

FYS-STK4155 – Project 1

Multiple Polynomial Regression

Julie Thingwall (juliethi)
Jonas Gahr Sturtzel Lunde (jonassl)
Jakob Borg (jakobbor)

also known as the three Js of the apocalypse

duh duh duh....duh duh...duhduh...

Morten Hjorth-Jensen's phone

Abstract

Regression is perhaps the most fundamental form of machine learning, describing a direct relation between sets of data through simple functions. In this report, we study multiple polynomial regression on synthetic and real data, with two explanatory variables. We employ three regression methods, namely ordinary least squares (OLS), Ridge regression, and Lasso regression, where we study the unique benefits and limitations of each method. A large emphasis is put on the study of fundamental concepts which carry over to more advanced machine learning methods, such as bias-variance tradeoff, resampling techniques, and cross-validation. We find that all methods are quite capable of predicting both datasets, but that OLS and partially Ridge suffers from overfitting, the effect of which is reduced by an increased size of the datasets. Lasso was found to be superbly overfit resistant, but at the cost of not achieving as good an optimal solution, and being very computationally costly. The Hessian matrix of high-order models was found to be very ill-conditioned. This posed practical limitations on the complexity of the model, where higher order columns of the design matrix are suppressed to practically zero, hindering the fit of high-order polynomials. For the synthetic Franke data, with $\sigma = 1$ gaussian noise, we were able to achieve an MSE of (...) and R2 score of (...) compared to the noisy data, and a MSE of (...) and R2 score of (...) compared to the Franke function, using a 5th order polynomial, and OLS. For the real terrain data, we were able to achieve an R2 score of 0.761, and an MSE of 0.0184 km, using a 80th order polynomial, and OLS.

Contents

1	Introduction	3
2	Theory	3
2.1	Linear regression	3
2.1.1	The design matrix	3
2.2	Ordinary Least Squares	3
2.2.1	Singular and ill-conditioned data	4
2.3	Ridge Regression	4
2.4	Lasso regression	4
2.5	Error Metrics	5
2.6	Variance and Confidence Intervals	5
2.7	Resampling techniques	5
2.7.1	Cross validation techniques	5
2.7.2	K-folding	5
2.7.3	Bootstrap	6
2.8	Bias variance tradeoff	6
3	Method	6
3.1	Data	6
3.1.1	Franke's function with Gaussian noise	6
3.1.2	Terrain data	7
3.2	Code structure	7
3.3	Multivariate polynomial regression	7
3.4	Basic analysis	7
3.4.1	K Fold Cross Validation	7
3.4.2	Ridge Regression	8
3.4.3	Lasso Regression	8
3.4.4	Bias Variance Tradeoff	8
3.4.5	Calculating the Bias and Variance	8
3.5	Numerical Limitations	8
4	Results	8
4.1	Ordinary Least Squares analysis	8
4.1.1	Franke data without cross-validation	8
4.1.2	Franke data with cross-validation	9
4.1.3	Terrain data	9
4.2	Confidence Intervals	10
4.3	Ridge, Lasso and λ -dependence	10
4.3.1	Franke data	10
4.3.2	Terrain data	11
4.4	Bias Variance Tradeoff	11
4.4.1	Franke Data	11
4.4.2	Terrain Data	12
4.5	Best Fit Terrain	14
5	Conclusion	14
	Appendices	15
A	Bias-Variance derivation	15
B	Additional Bias-Variance Results	16

1 Introduction

The simple nature of linear regression gives it some advantages to its more complicated cousins, like being very interpretable, often offering direct insights into relations between the data. Linear regression is usually also reducible to very simple, explicit, and relatively computationally cheap problem. Properties and concepts applicable to regression also show up in more sophisticated machine learning algorithms, making it a great gateway method to understanding important concepts. A lot of this project will be dedicated to studying these concepts, like resampling techniques, bias-variance tradeoff, under/overfitting, as well as stability and convergence properties.

We will focus on polynomial regression in two dimensions (two explanatory variables), employing three regression methods, namely Ordinary Least Squares (OLS), Ridge regression, and Lasso regression. The strength and limitations of these methods, as well as regression in general, will be discussed. We will use two different datasets for our explorations. Firstly, we will generate a dataset using the Franke function, a function commonly used to test interpolation algorithms. Normally distributed noise will be added to the generic data. Secondly, we will use digital terrain data of Norway, which should pose a much larger challenge to our methods.

2 Theory

2.1 Linear regression

Regression aims to explain some output data \mathbf{f}^1 as a function of some input data $\{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{p-1}\}$, plus some unknown noise term, ϵ . In other words,

$$\mathbf{f} = h(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{p-1}) + \epsilon \quad (1)$$

Here, \mathbf{f} and \mathbf{x}_i are in \mathbb{R}^n , meaning that there are n corresponding observations in each dataset, while there are p different input sets.

Linear regression assumes h is a linear function of the input variables, giving

$$\hat{\mathbf{f}} = h(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{p-1}) = \beta_0 + \beta_1 \mathbf{x}_0 + \dots + \beta_p \mathbf{x}_{p-1} \quad (2)$$

where $\hat{\mathbf{f}}$ denotes the regression fit of the true data \mathbf{f}

2.1.1 The design matrix

The linear nature of equation eq. (2) means it can be written as a matrix equation, on the form

$$\hat{\mathbf{f}} = \mathbf{X} \cdot \beta \quad (3)$$

where β is an unknown \mathbb{R}^p vector of the polynomial coefficients, and \mathbf{X} is a known $\mathbb{R}^{n \times p}$ matrix of the input data, the columns of which are simply the input data (and a vector consisting solely of ones, matching the β_0 terms):

$$\mathbf{X} = [\mathbf{1} \ \mathbf{x}_0 \ \mathbf{x}_1 \ \dots \ \mathbf{x}_{p-1}] \quad (4)$$

The purpose of linear regression is to solve equation eq. (3) for β . This equation does normally not have a solution, as it would require our function to perfectly fit every datapoint. Even if such a solution should exist, it would not even be desirable, as it would probably be a result of overfitting the model to the data, as all (real) data have some inherent noise and uncertainty in them. Instead, one employs some sort of *cost function* $C(\beta)$, which quantifies how good a fit a certain β gives. Regression then reduces to minimizing this cost function, resulting some optimal β . Choosing the right cost function is important. As a matter of fact, regression methods, are actually uniquely defined by their cost functions.

2.2 Ordinary Least Squares

The most common and simple linear regression method is the Ordinary Least Squares method. Its cost function is defined as

$$C^{\text{OLS}}(\beta) = \frac{1}{n} \sum_{i=0}^{n-1} |f_i - \mathbf{x}_i \cdot \beta|^2 \quad (5)$$

or, written on matrix form,

$$C^{\text{OLS}}(\beta) = \|\mathbf{X}\mathbf{f} - \beta\|_2^2 \quad (6)$$

where $\|\cdot\|_2$ indicates the L2 norm.

One of the advantages of OLS is its simple interpretation. Its cost function is the sum of the squared difference between our predicted value, and the actual value of the data.

It can be shown (ref plz) that minimizing the cost function of OLS actually gives an explicit formula for β , as

$$\min_{\beta} (C^{\text{OLS}}(\beta)) \Rightarrow \hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{f} \quad (7)$$

From a linear algebra standpoint, minimizing the OLS cost function is simply the orthogonal projection of \mathbf{y} onto the column space of \mathbf{X} . This can be seen from the fact that

$$\hat{\mathbf{f}} = \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{f} \quad (8)$$

where $P = \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is recognized as the projection matrix of the space spanned by the columns of \mathbf{X} .

¹The output is usually denoted \mathbf{y} , but we reserve this as the second input parameter.

2.2.1 Singular and ill-conditioned data

One obvious problem with OLS regression is that there is no guarantee that $(\mathbf{X}^T \mathbf{X})$ is invertible, giving no explicit solution to eq. (7). A singular matrix means there is no single best fit for the data.

Another related problem is that the matrix might analytically be invertible, but be very *ill-conditioned*. The *condition number* of a matrix A is defined as the ratio of the highest to the lowest singular value:

$$\text{condition number} = \frac{s_0}{s_{n-1}} \quad (9)$$

A high condition number makes the matrix almost impossible to row-reduce (and therefore invert), due to its very high sensitivity to numerical noise. This might be even worse than a singular matrix, as a numerical method would gladly "invert" the matrix, giving a totally wrong matrix.

One solution is simply employing a different regression method which circumvents this problem, like Ridge or Lasso regression. Another solution is to approximate the inverse of the matrix, using something known as a pseudo inverse, or a Moore–Penrose inverse. The pseudo inverse of a matrix, usually denoted A^+ , can be shown to give the best OLS fit of a model (REF PLZ), and is defined by the SVD composition of A , as

$$A^+ = U D^+ V^T \quad (10)$$

where $A = V D U^T$ is the SVD decomposition of A , and D^+ is defined as a diagonal matrix with the inverse of the diagonal elements of D , σ_i on the diagonal

$$D^+ = \begin{pmatrix} 1/\sigma_0 & & & \\ & 1/\sigma_1 & & \\ & & \ddots & \\ & & & 1/\sigma_{n-1} \end{pmatrix} \quad (11)$$

except for any $\sigma_i = 0$, in which case the diagonal element of D^+ is also set to 0.

