

Regression analysis and resampling methods

Nicholas Karlsen and Thore Espedal Moe

University of Oslo

(Dated: October 9, 2020)

In this paper we aim to explore the application of three distinct regression methods to both a sampling of the Franke function, as well as a set of terrain data. We utilize the ordinary-least-square (OLS) regression, Ridge regression and LASSO regression along with bootstrap and k-fold cross-validation re-sampling techniques in order to create and asses predictive polynomial models. We analyze the performance of the different regression models on the two data sets, with the ultimate goal of determining the best regression method and the best predictive polynomial model for each data set.

I. INTRODUCTION

In essence, Linear Regression is the process of taking points from a function, or a set of measurements, and mapping them linearly to coordinates in a chosen basis in order to create an approximation, or model, explaining the original dataset and estimating unseen points in the domain spanned by the original data set. That is, from a limited set of data try to infer a linear functional relationship between some chosen prediction variables and the measured/sampled responses, and use this functional relationship to predict new data points in the domain. Typically, one is either interested in extracting the functional relationship between the measured responses and the prediction variables in order to say something about the relative importance of the conjectured predictors; or one wants to create a "best possible" predictive map of the (mostly unmeasured) response variables as a function of the (mostly unmeasured) prediction variables. In both cases what one really wants is some form of optimal mapping from prediction variables to response variables, and that mapping can be estimated by the use of linear regression methods.

Obviously, this way of doing predictions and inferences has an enormous range of possible applications; in basically every instance where one expects some linear relationship between selected prediction variables and a response variable one can, if one has a data set containing measurements of the variables, try to use regression methods to create an optimal model. Examples might be the construction of a model for housing prices as a function of distance from city center, living area, building year and whatever else one could imagine might affect the prices. Or one could wish to determine the best coefficients in a model for nuclear binding energy as a function of nucleon numbers.

In this paper we investigate two illustrative and similar applications: The reconstruction of the Franke function from sampled values, and the (re-)construction of an elevation map from a downsampled set of terrain data. In a way, our data sets are "low-resolution images" of the Franke function and the terrain. Our goal is then to find the best possible "high-resolution images" by using different regression methods and resampling techniques on the sampled data, using tow-dimensional polynomi-

als as basis functions. In other words, we want to find the most faithful recreation of the original terrain data and the Franke function, in the basis of two-dimensional polynomials.

We begin with a theoretical discussion of the regression and resampling methods that will be used, as well as some comments on the significance of the so-called Bias-Variance-Tradeoff in the context of linear regression. We then proceed to investigate the application of ordinary-least-squares (OLS), Ridge regression and LASSO regression to samplings of the Franke function. We use k-fold cross-validation to assess the mean-squared-error (MSE) of our different models, and we use bootstrap resampling in order to estimate the bias and variance of our models. We discuss the influence of the number of points sampled, the polynomial degree we try to fit, the values of the hyper-parameter lambda for the Ridge and LASSO regressions and the presence/absence of noise in our samplings. Afterwards we will move on to the terrain data, where we perform much of the same analysis as for the Franke function. We then try to evaluate which regression methods perform the best for different use cases, and which resulting models might be considered the 'best'.

II. THEORY

A. Linear Regression

Consider a set of data points $\{(x_1, y_1), \dots, (x_N, y_N)\}$ which we wish to fit to some linear model $\tilde{\mathbf{y}}(\boldsymbol{\beta})$ where $\boldsymbol{\beta}$ is a vector containing the free parameters of the model. The model $\tilde{\mathbf{y}}$ will then be related to the true data points \mathbf{y} by

$$\mathbf{y} = \tilde{\mathbf{y}}(\boldsymbol{\beta}) + \boldsymbol{\varepsilon} \quad (1)$$

where $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_N)$ represents the error of the model. The aim of Linear regression models is thus to find the optimal parameters $\boldsymbol{\beta}$ such that not only the error of the model on some particular dataset is minimized, but also the error of the model applied to different data sets sampled in the same domain is minimized such that we may use our model to make predictions about new datasets.

1. The Design Matrix

When constructing a model, we first have to choose a set of basis functions. A popular choice for which is \mathbb{P}_n , which models a wide range of different phenomena and is also the choice we will make for this paper. But in principle, any set of linear basis functions may be chosen. For a single variate polynomial, a model would then take the form

$$\tilde{y}(x) = \beta_0 + \beta_1 x + \dots + \beta_n x^n \quad (2)$$

When trying to find the optimal β , a useful construct is the so-called Design Matrix, which for a polynomial basis takes the form

$$X = \begin{bmatrix} 1 & x_1 & \dots & x_1^n \\ 1 & x_2 & \dots & x_2^n \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_N & \dots & x_N^n \end{bmatrix} \quad (3)$$

where we see that entry X_{ij} contains the x_i data point evaluated in the j -th basis function. This matrix has the particularly useful property that we may then multiply it with β to obtain the corresponding set of predictions $\tilde{\mathbf{y}}$ like

$$\tilde{\mathbf{y}} = X\beta \quad (4)$$

leaving us with a linear algebra problem to solve in finding the β which optimizes $\tilde{\mathbf{y}}$.

This also extends directly to multivariate cases where we have more than one independent variable. In particular for two dimensional polynomials, our choice of basis for this article, we have that an n -th degree polynomial is constructed by all permutations of $x^p y^q$ where $p + q \leq n$ which means that the total number of predictors for degree n is given by

$$\binom{2+n}{n} = \frac{(2+n)!}{n!(2+n-n)!} = \frac{(n+1)(n+2)}{2} \quad (5)$$

2. Model assessment

Once we have constructed a model, we also need a way to quantitatively evaluate its performance. Here, we will focus on two statistical measures; First, we have the *mean squared error* (MSE) defined as

$$\text{MSE}(\mathbf{y}, \tilde{\mathbf{y}}) = \frac{1}{n} \sum_i (y_i - \tilde{y}_i)^2 \quad (6)$$

which as the name suggests is a measure of the mean square distance the model $\tilde{\mathbf{y}}$ has from the data set \mathbf{y} . We also have the *coefficient of determination*, defined as

$$R^2(\mathbf{y}, \tilde{\mathbf{y}}) = 1 - \frac{\sum_i (y_i - \tilde{y}_i)^2}{\sum_i (y_i - \mathbb{E}[\mathbf{y}])^2} \quad (7)$$

which loosely speaking is a measure of well the proposed model predicts the variance of the data set, where $R^2 = 1$ corresponds to a perfect fit and $R^2 < 1$ an increasingly worse fit. It is worth to note that the MSE and R^2 are essentially equivalent when measuring the performance of a model on a particular sample, which is ultimately determined by the $\sum_i (y_i - \tilde{y}_i)^2$ term in both cases. Thus they differ mainly in scaling & centering. However, R^2 will notably yield results which translates somewhat more directly across different models.

