

Project 1: Regression Analysis and Resampling Methods

Author:
Håkon FOSSEIM
October 27, 2021

In this project, we have assessed the performance of three different linear regression methods, namely Ordinary Least Squares (OLS), Ridge- and Lasso regression. These models were tasked with fitting a bivariate polynomial to three dimensional data points originating from either the Franke function or real terrain.

The Bootstrap and Cross-Validation resampling techniques were used to aid in our model assessment, the results of which were used to select the optimal regression method and its parameters.

For the generated Franke function data, the optimal OLS and Ridge models were found to perform very similarly in the parameter space of our search. The optimal OLS model was found to be at $p = 8$, while Ridge gave very similar results in a range of values. Both of these methods outperformed Lasso regression for our chosen level of noise.

For the terrain data, we found that Ridge regression using a polynomial of degree $25 \leq p \leq 35$ and $\lambda \approx -11$ typically slightly outperforms the optimal OLS model given by $p = 18$, but at approximately four times the computational cost. For this reason, this OLS model was chosen as the most suited to fitting the terrain data.

I. Introduction

Our goal in this report is to apply three different linear regression models, namely OLS, Ridge and Lasso, to two different data sets and assess their performance. The first dataset is generated by adding noise to the Franke function. The second uses real terrain data¹ in Norway. To perform this assessment, we calculate their bias, variance, MSPE and R2test score. These estimators will also allow us to gain deeper insights into our regression mod-

els.

However, due to the noise in our data, they will also possess some variance. We therefore use the bootstrap resampling technique to get more precise measurements as well as cross validation, which is a resampling technique better suited for small data sets.

¹ terrain 1, found at [3]

II. Theory

A. Linear Regression models

We are interested in approximating the relationship between a set of explanatory variables² $[x_0, x_1, \dots, x_{n-1}]$ and an unknown target function f . We are in possession of a noisy dataset given by $y(x_i) = f(x_i) + \varepsilon_i$ where ε is the noise in our data.

Linear regression tries to model the target function as a linear combination of the explanatory variables raised to various powers. For a given data point y_i , it creates a model from the corresponding x_i in the following manner

$$\tilde{y}(x_i) = \beta_0 + 0x_i^1\beta_1 + x_i^2\beta_2 + \dots + x_i^n\beta_n \quad (1)$$

Doing this for every y_i leads to the following matrix equation

$$\begin{bmatrix} \tilde{y}_0 \\ \tilde{y}_1 \\ \vdots \\ \tilde{y}_{n-1} \end{bmatrix} = \begin{bmatrix} 1 & x_0 & x_0^2 & \dots & x_0^n \\ 1 & x_1 & x_1^2 & \dots & x_1^n \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n-1} & x_{n-1}^2 & \dots & x_{n-1}^n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{bmatrix} \quad (2)$$

Which we denote as

$$\tilde{\mathbf{y}} = \mathbf{X}\boldsymbol{\beta} \quad (3)$$

Where \mathbf{X} is called our design matrix.

To find the optimal values of the regression coefficients, $\boldsymbol{\beta}$, a cost function $C(\mathbf{y}, \tilde{\mathbf{y}})$ is used. This cost function measures in one way or another how far off a hypothetical model $\tilde{\mathbf{y}}(\boldsymbol{\beta})$ is from the target function, and so the linear regression picks the model whose regression coefficients minimize this cost function, i.e.:

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

Where

$$\hat{\boldsymbol{\beta}} = \operatorname{argmin}_{\boldsymbol{\beta}} C(\mathbf{y}, \tilde{\mathbf{y}}(\boldsymbol{\beta})) \quad (5)$$

What defines a specific linear regression model is the way its cost function is implemented.

Without knowing better, one might suggest a cost function which measures the absolute distance to from the model to the target function, i.e. the mean absolute error(MSA), given by

$$\text{MSA}(\mathbf{X}, \boldsymbol{\beta}) \equiv \frac{1}{n} \sum_{i=0}^{n-1} |y_i - \tilde{y}_i| = \frac{1}{n} \sum_{i=0}^{n-1} |y_i - \mathbf{X}_{i,*}\boldsymbol{\beta}| \quad (6)$$

However, this function suffers from the drawback that it is not differentiable with respect to $\boldsymbol{\beta}$ and so we cannot find a closed form expression for $\hat{\boldsymbol{\beta}}$. A better choice is the mean squared error (MSE), defined as

OLS :

$$C(\boldsymbol{\beta}) = \text{MSE} \equiv \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2 = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \mathbf{X}_{i,*}\boldsymbol{\beta})^2 \quad (7)$$

This is the cost function used by the OLS. Since we are using the square distance, it will see models with outliers as particularly unfavorable, which is good because these outliers usually stem from spikes in the noise of our data.

Furthermore, unlike the MSA, the MSE is differentiable w.r.t. $\boldsymbol{\beta}$, meaning we can get a closed form expression for the optimal regression parameters, $\hat{\boldsymbol{\beta}}$.

