times, we have discovered tools that could have been seen as magic for only a century ago. The pure combination of programming and statistics form this magnificent basis as we dive into three different regression types while trying to avoid dangers as getting a too much biased model or using too few data points. The purpose of the regression analysis is to try and model the infamous Franke-function and, later on, a real topographic area in the north named Norway. The possibilities of this application are limitless and it inspires a whole new generation of data scientists to achieve what we did not dare to dream about in the past.

## II Theoretical Methods and Technicalities

### i Ordinary Least Squares

$$f(X) = \beta_0 + \sum_{j=1}^p X_j \beta_j \quad (1)$$

This linear regression model assumes that the regression function  $f(X)$  is a reasonable approximation to the linear model. Here the  $\beta_j$  are unknown parameters and the variables  $X_j$  can come from different sources.

The residual sum of squared errors (RSS) is deviations predicted from actual empirical values of data. It is a measure of the discrepancy between the data and an estimation model. A small RSS indicates a tight fit of the model to the data.

It is used as an optimality criterion in parameter selection and model selection.

$$RSS = \sum_{i=1}^N (y_i - f(x_i))^2 \quad (2)$$

$$= \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p X_j \beta_j)^2 \quad (3)$$

One can minimize RSS to obtain the unique solution

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (4)$$

with the design matrix

$$\mathbf{X} = \begin{bmatrix} x_{00} & x_{01} & x_{02} & \dots & x_{1m} \\ x_{10} & x_{11} & x_{12} & \dots & x_{2m} \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & x_{n2} & x_{n3} & \dots & x_{nm} \end{bmatrix} \quad (5)$$

giving the fitted values from the training inputs

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

This is indeed the solution we are looking for in ordinary least squares.

## ii Ridge regression

Ridge regression shrinks the regression coefficients  $\beta$  by imposing a penalty on their size.

$$\hat{\beta}^{\text{ridge}} = \operatorname{argmin}_{\beta} \left\{ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p X_{ij}\beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\} \quad (7)$$

Here  $\lambda$  is a hyperparameter that controls the amount of shrinkage. The ridge solution are not equivariant (acts symmetrically on another symmetric space), so one needs to standardize the inputs before solving this equation.

This can also be written in the RSE-form, but with a size constraint on the parameters.

$$\hat{\beta}^{\text{ridge}} = \operatorname{argmin}_{\beta} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p X_{ij}\beta_j \right)^2 \quad (8)$$

$$\text{subject to } \sum_{j=1}^p \beta_j^2 \leq t \quad (9)$$

There is a one to one correspondence between  $\lambda$  and  $t$ .

When there are too many correlated variables in a linear regression model, their coefficients can become poorly determined and exhibit high variance. A wildly large positive coefficient on one variable can be canceled by a similarly large negative coefficient on its correlated cousin. But by imposing a size constraint on the coefficients, as in (9), the problem decreases in importance.

The ridge solution can also be written in matrix form.

$$RSS(\lambda) = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda \beta^T \beta \quad (10)$$

with the ridge regression solution as

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

## iii Lasso regression

The LASSO (Least Absolute Shrinkage and Selection Operator) regression is also a shrinkage method like Ridge, but there is a small difference. Ridge regression wants to minimize the loss function, and to do so we add a hyperparameter  $\lambda$  to minimize the parameters that does not matter that much. LASSO has, however, the ability to make the coefficients zero and totally exclude them if they are not relevant for the model.

The implementation of LASSO was done by using a module from the scikit-learn package.

In addition, LASSO also depends on another parameter called the learning rate  $\gamma$ , and this parameter will be optimized by the scikit learn package using various methods, but as this is not an important subject of this project, we will not elaborate further.

#### iv Confidence interval

In statistics, a confidence interval is an estimate of an interval to a certain degree, where it might contain the true value of a population parameter.

To calculate the confidence interval of the parameter  $\beta$  one would have to first calculate the variance of  $\beta$ :

$$Var(\beta) = (\mathbf{X}\mathbf{X}^T)^{-1}\sigma_\epsilon^2$$

The standard deviation is defined as the square root of the variance, so the standard deviation of  $\beta$  would then be:

$$\sigma_\beta = \sqrt{(\mathbf{X}\mathbf{X}^T)^{-1}\sigma_\epsilon^2} = \sqrt{(\mathbf{X}\mathbf{X}^T)^{-1}\sigma_\epsilon}$$

This is a special case where we know the variance of  $\beta$ , but for the general case, this may not be known. One approach is then to approximate the variance with the mean squared error for  $\beta$ .

$$\begin{aligned} Var(\beta) &\approx MSE_\beta \\ \sigma_\beta &\approx \sqrt{MSE_\beta} \end{aligned}$$

Then to find a 95% confidence interval for  $\beta$ , one needs to use the formula:

$$\beta \pm z \times \sigma_\beta$$

Where  $z$  is the 95th percentile of the normal distribution. This corresponds to  $z = 1.96$ .

The theory can be found in the textbook of Hastie et al. [2].

#### v Bias-Variance Tradeoff

We have assumed that our data can be represented by an unknown function with the addition of some noise  $\epsilon$ .

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

Where  $\epsilon$  is normally distributed with a mean equal to zero and standard deviation equal to  $\sigma^2$ .