A particularly important point to make when considering these measures is that they will both generally favor increasingly higher degree polynomials when looking at some particular data set. Whilst yielding good MSE & R^2 scores this leads to over-fitting, where the model becomes too specialized to one particular sample which means that the overall predictive capabilities of the model actually decreases. We will look closer at the details of this in a later section. For now, we simply note that this problem motivates us to split our data sets into training and test sets. Where we train our model using the training set, and subsequently validate the model by computing the MSE & R^2 scores on the test sets. Thus, if our model becomes too specialized to the training set, this will lower MSE & R^2 scores in the test set, giving us a much better indication on the predictive capabilities of our model.

3. Ordinary Least Squares

In ordinary least squares (OLS), we aim to find an optimal set of parameters $\hat{\beta} = [\hat{\beta}_0, \dots, \hat{\beta}_n]^T$ such that the L^2 norm $\|\mathbf{y} - X\beta\|_2$ is minimal, with the L^2 norm being induced by the inner product

$$\|\mathbf{u}\|_2^2 = \sum_i u_i^2 = \mathbf{u}^T \mathbf{u} \quad (8)$$

This defines the cost function for OLS, which may be written as

$$C_{OLS}(\beta) = (\mathbf{y} - X\beta)^T (\mathbf{y} - X\beta) \quad (9)$$

In order to find the minima, we differentiate wrt to β and assert that $\partial_{\beta} C_{OLS} = 0$ for the optimal predictor. Taking the partial derivative yields

$$\frac{\partial}{\partial \beta} C_{OLS}(\beta) = -2X^T (\mathbf{y} - X\beta) \quad (10)$$

We assert this is zero at the minima, which yields

$$X^T \mathbf{y} = X^T X \beta \quad (11)$$

then taking the inverse of $X^T X$ on both sides then gives the optimal β as

$$\hat{\beta}_{OLS} = (X^T X)^{-1} X^T \mathbf{y} \quad (12)$$

which mathematically is the best projection of the data-points to our model. It can also be shown that the variance of these optimal predictors are given by [1]

$$\text{Var}(\hat{\beta}_{\text{OLS}}) = \sigma^2 (X^T X)^{-1} \quad (13)$$

Where σ^2 denotes the variance of the underlying noise of the sample.

While OLS yields the mathematically best model for some particular sample, it does not necessarily yield the best predictive model. OLS may suffer from over-fitting to the particular sample, and will therefore not always perform very well on new data points for which the model has not been trained. As such, we instead look at alternate regression methods with additional tuning parameters which allow us to account for the variance between different sets of data as to yield a model with better predictive abilities compared to the OLS.

4. Ridge Regression

One such method is the Ridge regression where we instead try to minimize $\|\mathbf{y} - X\beta\|_2 + \lambda\|\beta\|_2$, from which we define our cost function as

$$C_R(\beta) = (\mathbf{y} - X\beta)^T (\mathbf{y} - X\beta) + \lambda\beta^T \beta \quad (14)$$

Similar for what we did for OLS, we take the partial derivative wrt. β , where the only difference in finding the minima resides in the additional term

$$\frac{\partial}{\partial \beta} \lambda\beta^T \beta = 2\lambda\beta \quad (15)$$

if we again assert that the optimal β is given by $\partial_{\beta} C_R(\beta)$ and rearrange we get

$$\hat{\beta}_R = (X^T X + \lambda)^{-1} X^T \tilde{\mathbf{y}} \quad (16)$$

As the optimal β for Ridge regression. So we see that the only difference between OLS and Ridge regression is the addition of the "penalty" parameter λ to the L^2 norm of the coefficients β . The main point of introducing this parameter is to reduce the variance of the regression coefficients β . It introduces a constraint on the allowable values of β , which means that the method no longer is unbiased. The aim is that the reduction of the variance outweighs the increase of the method's bias, leading to an overall lower test error, cf. the later discussion of the Bias-Variance Trade-off. As a technical note, during our computations we center our response-variable and scale plus center our predictors. This has the effect of removing the constant column in the design matrix, so that penalty term doesn't apply to β_0 .

5. LASSO Regression

The cost function for the *least absolute shrinkage and selection operator* (LASSO) method is defined as

$$C_L(\beta) = \|\mathbf{y} - X\beta\|_2^2 + \lambda\|\beta\|_1 \quad (17)$$

Which as opposed to OLS and Ridge has no analytical solution for finding the optimal β , which is instead found by optimization methods such as gradient descent. The specifics of this is beyond the scope of this article, and we will simply use the existing libraries provided by scikit-learn [2] to perform the regression.

Qualitatively the LASSO method is distinguished from the Ridge regression by how the norm of β is penalized. While the Ridge regression applies the penalty to the L^2 norm, the Lasso aims to shrink the L^1 norm. Two important consequences arise from this: firstly, it is no longer possible to obtain a closed-form solution to the minimization problem; secondly, the regression coefficients may be shrunk far more unevenly. In contrast to the Ridge regression, LASSO regression may shrink individual β_i to zero. This is an extremely useful property of the method, since it allows unimportant predictors to be identified and discarded. The drawback, of course, is that the bias of the Lasso method is expected to exceed the bias of the Ridge method. Furthermore, for highly correlated prediction variables, like our two-dimensional polynomials, the LASSO method might struggle with which of those prediction variables to keep. Once again a technical note; by centering we remove the constant column from the design matrix to avoid penalizing the constant term through β_0 .

B. Singular Value Decomposition

Consider an $m \times n$ matrix X of rank r . The singular values of X are then defined as the root of the eigenvalues of the diagonal matrix $X^T X$. We may also express X in the so-called singular value decomposition (SVD)

$$X = U \Sigma V^T \quad (18)$$

where U , V are unitary matrices and

$$\Sigma = \begin{bmatrix} D & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix} \quad (19)$$

an $m \times n$ matrix, where the matrix D is a diagonal $r \times r$ matrix containing the singular values of $X^T X$

$$D = \begin{bmatrix} \sigma_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_n \end{bmatrix} \quad (20)$$

which by convention is ordered such that $\sigma_1 \geq \dots \geq \sigma_n$.

1. Application to OLS

We may use the SVD to re-express the expression for $\hat{\beta}$ in OLS given by Eqn. 12 by writing

$$X^T X = (U \Sigma V^T)^T (U \Sigma V^T) = V \Sigma^T U^T U \Sigma V^T \quad (21)$$

Since U is unitary, it follows that $U^T U = \mathbb{I}$, further we have that $\Sigma^T = \Sigma$ since it is a diagonal matrix. We therefore end up with

$$X^T X = V \Sigma^2 V^T \quad (22)$$

inserting this into Eqn. 12 gives us

$$\hat{\beta}_{OLS} = (V \Sigma^2 V^T)^{-1} V \Sigma U^T y \quad (23)$$

This gives an alternate way of solving the OLS problem, which becomes particularly useful in situations where we have a large amount of data sampled from a relatively small domain, increasing the chance of X having linearly dependent columns which in turns leads to X , and by extension $X^T X$ no longer being invertible.

