

Project 1 in FYS-STK4155

Adrian Martinsen Kleven
Simon Elias Schrader

Autumn 2020

Contents

1	Abstract	3
2	Introduction	3
3	Methods	3
3.1	Linear Regression	3
3.1.1	Ordinary Least Squares Regression	4
3.1.2	Ridge Regression	5
3.1.3	LASSO Regression	6
3.2	Resampling Methods and Data Splitting	6
3.3	Bias-variance relationship	6
3.3.1	Bootstrap method	7
3.3.2	K-fold Cross validation	7
3.4	Fitting the Franke function	8
3.5	Geographic Data	8
4	Computational Implementation	9
4.1	Finding the ideal parameter λ	9
4.2	Error in LASSO-Regression	9
4.3	Sampling Franke Function	9
4.4	The Design Matrix	10
4.5	Scaling	10
5	Results	10
5.1	β standard deviation in OLS and Ridge regression	10
5.2	Error comparison	10
5.2.1	Bootstrap vs. Cross validation	10
5.2.2	Franke Function	12
5.2.3	Geographic data	13
5.3	Bias-Variance tradeoff	15
5.3.1	Franke Function	15
5.3.2	Geographic data	16
5.4	Evaluation of the Geographic Fit	17
6	Conclusion	18
7	Appendix	19
7.1	Proof that the MSE can be decomposed into Bias and Variance	19

List of Figures

1	Optimal regularization parameter to noise	5
2	K-fold Cross Validation operation structure	7
3	Estimated MSE from OLS using different number of splits in k-fold CV	8
4	Taebbaek Mountain	9
5	Comparing Error estimates from Bootstrap vs. K- fold CV	11
6	Comparing Error estimates from Bootstrap vs. K- fold CV -1000 data points	11
7	Comparing time elapsed from Bootstrap vs. K- fold CV	12
8	Test MSE for LASSO, Ridge and OLS regression	12
9	Test R^2 -value for LASSO, Ridge and OLS regression	13

10	Test MSE for LASSO, Ridge and OLS regression	14
11	Bias, Variance, MSE & training MSE as a function of the polynomial degree using LASSO regression	15
12	Bias, Variance, MSE & training MSE as a function of the polynomial degree using OLS regression	15
13	Bias, Variance, MSE & training MSE as a function of the polynomial degree using Ridge regression	16
14	Bias, Variance, MSE & training MSE as a function of the polynomial degree for Geodata	17
15	Fitted Geodata	18

List of Tables

1	Confidence interval of estimators in OLS vs. Ridge	10
---	--	----

1 Abstract

In almost all fields of science, extrapolating data and creating a fitting function in a reliable way, is tremendously important, and thanks to Computers, doing so is increasingly possible using automatized methods. Using existing OLS, Ridge, and LASSO Regression, we fitted the Franke Function (2) with added noise and Geographic Data with 2D polynomials of varying degrees and analysed different statistical values of interest. We found that OLS is prone to overfitting due to a very large variance at high complexities, while Ridge and LASSO remain stable at these higher degrees. With our parameters & data, Ridge and OLS regression gave about the same Mean Square Error (MSE), while LASSO did not converge to a desirable value at higher polynomial degrees and hence gave larger errors. For the geographical data, we found that, with 5000 randomly chosen data points, using OLS Regression, the MSE is around 22,500 at the ideal value, while using Ridge Regression, it is 23,000. Hence, regression does not seem to be suitable for predicting geographical data. We also compared k-fold Cross Validation to Bootstrap and found out that these two methods do not give the same error estimates.

2 Introduction

The purpose of this article is to explore three different linear regression algorithms, their utility, application and some of the theory behind them. These are Ordinary Least Squares (OLS), as well as the two Shrinkage methods Ridge regression and LASSO regression. In this article, we compare these methods to one another, as well as describe their applicability and their peculiarities. These methods are tested on a model dataset constructed using the Franke function (2) with added normally distributed noise, then later applied to digital terrain data over a part of South Korea in an attempt to accurately recreate the image. We explore the two resampling techniques K- fold cross validation and Bootstrapping. Studying the bias- variance trade off in model selection. Resampling is usually used for models where the amount of data is restricted. Seeing as terrain data has a finite resolution, its application should be considered for error estimation or model selection although we don't cover this. The models we consider are not limited by lack of data, but resampling is nevertheless a valid method for predicting various statistical values.

All methods are implemented using the Python programming language. The programs are available at the [Github address](#).

3 Methods

3.1 Linear Regression

The aim of Linear Regression is to find the conditional distribution $p(y|\mathbf{x})$ for y given \mathbf{x} . In Linear Regression, a linear functional dependence is assumed, and the aim is hence to create a linear function that fits a set of p predictor variables \mathbf{x} to a target variable \mathbf{y} , where there are n discrete observations for each of the p predictor variables and the target variable. That is, the aim is to find a functional relationship of the form $f(\mathbf{x}) = y$.

We assume hence that each target variable y_i has a functional dependence of the predictor variables:

$$y_i = \tilde{y}_i + \epsilon_i = \beta_0 + \beta_1 x_{i0} + \beta_2 x_{i1} + \dots + \beta_p x_{i(p-1)} + \epsilon_i$$

As this is true for all y , this can be written in matrix form:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

where \mathbf{y} , and $\boldsymbol{\epsilon}$ are vectors of length n , $\boldsymbol{\beta}$ is of length p and \mathbf{X} is a matrix of size $n * p$. \mathbf{X} is called the design matrix. We now want to find fitting parameters $\boldsymbol{\beta}$ that represent the best linear fit of y to the input data \mathbf{x} . Hence, we want the fit

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

to be as close to the target data \mathbf{y} as possible. There are several approaches to define the best linear fit, and in this article, we looked at three different methods that define the best fit in a different way each. All methods have in common that they define a loss function, which describes the deviation from the fit to the real data, that needs to be minimized.

3.1.1 Ordinary Least Squares Regression

In Ordinary Least Square regression, the loss function is defined as the following:

$$C(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=0}^n (y_i - \tilde{y}_i)^2 = \frac{1}{n} (\mathbf{y} - \tilde{\mathbf{y}})(\mathbf{y} - \tilde{\mathbf{y}})^T = \frac{1}{n} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T$$

The solution minimizing this problem is found by deriving $C(\boldsymbol{\beta})$ with regards to each β_i and setting each derivative equal to zero. Doing this, we end up with the following equation (in matrix notation)

$$\frac{\partial C(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = 0 = \mathbf{X}^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

which gives us

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

which can be solved when $(\mathbf{X}^T \mathbf{X})$ is invertible. In the case of strongly correlated data and large matrices, There might be issues in inverting $(\mathbf{X}^T \mathbf{X})$ due to the matrix being singular. Also, for very large matrices, there might be issues caused by the fact that many mathematical operations can introduce a numerical error. Hence, an equivalent, but more stable way to implement the solution for $\boldsymbol{\beta}$ is to use Singular Value Decomposition (SVD). As every matrix has a SVD decomposition, we can write

$$\begin{aligned} \boldsymbol{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= (\mathbf{V} \mathbf{D}^T \mathbf{U}^T \mathbf{U} \mathbf{D} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{D}^T \mathbf{U}^T \mathbf{y} \\ &= (\mathbf{V} \mathbf{D}^T \mathbf{I}_p \mathbf{D} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{D}^T \mathbf{U}^T \mathbf{y} \\ &= (\mathbf{V} \mathbf{D}^2 \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{D}^T \mathbf{U}^T \mathbf{y} \\ &= \mathbf{V} \mathbf{D}^{-2} \mathbf{V}^T \mathbf{V} \mathbf{D}^T \mathbf{U}^T \mathbf{y} \\ &= \mathbf{V} \mathbf{D}^{-1} \mathbf{U}^T \mathbf{y} \end{aligned} \tag{1}$$