Furthermore, we also assume that the function  $f(\mathbf{x})$  can be approximated to a model  $\tilde{\mathbf{y}}$ , where the model is defined by a design matrix  $\mathbf{X}$  and parameters  $\beta$ .

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

The parameters  $\beta$  are in turn found by optimizing the mean squared error (MSE) via the so-called cost function

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

One can show that the cost function can be rewritten as

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

By simply using these relations

$$\begin{aligned} E[\mathbf{y}] &= \mathbf{f} \\ E[\epsilon] &= 0 \\ Var(\mathbf{y}) &= Var(\epsilon) = \sigma_\epsilon^2 \end{aligned}$$

Since the variance of  $\mathbf{y}$  and  $\epsilon$  are both equal to  $\sigma^2$ . The mean value of  $\epsilon$  is by definition equal to zero. In addition, the function  $\mathbf{f}$  is not a stochastic variable, and the same argument can also be used for  $\tilde{\mathbf{y}}$ . Then we plug in the definition for  $\mathbf{y}$  into the cost function:

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

Adding and subtracting  $E[\mathbf{y}]$  we get:

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

And simply by using the relations mentioned above concerning the expectation value for  $\mathbf{y}$  and the variances for both  $\mathbf{y}$  and  $\epsilon$ , the cost function can be rewritten to:

$$\begin{aligned} E[(\mathbf{y} - \tilde{\mathbf{y}})^2] &= E[(\mathbf{f} - E[\tilde{\mathbf{y}}])^2] + Var(\tilde{\mathbf{y}}) + \sigma_\epsilon^2 \\ E[(\mathbf{y} - \tilde{\mathbf{y}})^2] &= \frac{1}{n} \sum_i (f_i - E[\tilde{\mathbf{y}}])^2 + \frac{1}{n} \sum_i (\tilde{y}_i - E[\tilde{\mathbf{y}}])^2 + \sigma_\epsilon^2. \end{aligned}$$

The full derivation can be found in Appendix A.

The first term on the right hand side is the squared bias, the amount by which the average of our estimate differs from the true mean. The second term represents the variance of the chosen model and finally the last term is the variance of the error  $\epsilon$ .

The diagram above illustrates the effect the complexity of a model has on the MSE, bias and variance. The more complex model one has, the lower the bias becomes, and the higher the variance becomes. An ideal model would be one that simultaneously achieves low variance and low bias.

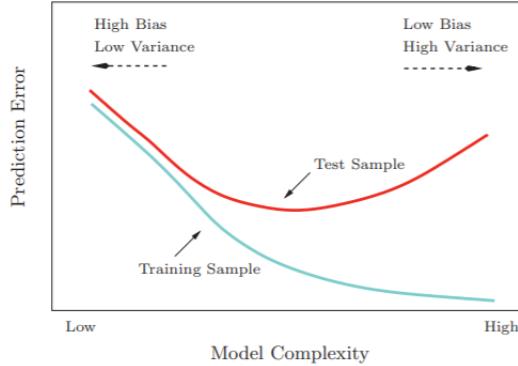


Figure 1: Test and training error as a function of model complexity [2]

## vi MSE and R2-score

Two of the main methods for determining how well our model fits the data, will be evaluating the MSE (mean squared error) and the R2-score. They are defined as

$$MSE(\hat{y}, \tilde{y}) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2,$$

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

For a perfect fit, these would be valued at 0 for the MSE, and 1 for the R2-score. These are the values our model should approach in order to be a good approximation of the data, and our goal is to make sure this is the case.

## vii Resampling

Resampling methods are an indispensable tool in modern statistics. They involve repeatedly drawing samples from a training set and refitting a model of interest on each sample in order to obtain additional information about the fitted model [1].

### i K-Fold Cross Validation

K-fold cross validation is a resampling technique where a given data set is divided into a K number of subsets/folds. One after another, each fold will play the role of the test set while the rest are used to train the model, as visualised in the diagram below where  $K = 5$ . This procedure ensures that the entire data set is being used both as training and testing data in a balanced manner.

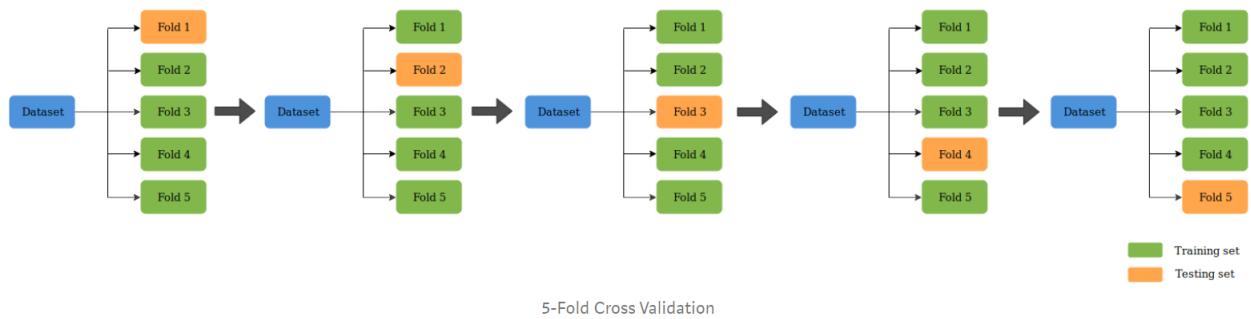


Figure 2: Diagram visualising the 5-fold Cross Validation [3]

## III Data Sets

### i Franke's Function

Our first data set is generated from [Franke's function](#). It is defined as

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

This is a function which has been widely used when testing various interpolation and fitting algorithms.