C. Re-sampling

Re-sampling methods are ways in which we can generate new statistics from our existing data, which as the name suggests implies sampling new data sets from our already existing data. By doing so, we may gain new insights about our data which may not be available through regular analysis, particularly in situations where we are limited by the number of data points.

Here, we will focus on two of many such techniques.

1. Cross Validation

In the cross-validation re-sampling method, we split our data set S into k equally sized subsets s_1, \dots, s_k . We then for each $i = 1, \dots, k$ assign the i -th subset as the test set and the remaining $k-1$ subsets as the training set and compute the statistics in the usual way. Then at the end, we compute the mean value of the k sets of statistics. A visual representation of this process can be seen in Fig. 1. When doing cross-validation, typical choices of k are 5 and 10 [1]. Which one is better will depend on how the error scales with the size of the training set, as such, choosing a suitable k requires some analysis. The cross-validation re-sampling provides a good estimate for the mean error of our estimates.

- Discuss this in more detail $\rightarrow \Delta Err$ wrt to number of data points

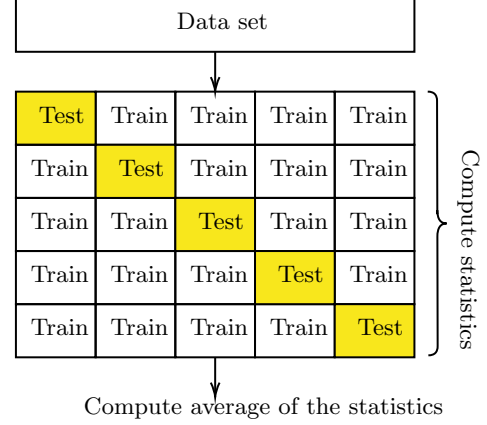


FIG. 1: Visual representation of k -fold cross-sampling for $k = 5$

2. Bootstrap

In the bootstrap re-sampling method, we sample our data set $S = \{s_1, \dots, s_N\}$ N -number of times, in particular, we allow sampling the same s_i multiple times. In this way, we generate new datasets in which some points are under-weighted and others overweighted with respect to the original dataset S . Loosely speaking, the concept corresponds to evenly sampling the ensemble of all possible input data sets to our regression method, in order to estimate the variance and expectation values of our model predictions over the ensemble. The idea is that if our training set representatively covers the domain of inputs, each re-sampled regression computation will give a certain prediction \tilde{y}_B and these \tilde{y}_B will appear with a frequency corresponding to their underlying probability distribution. When we later try to compute statistics of the predictions, we will then be able to use point-estimators on the computed \tilde{y}_B to obtain estimates of, for instance, the variance of the predictions or the mean value of the predictions over the ensemble of all possible inputs in the domain spanned by our data.

D. The Bias-Variance Trade-off

A key concept in much of machine learning is the so-called Bias-Variance Trade-off. Simply put the test error of a learned prediction method can be viewed as being partly the sum of two distinct quantities, the bias of the method and the variance of the method. The bias can be thought of as errors due to inflexible assumptions and constraints on the model, while the variance, quite opposingly, is the actual statistical variance of the procedure. Crucially, the variance is connected with the model having too much flexibility, over-fitting on the training data while giving very spread results on the test data. Naturally, the relative sizes of the errors due to bias and

the errors due to variance change with the flexibility/complexity of the model. For an inflexible model, e.g. a polynomial regression model of low degree, the bias is expected to be the dominant contribution to the test error. On the other side of the spectrum, for a very flexible model, e.g. a polynomial regression model of high degree, the variance would be anticipated to be the dominating source of the test error. Importantly, the relative effects of the bias and the variance won't necessarily change at the same rate. Thus, one could hope to find an optimally flexible model which minimizes the combined error due to variance and bias. This is the essence of the Bias-Variance Trade-off; one could, for instance, try to slightly increase the bias of a method in the hopes of reducing the variance far more.

Concretely for the problem at hand in this article, we use the expected MSE of our models on the test data to assess our models. This expected MSE can be explicitly decomposed into a bias term, a variance term and an irreducible error [1]. We assume that our response values \mathbf{y} are given by a "true" function \mathbf{f} with the addition of normally distributed noise $\boldsymbol{\varepsilon}$ with zero mean:

$$\mathbf{y} = \mathbf{f} + \boldsymbol{\varepsilon} \quad (24)$$

Denoting our model predictions by $\tilde{\mathbf{y}}$ and keeping our "true" test values \mathbf{f} fixed, we can then decompose the expected squared test error, the expectations being taken over the ensemble of possible predictions and noise values, as:

$$\begin{aligned} & \mathbb{E}[(\mathbf{y} - \tilde{\mathbf{y}})^2] \\ &= \mathbb{E}[(\mathbf{y} - \mathbb{E}[\tilde{\mathbf{y}}]) - (\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])]^2 \\ &= \mathbb{E}[(\mathbf{y} - \mathbb{E}[\tilde{\mathbf{y}}])^2] + \mathbb{E}[(\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])^2] \\ &\quad - \mathbb{E}[2(\mathbf{y} - \mathbb{E}[\tilde{\mathbf{y}}])(\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])] \\ &= \mathbb{E}[(\mathbf{y} - \mathbb{E}[\tilde{\mathbf{y}}])^2] + \text{Var}[\tilde{\mathbf{y}}] - 2(\mathbf{y} - \mathbb{E}[\tilde{\mathbf{y}}])\mathbb{E}[(\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])] \\ &= \mathbb{E}[(\mathbf{y} - \mathbb{E}[\tilde{\mathbf{y}}])^2] + \text{Var}[\tilde{\mathbf{y}}] - 2(\mathbf{y} - \mathbb{E}[\tilde{\mathbf{y}}])(\mathbb{E}[\tilde{\mathbf{y}}] - \mathbb{E}[\tilde{\mathbf{y}}]) \\ &= \mathbb{E}[(\mathbf{y} - \mathbb{E}[\tilde{\mathbf{y}}])^2] + \text{Var}[\tilde{\mathbf{y}}] = \text{Bias}^2[\mathbf{y}, \tilde{\mathbf{y}}] + \text{Var}[\tilde{\mathbf{y}}] \\ &= \mathbb{E}[(\mathbf{f} + \boldsymbol{\varepsilon} - \mathbb{E}[\tilde{\mathbf{y}}])^2] + \text{Var}[\tilde{\mathbf{y}}] \\ &= \mathbb{E}[(\mathbf{f} - \mathbb{E}[\tilde{\mathbf{y}}])^2] + \text{Var}[\tilde{\mathbf{y}}] + \text{Var}[\boldsymbol{\varepsilon}] \\ &= \text{Bias}^2[\mathbf{f}, \tilde{\mathbf{y}}] + \text{Var}[\tilde{\mathbf{y}}] + \sigma^2 \end{aligned} \quad (25)$$