Under the assumption that true data y is given as

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

where $\boldsymbol{\epsilon}$ follows a normal distribution $N(0, \sigma^2)$, information about the variance and the expectation value can be deduced. Assuming that $f(\mathbf{x}) = \mathbf{X}\boldsymbol{\beta}$, it is easily shown that

$$\mathbb{E}(y_i) = \mathbf{X}_{i,*} \boldsymbol{\beta}$$

$$\text{Var}(y_i) = \sigma^2$$

where $\mathbf{X}_{i,*}$ is the i -th row of the design matrix.

For $\boldsymbol{\beta}$, it can be shown that

$$\begin{aligned} \mathbb{E}(\boldsymbol{\beta}) &= \boldsymbol{\beta} \\ \text{Var}(\boldsymbol{\beta}) &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \end{aligned}$$

3.1.2 Ridge Regression

The above mentioned problems with OLS regression, namely those of correlated data and near singular matrices, can be amended by implementing ridge regression. This is done by introducing the regularization parameter λ in the cost function:

$$C(\mathbf{X}, \boldsymbol{\beta}) = \frac{1}{n} \{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\} + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta}$$

effectively imposing a penalty on the size of $\boldsymbol{\beta}$. One of the consequences of Ridge Regression is that it "shrinks" the singular values σ of the Design matrix \mathbf{X} : It can be shown that [1]

$$\begin{aligned} \mathbf{y}^{OLS} &= \mathbf{X} \hat{\boldsymbol{\beta}}^{OLS} = \mathbf{U} \mathbf{U}^T \mathbf{y} \\ \mathbf{y}^{Ridge} &= \mathbf{X} \hat{\boldsymbol{\beta}}^{Ridge} = \sum_{j=1}^p \mathbf{u}_j \frac{\sigma_j^2}{\sigma_j^2 + \lambda} \mathbf{u}_j^T \mathbf{y} \end{aligned}$$

where \mathbf{u}_j are the column of the matrix \mathbf{U} from the SVD of \mathbf{X} . It also follows directly from this formula that the smaller singular values are shrunk more.

This has the side effect of decreasing the variance of the parameters $\boldsymbol{\beta}$, as well as the values $\boldsymbol{\beta}$ themselves (it can be shown [2] that the variance for all $\lambda > 0$ is always lower than when using OLS), and thereby reducing the impact of overfitting.

The expression for the optimal value of $\boldsymbol{\beta}$ becomes

$$\boldsymbol{\beta}^{Ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}.$$

Where as before the matrix $\mathbf{X}^T \mathbf{X}$ could have been singular, by adding a non- zero, positive value to the diagonal, the eigenvalues become positive definite ensuring that the matrix is invertible.

Regularization parameter There are several factors that determine how different regularization parameters λ will minimize the cost function. A dataset with a lot of noise will cause more overfitting at higher model complexities. We expect that increasing the noise parameter σ in a model dataset will make models with a larger value of λ , perform better. We see an example of this in figure 1, where models with higher values of λ give the smallest MSE with increasing noise.

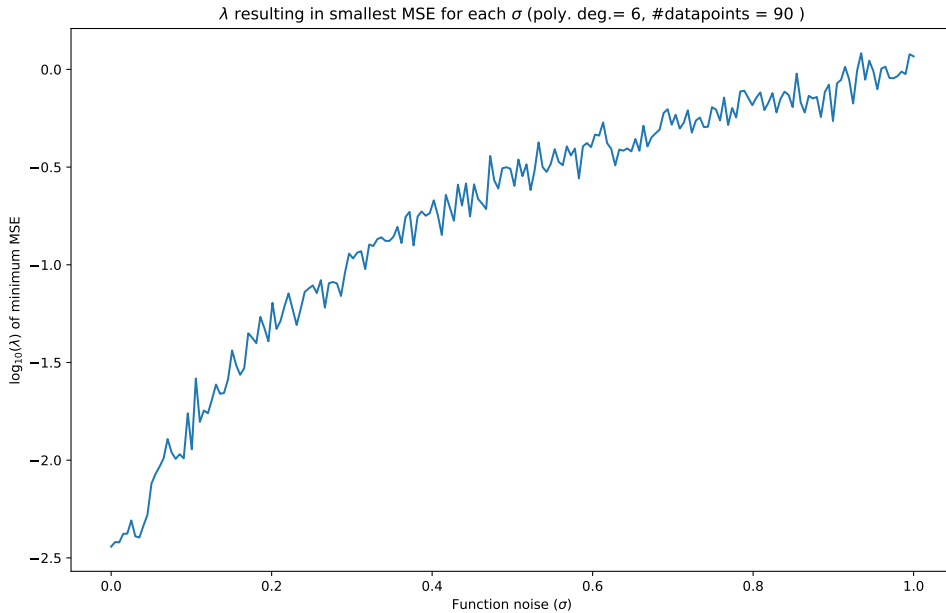


Figure 1: Performing Ridge regression on the Franke function (2), with a normally distributed error term $\sigma\epsilon$.

3.1.3 LASSO Regression

As with ridge, LASSO regression implements a new component to the cost function. Only, in this case the penalty is proportional to the L1- norm $\|\beta\|$:

$$C(\mathbf{X}, \beta) = \frac{1}{n} \{(\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta)\} + \lambda\|\beta\|.$$

One of the problems with ridge regression is that rather than eliminating covariant features, it will tend to shrink parameters, reducing their influence but not removing them. This is due to the behaviour of the squaring term in the L2- norm, which penalizes few large numbers more than many small numbers. For example, $\|[\beta_1, \beta_2]\|_2 \leq \|[\beta_1 + \beta_2, 0]\|_2$. Whereas in LASSO, the area of the L1- norm in the parameter space is more likely to intersect the least square loss function along the orthogonal axes in the parameter space, along which, at least one of the parameters are zero. This has a greater effect on eliminating features in a model, although which of two collinear parameters are eliminated is arbitrary.

3.2 Resampling Methods and Data Splitting

In order to get an approximate estimate of the quality of a fit, it is necessary that the fit is tested on a set that it has not been trained on. As the training set determines the parameters β , we need a set from the same distribution as the training set, which is different from the training set. For that reason, the full data set of predictors and targets is split into a test set and a training set. The test usually contains 20 – 25% of the full set, while the training set contains 75 – 80%. It is also possible to use a further validation set to do a final estimate using data that has never been "seen" by the programmer or the code, however, this was not implemented in this project.