Differentiating Eq.(7) with respect to $\boldsymbol{\beta}$ and setting the expression equal to zero gives the optimal regression parameters:

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

Another useful measure to assess a model's performance is the R^2 score, also called the coefficient of determination. It is defined as

$$R^2 = 1 - \frac{\sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2}{\sum_{i=0}^{n-1} (y_i - \bar{y})^2} \quad (9)$$

Where

$$\bar{y} = \frac{1}{n} \sum_{i=0}^{n-1} y_i \quad (10)$$

We see from Eq.(9) that a perfect fit of our model to the data points yields an R^2 score of 1. An R^2 score of 0 indicates that our model performed no better than the average value of the data points, and a negative score means it did worse and simply taking the average of the data is a better model.

1. Ridge regression Ridge regression is defined by the following cost function:

$$\begin{aligned} C(\boldsymbol{\beta}) &= \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \mathbf{X}_{i,*}\boldsymbol{\beta})^2 + \lambda \sum_{j=1}^{p-1} \beta_j^2 \\ &= \frac{1}{n} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta} \\ &= \frac{1}{n} (\mathbf{y}^T - \boldsymbol{\beta}^T \mathbf{X}^T) (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta} \end{aligned} \quad (11)$$

The optimal $\boldsymbol{\beta}$ -parameters are given by $\operatorname{argmin}_{\boldsymbol{\beta}} C(\boldsymbol{\beta})$, i.e. solving for $\partial C(\boldsymbol{\beta}) / \partial \boldsymbol{\beta} = 0$

² Also called independent variables

Since \mathbf{y} is not a function of $\boldsymbol{\beta}$, we get

$$\begin{aligned} & -\frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{y}^T \mathbf{X} \boldsymbol{\beta} - \frac{\partial}{\partial \boldsymbol{\beta}} \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y} + \frac{\partial}{\partial \boldsymbol{\beta}} \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} + n\lambda \frac{\partial}{\partial \boldsymbol{\beta}} \boldsymbol{\beta}^T \boldsymbol{\beta} \\ &= 0 \end{aligned} \quad (12)$$

And by using the following identity[2],

$$\frac{\partial \mathbf{a}^T \mathbf{A} \mathbf{a}}{\partial \mathbf{a}} = (\mathbf{A} + \mathbf{A}^T) \mathbf{a} \quad (13)$$

the third term in Eq.(12) becomes

$$\frac{\partial}{\partial \boldsymbol{\beta}} \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = 2 \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} \quad (14)$$

furthermore, we also have that[2]

$$\frac{\partial \mathbf{b}^T \mathbf{a}}{\partial \mathbf{a}} = \mathbf{b} \quad (15)$$

Which can be applied to the first term to yield

$$\frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{y}^T \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^T \mathbf{y} \quad (16)$$

Now, since every term in Eq.(11) is a scalar, we have

$$(\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y}) = (\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y})^T = \mathbf{y}^T \mathbf{X} \boldsymbol{\beta} \quad (17)$$

meaning that the first and second terms in Eq.(12) are equal, leaving us with

$$-2 \mathbf{X}^T \mathbf{y} + 2 \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} + 2n\lambda \boldsymbol{\beta} = 0 \quad (18)$$

$$(\mathbf{X}^T \mathbf{X} + n\lambda \mathbf{I}) \boldsymbol{\beta} = \mathbf{X}^T \mathbf{y} \quad (19)$$

If we 'absorb' n into λ , we are left with the following analytical expression for the optimal regression coefficients

$$\hat{\boldsymbol{\beta}}_{Ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (20)$$

B. Lasso regression

Lasso regression is defined by the following cost function:

$$C(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \mathbf{X}_{i,*} \boldsymbol{\beta})^2 + \lambda \sum_{j=1}^{p-1} |\beta_j| \quad (21)$$

LASSO is an acronym for Least Absolute Shrinkage and Selection Operator and it performs best for data with many features. This is because, like Ridge, it is designed

to prevent overfitting by suppressing the regression coefficients $\hat{\boldsymbol{\beta}}$. However, unlike Ridge, one cannot find a closed form solution for the optimal regression coefficients. For this, one needs to use numerical methods.

It also has the ability to force the β values to exactly zero, which Ridge regression cannot do. We saw this happen for $\log \lambda > -1.21$ with $n = 20$ for the Franke function, where Lasso regression returns only a flat surface at 0. Varying n also varies at which λ this happens.