In a convenient abuse of notation, the above calculation is meant to be performed element-wise. What has been computed above is the expected squared test error for each data point in the test set separately; for notational convenience stacked up as vectors. To get a total MSE for the whole test set, one simply must take the mean

over the vector elements afterwards. This yields:

$$\begin{aligned} \text{MSE} &= \frac{1}{n} \sum_i^n \mathbb{E}[(y_i - \tilde{y}_i)^2] \\ &= \frac{1}{n} \sum_i^n (\mathbb{E}[(y_i - \mathbb{E}[\tilde{y}_i])^2] + \text{Var}[\tilde{y}_i]) \\ &= \frac{1}{n} \sum_i^n (\mathbb{E}[(f_i + e_i - \mathbb{E}[\tilde{y}_i])^2] + \text{Var}[\tilde{y}_i]) \\ &= \frac{1}{n} \sum_i^n (\mathbb{E}[(f_i - \mathbb{E}[\tilde{y}_i])^2] + \text{Var}[\tilde{y}_i] + \text{Var}[e_i]) \\ &= \frac{1}{n} \sum_i^n ((\mathbb{E}[(f_i - \mathbb{E}[\tilde{y}_i])^2] + \text{Var}[\tilde{y}_i] + \sigma_i^2)) \\ &= \frac{1}{n} \sum_i^n (\text{Bias}^2[f_i, \tilde{y}_i] + \text{Var}[\tilde{y}_i] + \sigma^2) \end{aligned} \quad (26)$$

Here σ_i^2 is the variance of the noise for each test point. Since all test points are assumed to have the same noise distribution, their variances will be equal. To be perfectly clear, the random variables for both of the preceding calculations are the model predictions \tilde{y}_i and the test point noises ε_i , which are considered to be independent of each other. As an annotation aside; similar derivations may be found in most, if indeed not all, textbooks on machine learning, though they are often less than explicit in defining what domains the expectation values are taken over. For a refreshingly clear presentation, complete with pointing arrows and textboxes, the reader is referred to the set of lecture notes by Professor W. Cohen [3].

In practical terms, we use the bootstrap re-sampling to estimate the bias and variance of the model. This is achieved by replacing \tilde{y}_i with $\tilde{y}_{i,B}$ in Eqn. 26. The expectation values will then be computed over all the bootstrap iterations. Furthermore, we are not able to extract the noise from the response values, and must then use the bias estimate $\text{Bias}^2[y_i, \tilde{y}_{i,B}]$. Explicitly, the MSE, the bias and the variance estimates obtained from the bootstrap re-sampling become, for a test set with size n and a number M bootstrap iterations:

$$\text{MSE}_B = \frac{1}{n} \sum_i^n \left(\frac{1}{M} \sum_B^M (y_i - \tilde{y}_{i,B})^2 \right) \quad (27)$$

$$\text{Bias}_B^2 = \frac{1}{n} \sum_i^n \left(y_i - \frac{1}{M} \sum_B^M \tilde{y}_{i,B} \right)^2 \quad (28)$$

$$\text{Var}_B = \frac{1}{n} \sum_i^n \left(\frac{1}{M} \sum_B^M \left(\tilde{y}_{i,B} - \frac{1}{M} \sum_B^M \tilde{y}_{i,B} \right)^2 \right) \quad (29)$$

E. The Franke Function

The Franke function is defined as

$$\begin{aligned}
 f(x, y) = & \frac{3}{4} \exp \left(-\frac{(9x-2)^2}{4} - \frac{(9y-2)}{4} \right) \\
 & + \frac{3}{4} \exp \left(-\frac{(9x+1)^2}{49} - \frac{(9y-21)}{10} \right) \\
 & + \frac{1}{2} \exp \left(-\frac{(9x-7)^2}{4} - \frac{(9y-3)}{4} \right) \\
 & - \frac{1}{5} \exp \left(-(9x-2)^2 - (9y-2) \right)
 \end{aligned} \tag{30}$$

and will serve as function on which we test the performance of our models and build an understanding before we tackle the real terrain data, which we do not have direct control over in the same way. For reference sake, we have plotted the Franke function for values $[0, 1] \times [0, 1]$ which can be seen in Fig. 2.

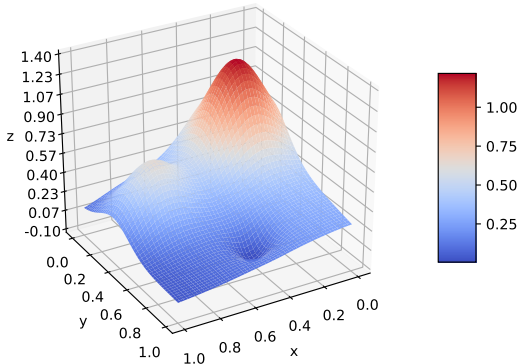


FIG. 2: The Franke Function

III. RESULTS & DISCUSSION

We aim to model both the Franke function and the terrain data as sums of polynomials in x and y . Our overarching goal is to determine which regression method performs the best (i.e. most faithfully reproduces the full data-set on which they are sampled) in both cases, and for which degree of the polynomial model and which value of the penalty parameter λ . The Franke function is given as in Eqn. 30 over the domain $x, y \in [0, 1]$.

The terrain data, which has been downloaded from [4], is a collection of measured elevation heights at a discrete set of 3601×1801 pixel coordinates. For ease of interpretation and computation, we parameterize the pixel indices as uniformly spaced $x, y \in [0, 1]$, and restrict ourselves to the first square of the terrain map, i.e. the first $[1801, 1801]$ pixels.

We will first consider random samples of the Franke function. We will analyze the performance of the different regression methods for various models (polynomial degrees) and penalty parameters. We will investigate the bias-variance behaviour of the methods, and the effect of the sampling size. Finally we will compare various 'best predictions' for the different regression methods, in order to visualize the significance of their different approaches. These predictions will be learned on the same randomly sampled set of points in $x, y \in [0, 1]$, and applied to a grid of 2000 points in each direction, with their resulting predictions being compared amongst the methods and with the ground truth of the Franke function evaluated on that grid.

Armed with insights from the investigation of the Franke function, we will subsequently do much of the same analysis for the terrain data. Ultimately we will compare the best predictions for each method (learned on the same down-sampled data set) with each other and with the ground truth of the full $[1801, 1801]$ terrain map we down-sampled from.

A. The Franke function

Throughout our analysis of the Franke function our focus was primarily on a data set consisting of $N = 500$ uniformly distributed random points, though we also made a brief analysis of the effects of varying the sample size. In particular, we present a limited analysis of a $N = 200$ sample in order to infer about the effects on which N has on the various regression methods.