3.3 Bias-variance relationship

In order to evaluate the quality of a fit, it is customary to evaluate the Mean Square Error (MSE) of the test data to the prediction. The Mean Square Error is defined as

$$MSE(\mathbf{y}, \tilde{\mathbf{y}}) = n^{-1} \sum_i^N (y_i - \tilde{y}_i)^2 = \mathbb{E}[(\mathbf{y} - \tilde{\mathbf{y}})^2]$$

where the left formula is the within-sample mean square error. The Mean Square Error can be decomposed into three parts: The unavoidable Error σ^2 , the squared Bias (which is the deviation of the mean value of an estimator and the true value), and the Variance (which is a measure of the spread of the estimator). Hence

$$MSE(\mathbf{y}, \tilde{\mathbf{y}}) = Bias[\tilde{\mathbf{y}}]^2 + Var[\tilde{\mathbf{y}}] + \sigma^2 = \left(f - \mathbb{E}(\tilde{f})\right)^2 + \mathbb{E}\left(\left[\mathbb{E}[\tilde{f}] - f\right]^2\right) + \sigma^2$$

A proof of this decomposition can be found in the appendix (3). In this article, we will however calculate the following values:

$$Bias[\tilde{\mathbf{y}}]^2 = n^{-1} \sum_i^N (y_i - avg(\tilde{\mathbf{y}}))^2$$

$$Var[\tilde{\mathbf{y}}]^2 = n^{-1} \sum_i^N (\tilde{y}_i - avg(\tilde{\mathbf{y}}))^2$$

Where $avg()$ stands for the average value. This is because the true values f_i are not known and thus impossible to obtain.

An important observation when fitting the data is that the Bias decreases as the model gets more complex, while the variance increases. When the Bias is large, important relations of the underlying data are not considered. When the Variance is high, small, unimportant features in the training set have too high of an impact. This is called overfitting. The aim is hence to find the model complexity that gives the lowest MSE. This is usually the case around the point where the Bias has the same magnitude as the Variance.

3.3.1 Bootstrap method

The Bootstrap method is a method that can be used to calculate statistical quantities. The main idea is the following: Let n be the size of the training set $Z = (z_1, z_2, \dots, z_n)$, where each z_i contains p predictors and one target value, that is $z_i = (x_{i,1}, x_{i,2}, \dots, x_{i,p}, y_i)$. Now, n data points z_i are drawn randomly from Z with replacement. This is done B times, creating B new sets Z'_j . We can now create a fitting model for each of the sets Z'_j . For each of these data sets Z'_j , we can now calculate any variable $\hat{\theta}_j$ given the data. Now, the distribution of $\hat{\theta}_j$ can be used to estimate the distribution of θ and get values of interest, such as the mean value or the variance. In this article, Bootstrap is used to get information about the test MSE, variance and bias in order to evaluate the quality of the fit. This is done by using the B bootstrap sets Z'_j to calculate the fitting parameters β and calculate B fits to the test set. \tilde{Y}_j . From these B fits, we can then approximate the test MSE, variance and bias.

TODO: Talk some more about the mathematical validity of Bootstrap.

3.3.2 K-fold Cross validation

K- fold Cross Validation is a useful technique for when you have a limited data set, and want to make sure that all features are captured in a training set. K- fold Cross Validation does this by partitioning the data set into k so called folds, assigning $k - 1$ folds as the training set and 1 of the folds as the testing set. It does this k times, making sure to utilize each fold as the testing set once, as seen in figure 2. In this way, we end up generating k different models, whose performance can be measured with respect to the testing fold or additionally, a separate testing set withheld from the partitioning process.

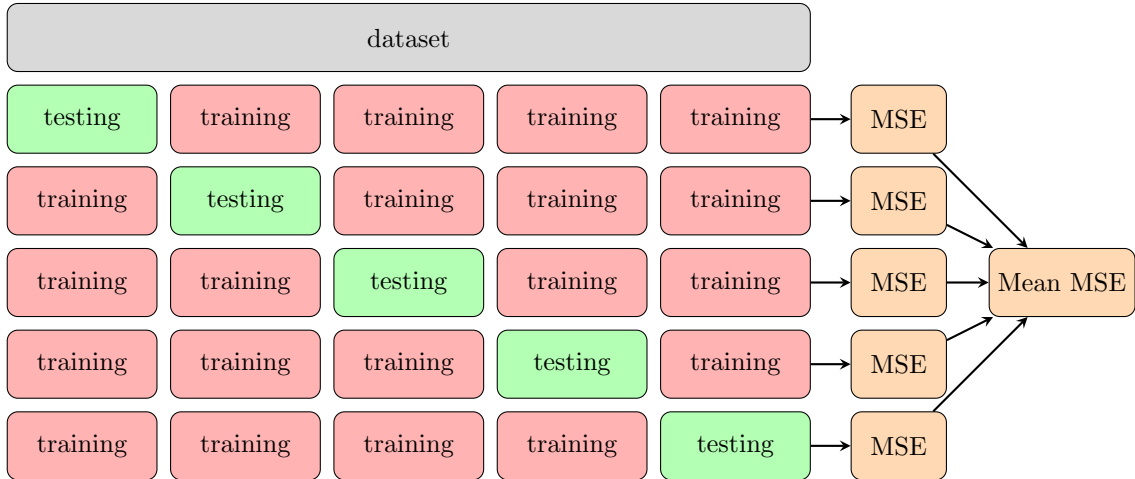


Figure 2: In this example, 5- fold Cross Validation is used on the entire dataset, with one model being generated for each row. The MSE on the test set is then calculated and averaged with the MSEs of every other model.

The value we choose for k will impact the accuracy of our model. To see this, we performed OLS regression on the same function (2) using different values for k resulting in figure 3.