C. Confidence interval of $\hat{\boldsymbol{\beta}}$

Assuming the estimators $\hat{\beta}_i$ are normally distributed with mean $\mathbb{E}[\hat{\beta}_i]$ and variance $\sigma_{\beta_i}^2$, i.e.

$$\hat{\beta}_i \sim \mathcal{N}(\mathbb{E}[\hat{\beta}_i], \sigma_{\beta_i}^2) \quad (22)$$

, then we can say that there is a $(1 - \alpha)100\%$ likelihood that $E[\hat{\beta}_i]$ lies in the interval $[\hat{\beta}_i - Z\sigma_{\beta_i}, \hat{\beta}_i + Z\sigma_{\beta_i}]$. In other words

$$P(\hat{\beta}_i - Z\sigma_{\beta_i} \leq \mathbb{E}[\hat{\beta}_i] \leq \hat{\beta}_i + Z\sigma_{\beta_i}) = 1 - \alpha \quad (23)$$

Where Z is how many standard deviations you have to go to the right or left of the mean of the standard normal to get an area of $(1 - \alpha/2)$. It is therefore given by

$$Z = \phi^{-1}\left(1 - \frac{\alpha}{2}\right) \quad (24)$$

Where φ^{-1} is the inverse cumulative distribution function of the standard normal distribution. For a 95% CI, i.e. $\alpha = 0.05$, we get $Z \approx 1.96$.

It can be shown[1] that the variance of the estimator $\hat{\boldsymbol{\beta}}$ is given by

$$\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2 \quad (25)$$

Where \mathbf{X} is the design matrix and $\sigma^2 = \text{var}(\epsilon_i)$. Furthermore, $\text{var}(\hat{\boldsymbol{\beta}})$ is defined as

$$\text{var}(\hat{\boldsymbol{\beta}}) = \begin{bmatrix} \text{var}(\beta_0) & \text{cov}(\beta_0, \beta_1) & \dots & \text{cov}(\beta_0, \beta_{p-1}) \\ \text{cov}(\beta_1, \beta_0) & \text{var}(\beta_1) & \dots & \text{cov}(\beta_1, \beta_{p-1}) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(\beta_{p-1}, \beta_0) & \text{cov}(\beta_{p-1}, \beta_1) & \dots & \text{var}(\beta_{p-1}) \end{bmatrix} \quad (26)$$

And so $\text{var}(\hat{\beta}_i) = \sigma_{\hat{\beta}_i}^2$ is given by the i -th diagonal entry of $(\mathbf{X}^T \mathbf{X})^{-1} \sigma^2$.

D. Bias-variance trade-off

We assume that our data is distributed according to

$$y_i = f(x_i) + \varepsilon_i \quad (27)$$

Where f is the target function to be estimated, and ε_i are i.i.d. according to

$$\varepsilon_i \sim \mathcal{N}(0, \sigma_\varepsilon^2) \quad (28)$$

Here y_i , x_i and ε_i are instances of the variables \mathbf{y} , \mathbf{x} and $\boldsymbol{\varepsilon}$, of which \mathbf{y} and $\boldsymbol{\varepsilon}$ are stochastic. The expression we want to decompose is given by

$$C(\mathbf{X}, \boldsymbol{\beta}) = \text{MSPE} \equiv \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2 \quad (29)$$

Now, since we are in fact dealing with continuous variables, the **MSPE** will be an approximation of $\mathbb{E}[(\mathbf{y} - \tilde{\mathbf{y}})^2]$. This approximation should improve as we increase n .

$$\begin{aligned} \text{MSPE} &\equiv \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2 \approx \mathbb{E}[(\mathbf{y} - \tilde{\mathbf{y}})^2] \quad (30) \\ &= \mathbb{E}[(f(\mathbf{x}) + \boldsymbol{\varepsilon} - \tilde{\mathbf{y}})^2] \quad (31) \end{aligned}$$

Adding and subtracting $\mathbb{E}[\tilde{\mathbf{y}}]$ yields

$$\begin{aligned} \mathbb{E}[(\mathbf{y} - \tilde{\mathbf{y}})^2] &= \mathbb{E}\left[\left((f(\mathbf{x}) - \mathbb{E}[\tilde{\mathbf{y}}]) - (\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}]) + \boldsymbol{\varepsilon}\right)^2\right] \\ &= \mathbb{E}[(f(\mathbf{x}) - \mathbb{E}[\tilde{\mathbf{y}}])^2] - \mathbb{E}\left[(f(\mathbf{x}) - \mathbb{E}[\tilde{\mathbf{y}}])(\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])\right] \\ &\quad + \mathbb{E}\left[\boldsymbol{\varepsilon}(f(\mathbf{x}) - \mathbb{E}[\tilde{\mathbf{y}}])\right] - \mathbb{E}\left[(\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])(f(\mathbf{x}) - \mathbb{E}[\tilde{\mathbf{y}}])\right] \\ &\quad + \mathbb{E}[(\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])^2] - \mathbb{E}\left[\boldsymbol{\varepsilon}(\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])\right] \\ &\quad + \mathbb{E}\left[\boldsymbol{\varepsilon}(f(\mathbf{x}) - \mathbb{E}[\tilde{\mathbf{y}}])\right] - \mathbb{E}\left[(\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])^2\right] + \mathbb{E}[\boldsymbol{\varepsilon}^2] \quad (32) \end{aligned}$$

If two stochastic variables \mathbf{x} and \mathbf{y} are independent, then the expectation value of their product is equal to the product of their expectation values, i.e.

$$\mathbb{E}[\mathbf{x}\mathbf{y}] = \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{y}] \quad (33)$$

In general, $\tilde{\mathbf{y}}$ and $\boldsymbol{\varepsilon}$ won't be independent and so (33) would not hold for these variables. This is because $\tilde{\mathbf{y}}$ depends on \mathbf{y} which again depends on the stochastic noise, $\boldsymbol{\varepsilon}$.