1. Noise & confidence intervals

We start by looking at the behaviour of OLS limited to relatively modest complexities of polynomials up to the 5th order modeled from data sets consisting of $N = 500$ random samples of the Franke function with added random noise in the form $0.2\mathcal{N}(0, 1)$. We also chose our confidence interval for the predictors as $\beta_i \pm 2\sigma_i$ in adherence with the literature as a standard choice [1], which corresponds to 95% of similarly sampled models falling within our range of uncertainty.

In Fig. 4 we see the size of each predictor β_i for a 5th order polynomial. Further, we also see the full size of confidence interval for each term for all the polynomial degrees $p = 1, \dots, 5$.

We observe that the confidence intervals markedly increases as a function of the models complexity, as shown more explicitly in Fig. 5, which seems to suggest a \sim quadratic trend. This implies that as we try to model our data on higher order polynomials, we will at some point run into the issue of the confidence interval of our predictors becoming an increasingly significant problem for the regression.

We then investigated the effects that different magnitudes of noise had on OLS. In particular, we added random noise sampled from a weighted normal distributions $\delta \cdot \mathcal{N}(0, 1)$ for different weights δ , where the R^2 score for each case over the different complexities are shown in Fig 6. We observe here that the regression becomes increasingly worse as δ approaches 0.5, and the noise begins to dominate.

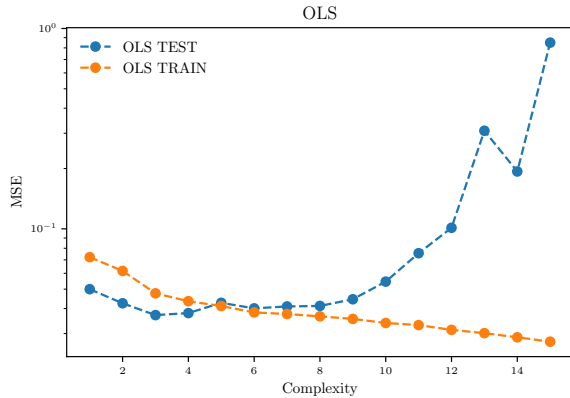


FIG. 3: The MSE for both the testing and the training data performed with OLS regression on $N = 500$ random samples of the Franke Function

2. Resampling

Keeping the noise at 0.2, we will move on trying to find an optimal model complexity for our $n=500$ data set. Splitting our data into a training and test set, with a test/training ratio 0.2, we can see how the MSE evolves as a function of polynomial degree in Fig. 3.

We may anticipate the rapid increase in the MSE of the test set being due to the larger variance at higher complexity. This motivates an application of Ridge and Lasso regression, which aim to reduce that variance at the cost of an increased bias.

Indeed, performing a bootstrap analysis (with $M=100$ bootstraps) shows, in Fig. 26, estimates for the test-MSE, the bias squared and the variance for the three different methods. The bottom row of the figure shows the situation for our original $n=500$ data set. We do confirm that the explosion in the test-MSE is due to the increased variance for higher polynomial degrees. Comparatively, we see that both Ridge and Lasso regression successfully manage to keep the variance down, with seemingly little increase in bias. On the top row we see the same type of bootstrap analysis for a data-set of $n=200$ points. There, the behavior is qualitatively similar, but the variance is higher for all methods, and it explodes the OLS at a lower polynomial degree.

We want to have a high enough complexity to be able to accommodate the peaks and valleys of the Franke func-

tion, while still not over-fitting to the noise. In order to find a suitable polynomial degree, we apply k-fold cross-validation, with $k=5$, to our full data set in order to get a more reliable estimate of the test-MSE than that provided by the bootstrap. For the Ridge and Lasso regressions, the k-fold cross-validation is also performed for a grid (with size 30) of possible values for the regularization parameter $\lambda \in [10^{-6}, 10^0]$ in order to select the best value (giving the lowest test-MSE) for a given degree.

We see, from Fig. 7, that while the OLS solution still blows up past a certain complexity, the Ridge and Lasso solutions keep a very stable MSE-score. Looking at the contour-plots of MSE versus λ and complexity, we find that there is little change in the test-MSE for Ridge and Lasso over a wide range of λ values and polynomial degrees. Though not shown here, we have investigated and found that, with the exception of for very high λ -values, the bias and variance versus complexity for both ridge and lambda does not vary appreciably with λ , nor do the actual predictions resulting from those methods trained on our 500-point data set.

Just looking at the MSE-scores one might be tempted to simply pick the method and complexity with the lowest MSE-score. This is not always a good idea. As we will see, a low MSE-score does not guarantee that predictions will be good. In fact, low complexity models, while having small MSE-scores, fail to reproduce more than the main peak of the Franke function. They simply do not have the necessary degrees of freedom to incorporate the finer qualitative details. In Fig. ?? we show the Franke function and the predictions learned on our $n=500$ points for OLS, Ridge and Lasso, for the highest complexity OLS manages before the variance becomes large enough to ruin it. This is notably NOT the highest degree that still keeps a low MSE-score, but rather the highest degree before the variance becomes appreciable in our bootstrap plots.

We see that the OLS solution clearly most resembles the actual Franke function, even though our error estimates give it a larger MSE than the Ridge and Lasso methods. We find that, while the Ridge and Lasso methods reproduce the main feature (the peak), their attempts at killing of the variance also kills of the actual variations in the underlying Franke function! We do not show it here, but we have tested these predictions for higher polynomial degrees on the exact same data. Those tests reveal that rapidly the OLS solution becomes dominated by the noise, even while keeping a low MSE, while the Ridge and Lasso solutions remain more or less unchanged.

This leads us then to conclude that the "best" fit to this kind of "terrain" data, if the goal is to actually reproduce the qualitative features of the "terrain", is the highest degree of OLS that is not overpowered by the variance. If one insists on using a higher complexity, Ridge and Lasso will remain stable, but might be markedly worse at reproducing the underlying features than a lower degree OLS. As might be imagined after

our earlier discussion regarding the effect of the sample size, a larger sample size allows for the use of higher degrees before OLS meets a variance barrier; conversely, with fewer data points the highest "safe" complexity decreases. In the unfortunate case that one has so few data-points that variance becomes an issue even for say polynomial degree 4, then one might actually prefer a higher order Ridge or Lasso fit, to at least capture the main "feel" of the terrain.

B. The terrain data

We then proceed to treat the terrain data in an almost identical fashion. We began by using a data set containing 2116 pixels/points created by picking only every 40th point in x and y from the original 1801×1801 map. That is, we down-sample into a map consisting of 46×46 pixels.

Once again using k -fold cross-validation, with $k = 5$, on our complete data set, we try to estimate what the optimal complexities and penalty parameters might be for the three methods. The results are shown in Fig. 12. We see, once again, that the MSE-scores remains quite consistent for Ridge and Lasso, and that there is little dependence on λ outside the very largest values. Here, the MSE for OLS is actually decreasing with complexity, until it blows up past a certain point.