In the case of a non-singular matrix, the pseudo inverse matrix is simply the inverse matrix. In the world of numerical inaccuracy, the pseudo inverse is also much more well behaved than the inverse for ill-conditioned matrices.

2.3 Ridge Regression

Ridge regression is as a slight modification to OLS regression. Ridge regression keeps the nice analytical properties of OLS (beta can be calculated explicitly), while outright solving the issues off ill-conditioned and singular matrices. It also holds other properties which, making it superior to OLS in some ways.

The cost function of Ridge regression is

$$C^{\text{Ridge}}(\beta) = \frac{1}{n} \sum_{i=0}^{n-1} |f_i - \mathbf{X}_{i*} \cdot \beta|^2 + \lambda \sum_{j=0}^{p-1} \beta_j^2 \quad (12)$$

or, on matrix form

$$C^{\text{Ridge}}(\beta) = \|\mathbf{X}\mathbf{f} - \beta\|_2^2 + \lambda \|\beta\|_2^2 \quad (13)$$

$$= C^{\text{OLS}}(\beta) + \lambda \|\beta\|_2^2 \quad (14)$$

The cost function of Ridge is simply that of OLS, with an added term. The new term is the squared size of the beta-vector, times some hyperparameter λ . The size of the new term is controlled by a hyperparameter λ . As we can see, this new term serves to penalizes the size of the β s.

It is common to rescale the output data around its mean, resulting in a small or zero intercept. This means Ridge won't penalize the intercept β_0 . The purpose of Ridge regression, as opposition to OLS, is to penalize the *slopes*, β_i , $i \neq 0$ of the model. This penalty is proportional to the hyperparameter λ . Ridge will therefore naturally have lower beta values than OLS. As the beta values define the correlation between the input and output data, increasing lambda can be interpreted as an increased skepticism to the input variables ability to predict the output. And, naturally, the steeper the minimum of the cost function was around a certain β , the less it will be reduced by Ridge. In other words, less confident β s get moved more.

Like OLS, Ridge has an analytical expression for β , derived from the minima of the cost function, given as

$$\beta = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{f} \quad (15)$$

the only difference from OLS being the addition of λ along the diagonal of the matrix being inverted.

This addition has another nice consequence. If $\mathbf{X}^T \mathbf{X}$ is singular, adding a small value to the diagonal will make it non-singular. This is an alternate solution to the singular problem of OLS.

2.4 Lasso regression

Lasso regression reminds a lot of Ridge regression. It's cost function is

$$C^{\text{Lasso}}(\beta) = \frac{1}{n} \sum_{i=0}^{n-1} |f_i - \mathbf{X}_{i*} \cdot \beta|^2 + \lambda \sum_{j=0}^{p-1} |\beta_j| \quad (16)$$

or, on matrix form

$$C^{\text{Lasso}}(\beta) = \|\mathbf{X}\mathbf{f} - \beta\|_2^2 + \lambda \|\beta\|_1 \quad (17)$$

$$= C^{\text{OLS}}(\beta) + \lambda \|\beta\|_1 \quad (18)$$

Instead of penalizing with the squared of the coefficients, Lassos cost function only employs the absolute value of β . This might not look like a large change, but it has a few consequences. Firstly, it sadly means that β is not longer analytical. The cost function of Lasso must be minimized using some sort of optimization algorithm, which is often more costly, and introduces problems like convergence and exactness.

The main difference in the results of Ridge and Lasso is that Lasso has a tendency to suppress some β s all the way to zero, while Ridge usually just reduce them somewhat. In Ridge, the

2. THEORY

penalty of having a small β is very small, since we square it. It will usually reach a point where reducing it further will increase the OLS cost term more than it will decrease the λ penalty term. In Lasso, the penalty to the β s isn't squared, and reducing the β to zero might very well give a reduction to the cost function.

Since Lasso tends to reduce some β s to zero, and is preferred when one might have reason to believe some of the explanatory variables are bad predictors.

2.5 Error Metrics

Different error metrics are often employed to quantify the quality of a given model. If a model gives a prediction $\hat{\mathbf{f}}$ for a data set with true values \mathbf{f} , we propose the following metrics.

Mean Squared Error (MSE)

$$\text{MSE}(\mathbf{f}, \hat{\mathbf{f}}) = \frac{1}{n} \sum_{i=0}^{n-1} (f_i - \hat{f}_i)^2 \quad (19)$$

Being simply the squared difference between prediction and true data, the mean squared error is a popular error metric. The MSE is exactly what OLS uses as its cost function, and MSE will therefore be guaranteed to minimize the MSE in respect to the data it was trained on.

MSE has the disadvantage having units squared of whatever the units of the data was. This might reduce the direct interpretability of the numbers.

R2 score

The R2 score is another common metric to quantify the success of our model. It is defined as

$$R^2(\mathbf{f}, \hat{\mathbf{f}}) = 1 - \frac{\sum_{i=0}^{n-1} (f_i - \hat{f}_i)^2}{\sum_{i=0}^{n-1} (f_i - \bar{f})^2}. \quad (20)$$

The R2 score looks weird, but has some nice interpretations. The numerator is simply the MSE of the model, while the denominator is the MSE of a model which just predicts all the values to be the average of the dataset, in other words, assuming there are no correlations. The fraction is therefore a metric for how bad the model is compared to an entirely "blank" model. Since R2 is one minus this quantity, an R2 score closer to 1 means a better model.

2.6 Variance and Confidence Intervals

Even after an optimal β is found, it is important to establish the level of confidence in the derived values. The *variance* of a variable θ is defined as

$$\text{Var}[\theta] = E[\theta^2] - E[\theta]^2 \quad (21)$$

and measures the predicted spread in the variable.

((ref)) shows that the variance-covariance matrix of our estimators β can be written as

$$\text{Var}[\beta] = (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2 \quad (22)$$

where σ^2 is the variance of the predicted data. This is typically estimated by

$$\hat{\sigma}^2 = \frac{1}{N - p - 1} \sum_{i=0}^{N-1} (y_i - \hat{y}_i)^2 \quad (23)$$

A related metric is the *confidence interval* for $\hat{\beta}$. A confidence interval of level α is defined as an interval around the predicted β s, where, if we were to be given many different data sets, and estimated a series of β s, $100(1 - \alpha)\%$ of them would fall in our chosen interval. It can be interpreted as a $100(1 - \alpha)\%$ chance that our chosen β s are correct, even though this interpretation is loathed by statisticians.

Assuming that the β s are drawn from a normal distribution, which is a standard and reasonable approximation, the confidence interval for $\hat{\beta}$ is

$$\beta \in [\hat{\beta} \pm z^{1-\alpha} \text{Var}[\hat{\beta}]], \quad \text{Var}[\hat{\beta}] = \mathbf{v}^{1/2} \hat{\sigma} \quad (24)$$

where $\mathbf{v} = \text{Diag}[(\mathbf{X}^T \mathbf{X})^{-1}]$, and $z^{1-\alpha}$ is the cumulative distribution function of the normal distribution, at a level $1 - \alpha$.

2.7 Resampling techniques

The primary limitation of all machine learning models is the lack of data. Data is needed to train the model, test the validity of the model, and infer important properties of the model, like its bias, variance, or the confidence interval of its parameters. *Resampling techniques* is a collective description of ways to split, reuse, and augment the data in clever ways, to get the most out of it.

2.7.1 Cross validation techniques

Cross validation techniques splits the data in a training and testing set, in order to validate the quality of the model. The model is trained only on the training set, and is then asked to predict the values of the testing set. The quality of the model can then be inferred from its success in predicting the values of the testing set, which the model has never seen before. In the simplest version of cross validation, one would simply split the data in two parts, train on one, and test on the other. Usually, the testing set contains around 1/3 to 1/5 of the data.

2.7.2 K-folding

K-folding is a type of cross validation technique, where the data is split in K random parts, or folds. One of the folds is then selected as a testing set, while the remaining $K - 1$ folds are used as training data. This is repeated K times, until all the folds have been used as testing data once. We're then left with a prediction on all the data, without the model ever "cheating" and seeing the data it's supposed to predict. Pretty cool stuff.

2.7.3 Bootstrap

Bootstrap is a popular resampling technique used to estimate statistical quantities of distributions where the proper probability distribution is unknown. This is accomplished through an iterative process, producing better and better estimators for the distribution’s statistical properties through a frequentist approach. Using the central limit theorem, (Devore und Berk, 2012), one can show that by performing multiple experiments and averaging the resulting estimators will approach their true value for sufficiently large number of experiments. This is easily generalised to multidimensional distributions given that the stochastic variables are independent and identically distributed (i.i.d). If not one would get an additional contribution from the covariance of the variables.