However, if we train the model $\tilde{\mathbf{y}}$ on a training set and then calculate the cost function $C(\mathbf{X}, \boldsymbol{\beta})$ on the test set only, then the noise in this set will not have affected our model, in which case we have

$$\mathbb{E}[\boldsymbol{\varepsilon}\tilde{\mathbf{y}}] = \underbrace{\mathbb{E}[\boldsymbol{\varepsilon}]}_{=0} \mathbb{E}[\tilde{\mathbf{y}}] = 0 \quad (34)$$

Since $f(\mathbf{x})$ and $\boldsymbol{\varepsilon}$ also are independent, the third, sixth and seventh terms in (32) are zero.

Expanding the third term yields

$$\begin{aligned} \mathbb{E}\left[(f(\mathbf{x}) - \mathbb{E}[\tilde{\mathbf{y}}])(\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])\right] &= \mathbb{E}[f(\mathbf{x})\tilde{\mathbf{y}}] - \mathbb{E}[f(\mathbf{x})]\mathbb{E}[\tilde{\mathbf{y}}] \\ &\quad - \mathbb{E}[\tilde{\mathbf{y}}]\mathbb{E}[\tilde{\mathbf{y}}] + \mathbb{E}[\tilde{\mathbf{y}}]^2 = \mathbb{E}[f(\mathbf{x})\tilde{\mathbf{y}}] - \mathbb{E}[f(\mathbf{x})]\mathbb{E}[\tilde{\mathbf{y}}] \\ &= \text{cov}(f(\mathbf{x}), \tilde{\mathbf{y}}) \quad (35) \end{aligned}$$

However, since \mathbf{x} is not a stochastic variable³, we can set it outside the expectation value operator:

$$\begin{aligned} \text{cov}(f(\mathbf{x}), \tilde{\mathbf{y}}) &= \mathbb{E}[f(\mathbf{x})\tilde{\mathbf{y}}] - \mathbb{E}[f(\mathbf{x})]\mathbb{E}[\tilde{\mathbf{y}}] \\ &= f(\mathbf{x})\mathbb{E}[\tilde{\mathbf{y}}] - f(\mathbf{x})\mathbb{E}[\tilde{\mathbf{y}}] = 0 \quad (36) \end{aligned}$$

The final expression for the MSPE then becomes

$$\begin{aligned} \text{MSPE} &\approx \mathbb{E}[(\mathbf{y} - \tilde{\mathbf{y}})^2] = \underbrace{\mathbb{E}\left[(f(\mathbf{x}) - \mathbb{E}[\tilde{\mathbf{y}}])^2\right]}_{\text{bias}^2} \\ &\quad + \underbrace{\mathbb{E}\left[(\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])^2\right]}_{\text{variance}} + \underbrace{\mathbb{E}[\boldsymbol{\varepsilon}^2]}_{\text{irr.err.}} \quad (37) \end{aligned}$$

The first term is the square of the bias, which is a measure of how well the interpolated model fits the target function. Low complexity models often have high bias, as they lack the flexibility to accurately capture the characteristics of the target function. We then say that the model suffers from *underfitting*. Note that the target function $f(\mathbf{x})$ is usually not available, since this is what we are trying to interpolate. In its place, one uses the response variable \mathbf{y} to get an estimate of the bias.

The second term measures the variance of the model, i.e. how much it varies around its mean. very high complexity models can fit the training data much better, and thus give low bias. However, it will also fit the stochastic noise in the training data, and thus pick up characteristics that aren't really existent in the target function. When the model is then evaluated on the test data, it performs poorly because it failed to adequately capture the relationship between \mathbf{f} and \mathbf{x} during training. We say that the model suffers from *overfitting*. In this case, the variance will typically be high, since the high complexity has allowed the model to fluctuate a lot in order to fit the noisy training data.

The last term is the irreducible error, and is the variance of the stochastic noise in the data. No matter how good of a model we generate, it can never predict the noise of yet-unseen data samples, and so this forms a lower bound on the prediction error.

³ i.e. \mathbf{x} is deterministic

E. Resampling methods

1. Bootstrap Each time a specific model is fitted to the data, we evaluate its performance through various estimators, $\hat{\theta}$. However, the data will likely contain some noise, which means these estimate will vary around their true values each time the model is fitted with data containing new noise.

The central limit theorem tells us that in the case that the noise is i.i.d., we can improve the precision of an estimate by replacing the estimator $\hat{\theta}$ with

$$\tilde{\hat{\theta}} = \frac{\hat{\theta}_0 + \hat{\theta}_1 + \dots + \hat{\theta}_B}{B} \quad (38)$$

Where $\tilde{\hat{\theta}}$ is the mean of $\hat{\theta}$ calculated on B data samples, each containing different noise.