A bootstrapped bias-variance analysis of a train-test split, with test/train ratio 0.2, is then performed, in order to see whether the exploding OLS-MSE is due to variance and whether Ridge and Lasso manage to keep the variance down. Fig. 11 shows this to be the case. Ridge and Lasso do keep the variance in check, so that most of their errors are due to bias. For OLS we see that it performs quite well, according to MSE, until the variance kicks in and poisons the solutions.

Keeping in mind our previous lesson about MSE-scores not being the whole picture, we try to recreate our original terrain data with predictions learned on our 2116 points. We try the highest polynomial degree before the OLS solution blows up, and find that the variance is still

significantly high to severely disfigure our predictions, even though the MSE-score is actually at its lowest. Going down in complexity until the variance becomes a sufficiently minor component of the error leads us to what we consider the best prediction for the given number of 2116 data points. This is shown in Fig. ???. We see that only the OLS method really manages to capture the variability of the terrain, the other two mainly identify a peak and a bottom and smooth out the paths between.

IV. CONCLUSION

We have analyzed the performance of OLS, Ridge and LASSO regression on 'terrain-like' types of data, where the regression model was two-dimensional polynomials of varying degrees. We find that in order to qualitatively predict finer features of the terrain, only OLS is really useful; the other two methods smooth out the actual terrain variations in an effort to be stable against variance. The best qualitative description, for a given number of input sample points, is given by the highest degree OLS-solution where the variance remains insignificant. At higher degrees the variance becomes strong enough to completely mar the predictions, even though the MSE-score may remain deceptively low. Should the number of input data-points be so low that the highest "safe" OLS solution is such an exceedingly low degree as to be useless, then Ridge or Lasso at higher complexities may at least predict the main features of the terrain.

As a concluding observation, the authors conjecture that the use of a different type of basis functions might be appropriate for these kind of 'terrain-like' data sets. For instance orthogonal polynomials, or a Fourier-type basis set, could be global basis sets that allow Ridge and Lasso to be more discerning in their regularizations, than with the highly correlated polynomials that have been used in the present study. Or basis functions with compact support, which are only non-zero on limited domains, such as B-splines could provide more freedom to adapt to the actual variations in the terrains, without having over-fitting ramifications outside their support.

-
- [1] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*, 2nd ed. (2009).
 - [2] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, Scikit-learn: Machine learning in Python, *Journal of Machine Learning Research* **12**, 2825 (2011).
 - [3] W. W. Cohen, [Bias-variance in machine learning](#).
 - [4] [Source of terrain data](#).

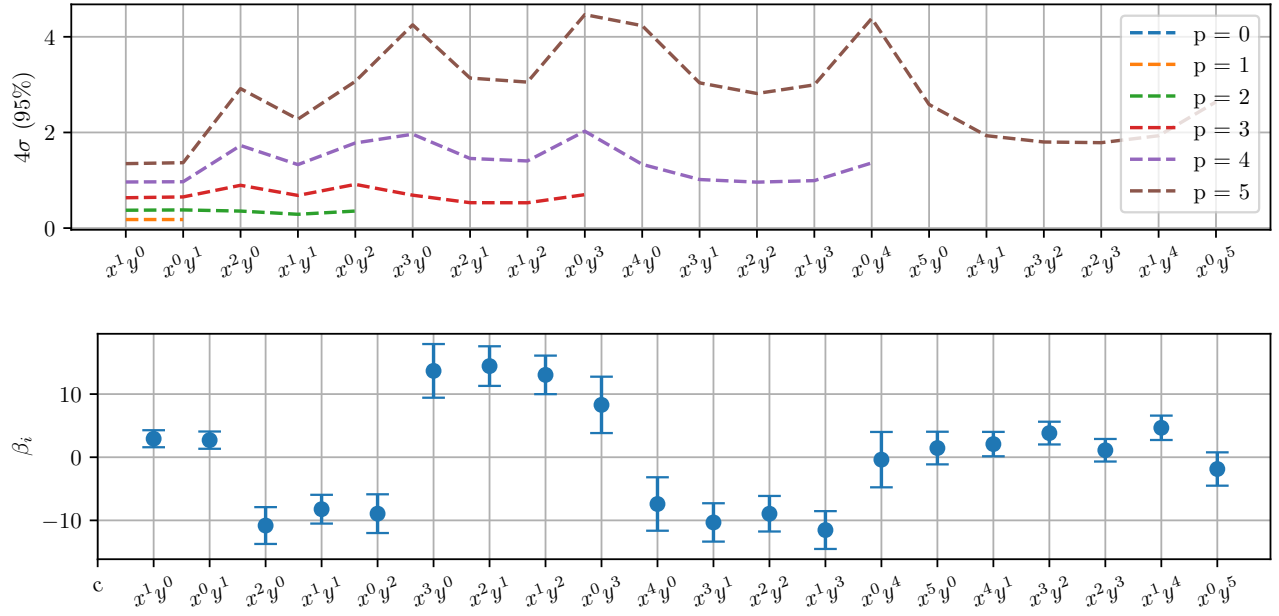


FIG. 4: Predictors β_i for a 5th degree polynomial modeled with OLS on random samples of the Franke Function with noise $0.2 \cdot \mathcal{N}(0, 1)$

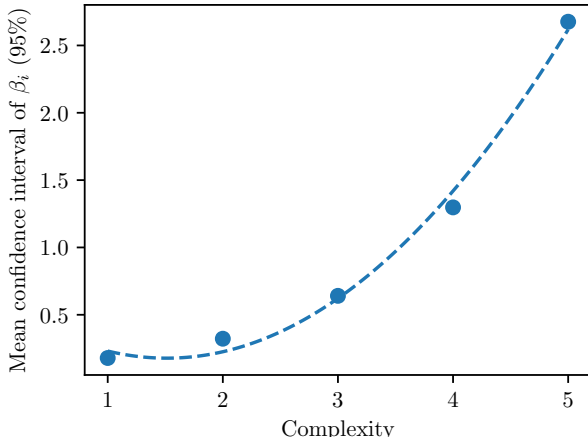


FIG. 5: Mean size of the 4σ (95%) confidence interval across all the predictors for polynomials up to the 5th order with OLS for $N = 500$ random samples of the Franke Function. Where the dashed line indicates 2nd order polynomial regression of the points

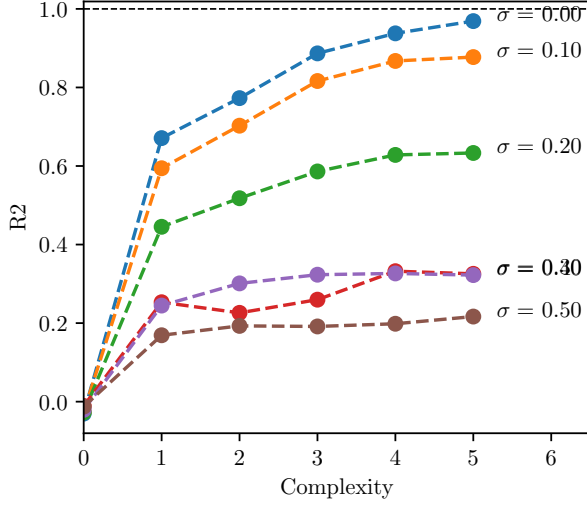


FIG. 6: R^2 measured for $N = 500$ random samples of the Franke function with added noise in the form $\delta \cdot \mathcal{N}(0, 1)$ modeled with complexities 0 to 5.