With a given data set of values $\mathbf{X} = (x_1, x_2, \dots, x_n)$, pick n values at random with replacement from \mathbf{X} and treat this subset $\tilde{\mathbf{X}}$ as the distribution in question. From this subset calculate the required estimators. Repeat this process for a given number of bootstraps, say k . By the central limit theorem, these estimators will get increasingly better for larger values of k

$$\mu_{\tilde{\mathbf{X}}} = \frac{1}{k} \sum_{i=0}^{k-1} \mu_i \approx \mu_{\text{True}}$$

$$\sigma_{\tilde{\mathbf{X}}}^2 = \frac{1}{k} \sum_{i=0}^{k-1} \sigma_i^2 \approx \sigma_{\text{True}}^2$$

In the setting of machine learning this could be used to estimate how well a given algorithm learns. This will be discussed further in sections 2.8 and 3.4.4.

2.8 Bias variance tradeoff

When assembling the design matrix from section 2.1.1 some choice is left to the designer. If the model is, for example, a polynomial fit, choosing the order of the polynomial will have a great effect on the error of the fit. This is often called the bias-variance tradeoff. It may be expressed as

$$C(\mathbf{X}, \beta) = \mathbb{E} \left[(\mathbf{f} - \mathbb{E}[\tilde{\mathbf{f}}])^2 \right] + \mathbb{E} \left[(\tilde{\mathbf{f}} - \mathbb{E}[\tilde{\mathbf{f}}])^2 \right] + \sigma^2$$

$$= \frac{1}{n} \sum_{i=0}^{n-1} \left(f_i - \mathbb{E}[\hat{f}_i] \right)^2 + \frac{1}{n} \sum_{i=0}^{n-1} \left(\hat{f}_i - \mathbb{E}[\hat{f}_i] \right)^2 + \sigma^2, \quad (25)$$

where the first term is the bias squared, the second is the variance and the third term is the irreducible error. This tradeoff expresses one of the central problem in any machine learning algorithm.

With a low model complexity, the fit may miss important regularities and features in the training data. When the model then is applied to the test data the expectation value from the fit, $\mathbb{E}[\tilde{\mathbf{f}}]$ will generally be quite different from the actual data, \mathbf{f} . This is a model with high bias, and is often called underfitting. As

the complexity is low, the variance will generally be low when applied beyond the training set, or in combination with resampling techniques.

With increasing complexity, the fit is more likely to find small and detailed features present in the training set, as well as consistencies which is a consequence of noisy or limited data. This way it represent the training data really well, but will tend to find regularities in the test data which should not be there. This is usually called overfitting, where a small change in the training data will impact the model in a great way. Over multiple iterations, through for example the bootstrap method, the predicted data will vary wildly so that $\tilde{\mathbf{f}} \neq \mathbb{E}[\tilde{\mathbf{f}}]$ and we get a high variance.

3 Method

3.1 Data

3.1.1 Franke’s function with Gaussian noise

As a testing ground for our regression methods, we first employ some generic data, generated by Franke’s function. Franke’s function is a smooth exponential function, often used to test interpolation methods. It is defined as

$$f(x, y) = \frac{3}{4} \exp \left(-\frac{(9x-2)^2}{4} - \frac{(9y-2)^2}{4} \right) + \frac{3}{4} \exp \left(-\frac{(9x+1)^2}{49} - \frac{(9y+1)^2}{10} \right) + \frac{1}{2} \exp \left(-\frac{(9x-7)^2}{4} - \frac{(9y-3)^2}{4} \right) - \frac{1}{5} \exp \left(-(9x-4)^2 - (9y-7)^2 \right). \quad (26)$$

In addition, we apply some Gaussian noise with $\sigma = 1$, to better simulate real world data. Our data generating function is therefore

$$f_{\text{noisy}}(x, y, \sigma) = f(x, y) + N(0, 1) \quad (27)$$

We use the interval $x, y \in [0, 1]$ of Franke’s function, where it contains two gaussian peaks and one dip (see left of fig. 1 below). We use a meshgrid of 101 points in each direction, gives us a total of 10201 datapoints across the mesh. After applying the noise, we see from the right of fig. 1 that the shapes are hardly recognizable.

The axis will in practice be shifted and scaled to the interval $x, y \in [-1, 1]$, such that the explanatory variables are centered around 0. The output data remains unchanged. This is further discussed in section 3.5.

3. METHOD

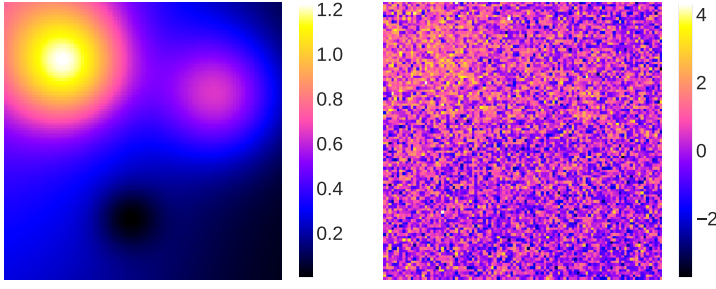


Figure 1 – Left: Generic data from Franke’s function. Right: Data from Franke’s function with added gaussian noise $N(0, 1)$.

3.1.2 Terrain data

As a more serious test of our methods, we use elevation data of the south-western part of Norway, gathered from <https://earthexplorer.usgs.gov/>, and shown in fig. 2. The elevation is in units of km, and we use explanatory variables $x, y \in [-1, 1]$, just like for the Franke data. The original image is very large, which might present problems with both the runtimes, and exploring the over- and underfitting properties of the different methods. We therefore employ downsampling of the image, usually by 16×16 pixels if nothing else is stated. We average pool over these pixels, resulting in a dataset of $113 \cdot 226 = 25538$ datapoints. The downsampled image is shown on the right of fig. 2.

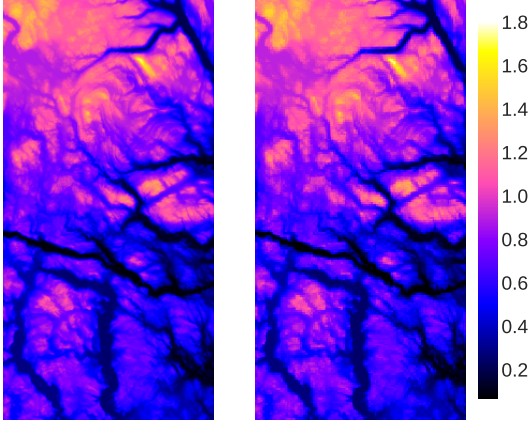


Figure 2 – Left: Elevation map of south-western Norway, 1801×3601 pixels, in units of km. Right: Downsampled version of elevation map on left, 113×226 pixels, with average pooling of 16×16 pixels.

3.2 Code structure

For our implementation we wished to utilize the modular nature of Jupyter Notebook. Thus, our code is structured around a python class called Regression, and several Jupyter Notebook files corresponding to solving different parts of the project.

The Regression class is then the back-bone of our whole implementation. This is where we have written methods for every regression method we need, namely OLS, Ridge- and Lasso regression, with and without cross-validation algorithms.

3.3 Multivariate polynomial regression

We are specifically dealing with polynomial regression with two explanatory variables. This means we have one set of output data, \mathbf{f} , which we attempt to predict as a polynomial of two sets of input data, \mathbf{x} and \mathbf{y} .

If we think of \mathbf{f} as a function of the input variables, $\mathbf{f} = h(\mathbf{x}, \mathbf{y})$, \mathbf{f} becomes a surface in two dimensions (or simply a series of heights, since all three are discrete). The polynomial fit will, either way, be a smooth surface in two dimensional space.

Now, assuming a polynomial of some order m , we get a design matrix on the form

$$\mathbf{X} = [\mathbf{1} \quad \mathbf{x} \quad \mathbf{y} \quad \mathbf{xy} \quad \dots \quad \mathbf{xy}^{m-1} \quad \mathbf{x}^m \quad \mathbf{y}^m] \quad (28)$$

Where \mathbf{X} will be a $n \times p$ matrix, where $p = \frac{(m+1)(m+2)}{2}$, because of all the crossterms between \mathbf{x} and \mathbf{y} . Similarly, we have β as a vector of length p .

3.4 Basic analysis

First and foremost, we performed regular OLS regression with a 5th order polynomial and without cross-validation the generated data as described by eq. (27). This was done mainly as a sanity check, given that we have the luxury of knowing the underlying properties of the data we’re trying to fit, making it easy to check if everything yields expected results. An important fact to be aware of when not using any cross validation algorithms is that we are more likely to end up with over fitting our data set.

This basic analysis consisted of checking the error metrics, namely the MSE and R2 score as described in section 2.5. This was done using simply by using the method found in the Scikit-Learn python package: `sklearn.errormetrics`. We also explored the variance and confidence interval of the β s, following the logic described in section 2.6.

Moving on, with confidence in our implementation this far, we performed the same basic analysis on the terrain data. To minimize computational time we downsample the data set as explained in fig. 2.