In this project added stochastic noise was to Franke's function before fitting to it. The included noise was normally distributed with mean equal to zero and standard deviation set to one,  $N(0, 1)$ .

A plot of the Franke function without the noise can be seen in the figure 3.

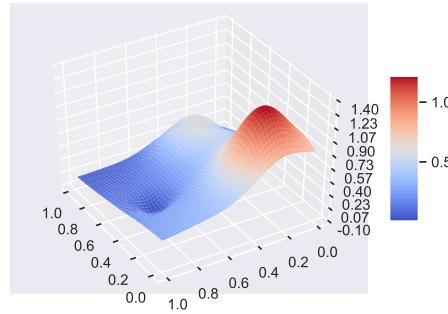


Figure 3: The Franke function plotted for  $0 \leq x, y \leq 1$

## ii Terrain Data

Our second data set is terrain data from an area on the south-eastern part of Norway, not far from the city of Stavanger. The data is sourced from the U.S. Geological Survey [4]. A visual representation of the data is shown in the figure 4.

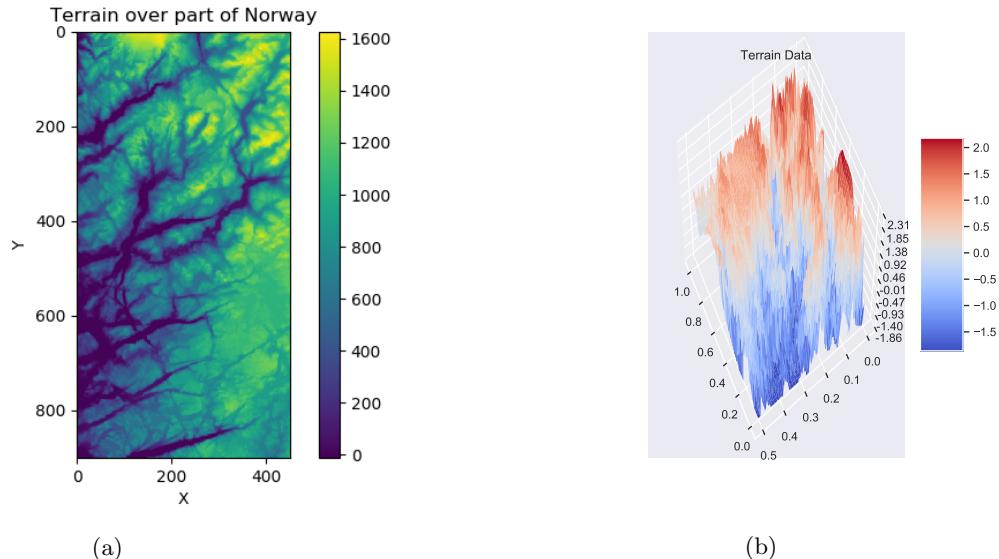


Figure 4: The terrain data visualized in both two and three dimensions. Figure (b) has been normalised while the z-value in (a) has real values. The data is gathered just outside of a place in the south-west part of Norway called Stavanger. Do you notice the steep fjords and the calm waters?

## IV Implementation

Programs used in this project can be found on our github repository [5]. These programs are inspired by the codes found in the course's github repository [1].

## V Results

### i Data from Franke's function

All calculations on Franke's function has been with a total amount of data points as 400.

#### i Choosing hyperparameter $\lambda$

From figure 5 it is clear that the choice of  $\lambda$  is not as crucial as choosing the correct polynomial degree to minimise the mean square error. We can see that several  $\lambda$ -values give satisfactorily small MSE, thus we choose  $\lambda = 10^{-3}$  &  $\lambda = 10^{-4}$  for further analysis with Ridge- and Lasso-regression.

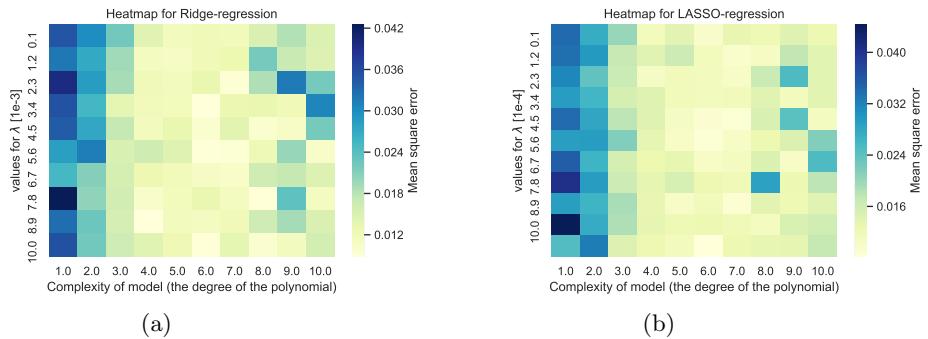


Figure 5: Heatmaps showing different hyperparameters  $\lambda$  as a function of model complexity and mean squared error for (a) Ridge-regression, and (b) LASSO