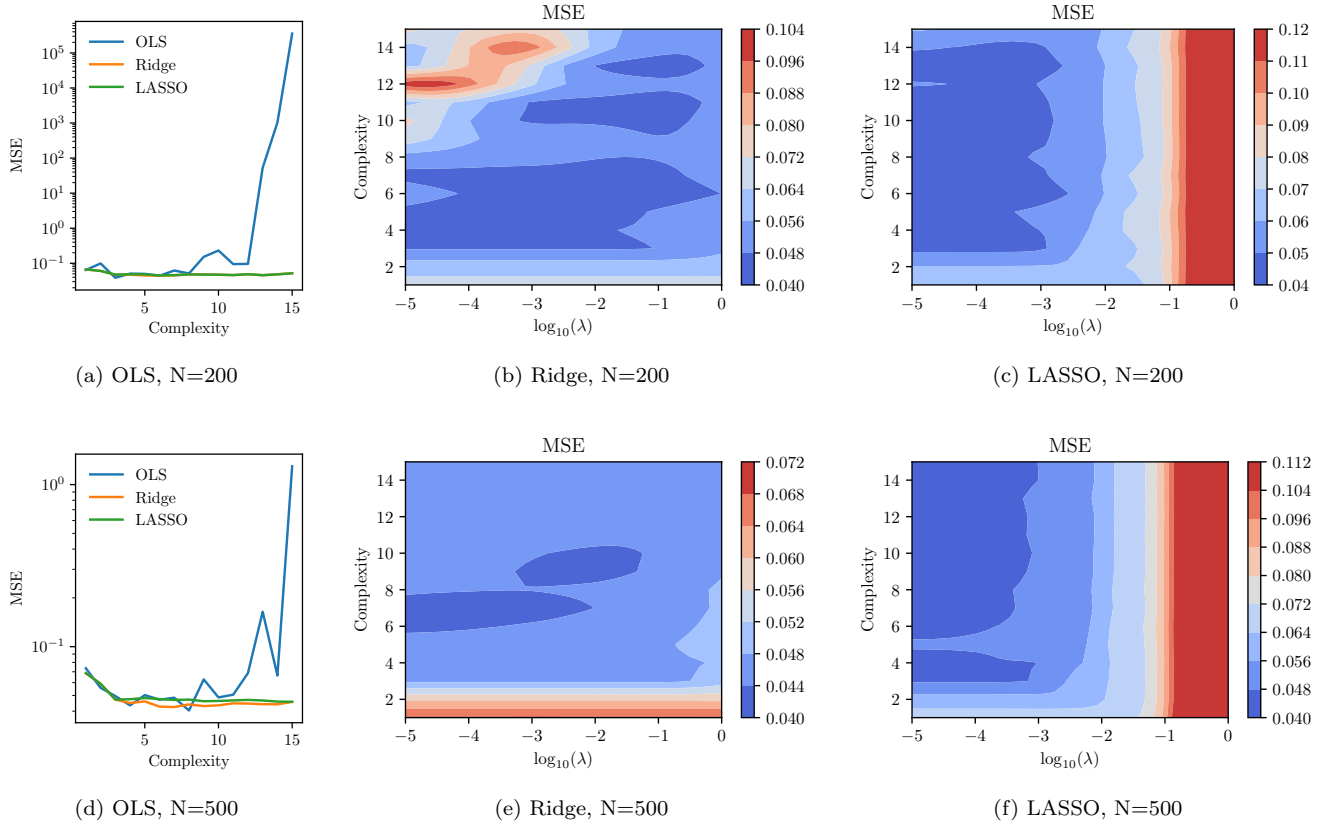


FIG. 7: MSE for $k = 5$ cross-validated Ridge & LASSO regression performed on N samples of the Franke Function with added random noise in the form $0.2 \cdot \mathcal{N}(0, 1)$