3.4.1 K Fold Cross Validation

The next natural step is to implement some form of cross-validation algorithm. We implemented the K fold resampling technique as explained in section 2.7.2. Further on, unless stated otherwise, we will always utilize cross-validation when performing any type of regression, and we will by default use $K = 10$.

After implementing K-folding, we tested it by looking at the different error metrics and compare them to the results from regular OLS without K-folding, both on the generated data and the terrain data. We expect that any signs of over fitting should disappear, as we now train and test on different parts of the data sets.

4. RESULTS

3.4.2 Ridge Regression

Furthermore, we implemented Ridge regression. Moving to Ridge regression, the β s are explicitly calculatable as

$$\beta = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^+ \mathbf{X} \mathbf{f}$$

as discussed in section 2.3. λ is a hyperparameter which needs to be tuned. We will explore the implications of λ , and tune it for each dataset.

3.4.3 Lasso Regression

Lasso, while very similar to Ridge, doesn't have an explicit expression for its optimal β . A optimization scheme will therefore need to be employed, to find the β which minimizes the cost function of Lasso, as discussed in section 2.4. Optimization is beyond the scope of this report, and we will employ the coordinate descent optimizer found in the Scikit-Learn Python package `sklearn.linear_model.Lasso`.

Just as with Ridge, the impact of the λ hyperparameter of Lasso will be studied.

3.4.4 Bias Variance Tradeoff

As discussed in section 2.8, overfitting, resulting in high model variance, can be a serious problem for highly complex models. We therefore studied the training and testing MSE of the models for a varying set of complexities. We split the dataset in a training set of 75% and a testing set of 25%. The models were thereafter run for polynomial order 0 through 20 for the Franke data, and 0 through 110 for the Terrain data, downsampled by 8x8. A $\lambda = 10^{-4}$ used for Ridge and Lasso. We also downsampled the terrain data to 32x32, from polynomial order 0 through 35, to see how a reduction in dataset size would impact the overfitting.

We also studied the numerical stability of the model. If the matrix $\mathbf{A} = \mathbf{X}^T \mathbf{X}$ is ill-conditioned, we might observe unstable or suppressed results. The test and train MSE of Ridge and OLS was run for polynomial orders 0 through 65, for a shifted set of explanatory variables, $x, y \in [0, 2]$. The condition number of the matrix \mathbf{A} for OLS and Ridge for $\lambda = 10^{-2}$ and $\lambda = 10^{-6}$ was plotted against polynomial orders 0 through 60. This was repeated for the shifted explanatory variables.

3.4.5 Calculating the Bias and Variance

In order to separate the contributions to the MSE into bias and variance as discussed in section 2.8, we need to some way to calculate the expectation value of our fit given slightly different training sets. To do this we implemented the Bootstrap method from section 2.7.3. First we split the data, making sure we test the fits on the same test data each bootstrap. Then the remaining training data is used as the set we pick from with replacement. After k bootstraps we are left with multiple slightly different fits, which we used to calculate the quantities involved in the bias and

variance terms of the MSE. In this way we usually can't quantify the irreducible error σ^2 , as this is generally unknown. So the sum of the bias and the variance should be less than the actual MSE.

3.5 Numerical Limitations

There might be limitations posed upon the complexity of our model by our numerical accuracy. Our design matrix will hold polynomials of increasing order, up until some last column of order \mathbf{x}^m . This column might differ hugely in magnitude to the first column, which is of order 1. The decision of how to scale the explanatory variables \mathbf{x} and \mathbf{y} might have a lot of impact on the numerical stability of the design matrix. We have chosen to space our explanatory variables in $[-1, 1]$ to limit such effects.

If an interval for \mathbf{x} and \mathbf{y} smaller than this is chosen, the entire columns will be largely suppressed. Say an interval of $[-0.5, 0.5]$ was chosen, it would have its largest value suppressed to 0.5^m . Around $m = 50$, it would lose approximately all its significant digits in an arithmetic operation with 1, because $0.5^{50} \approx 10^{-15}$, approaching double float precision compared to 1. The columns of polynomial order 50 or higher in the design matrix would for all intents and purposes be 0. Choosing an interval exceeding 1, like $[-2, 2]$, might be even worse, as it would render the *first* columns practically zero instead of the last, ruining the lower order polynomials.

Choosing an interval $[-1, 1]$ in no way solves the problem, as much of the column will be suppressed to zero at higher orders, but neither does it entirely blow up or down, and seemed like the best choice. For OLS especially, the suppressed later columns might lead to a singular or extremely ill-conditioned matrix. Pseudo-invertation, as discussed in [...] will be employed to mitigate this problem.

4 Results

4.1 Ordinary Least Squares analysis

4.1.1 Franke data without cross-validation

In table 1 we see the results of performing an ordinary polynomial regression on eq. (27) using a polynomial of order 5. Since we have the luxury of working with a data set where we know what the underlying function should look like, we compare our fit both to the noisy data, and also the clean Franke function. First and foremost, we notice that the MSE when compared to the noisy data is less than 1. As we have full control of the generated noise in our data set, and we have set this noise to be Gaussian with a standard deviation $\sigma = 1$, we know that the MSE should not go lower than 1. Thus, what we observe here seems to be a classic case of over-fitting, which is not surprising considering we have not yet implemented any type of resampling or cross-validation, meaning we both train and test on the same data set.

4. RESULTS

Aside from that, we see that when comparing to the underlying function, the MSE score does indeed approach 0, and the R2 score is creeping up towards 1.

	MSE	R2
Franke + noise	0.997	0.047
Noiseless	0.004	0.843

Table 1 – Table showing the MSE- and R2 score when performing ordinary polynomial regression on the generated noisy data (eq. (27)) without any resampling techniques. The upper row contains the errors calculated with respect to the noisy data, while the lower row shows the errors with respect to the underlying Franke function with no added noise.

4.1.2 Franke data with cross-validation

In fig. 3, we see the result of fitting eq. (27) with a 5th order polynomial using OLS regression and the K-fold cross-validation algorithm compared with the underlying noiseless Franke function. Here we observe that the MSE is no longer lower than 1 when compared to the noisy data, which points to cross-validation minimizing the chance of over-fitting.

	MSE	R2
Franke + noise	1.001	0.047
Noiseless	0.004	0.820

Table 2 – Table showing the MSE- and R2 score when performing OLS regression with K-folding. The upper row shows the error metrics when compared with the noisy data, while the lower row contains the error when comparing to the clean Franke function.

A visual comparison of the underlying Franke function and our polynomial fit can be seen in fig. 3. While stating that "they clearly look alike" is not the most scientific way of deciding whether a method yields sufficient results or not, they **do** clearly look similar!

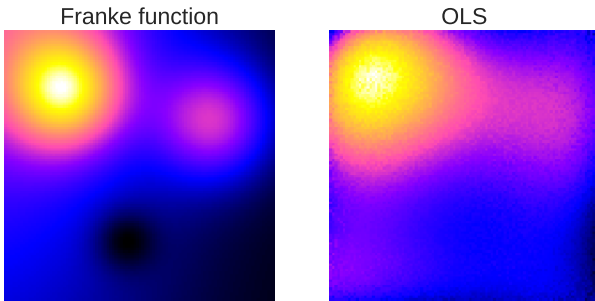


Figure 3 – Figures comparing the Franke function without noise (left) and the prediction made using OLS with a 5th order polynomial and K-folding with $K=10$ (right).

4.1.3 Terrain data

In table 3 we see the MSE and R2 score when trying to fit the terrain data to a polynomial of order 5 using regular OLS and

k-folding with $k=10$. Given the low R2-score we can conclude that a 5th order polynomial might not be the best fit for this data-set. Intuitively, this makes sense, as a terrain will be made up of a high number of nooks and crannies which is not possible to represent with a polynomial with such a low order.

	MSE	R2
Terrain data	0.042	0.456

Table 3 – Table showing the error metrics when fitting the terrain data using OLS regression with a polynomial of order 5 and k-folding with $K=10$.

The quality of the fit can also be seen in fig. 4. Here we see that, while a 5th order polynomial does preserve the overall shape of the data, it does not do a good job at replicating the finer details. How the error metrics behave for a model with higher complexity is explored more closely later on.

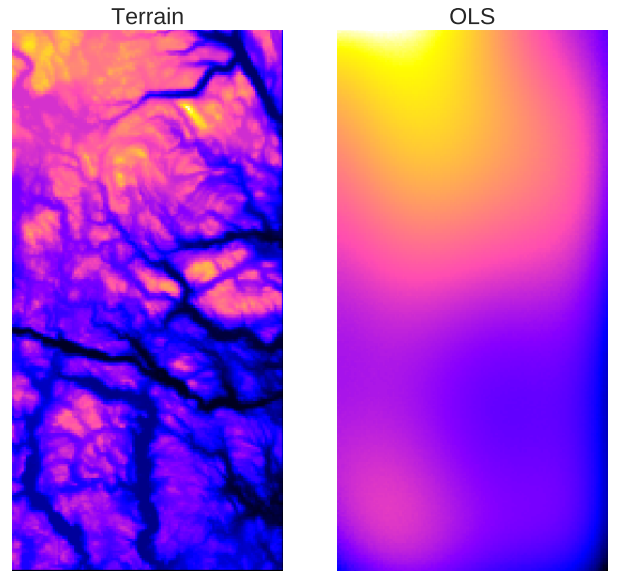


Figure 4 – Figures comparing the terrain data (left) and the prediction made using OLS with a 5th order polynomial, and K-folding with $K=10$ (right).