The variance of this new estimator is then given by[1]

$$\text{Var} [\tilde{\hat{\theta}}] = \frac{\text{Var} [\hat{\theta}]}{B} \quad (39)$$

One time consuming and expensive way to calculate our new estimates would be to repeat the experiment B times with completely new data sets each time.

A much more practical way of generating "different" data sets is to randomly sample with replacement N data points from the original data set, from which we can calculate $\hat{\theta}_i$.

2. K-fold Cross-Validation As alluded to in IID, in order to assess the performance of a model on unseen data, one splits the original data set into a training set and a test set. The training data is then used to produce the model and its performance is gauged using the test data. This way we can assess how well it learned the features of the data rather than how good it is at fitting noisy data, which in the case of high amplitude noise will merely be a measure of its flexibility.

For small data samples, there is a higher chance that the test set given to the model will not be representative of the full data set. For example, the model might be assigned a test set which is a lot less or much more noisy than the training data, resulting in us either over- or underestimating how the model would perform on new data.

To remedy this, the cross validation algorithm splits the original data set into K samples, where 1 of these is to be used for testing and the remaining $K - 1$ samples are to be collectively used as training data.

This is done K times, but each time we pick one of the previously unchosen K samples as our test set. In this way, all the original data points will have been included once in the test and $K - 1$ times in the training set. The expected value and variance of the estimators are

calculated for each fold, and the mean of these quantities is returned. This averages out any anomalous estimates that any one of the splits might produce due to being assigned an unrepresentative test set.

F. Feature Scaling

Some machine learning algorithms may have reduced performance unless the features are distributed according to the standard normal distribution. An example of this is the K-nearest neighbours algorithm, which counts the K nearest neighbours in the feature space and performs classification according to which of the categories is most represented. To do this requires it to calculate the euclidean norm to its neighbours, and unless the features are scaled, this yields suboptimal results.

Furthermore, in using gradient descent to find the minimum of the cost function, it is important to scale the data. This is because if e.g. one of the features has a large magnitude compared with the others, the surface (in the case of two features) will be elongated along that feature's axis, making gradient descent zigzag its way to the minimum rather than taking a straighter path, thereby decreases its performance.

In our case the features have the same mean and standard deviation and so scaling them is somewhat redundant. We still scale our data⁴, since it is good machine learning practice.

G. The Franke function

The Franke's bivariate test function is commonly used as a test function in interpolation problems, and is given by

$$\begin{aligned} f(x, y) = & \frac{3}{4} \exp \left[-\frac{1}{4}(9x - 2)^2 - \frac{1}{4}(9y - 2)^2 \right] \\ & + \frac{3}{4} \exp \left[-\frac{1}{49}(9x + 1)^2 - \frac{1}{10}(9y + 1) \right] \\ & + \frac{1}{2} \exp \left[-\frac{1}{4}(9x - 7)^2 - \frac{1}{4}(9y - 3)^2 \right] \\ & - \frac{1}{5} \exp \left[-(9x - 4)^2 - (9y - 7)^2 \right] \end{aligned} \quad (40)$$

How to set up this interpolation problem given our two dimensional target function is described in the appendix.

⁴ by subtracting the mean of each column in the design matrix from every element in that column.

III. Methods

A. Program description

Our current program stems from an initial desire to be able to control a messy code and produce the desired plots with more ease. To do this, we represented all the variables needed in the form of bools and ints in the main() function, in what we will refer to as our control panel⁵. Doing so makes it easy to adjust them and re-run main.py in order to produce a different plot, testing a different part of our program in the process. It also makes it necessary to "reformat" these variables before usage in the rest of the program for the sake of readability and sanity. For example, the regression method is set as an int in the control panel;

```
regInt      = 0 #0=OLS,           1=ridge,           2=lasso
```

This int gets converted to the string regMeth before being passed on to the rest of the program. regMeth gets used by all three functions in resampling_methods.py to pick the regression method:

```
if(regMeth=='OLS'):
    beta_hat=OLS(X_train,z_train,skOLS)
if(regMeth=='ridge'):
    beta_hat=ridge(X_train,z_train,lmd)
if(regMeth=='lasso'):
    beta_hat=lasso(X_train,z_train,lmd)
```

It is also used repeatedly in plotters.py to provide descriptive plot titles and filenames. It would be much more tedious and error prone to remember which ints correspond to which methods etc., in addition to the code being unreadable and containing functions with tons of arguments.

And so, below the control panel, we have included the function reformatVariables(original variables), which lies in the scope of the control panel variables, and automatically takes its variables to be reformatted as lists and dictionaries etc. in reformat_variables.py.

It also generates the data to be interpolated based on the terrainBool and n , the number of points along the x- and y axis, which are specified in the control panel.