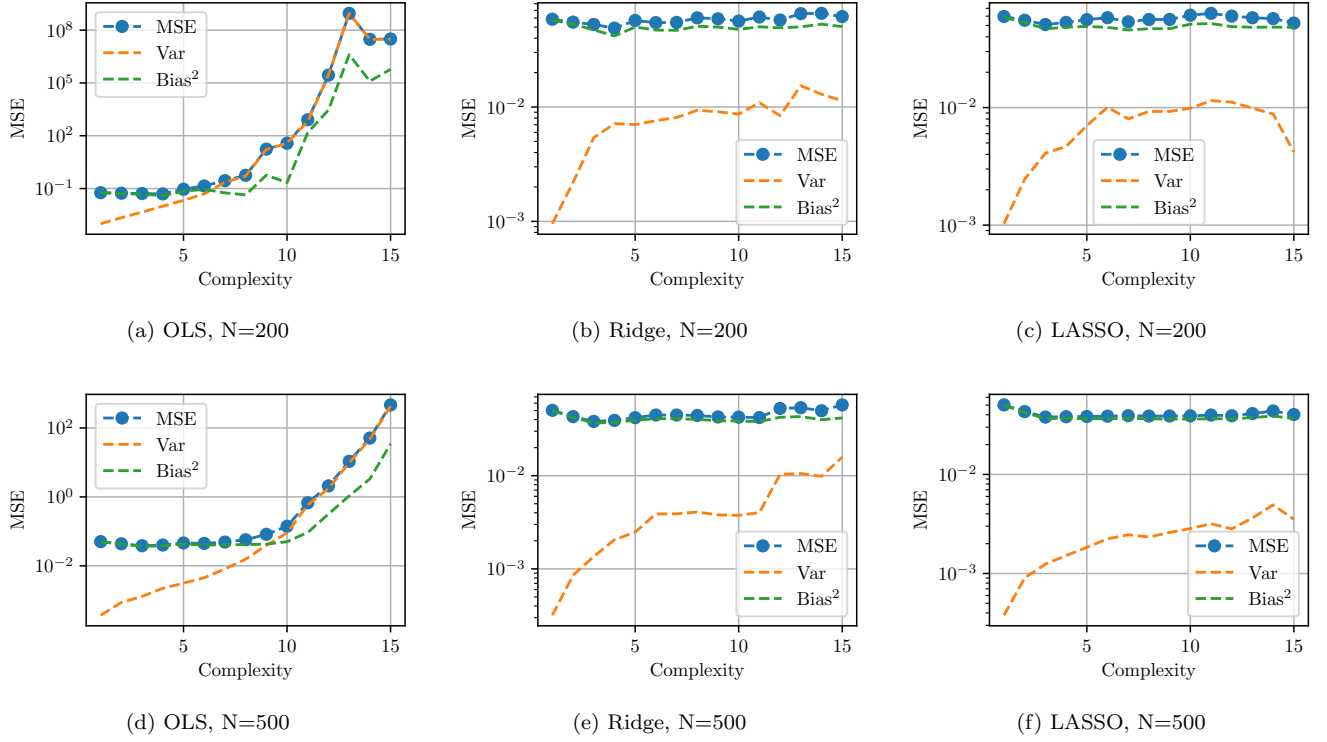


FIG. 8: $M = 100$ Bootstrapped MSE, Bias² and Variance for OLS, Ridge and LASSO regression on N random samples of the Franke Function.

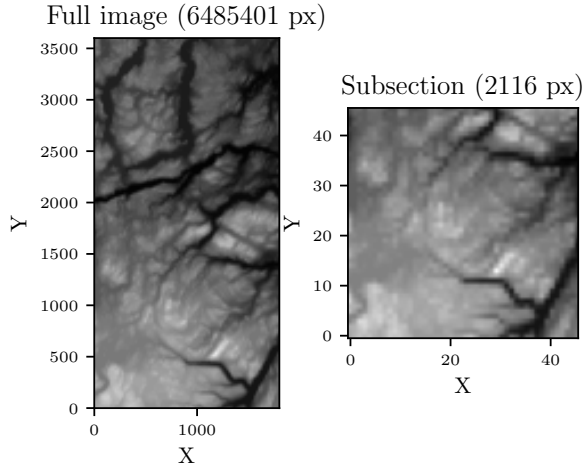


FIG. 9: The full terrain data as well as the down-sampled subsection which was used as the data set for our model.

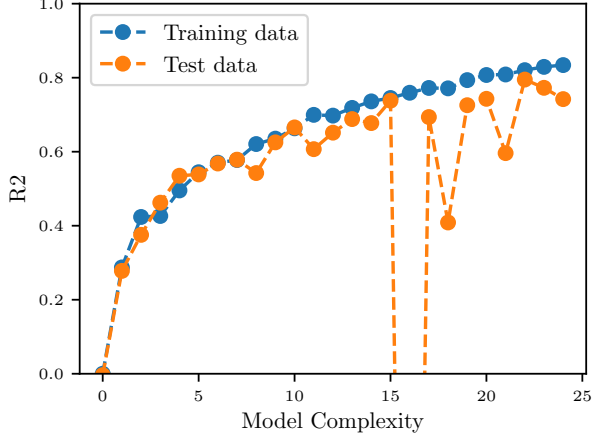
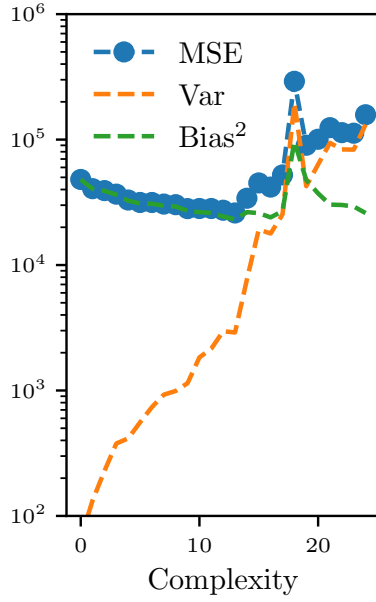
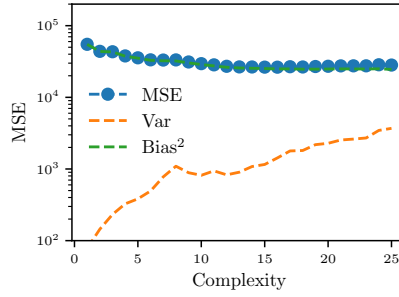


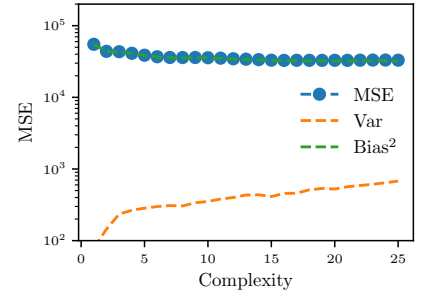
FIG. 10: R^2 for the Terrain data using OLS regression.



(a) OLS, $N=500$



(b) Ridge



(c) LASSO

FIG. 11: $M = 100$ bootstrapped MSE for LASSO and Ridge regression on the Terrain data for the λ parameters yielding the best results for $k = 5$ fold cross-validation along with the $M = 100$ bootstrapped OLS.

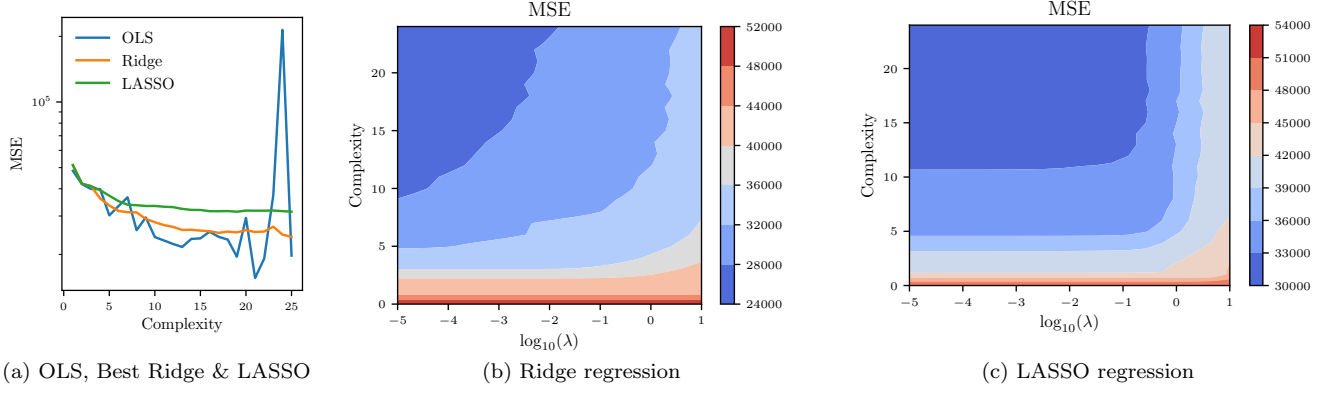


FIG. 12: MSE for $k = 5$ cross-validated Ridge & LASSO regression performed on the terrain data