4.2 Confidence Intervals

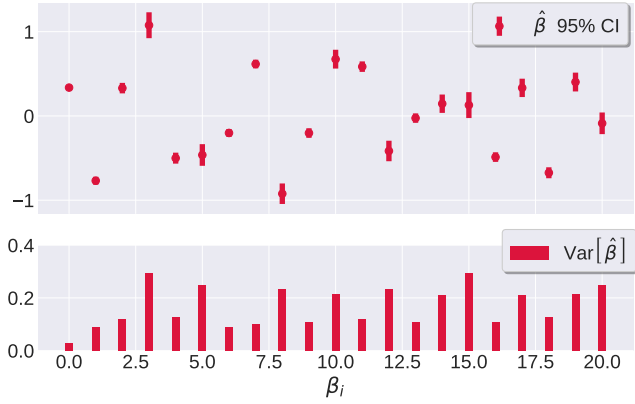


Figure 5 – Regression fit of 5th order polynomial on Franke data, using OLS with K-folding. Top: Predicted β s with 95% confidence intervals. Bottom: Variance of predicted β s.

4.3.1 Franke data

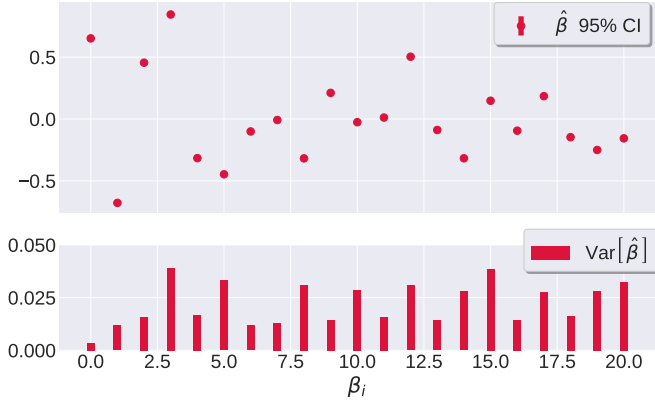


Figure 6 – Regression fit of 5th order polynomial on Terrain data, using OLS with K-folding. Top: Predicted β s with 95% confidence intervals. Bottom: Variance of predicted β s.

4.3 Ridge, Lasso and λ -dependence

In fig. 7 and fig. 8 we explore how the MSE changes as a function of the hyperparameter λ when fitting the data using Ridge regression and Lasso regression respectively.

First and foremost, we see that when using Lasso regression, we need a significantly lower value for λ when compared to Ridge regression if we wish to reach the same accuracy as with regular OLS. For Ridge, we see that for $\lambda < 10^1$, we get approximately the same precision as with regular OLS, while with Lasso we need a $\lambda < 10^{-4}$ to reach sufficient accuracy.

Error metrics for different values of λ using Ridge regression

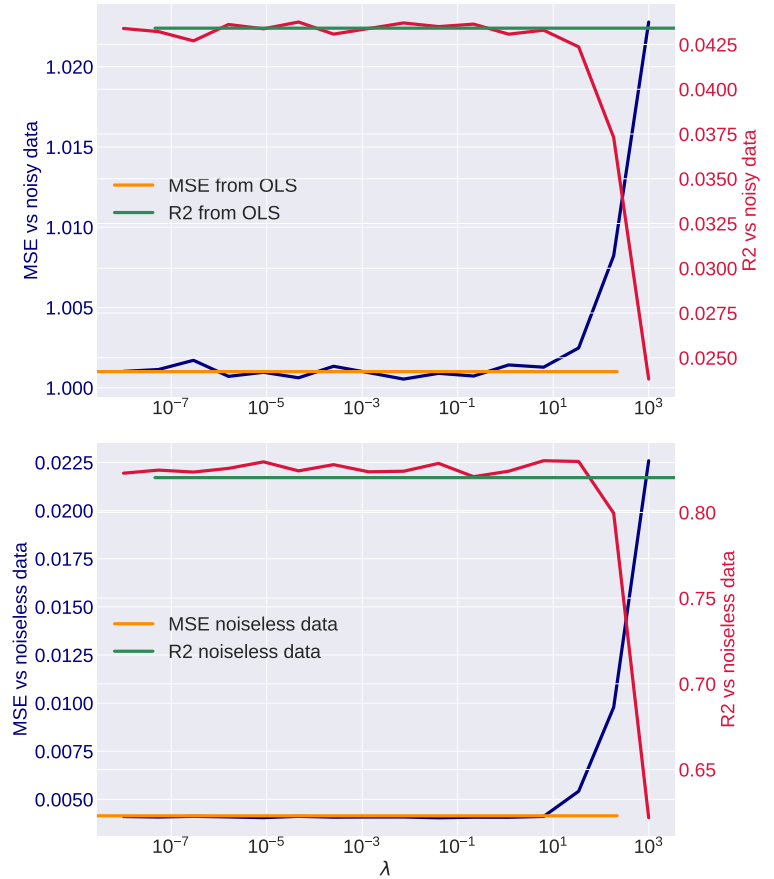


Figure 7 – Plots showing how the MSE- and R2 score behaves for increasing values of λ . The upper plot illustrates these metrics when comparing the predicted fit from Ridge regression to the Franke function with added noise. The lower plot shows the behavior of these metrics when comparing the fit to the noiseless Franke function. The lines represent the MSE and R2-score acquired from OLS regression.

4. RESULTS

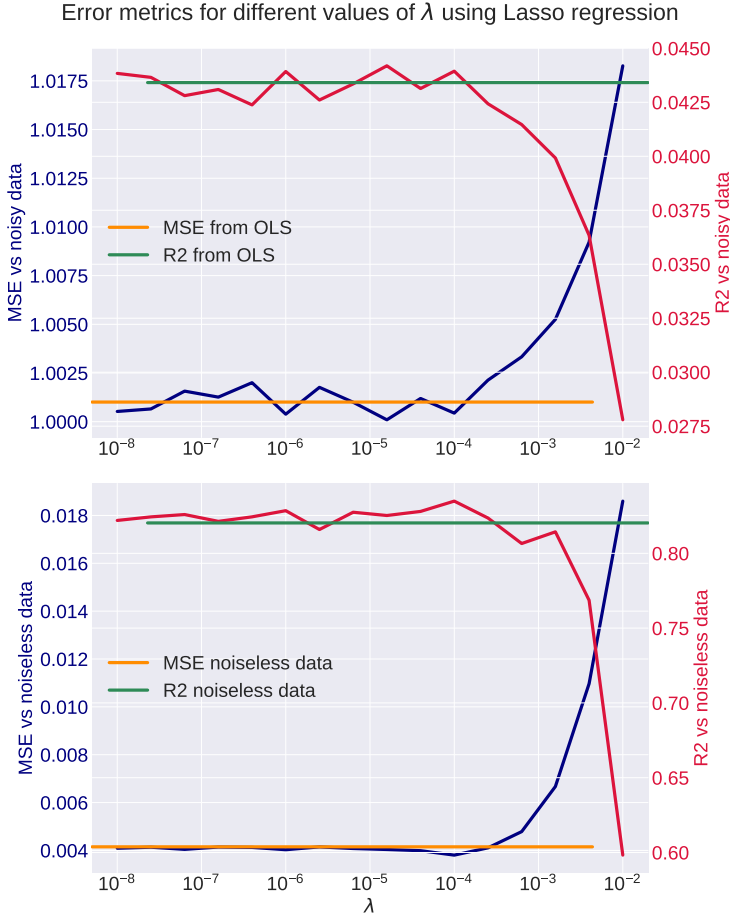


Figure 8 – Plots showing how the MSE- and R2 score behaves for increasing values of λ . The upper plot illustrates these metrics when comparing the predicted fit from Lasso regression to the Franke function with added noise. The lower plot shows the behavior of these metrics when comparing the fit to the noiseless Franke function. The lines represent the MSE and R2-score acquired from OLS regression.

4.3.2 Terrain data

As with the generated Franke data, we wish to explore how the error metrics behave as a function of λ for both Ridge and Lasso regression on the terrain data set. The results of this can be seen in fig. 9 and fig. 10. Unexpectedly, we see the same trend here, where we need a much lower value for λ to reach sufficient accuracy with Lasso compared to Ridge. When using Ridge regression we reach OLS error metrics when $\lambda < 10$, while Lasso doesn't converge until $\lambda < 10^{-7}$.