## ii Train and test model

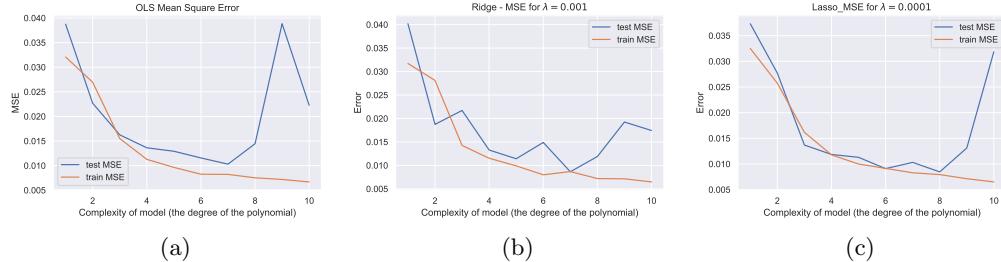


Figure 6: The mean squared error of the testing and training data as a function of model complexity on Franke's Function for (a) OLS, (b) Ridge and (c) LASSO-regression.

## iii Bias and variance

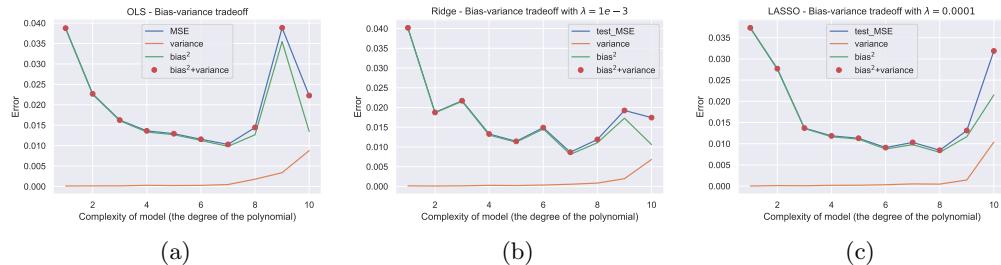


Figure 7: The mean squared error, bias and variance as a function of model complexity on Franke's Function for (a) OLS, (b) Ridge and (c) LASSO-regression.

#### iv Confidence Interval

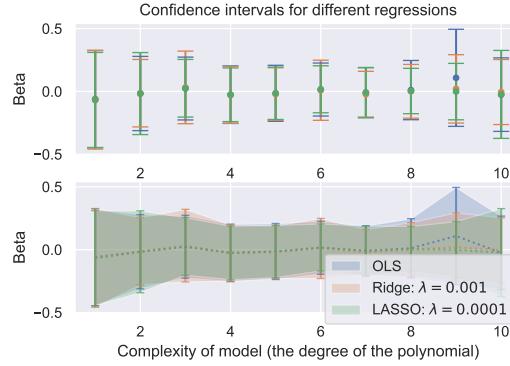


Figure 8: Confidence intervals for the regression coefficients  $\beta$  for all three regression methods. For polynomial degrees 4-7, the intervals are quite small and overlapping. Outside this range, the intervals are shown to be larger.

#### v Regression analysis

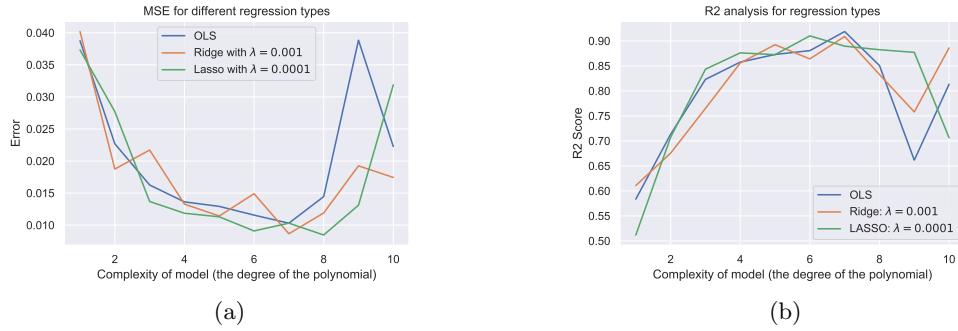


Figure 9: The mean squared error and the  $R^2$ -score as a function of the model complexity on Franke's function with OLS, Ridge & LASSO.

Table 1: Lowest mean squared error calculated on the testing data.

Regression method	Model complexity	MSE
OLS	7	0.010
Ridge	7	0.008
LASSO	8	0.007

Table 2: Highest  $R^2$ -score calculated on the testing data

Regression method	Model complexity	$R^2$ -score
OLS	7	0.92
Ridge	7	0.91
LASSO	6	0.91

## ii Terrain data

Here, we repeat the above analysis on the new data set from the terrain data.

### i Choosing hyperparameter $\lambda$

As mentioned before, choosing a hyperparameter  $\lambda$  is normally of great importance for Ridge and LASSO. However, figure 10 shows similar trends as figure 5, that the choice of  $\lambda$  is not very important, and thus we will choose the same values as was done with Franke's function for further analysis on the terrain data.