Each resampling method has its own list, which contains the name of the method as its first element along with the variables it needs to perform its task.

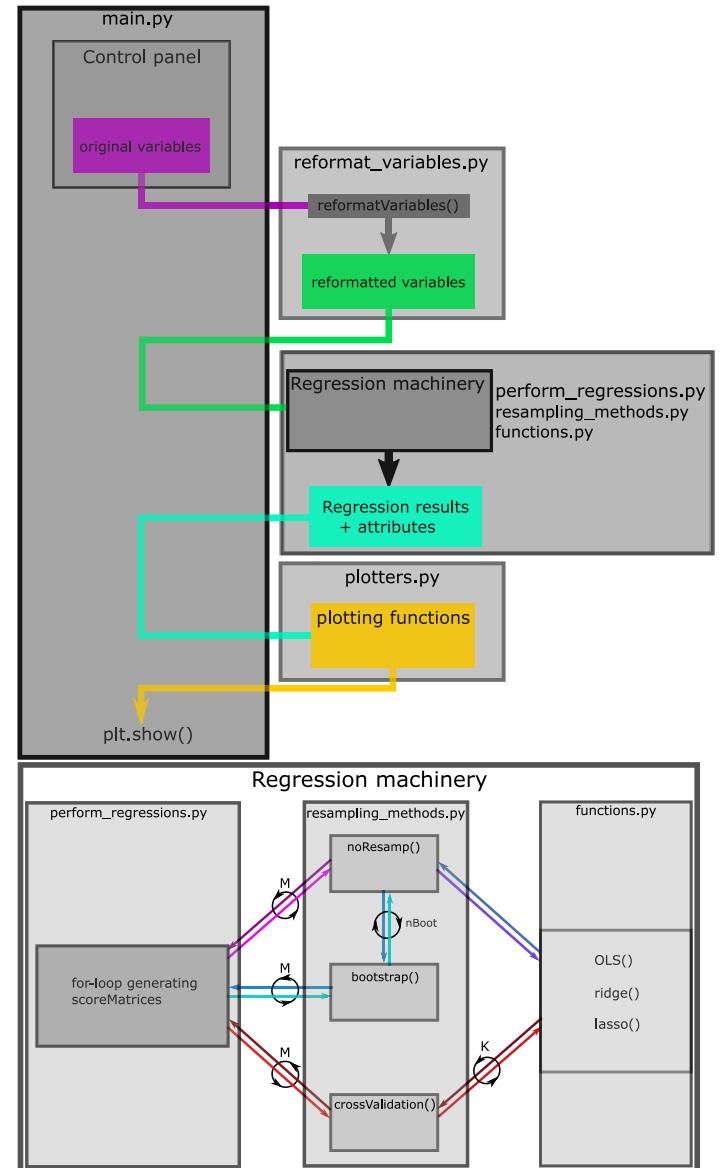


FIG. 1. Here we see the program flow of our code as well as the "regression machinery", which is tasked with returning the requested regression results. The outgoing values are represented by a lighter color than the returning results. The circles specify how many times the the function call occurs after we run main.py. Here M is the number of elements in each of the matrices contained in the scoreMatrices dict. K is the number of folds specified for Cross-Validaiton and $nBoot$ the number of bootstrap samples to be used. After each call to a regression function, the result is used to again call on the getScores() function in functions.py in order to calculate the errors/scores of that given regression iteration.

⁵ can be seen in the appendix

This is seen in reformat_variables.py⁶:

```
no_resamp=['no_resamp.',commPars]
bootstrap=['bootstrap', commPars,nBoot]
crossval=['crossval', commPars,skCV,shuffle,K]
resampMethods = [no_resamp, bootstrap, crossval]
resampMeth = resampMethods[resampInt]
```

Which takes resampInt from the control panel,

```
resampInt = 1
#0=no_resamp., 1=bootstrap, 2=crossval
```

, as an argument and determines which of the above lists to pass as an argument to perform_regression().

This list then determines the resampling method to be used, as well as providing the necessary variables to do so.

We can also specify which ranges of the hyperparameters λ , p and σ to iterate over. This determines what plot gets produced. To store these results, we initialize an "empty" dict containing only the scores⁷ specified by the user by passing in the integers corresponding to the scores seen in the dummyList below. This slightly reduces computational cost compared to calculating all scores each time.

```
def getScores(*args):
    dummyList = ['bias','variance','MSEtest',...]
    emptyScoreScalars = {}
    for i in range(0,len(args)):
        scoreName = dummyList[args[i]]
        emptyScoreScalars[scoreName]=0
    print("The following scores are being
calculated:\n", emptyScoreScalars,'\'\n' )
    return emptyScoreScalars
```

As seen, the values belonging to each score string are initially set to zero, but will be (repeatedly re-)filled when the regression method of choice calls the getScores() function in functions.py. After the scores are calculated, they get stored as an element in their corresponding score matrix in the dictionary scoreMatrices. This dictionary is created in perform_regression.py and will create entries with titles corresponding to those of emptyScoreScalars.

```
scoreMatrices = {}
for scoreName in emptyScoreScalars:
    scoreMatrices[scoreName] =\
    np.zeros((nLambdas,nOrders,nSigmas))
```

⁶ Here commPars refers to the variables xr,yr,z and the bools scaling and skOLS. Future improvements include implementing these common parameters in their own list.