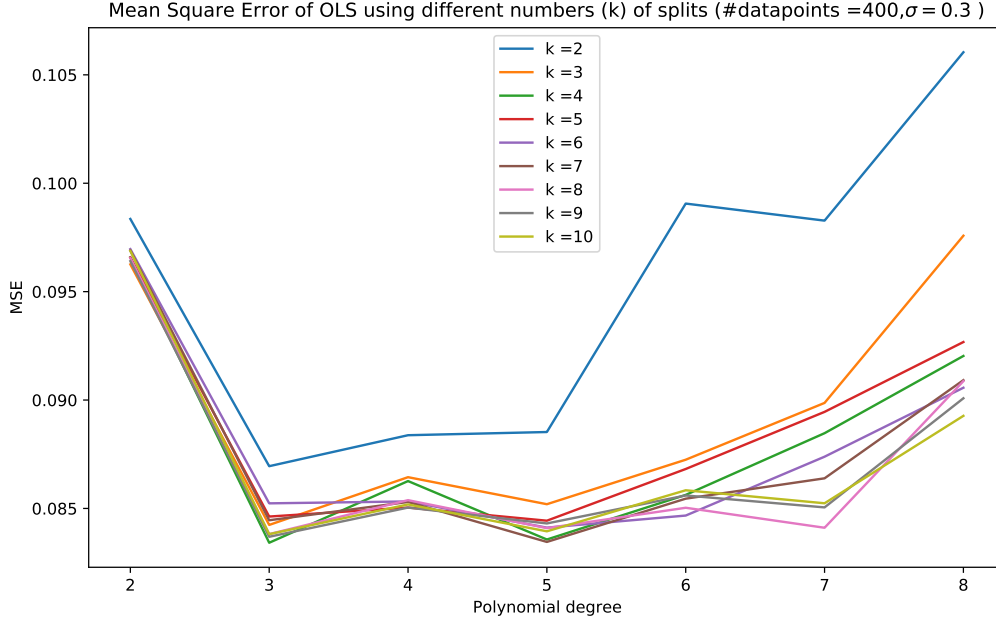


Figure 3: Estimated MSE as a function of model complexity (polynomial degree) for different values of k , using k - fold CV OLS on the Franke function (2). We can see that higher values of k consistently make lower error estimates though there is no clear correlation among the highest few values.

We expect to see this trend based on the increasing size of the training set from 50% to 75% aso. going from $k = 2$ to $k = 3$ and on.

3.4 Fitting the Franke function

The Franke Function is given by

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

and serves as a test function to check how well the implemented methods perform. We also add normally distributed noise in order to evaluate the behaviour of our fit.

3.5 Geographic Data

In order to test how well the Regression methods work with real data, we downloaded Geographic Data representing the Taebaek Mountains in South Korea. This has a surface area of approximately $3601 \times 3601 km^2$. We used a picture with resolution of 3601×3601 pixels, hence each square kilometer is represented as one pixel, where colour intensity represents the height (black equals height at sea level). This can be seen in figure 4.

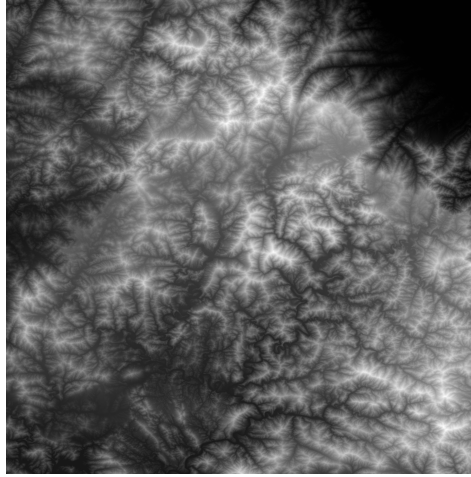


Figure 4: The Taebaek Mountain area in South Korea. The original image has a resolution of 3601×3601 pixels.

4 Computational Implementation

4.1 Finding the ideal parameter λ

In order to find the ideal parameter λ for both LASSO and Ridge Regression, we looped over a set of lambda-values in different magnitudes and estimated the MSE using 4-fold Cross Validation. We chose to use the parameter λ that gave the lowest MSE, averaged over two uncorrelated runs. This was repeated for each fit, that is, λ was recalculated for each new Design Matrix X . This is because the ideal λ value changes with model complexity.

4.2 Error in LASSO-Regression

As there is no way to find the parameters **beta** for LASSO-Regression analytically, numerical methods need to be used. Here, we used Scikit-Learn in order to solve the problem numerically. We have observed that the LASSO-method converges rather slowly. For this reason, we have chosen the tolerance to be equal to 0.03 for fitting the Franke function, while the number of maximum iterations is chosen to be 10^5 . We observed that decreasing the tolerance to 0.01 lead to no convergence even at 10^5 iterations. Similarly, we used XXX as tolerance at 10^5 maximum iterations for fitting the country data. As this is very time-consuming, we found it impossible to find the ideal parameters **beta** for LASSO-Regression. Each time the LASSO-fit was higher than the OLS fit, no ideal parameters have been found, as the ideal LASSO-fit cannot be no larger than the OLS-fit (because $\lambda = 0$ is the OLS fit).

4.3 Sampling Franke Function

We sample randomly points x_i and y_i between 0 and 1 and calculate the target points

$$z_i = f(x_i, y_i) + \epsilon_i$$

where ϵ_i follow a normal distribution with standard deviation σ . We used different values for σ , usually in the range $[0.01, 1]$, which has a profound impact on the fit.

4.4 The Design Matrix

The Design Matrix implemented is a 2D-version of the Vandermonde Matrix, hence, approximate z_i as, for a given degree k

$$\tilde{z}_i = \beta_0 + \beta_1 x_i + \beta_2 y_i + \beta_3 x_i^2 + \beta_4 x_i y_i + \dots + \beta_p y_i^k$$

In the practical implementation, we "drop" β_0 by removing the first column from the Design Matrix, as our data is scaled.

4.5 Scaling

The design matrix is scaled using Scikit-Learn's StandardScaler. While this is not necessarily necessary for OLS, scaling is needed in order to ascertain that Ridge Regression & LASSO Regression scale parameters equally, giving an overall better result.

5 Results

5.1 β standard deviation in OLS and Ridge regression

Using OLS and Ridge regression, for a given polynomial degree, the radius of the 95'th percentile confidence interval of the parameter β with the largest standard deviation was calculated and tabulated. The result is shown in table 1.

Table 1: Half the width of the 95'th percentile confidence interval for the estimators β in OLS and Ridge regression with the highest standard deviation, using the regularization parameter λ .

polynomial degree	OLS	Ridge ($\lambda = 0.1$)	Ridge ($\lambda = 1.0$)
1	0.005	0.014	0.013
2	0.023	0.062	0.047
3	0.128	0.207	0.077
4	0.632	0.250	0.086
5	2.977	0.278	0.091

It's clear that Ridge regression has the effect of tightening the confidence interval for the parameters β belonging to the highest order polynomials. It's also clear that increasing the regularization parameter λ has the expected effect of further tightening the confidence interval of the parameters. It's worth noting that while in OLS the values occupy four different orders of magnitude, but using Ridge, that decreases to two or one. The cost function imposing a penalty on the size of the β 's, forces every parameter to within a smaller subset of options.

5.2 Error comparison

5.2.1 Bootstrap vs. Cross validation

Bootstrap and K- fold cross validation is here used to estimate the error in our models. We're evaluating the error estimation of Bootstrap and K- fold cross validation in terms of MSE on the same dataset using a fixed value for σ , a fixed number of data points and fixed model complexity. The objective being to compare the two resampling methods, varying only the amount of folds used in K- fold cross validation. Using a dataset constructed from the Franke function, we get figure 5.