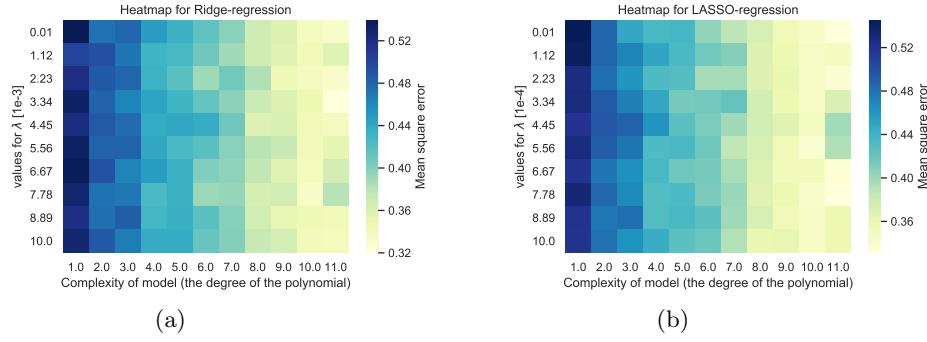


Figure 10: Heatmaps showing different hyperparameters  $\lambda$  as a function of model complexity and mean squared error for (a) Ridge-regression, and (b) LASSO for Franke's function.

## ii Train and test model

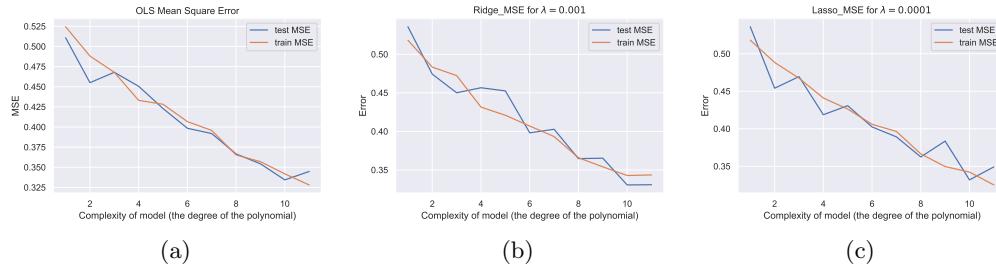


Figure 11: The mean squared error of the test and train data as a function of model complexity on the terrain data for (a) OLS, (b) Ridge and (c) LASSO-regression.

## iii Bias and variance

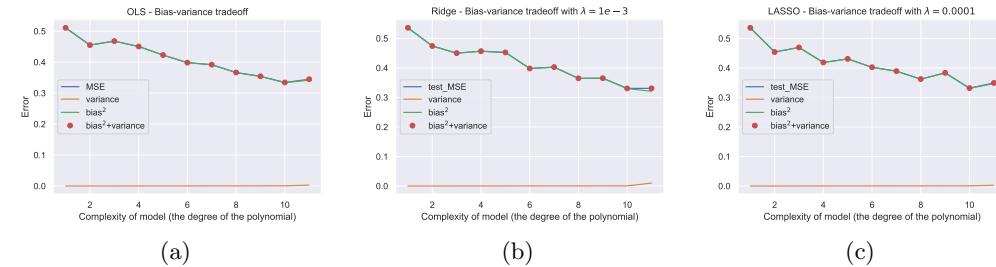


Figure 12: The mean squared error, bias and variance as a function of the model complexity on the terrain data for (a) OLS, (b) Ridge and (c) LASSO-regression.

#### iv Confidence Interval

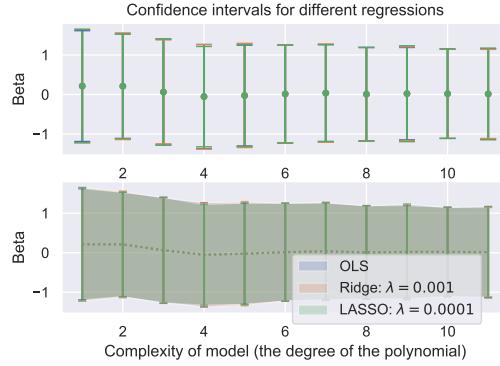


Figure 13: Confidence intervals for the regression coefficients  $\beta$  for the terrain data. We can see the intervals essentially completely overlap for all three regression methods.

#### v Regression analysis

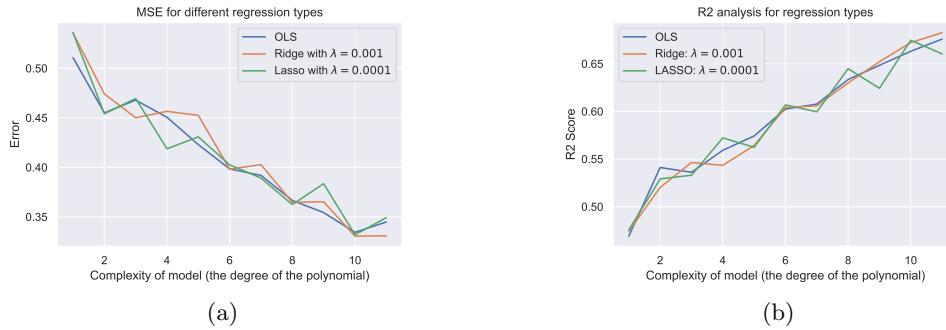


Figure 14: The mean squared error and  $R^2$ -score as a function of the model complexity on the terrain data with OLS, Ridge & LASSO.