⁷ We use the term score for any calculation done on the regression data. Due to the use of dicts, adding further scores to our program shouldn't be an issue.

The regression results from different lambdas are stored along the "y-axis" of this matrix, while the result from different orders and sigmas are stored along the x and z axes respectively. We store these λ , p and σ arrays as nested lists of the list hyperPars, which again is used as an entry in the dictionary calcAtts along with other things needed for plotting:

```
calcAtts = {'hyperPars': hyperPars,\n'plotTitleAtts':[regMeth,resampMeth,nBoot,K,n]}
```

As the for-loop over all hyperparameters is finished, and all 3D matrices are filled, the calculation results in the form of the dict scoreMatrices, is returned to main()⁸ and lastly sent to the plotting functions in plotters.py, which deduces the dimensions of our score-result matrices and generates figures with the relevant variables included in the title, so as to be reproducible by setting the control panel accordingly.

Note that we also included the script run_exercises.py which allows one to more easily reproduce our results. It is simply a series of if-statements with the variables set according to the exercise and "task", with the reformatter, regression machinery and plotter in their scope.

⁸ in addition to the calculation attributes, calcAtts.

IV. Results

A. The Franke function

a. Ordinary least squares

Throughout the project, our bias calculation were clearly off. They refused to go to zero, which in retrospect might make sense, since we were calculating it on the test data motivated by our "results" in the bias-variance section. Using the following expression

```
scoreVal = (np.mean((z_tilde-z_train)**2))
```

where $\text{scoreVal} = \text{bias}^2$ yields much better results, which also makes sense. As the model fits the training data better and better, this expression should go to zero. There is however no reason for why it should decrease with complexity if evaluated on the test data, since it has not seen the noise in that data, and won't do any better job of predicting the true test data plus noise than the true test data, which it does a horrible job of when overfitted.

We choose to divide the x- and y-axis into $n = 40$ points, yielding 1600 gridpoints. A larger n could be chosen to give better fits, at the cost of computational time. A fixed σ will be used throughout evaluating regression and resampling methods on the Franke function. We will be evaluating these methods by use of two dimensional polynomials of degrees $p = 0$ up to $p = 20$ to find which one is best suited to fit the Franke function at a specific level of noise. To find an appropriate noise level, we can study how the fit varies with σ for a polynomial fit of e.g. degree $p = 10$, as seen in Fig.2.

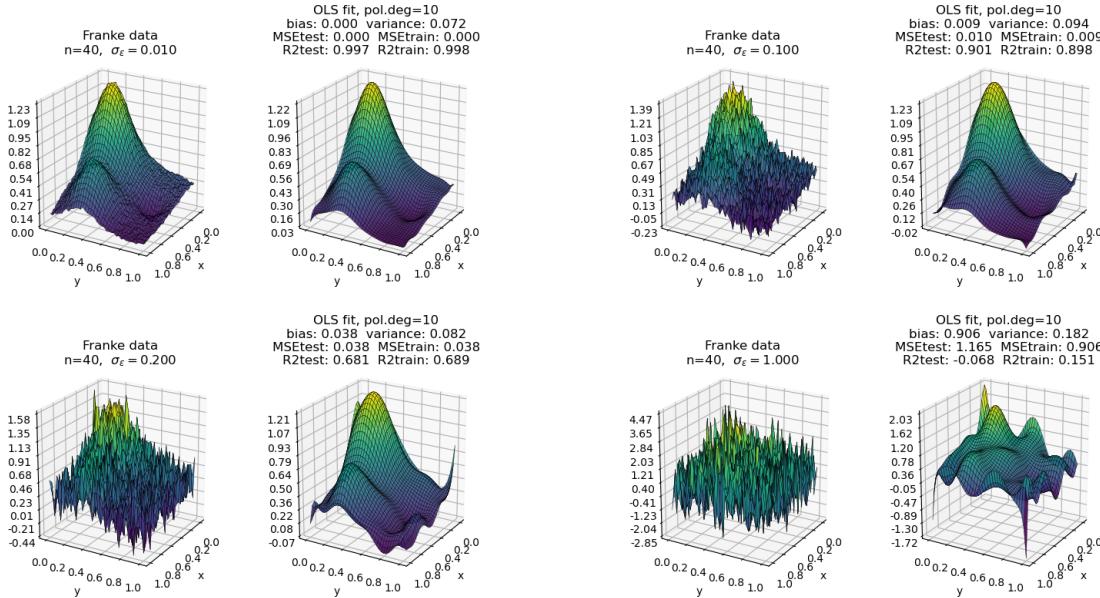


FIG. 2.