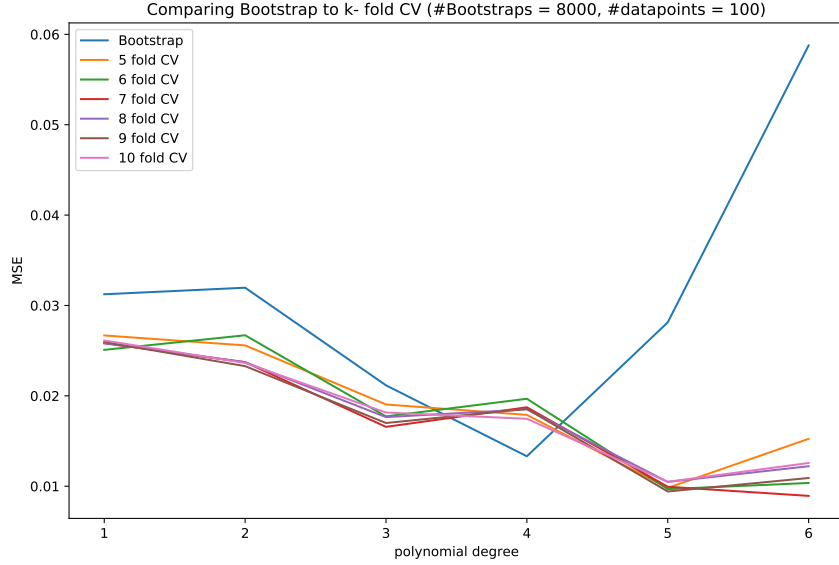


Figure 5: Estimated MSE as a function of model complexity (polynomial degree) for different values of k , using k - fold CV and Bootstrap OLS on the Franke function (2). We see that all iterations of K - fold cross validation makes a lower error estimate than 8000 bootstraps on this very limited data set of 100 entries

Trying again, this time with 10 times as many data points. We get figure 6

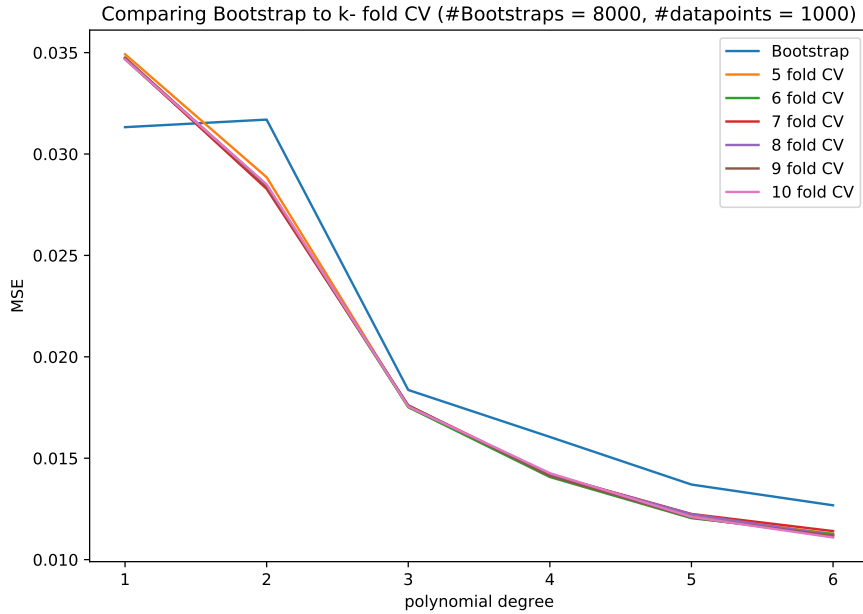


Figure 6: Estimated MSE as a function of model complexity (polynomial degree) for different values of k , using k - fold CV and Bootstrap OLS on the Franke function (2). Again, all iterations of K - fold cross validation makes a lower error estimate than 8000 bootstraps. However, increasing the number of data points made the difference significantly smaller, especially for the more complex models.

Figures 5 and 6 give the impression that on a sufficiently large dataset, Bootstrap gives comparable error estimates with K - fold cross validation. Bootstrap, however comes at a significantly higher computational cost which we've yet to show explicitly.

Using the same model as above, we run K- fold cross validation with $k = 10$ and compare the elapsed time with Bootstrap. The result is shown in figure 7.

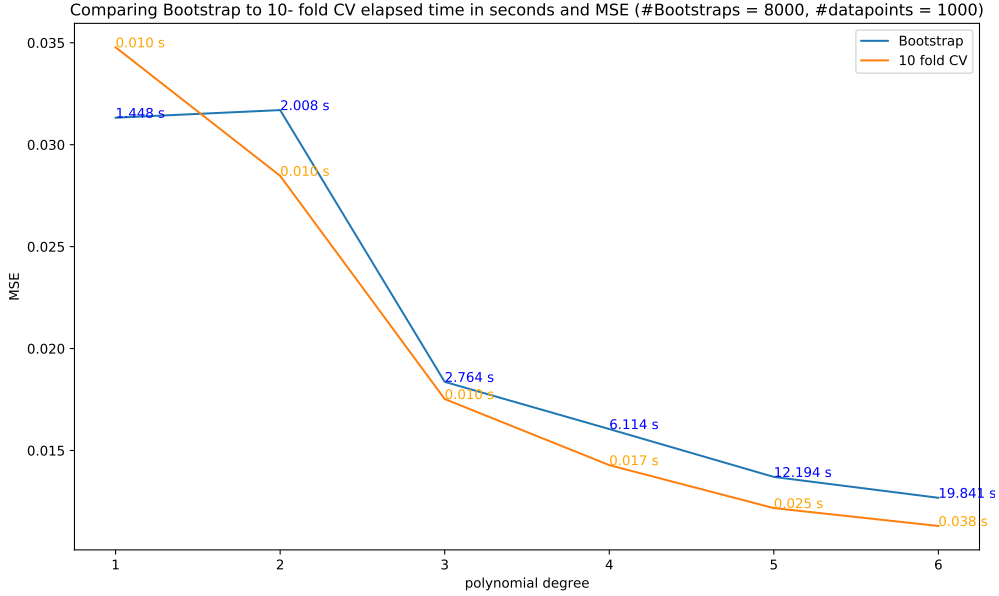


Figure 7: Estimated MSE as a function of model complexity for $k = 10$, using k- fold CV and Bootstrap OLS on the Franke function (2). The time elapsed for each model is plotted. Not only does K- fold CV take considerably less time, but the time increases much slower with increased model complexity compared to Bootstrap.

5.2.2 Franke Function

Using resampling methods, the MSE for OLS, Ridge- and LASSO regression on the Franke function were plotted against the model complexity.