Table 3: Lowest and highest mean squared error &  $R^2$ -score calculated on the testing data.

Regression method	Model complexity	MSE	$R^2$ -score
OLS	10	0.340	0.66
Ridge	11	0.325	0.68
LASSO	10	0.330	0.67

## VI Discussion

### i Franke's Data

As discussed earlier in the text, both Ridge and LASSO depend on the hyperparameter  $\lambda$  for minimising the mean squared error. And this is clearly illustrated for both Ridge and LASSO in figure 5, where the MSE has been plotted for various  $\lambda$  as a function of the complexity of the model. The figure also tells us that there are no specific  $\lambda$  that minimises the error the most, and that a lot of different values of  $\lambda$  minimise the error equally. For that reason we chose two of the best performing values of  $\lambda$ .

From figure 6 we clearly observe that the training data is performing better than the testing data for all the regression methods. As the complexity increases, the mean square error for the training data decreases, and will in theory decrease until we get an ideal fit. However, this is not the case for the testing data. As we train our model on the training data, we will make an increasingly tight fit to it which may cause overfitting. We do want to decrease our training mean square error, but we would also like to use our model on independent test sets and still have the error as low as possible. We will not be able to do this if we overfit the model to the training data, and this is shown in the figure where the training data gets a tighter fit, but the testing data does not experience the same for increasing polynomial degrees.

By reading off the results of figure 7, we can see that the bias is the main contribution to the MSE up until a certain complexity of the models for all the regression methods. At a certain point, however, the variance drastically increases. These results reflect the behaviour we discussed earlier in the report, where the complexity of the model affects the bias and variance. The complexity of the model that minimises the bias and the variance is not the same for all the regression models. OLS achieves this at polynomial degree 7, as well as Ridge with  $\lambda = 10^{-3}$ , while LASSO with  $\lambda = 10^{-4}$  does it at polynomial degree 8.

With an increasing number of data points for Franke's function, it was found that the model became more stable (see Appendix C) and did not explode as it did with lower amount of data points. The data was plotted while choosing the lowest mean square error and the highest  $R^2$  score for up to ten polynomial degrees. This defends the use of only 400 data points, as it does not mean that we will get any significantly better scores running computationally heavy algorithms with huge matrices.

Looking back at the figure 5 one can easily see deviations of the model. One example is the ordinary least square method in figure 7a where the mean square error peaks at polynomial degree 9. One can look at the corresponding confidence interval for beta in figure 8, where the interval peaks as well at polynomial degree 9. This was not the case when the number of data points were increased, but on the other hand we did not get an increase in either variance, or a substantial decrease in bias for a polynomial degree up to 10. When increasing the polynomial degree up to 15, the values were exploding, both with a small and large amount of data points. This can easily be visualised in the jupyter notebook in our GitHub repository [5]. One of the reasons behind this might be lack of numerical precision. Another reason can be instability in the implemented algorithms.

From figure 9a the total lowest mean square error is achieved at polynomial degree 8 for LASSO with a value of 0.007 (1). For an ideal fit, the MSE should be zero. Thus, all of our types of regressions do indeed have a good fit compared to an ideal fit. On the other hand it is the ordinary least square method that results in the total biggest  $R^2$ -score for the polynomial degree of 7 with a value of 0.92, where the ideal value should be at 1. At each time the program is executed, the data that becomes the testing data changes, and the mean squared error and  $R^2$ -score changes along with it. This can be observed in the jupyter notebook [5]. This also indicates that the best regression type might change every time the program is run.

## ii Terrain Data

The terrain data includes too many data points which would need a lot of computational power, thus we take every 25th point, giving a total of 10550 points.

Figure 10 is similar to 5, but with some differences. The complexity of the terrain data will demand a higher polynomial degree for low mean square errors, and once again, the hyperparameter  $\lambda$  is not as important as the complexity. The mean square error is a function of both bias and variance, but for the lambdas and degrees of polynomial that were simulated, we did not see a substantial increase in mean square errors. For higher polynomial degrees, it was expected that the bias-variance trade-off that was seen for lower polynomial degrees for Franke's function would present itself, but unfortunately, due to both numerical instability and lack of precision, the numbers were exploding.

There was not observed any clear signs of overfitting when observing the figure 11, due to an improvement of the mean square error for both training and testing data. This is further strengthened by looking at figure 12 where the only contribution for the mean square error is the bias. However, this might be a sign of underfitting, which could indicate either that the model has not trained enough, or it is not using a high enough polynomial degree.

In a highly varying landscape such as our terrain data, the calculated confidence interval for the regression coefficients  $\beta$  will behave accordingly. This behaviour is seen in figure 13 for all regression types where significant intervals indicates a not very good fit. This is expected because of the varying data.