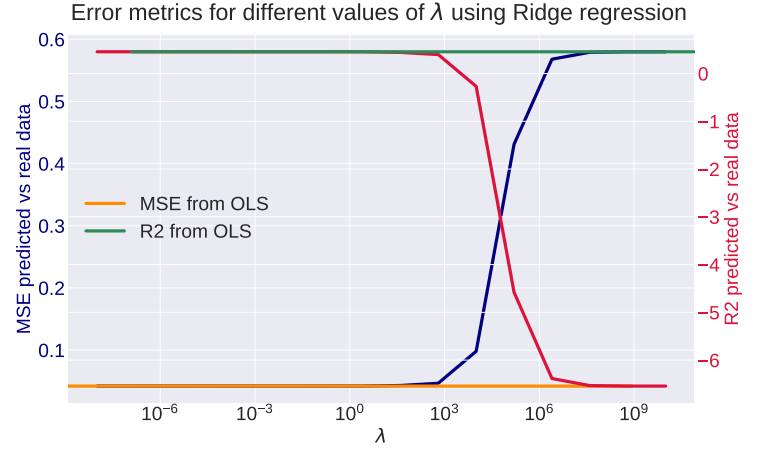


Figure 9 – Plot showing how the MSE and R2-score behaves for different values of λ when doing Ridge regression on the terrain data. The lines represent the error metrics obtained from regular OLS.

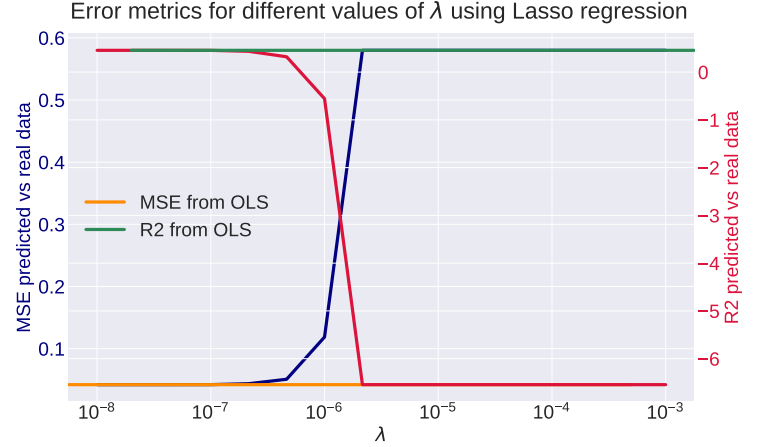


Figure 10 – Plot showing how the MSE and R2-score behaves for different values of λ when doing Lasso regression on the terrain data. The lines represent the error metrics obtained from regular OLS.

4.4 Bias Variance Tradeoff

4.4.1 Franke Data

As discussed in section (...), a high complexity model might lead to overfitting to the training data. In our case, complexity corresponds to the polynomial order of our model. In fig. 11, the train and test MSE of Franke data is shown, for all three regression methods. Clear signs of overfitting (high model variance) is shown for OLS at higher polynomial orders. This manifests itself as an increase in testing error, while the training error still decreases. As discussed in section 2.8, this is due to the model finding false trends in the training data, that doesn't actually exist in the underlying model. A high complexity model will fit to these trends, which won't exist in the testing set, and therefore lead to an increased error.

4. RESULTS

The trend is also visible in Ridge regression, but with a substantially lower slope in the error, indicating a stronger resistance to overfitting. Lasso regression is virtually immune to overfitting at higher complexities, having almost no slope at all. All models seem to reach a testing error minimum around polynomial order 5, increasing for higher orders, but at vastly different rates.

It's however still interesting to note that under an ideal complexity (i.e. poly order 5), OLS outperforms the two other models, with a slightly lower testing error.

Another interesting observation is that the testing error on OLS and Ridge fall below 1 around polynomial order 7.5. This itself is an indication of overfitting, as our data is generated with $MSE = 1$ by design. The noiseless Franke's function, the function we're trying to reproduce, has an MSE of 1 compared to our data. Gaining an MSE lower than this indicates that we're finding trends that doesn't exist in the underlying model. Lasso interestingly converges on a training error of exactly 1.

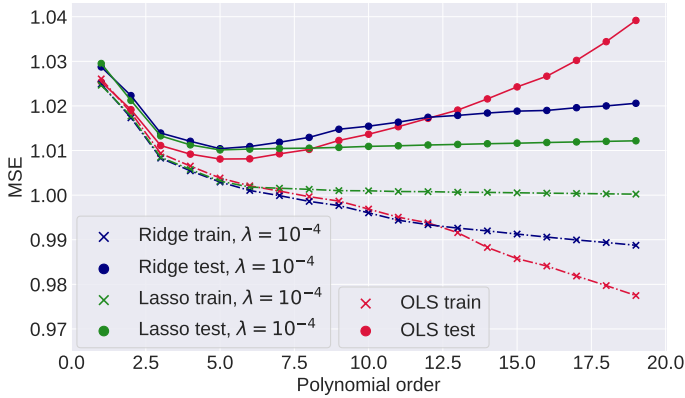


Figure 11 – MSE of testing and training data for OLS, Ridge, and Lasso as function of polynomial order, using K fold validation on the Franke data. $\lambda = 10^{-4}$ was used for Ridge and Lasso. We see that Lasso keeps a much more stable MSE for increased complexity, while the MSE on the OLS testing data rapidly diverges, indicating an overfitted model. Ridge falls somewhere in between.

Using the bootstrap method as discussed in section 3.4.4 we are able to calculate the separate components in the MSE to get a deeper insight into the evolution with polynomial order. Here we include only the result from OLS, while similar results for Ridge and Lasso may be found in appendix B. In fig. 12 are the results for the Franke data set, and the results for the terrain data is included in fig. 14. As expected we see how the bias is high for low complexity, as a result of underfitting. For increasing complexity the variance start increasing and the bias drops, which leads to a minimum in the MSE at polynomial order $\sim 4,5$ (for the Franke data). At higher complexity the model start overfitting, which is evident in the increase of the variance leading to an increase in the MSE. We note however that we also get a similar increase in the bias for higher complexity. This is not expected and the reason is not understood. Also, as in this case the irreducible error is known we have subtracted it from the MSE and the bias.

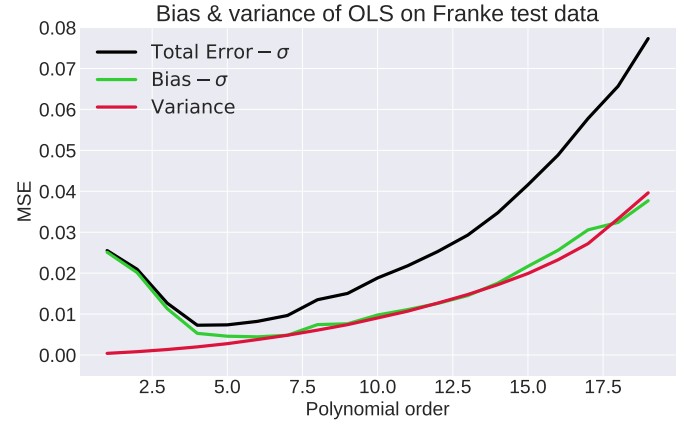


Figure 12 – The bias variance tradeoff as a function of polynomial order on the Franke test set. Here we can clearly see how the underfitting at low complexity is evident in the bias. While for higher complexity the model start overfitting and the variance of the model increase. However, we note that the bias also start increasing along with the variance for higher complexity which is unexpected. The reason for this is still unknown.

4.4.2 Terrain Data

Looking at the real terrain data, the picture becomes a bit more complicated due to a number of factors, mostly related to the fact that this dataset contains considerably more points, leading to numerical limitations on the complexity.

Figure 13 shows the train and test error of OLS and Ridge on the terrain data downsampled by 8×8 . Here we are, up to seemingly arbitrary complexity, unable to provoke overfitting on either methods, which both remain constant at higher orders. Our theory is that we have reached a level of complexity which is too numerically unstable and ill-conditions to give reliable results. Even if there in theory was a complexity which provoked overfitting in the methods, we might be unable to see them due to suppression of the higher order terms in the design matrix, as discussed in section (...). After a certain order, the columns of the design matrix become virtually zero-filled, compared to earlier columns. Increasing the complexity therefore has no effect on the solution.

Another observation is that OLS seems to substantially outperform Ridge in this scenario. This is probably due to the fact that the small λ added to the Hessian matrix ends up disturbing higher order terms, which are themselves very small. This is of course some of the idea behind Ridge, but in this scenario where overfitting isn't an issue, it actually ends up hurting the predictions.

Using K -folding it seems we are not able to provoke overfitting, but using our algorithm for calculating the variance and bias using the bootstrap method and OLS we can see the similar trend we saw for the Franke data in fig. 14.