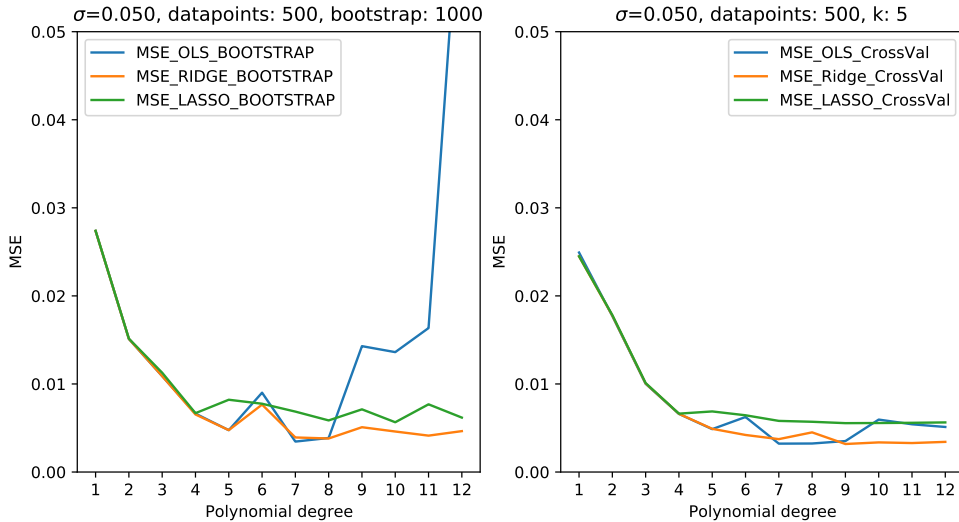


Figure 8: Test MSE for LASSO, Ridge and OLS regression. Estimated using Bootstrap (left) or k-Fold Cross validations. The exact parameters can be found in the graph's title.

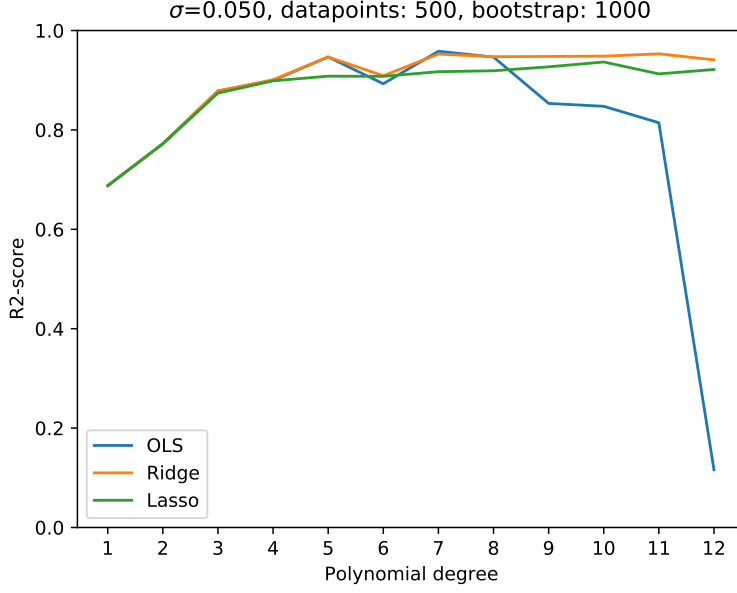
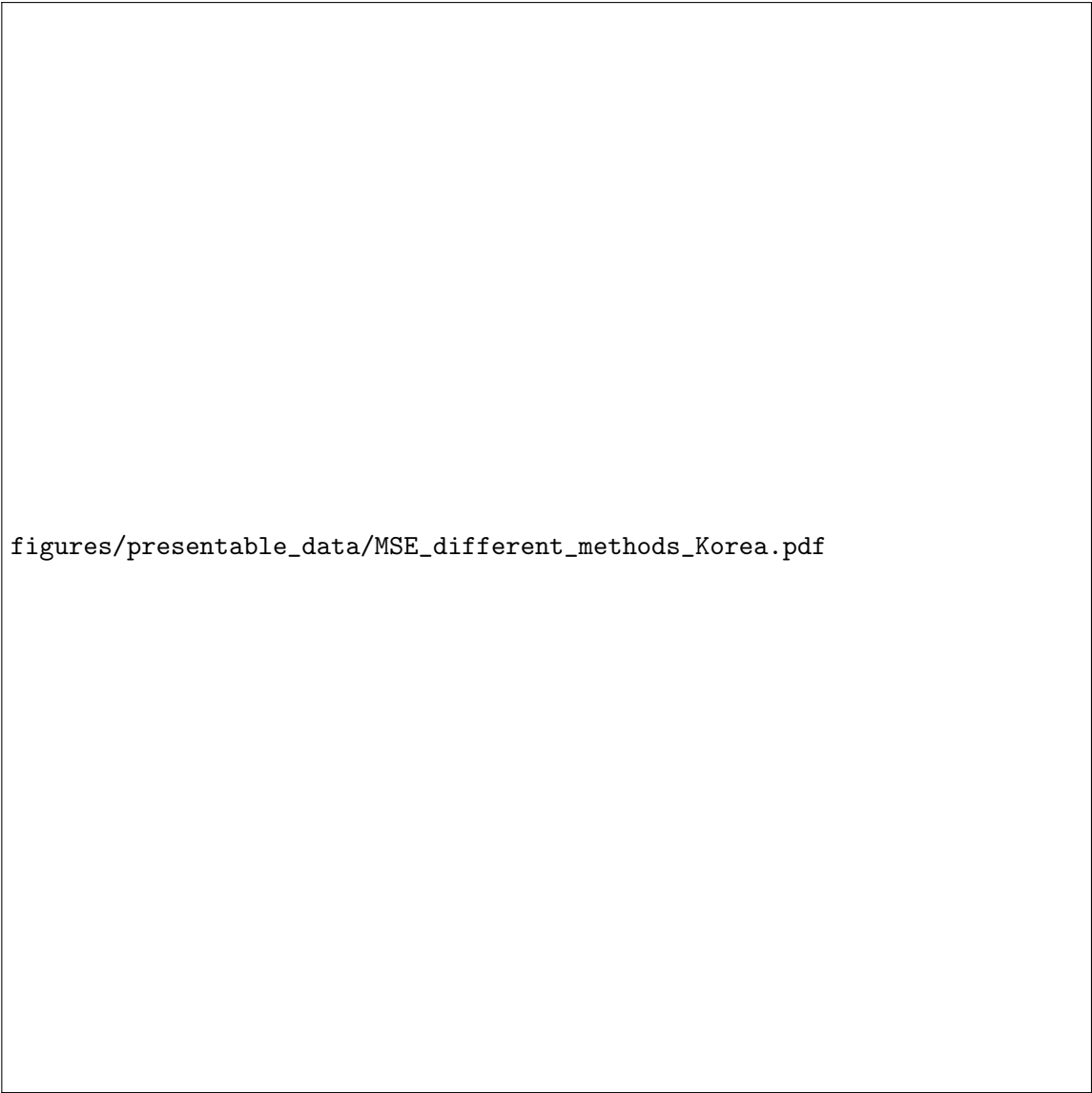


Figure 9: Test R^2 -value for LASSO, Ridge and OLS regression. Estimated using Bootstrap (left) or k-Fold Cross validations. The exact parameters can be found in the graph's title.