The  $R^2$ -scores and the mean square errors in figure 14 show that there is not necessarily one regression type that excels, but that every one has their pros and cons when they work in their best environment. It was expected that the regression type with the lowest mean square error would also produce the highest R2 score, but our results for both Franke's function and the terrain data contradict this hypothesis. Nonetheless, we are seeing a much better fit of the testing data when comparing figure 16 and figure 14. As for Franke's function, the mean square error is significantly smaller than the mean square error of the terrain data, and the  $R^2$  scores for Franke's function reach values above 0.90 against the terrain data's 0.67. The mean square error is not necessarily an indication of a model that will do well on numerous different data sets, because a model can be fantastic while still producing non-optimal mean square errors. The  $R^2$  score, however, will give a score based independently on the amount of data sets. Thus, a higher  $R^2$  score will correspond to a better model. Based on this argument, we can make a choice as to which method generated the best model for our data sets.

It is of interest to train on a large scale of data, and test on the remaining small parts of the

data. If there is not enough data points to train on, one will not get a general enough model that will work on data different from the training data. On the other hand, if there is too much data with high complexity, as is the case with the terrain data, the difficulty increases for developing and enhancing a model that will work on a given data set. This demands stronger optimization of code and more powerful machinery.

## VII Conclusion

In this project, we have implemented the regression methods OLS, Ridge and LASSO for data from Franke's function and terrain data from the south-western part of Norway. Observing the results on the terrain data, the Ridge-regression model is slightly better than the OLS and LASSO, when evaluating the mean squared error and the  $R^2$ -score. As for the Franke function, the results show no clear indication as to which model fit the data best, as LASSO produced the lowest MSE, and OLS produced the highest  $R^2$ -score. As we argued that the  $R^2$ -score is the best indicator of quality, we then conclude that OLS is the favoured model for this data set. Furthermore, these results may vary each time the program is executed, because of the random choice of test data, as mentioned before. This discrepancy may also have been caused by possible errors in our code which could have given us these contradicting results.

The art of balancing the eternal bias-variance trade-off while struggling with a huge amount of calculations, while only at an apprentice level, has yielded us great experience and a unique insight into data handling, regressions and resampling.

As an evaluation of our work through the project, we found some potential in writing down diaries during the weeks, since there were many times we spent a tiny amount of a lifetime on figuring out what we had done earlier. In addition, using time- and CPU-efficient packages such as scikit-learn has shown to provide more stable and robust machinery, yielding better results for even higher complexities.

## References

- [1] Morten Hjorth-Jensen, Github Repository, <https://github.com/CompPhysics/MachineLearning/tree/master/doc>
- [2] Trevor Hastie, Robert Tibshirani, Jerome H. Friedman, The Elements of Statistical Learning, Springer <https://www.springer.com/gp/book/9780387848570>.
- [3] Medium, Data Driven Investor, K-Fold Cross Validation  
<https://medium.com/datadriveninvestor/k-fold-cross-validation-6b8518070833>
- [4] U.S. Geological Survey, EarthExplorer <https://earthexplorer.usgs.gov/>
- [5] GitHub repository containing all code and plots, as well as a Jupyter notebook for easy reproducibiliy <https://github.com/ohebbi/FYS-STK4155/tree/master/project1>

# Appendices

## A Derivation of The Bias -Variance Tradeoff

$$\begin{aligned} E[(\mathbf{y} - \tilde{\mathbf{y}})^2] &= E[(\mathbf{f} + \epsilon - \tilde{\mathbf{y}})^2] \\ E[\mathbf{y}] &= \mathbf{f} \\ E[\epsilon] &= 0 \\ Var(\epsilon) &= \sigma_\epsilon^2 \end{aligned}$$

Adding and subtracting  $E[\mathbf{y}]$  we get

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

Then we need to find an expression for  $E[(\mathbf{f} - \tilde{\mathbf{y}})^2]$ . And we do that by adding and subtracting  $E[\tilde{\mathbf{y}}]$

$$\begin{aligned} E[(\mathbf{f} - \tilde{\mathbf{y}})^2] &= E[(\mathbf{f} - \tilde{\mathbf{y}} + E[\tilde{\mathbf{y}}] - E[\tilde{\mathbf{y}}])^2] \\ E[(\mathbf{f} - \tilde{\mathbf{y}})^2] &= E[(\mathbf{f} - E[\tilde{\mathbf{y}}])^2] + E[(E[\tilde{\mathbf{y}}] - \tilde{\mathbf{y}})^2], \quad 2E[(\mathbf{f} - E[\tilde{\mathbf{y}}])(E[\tilde{\mathbf{y}}] - \tilde{\mathbf{y}})] = 0 \\ E[(\mathbf{f} - \tilde{\mathbf{y}})^2] &= E[(\mathbf{f} - E[\tilde{\mathbf{y}}])^2] + E[(\tilde{\mathbf{y}} - E[\tilde{\mathbf{y}}])^2] \\ E[(\mathbf{f} - \tilde{\mathbf{y}})^2] &= E[(\mathbf{f} - E[\tilde{\mathbf{y}}])^2] + Var(\tilde{\mathbf{y}}) \end{aligned}$$

Then we put that expression back in above.