4. RESULTS

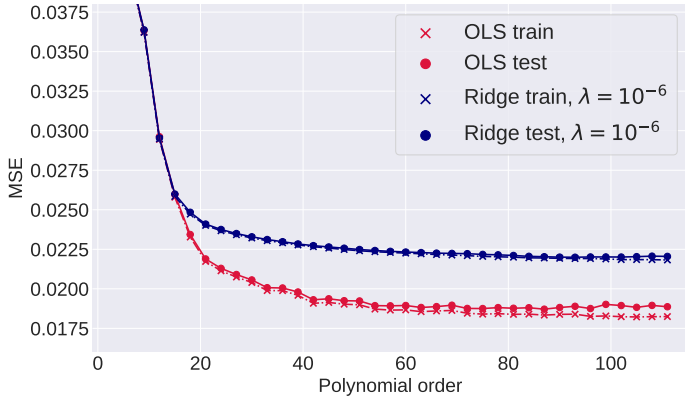


Figure 13 – MSE of testing and training data for OLS and Ridge as function of polynomial order, using K fold validation on Terrain data, downsampled by 8×8 . There is no sign of overfitting on either model, even up to very high orders, which might indicate that overfitting occurs at a level of complexity we can't numerically represent, due to suppression of high order terms in the design matrix.

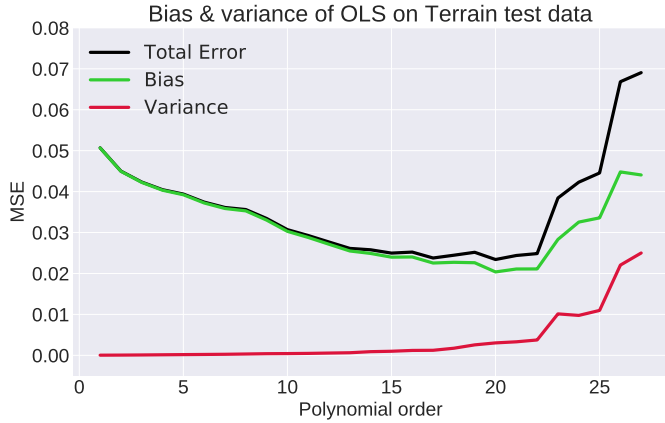


Figure 14 – The bias variance tradeoff as a function of polynomial order on the terrain test set. Like for the Franke data in fig. 12 we can see how the bias is prominent at low complexity, while for higher complexity the variance of the model increase. Here too the bias start increasing for higher complexity which is unexpected.

Like discussed in section (...), this would blow up reversely if our explanatory variables were in a range exceeding zero. In fig. 15, we have done exactly this. Analytically, this should leave the solutions exactly the same. We see, however, that, unlike fig. 13, both the train and test MSE blows up around order 25. If this was an overfit, the training error would remain low. The exploding training error highly indicates a numerical instability. Like discussed in section (...), the low order columns of the design matrix is now suppressed, instead of the latter, having an entirely different effect on the results.

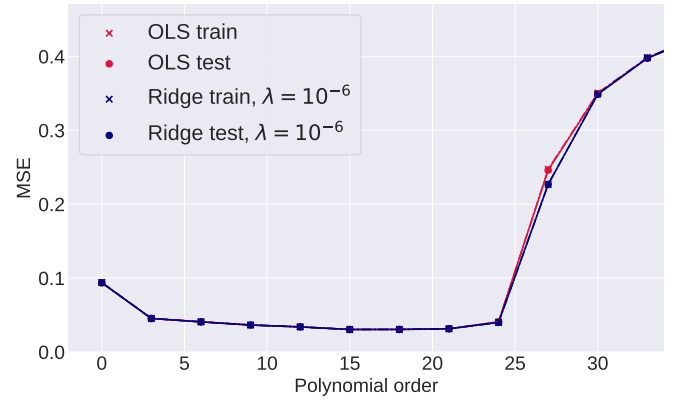


Figure 15 – Same setup as fig. 13, except explanatory variables are shifted to $x, y \in [0, 2]$. This causes the MSE to diverge at around polynomial order 25.

Another indication of the numerical instability can be seen from looking at the condition number (discussed in section (...)) of the design matrix, and the Hessian matrix (which is the one we are actually inverting). In fig. 16 the condition number of both matrices are plotted, showing the large numerical instability of higher order models.

An interesting correlation is that the conditioning number for $\lambda = 10^{-6}$ stops increasing exactly at the same point (order 25) that OLS starts outperforming Ridge, as we saw in figure fig. 13. Apparently, from order 25 and onward, the λ value dominates the size of the columns, causing further columns to add less to the predictive capabilities, but also produce less numerical instability.

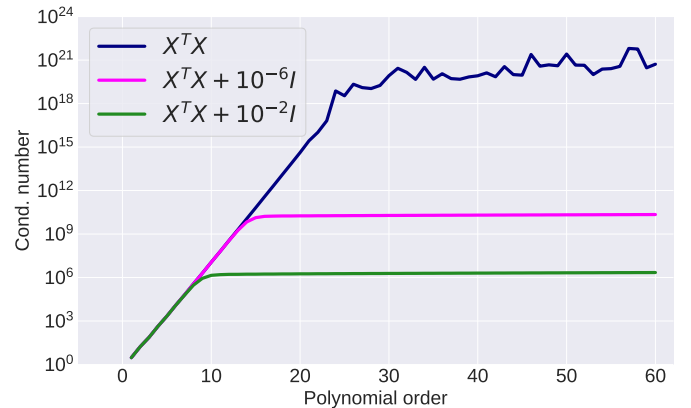


Figure 16 – Condition number of Hessian matrix as function of polynomial order. Hessian matrix for Ridge, with λ of 10^{-2} and 10^{-6} is also shown. The condition number, and thereby numerical instability, grows large for large polynomial orders, but stabilizes around 10^{20} . The lambdas have a clear suppression effect on the condition number, stabilizing the condition number much earlier.

Figure fig. 17 shows the same conditioning plot, but for the explanatory variables in the interval $[0, 2]$, just as we looked at in fig. 15. Here the reason for the diverging behavior becomes even more obvious. The conditioning number of the Hessian matrix not only increases a lot quicker (it is 10^{16} at poly order 10, to the 10^7 in the former case), but it also never stops diverging.

5. CONCLUSION

Adding a λ term now also have virtually no effect, which explains why Ridge also diverges in a similar fashion to OLS.

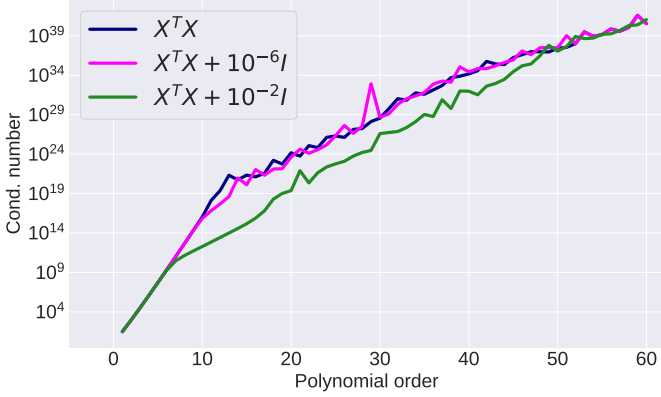


Figure 17 – Figure with same setup as fig. 16, except explanatory variables are shifted to $x, y \in [0, 2]$. The condition number now never stops increasing, and the lambdas no longer have the same suppression effect.

We are able to provoke overfitting without running into the aforementioned problems by reducing the number of datapoints. Down-sampling the image by 32×32 , we see from fig. 18 that the testing error of OLS increases from its lowest point around order 25. From there, the error diverges very rapidly, especially from order 45 and higher. This overfit increases much more rapidly than it did for the Franke data, which was more gradual. A good theory as to why is that, downsampled by 32×32 , the image contains only 112×56 pixels. The error diverges when the polynomial order starts approaching the number of points on the shortest axis. We know that a polynomial of degree N is capable of perfectly fitting a dataset of N points. This allows for a massive overfitting along the x-axis.

Ridge is a lot more resistant, but we see a slight increase in higher orders, which we didn't see for the less downsampled image.

References

- [Devore und Berk 2012] DEVORE, Jay L. ; BERK, Kenneth N.: *Modern Mathematical Statistics with Applications*. Second Edition. Springer, 2012

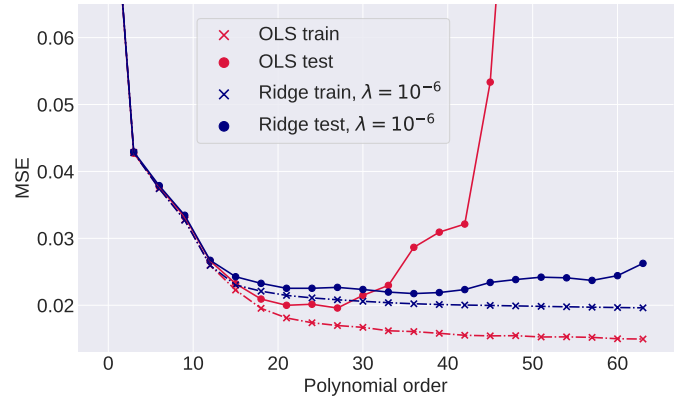


Figure 18 – Same setup as fig. 13, except image is downsampled by 32×32 . OLS heavily overfits for polynomial orders above 30. Ridge seems almost unaffected by the increased complexity, holding a relatively steady training MSE.