As expected from the theory, increasing the model complexity lead to a decreased Mean-Square-Error and an increased R2-Score for all three methods when fitting the Franke Function. However, as can be seen in figures 8 and 9, the mean square error for the Ordinary Least Square approximation suddenly raises sharply (and the R2 score drops sharply) after some complexity is reached: This behaviour cannot be observed with LASSO and Ridge. Neither LASSO nor Ridge give conceivably better results than OLS when looking at the complexity giving the best results, that is, with the chosen parameters, OLS, LASSO and Ridge regression are about equally fast. As OLS is the most basic method, easiest to implement and also fastest, we figured out that it is the most suitable method for this problem - adding more complexity while avoiding a high variance (Ridge or LASSO) does not improve the quality of the fit. Another interesting observation is that there is much more fluctuation in the MSE when using bootstrap than when using 5-fold Cross validation. We are not quite sure why this happens, but we assume that one possible cause is that when using bootstrap, some very "undesirable" combinations of values might be chosen, which leads to a very high Bootstrap. Thus, the graphs might look smoother when running more Bootstrap-samples. Another reason is that the training set in the Bootstrap method is not equally large as the training set when using 5-fold cross validation (75% vs 80%). We also observed that Bootstrap is more time-demanding than k-fold Cross Validation. As Bootstrap and Cross Validation are two very different methods to estimating a value of interest, the MSE in this case, different results are not surprising anyways. It would however be troublesome if the two methods gave different "ideal" models. The large deviation between LASSO and OLS is due to the aforementioned error.

5.2.3 Geographic data

Figure 10 shows the Mean Square Error of fitting the geographic data as a function of polynomial degree for both Bootstrap and K-Fold Cross Validation for all three methods. We have chosen a value λ which minimizes the MSE.



figures/presentable_data/MSE_different_methods_Korea.pdf

Figure 10: Test MSE for LASSO, Ridge and OLS regression. Estimated using Bootstrap (left) or k-Fold Cross validations. The exact parameters can be found in the graph's title.

The data looks similar to the data attained when fitting the Franke Function, just in a different magnitude of error: The error decreases sharply as the model complexity is increased, remains rather constant (or increases slowly) as it is increased more, and then raises again, sharply for OLS. The LASSO data is unreliable, as we have chosen to only include it for completeness reasons - we did not manage to reach a sufficiently converged value, hence the MSE is much larger than the minimal possible LASSO data.

5.3 Bias-Variance tradeoff

5.3.1 Franke Function

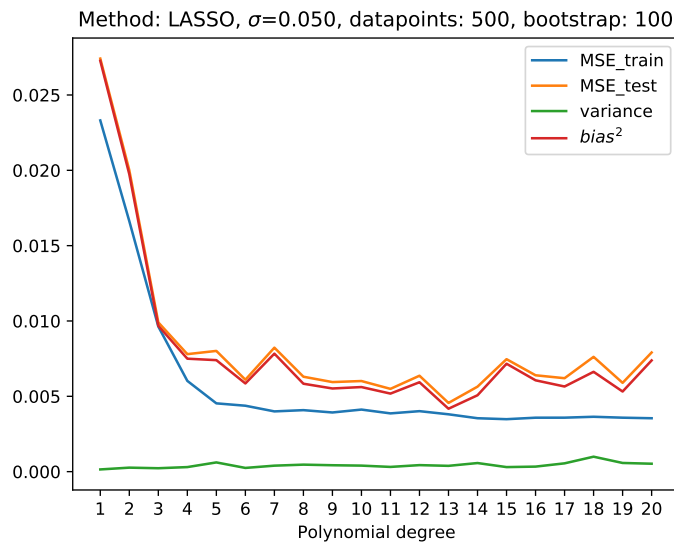


Figure 11: Bias, Variance, MSE & training MSE as a function of the polynomial degree using LASSO regression. The values for σ , the size of the dataset and the number of bootstraps can be found in the title.

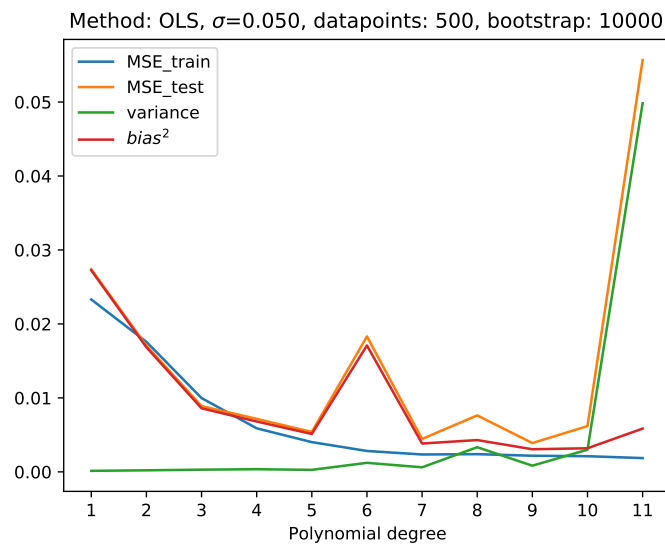


Figure 12: Bias, Variance, MSE & training MSE as a function of the polynomial degree using OLS regression. The values for σ , the size of the dataset and the number of bootstraps can be found in the title.

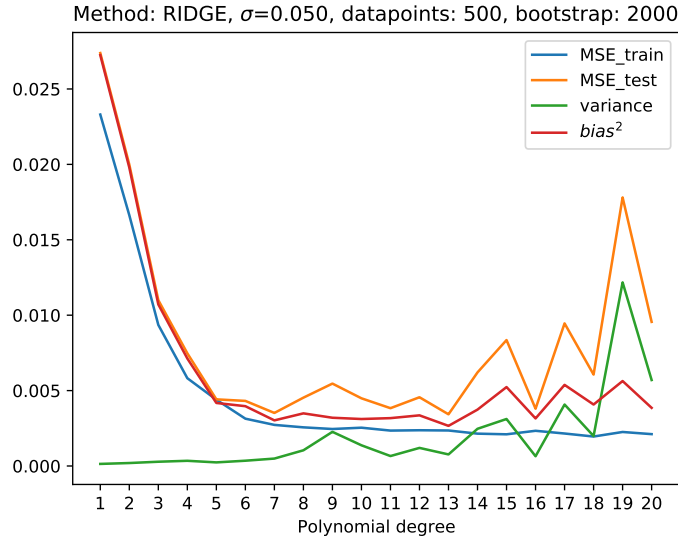


Figure 13: Bias, Variance, MSE & training MSE as a function of the polynomial degree using Ridge regression. The values for σ , the size of the dataset and the number of bootstraps can be found in the title.

Figures 11, 12 and 13 show the MSE, the Training-MSE, the Variance, and the bias for the Franke Function. In all figures, the aforementioned statement that the variance increases and the bias decreases with increasing model complexity, hold true. Especially for the OLS, we see that this resembles figure 2.11 on page 38 in Hastie et al. by a lot[1]. One can see that OLS is more prone to a higher variance at higher polynomial degrees, as is expected - shrinkage methods explicitly reduce the variance. There seems to be especially little variation in the LASSO approach, however, it is possible that this is due to the aforementioned issue with LASSO regression, as we have not found the ideal LASSO-coefficients due to the large error.