$$\begin{aligned} E[(\mathbf{y} - \tilde{\mathbf{y}})^2] &= \sigma_\epsilon^2 + E[(\mathbf{f} - \tilde{\mathbf{y}})^2] \\ E[(\mathbf{y} - \tilde{\mathbf{y}})^2] &= E[(\mathbf{f} - E[\tilde{\mathbf{y}}])^2] + Var(\tilde{\mathbf{y}}) + \sigma_\epsilon^2 \\ \mathbb{E}[(\mathbf{y} - \tilde{\mathbf{y}})^2] &= \frac{1}{n} \sum_i (f_i - \mathbb{E}[\tilde{\mathbf{y}}])^2 + \frac{1}{n} \sum_i (\tilde{y}_i - \mathbb{E}[\tilde{\mathbf{y}}])^2 + \sigma^2. \end{aligned}$$

## B Ordinary least square method fit on Franke's function.

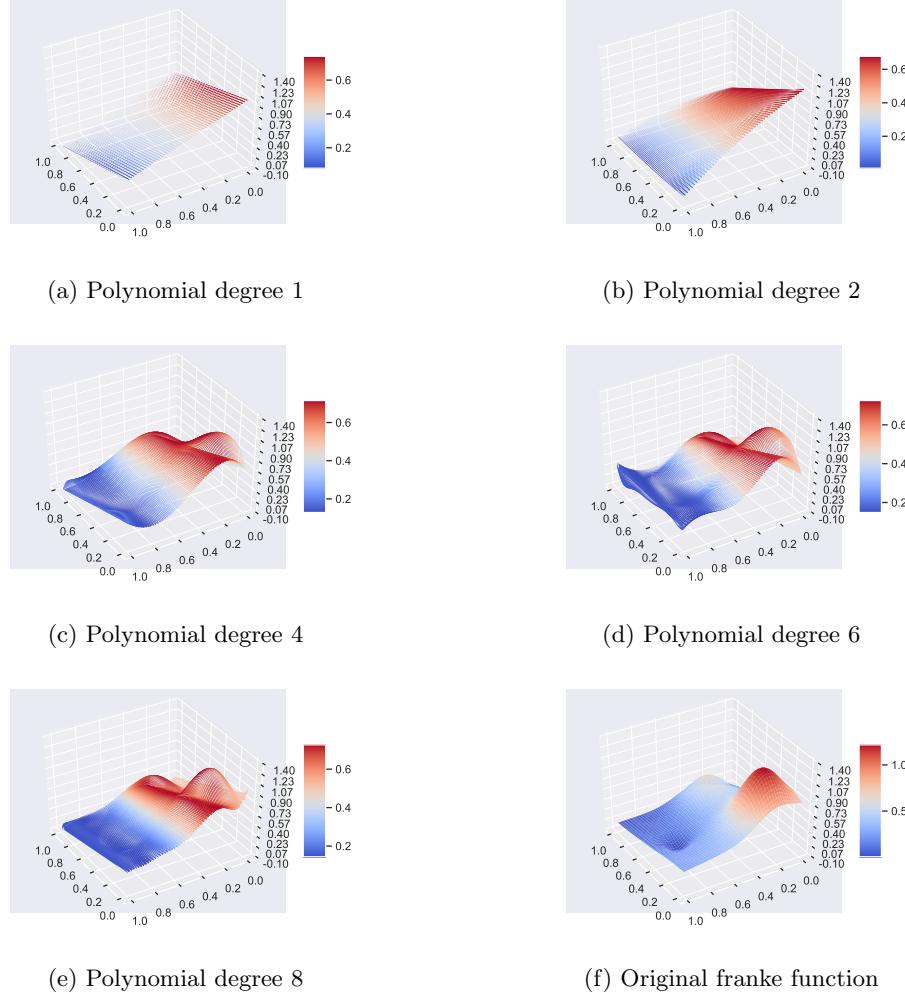


Figure 15: These figures show how the polynomial grad changes the contour of the landscape. The regression type used is least squared method.

## C Model fit as a function of data points

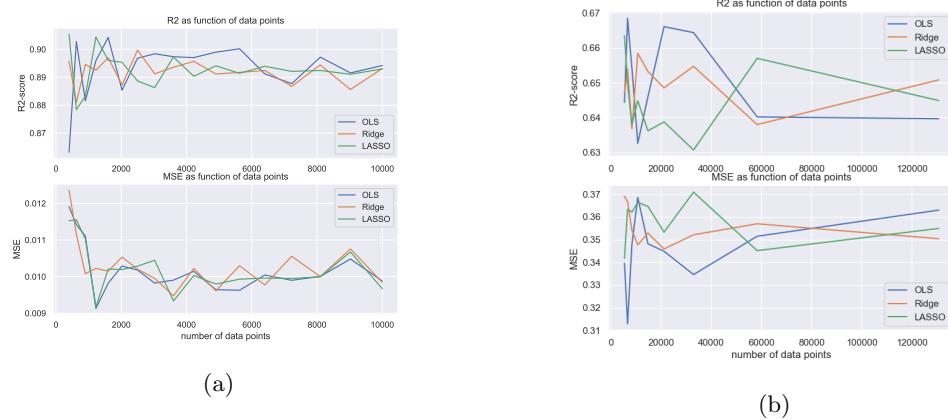


Figure 16: Franke's function and terrain data as function of data points. The values of  $R^2$  and MSE do not change significantly as the amount of data points increases.