4.5 Best Fit Terrain

Having studied the implications of both the model complexity and hyperparameters, we can safely conclude that, as long as we're dealing with a large amount of data, OLS with a high order polynomial gives the best fit. We therefore employ a 80 degree polynomial on the terrain data, downsampled only by 4×4 . In fig. 19 we see the result, together with an absolute difference plot. We see that the fit does a good job at capturing the large scale structures of the data, while struggling with the finer details. The final fit results in an MSE of 0.0184, and a R2 score of 0.761.

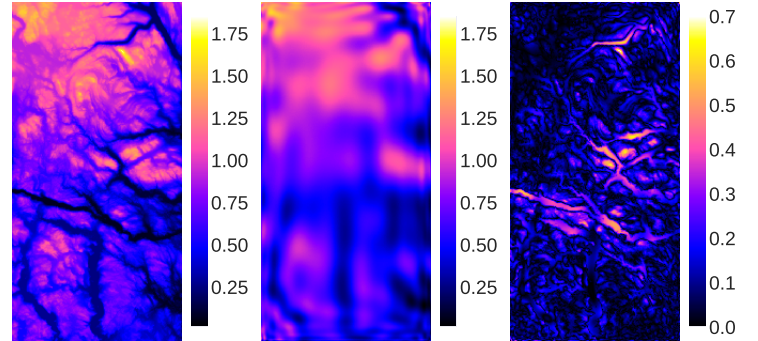


Figure 19 – Left: Original terrain image, downsampled by 4×4 . Center: Best fit OLS, polynomial order 80, with K fold validation. Right: Absolute difference between the former two. We clearly see that the fit retains the larger trends of the image, but struggles to capture the more fine-structured qualities.

5 Conclusion

Appendices

A Bias-Variance derivation

In this section we will show the derivation of eq. (25). We start by the definition of mean squared error, eq. (19). This is manipulated by adding and subtracting the expectation value of our fit inside the expression. By doing this we can order the equation in a way that decompose the MSE in to two contributions; the bias and the variance. In this derivation the following notation is used;

- The data is expressed as $\mathbf{y} = \mathbf{f} + \boldsymbol{\epsilon}$, where \mathbf{f} is deterministic and $\boldsymbol{\epsilon}$ is stochastic noise from the normal distribution $\boldsymbol{\epsilon} \sim N(0, \sigma^2)$. This means that $\mathbb{E}[\mathbf{f}] = \mathbf{f}$, $\mathbb{E}[\boldsymbol{\epsilon}] = 0$, $\text{Var}(\boldsymbol{\epsilon}) = \mathbb{E}[\boldsymbol{\epsilon}^2] = \sigma^2$
- The fit is expressed as $\tilde{\mathbf{y}} = \mathbf{X}\boldsymbol{\beta}$

$$\frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2 = \mathbb{E}[(\mathbf{y} - \tilde{\mathbf{y}})^2] = \mathbb{E}[(\mathbf{f} + \boldsymbol{\epsilon} - \tilde{\mathbf{y}} + \mathbb{E}[\tilde{\mathbf{y}}] - \mathbb{E}[\tilde{\mathbf{y}}])^2]$$

expanding the parenthesis:

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

using $\mathbb{E}[\boldsymbol{\epsilon}] = 0$:

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

$\mathbb{E}[x]$ is a linear operator:

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

Now we have separated out the three terms we want, so for this to be equal the expression in eq. (25), the last term must be equal to zero. To show this we factor out the deterministic terms out of the expectation value.

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

$$\begin{aligned} \text{where } \mathbb{E}[\mathbb{E}[\tilde{\mathbf{y}}]] &= \mathbb{E}[\tilde{\mathbf{y}}] \\ &= 0 \end{aligned}$$

So the MSE can be expressed, inserting the summation form of the expectation value

$$\begin{aligned} \text{MSE}(\mathbf{y}, \tilde{\mathbf{y}}) &= \mathbb{E}[(\mathbf{f} - \mathbb{E}[\tilde{\mathbf{y}}])^2] + \mathbb{E}[(\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])^2] + \sigma^2 \\ &= \frac{1}{n} \sum_{i=0}^{n-1} (\mathbf{f} - \mathbb{E}[\tilde{\mathbf{y}}])^2 + \frac{1}{n} \sum_{i=0}^{n-1} (\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])^2 + \sigma^2 \end{aligned}$$

B Additional Bias-Variance Results

Here we present additional results from our bias variance tradeoff experiments which we chose to not to include with the rest of the results to reduce cluttering. As discussed for the results from OLS, the bias is high for low complexity which leads to a high MSE. Increasing complexity reduces the bias which leads to an optimal polynomial order. In the case of OLS we could see how the variance started to increase when the model started overfitting, but using Ridge and Lasso we are not able to run calculations with high enough polynomial order to make this trend evident, at least not for the terrain data where the MSE and the bias is almost equal. For the Franke data we see how the variance start to increasing ever so slightly for higher complexity, and still the bias is also increasing which is unexpected.

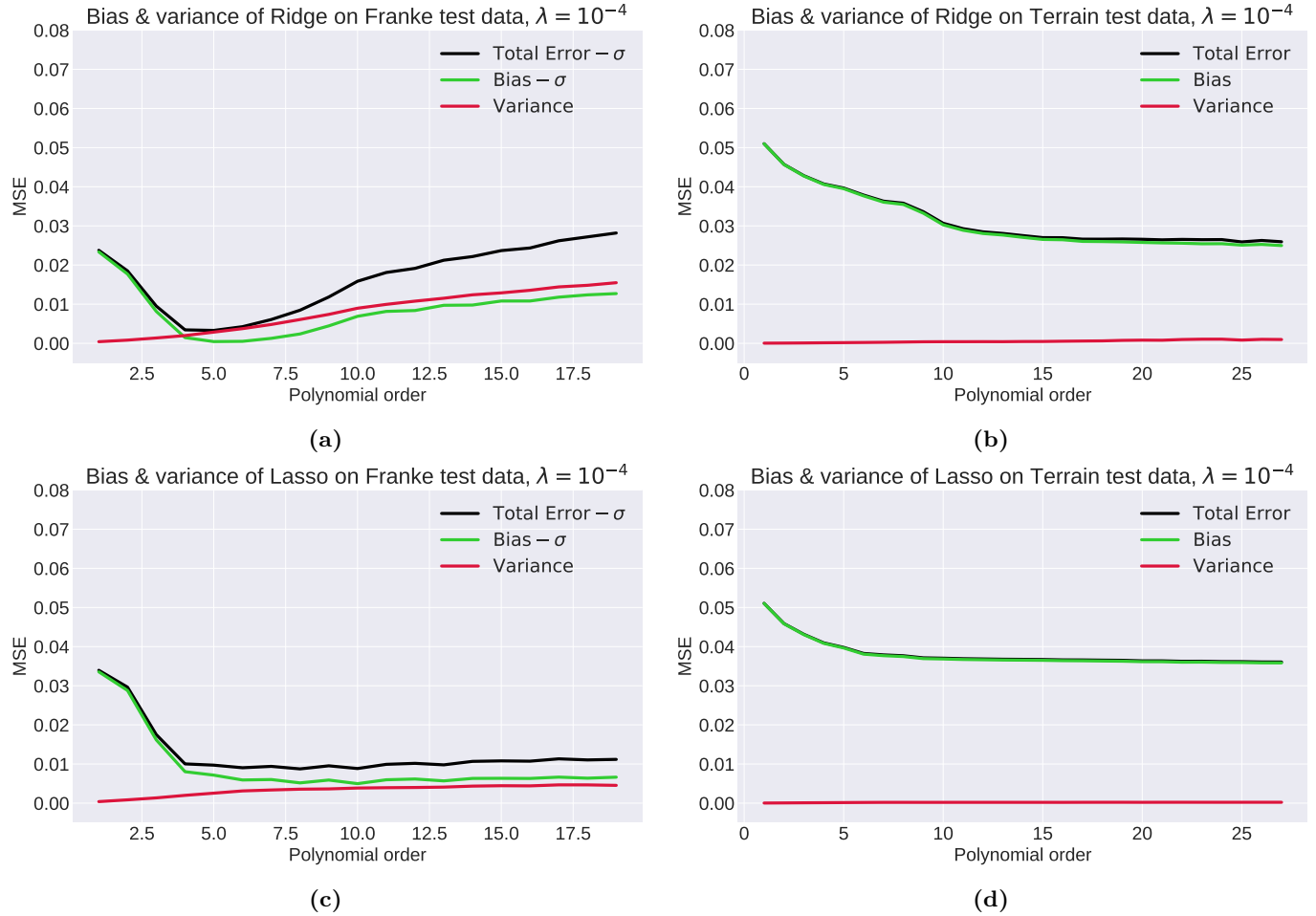


Figure 20 – Bias-variance results for Ridge in (a) and (b), Lasso in (c) and (d), all run with $\lambda = 10^{-4}$. For the Franke data the irreducible error is subtracted from the bias and the MSE.