5.3.2 Geographic data

Figure 14 shows the bias, the variance and the MSE as a function of polynomial degree for Ridge and OLS Regression. We have chosen a value λ which minimizes the MSE. We decided not to include LASSO as the LASSO data we attained is not the true LASSO data.



figures/presentable_data/Bias_variance_Korea.pdf

Figure 14: Bias, Variance, MSE & training MSE as a function of the polynomial degree using OLS and Ridge regression. The size of the dataset and the number of bootstraps can be found in the title.

Again, one can see how the Bias decreases more or less monotonously for both methods, until it eventually stagnates, while the variance increases later on. However, the increase in variance is not as drastic for Ridge regression, where it goes up, but never above 10.000, while it explodes for OLS, which fits with our expectations.

5.4 Evaluation of the Geographic Fit

Figure 15 shows scaled versions of the original data using the Ridge fit giving the minimal error, the OLS fit giving the minimal error, as well as the scaled original picture.

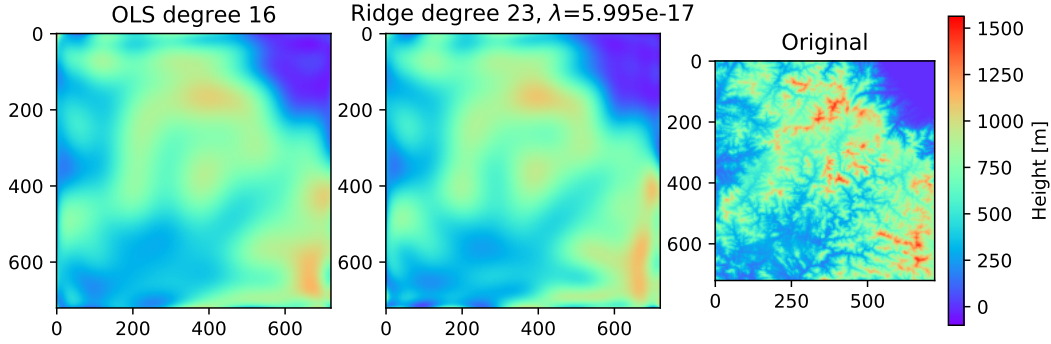


Figure 15: Fit of the Geodata with parameters that reduce the MSE using Ridge Regression and OLS Regression. The original (scaled) image is included as well. 5000 data points were used.

Just from this picture, we can see a lot. Both the OLS and the Ridge fit manage to fit the area in such a way that is generally correct - the large mountain ranges and the deep valleys are included, the sea is the "deepest" part, and of course, the dimensions are correct. However, there is a profound lack of detail both in the OLS and in the Ridge fit. Even the "easy" part (the dark area, which symbolizes the sea) is not very well fitted (it goes below zero), which lies in the nature of a polynomial fit. This is rather frustrating, as both creating these images and finding the ideal parameters is very time-consuming. Increasing the data points would most likely give a better result, but by the fact that the MSE does not go well below 20.000, it is unlikely that there is a huge improvement. A MSE of 20.000 means that the minimal root MSE is around 140, which, while different from the Mean Absolute Error, still gives an indication on the "average difference" between the original data and our fit. All of this combined shows that Regression is not the ideal way to match such complex data where there does not seem to exist a functional relationship with little error.

6 Conclusion

Our analysis of methods has given a lot of insight. First of all, fitting a "simple" function such as the Franke function, can be done with little loss using any of the methods, especially when the underlying error is small. Sadly, this cannot be said about complex data such as the geographic data, where it was not possible to predict more than the most basic traits due to the nature of the data, making it look like a picture from the late 18th century. OLS seems to be sufficient for both problems, as it reached an MSE on par with Ridge regression. If this is the case, it is the ideal method, as it is computationally cheap. However, OLS is prone to overfitting, which is not the case with Ridge regression. LASSO is a very expensive method due to the lack of an analytical solution, and in this analysis, it has been prohibitively expensive. Finding the parameter λ by iterating over all values is very expensive, and using more clever methods is necessary in large-scale projects.

Error estimation is no easy task, and two of the most used methods give different error estimates - while they are in the same magnitude, they do differ largely, especially where there are steep changes.

7 Appendix

7.1 Proof that the MSE can be decomposed into Bias and Variance

$$\begin{aligned} MSE(y, \tilde{y}) &= \mathbb{E} [(y - \tilde{y})^2] \\ &= \mathbb{E} [(f + \epsilon - \mathbb{E}[\tilde{y}] + \mathbb{E}[\tilde{y}] - \tilde{y})^2] \\ &= \mathbb{E} [(f + \epsilon - \mathbb{E}[\tilde{y}])^2 + (\mathbb{E}[\tilde{y}] - \tilde{y})^2 + 2\mathbb{E}[(f + \epsilon - \mathbb{E}[\tilde{y}])(\mathbb{E}[\tilde{y}] - \tilde{y})]] \\ &= \mathbb{E} [\epsilon^2] + \mathbb{E} [(f - \mathbb{E}[\tilde{y}])^2] + \mathbb{E} [\mathbb{E}[\tilde{y}] - \tilde{y})^2] \\ &\quad + 2\mathbb{E}[\epsilon] \mathbb{E}[f - \mathbb{E}[\tilde{y}]] + 2\mathbb{E}[(f - \mathbb{E}[\tilde{y}])(\mathbb{E}[\tilde{y}] - \tilde{y})] + 2\mathbb{E}[\epsilon(\mathbb{E}[\tilde{y}] - \tilde{y})] \\ &= \sigma^2 + Bias(\tilde{y})^2 + Var(\tilde{y}) \\ &\quad + 2\mathbb{E}[\epsilon] \mathbb{E}[f - \mathbb{E}[\tilde{y}]] + 2\mathbb{E}[(f - \mathbb{E}[\tilde{y}])(\mathbb{E}[\tilde{y}] - \tilde{y})] + 2\mathbb{E}[\epsilon(\mathbb{E}[\tilde{y}] - \tilde{y})] \\ &= \sigma^2 + Bias(\tilde{y})^2 + Var(\tilde{y}) \\ &\quad + 2\mathbb{E}[\epsilon] \mathbb{E}[f - \mathbb{E}[\tilde{y}]] + 2(f - \mathbb{E}[\tilde{y}])\mathbb{E}[(\mathbb{E}[\tilde{y}] - \tilde{y})] + 2\mathbb{E}[\epsilon] \mathbb{E}[(\mathbb{E}[\tilde{y}] - \tilde{y})] \\ &= \sigma^2 + Bias(\tilde{y})^2 + Var(\tilde{y}) \end{aligned} \tag{3}$$

where we used that $\mathbb{E}[\epsilon] = 0$, that $(f - \mathbb{E}[\tilde{y}])$ simply is a number (f is deterministic) and that ϵ and \tilde{y} are independent variables (hence the expected value can be separated), as well as the fact that $\mathbb{E}[\mathbb{E}[\tilde{y}] - \tilde{y}] = \mathbb{E}[\tilde{y}] - \mathbb{E}[\tilde{y}] = 0$

References

- [1] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. Springer Series in Statistics. Springer New York Inc., New York, NY, USA, 2001.
- [2] Wessel N. van Wieringen. Lecture notes on ridge regression, 2015.