We see that for $\sigma = 0.001$, the noise amplitude is so low that it yields an almost perfect fit to the target function, the shape of which can be seen here. However, this isn't of interest to us, since we want to test our regression methods on a more realistic data set. At around $\sigma_\epsilon = 0.2$, the model starts to overfit, as indicated by a higher bias and lower R2test score. For $\sigma_\epsilon = 1$ and above, the fit barely resembles the target function. The negative R2test-score indicates that even just using the mean of the noisy data is a better model. We therefore settle on $\sigma_\epsilon = 0.2$ for our Franke function interpolation.

We start by evaluating the 95% confidence intervals for $\hat{\beta}$ using OLS for $p = 1$ up to $p = 5$. The fitted function along with its regression coefficients β_i ⁹, are shown in Fig.3 and Fig.4, respectively.

⁹ of which there are $(p+2)(p+1)/2$, as discussed in the appendix.

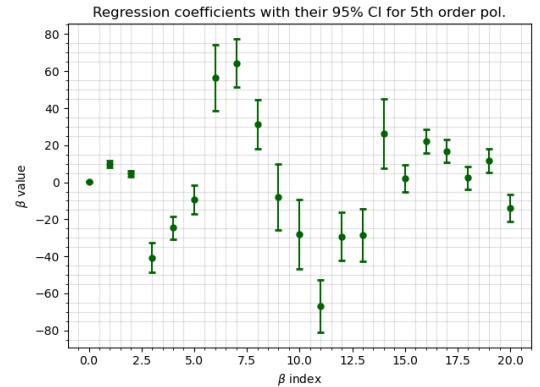
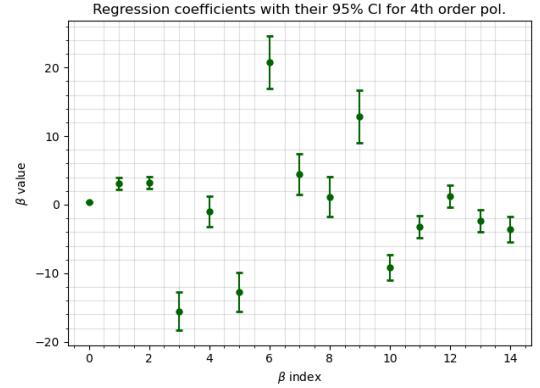
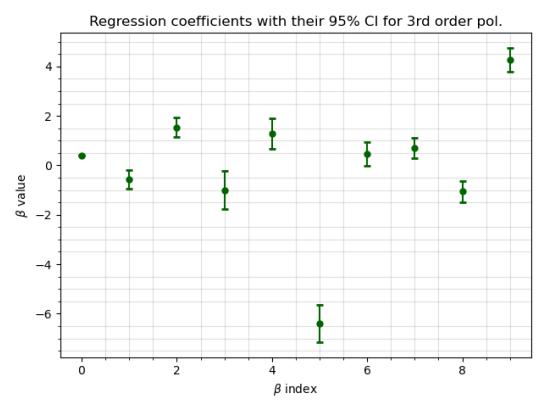
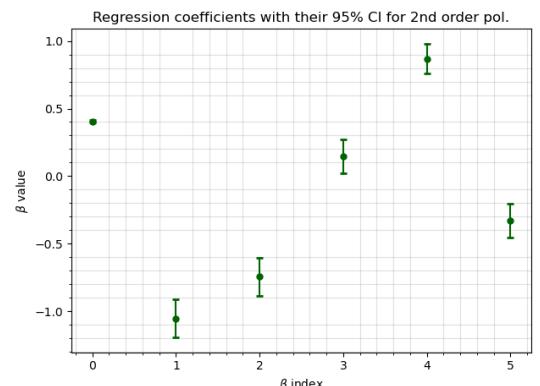
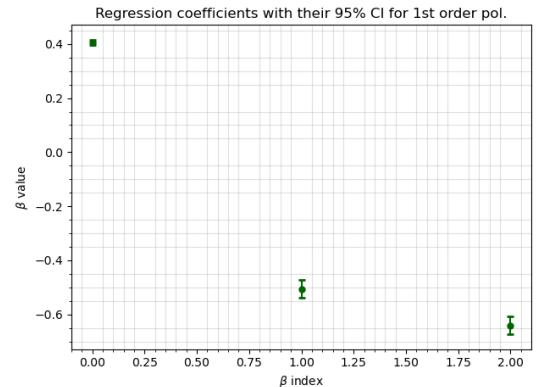
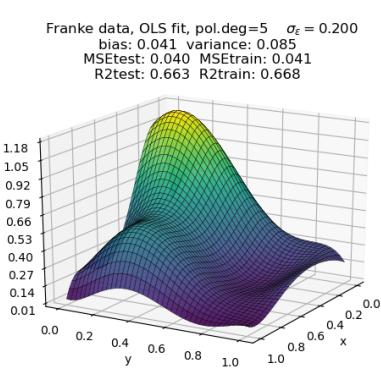
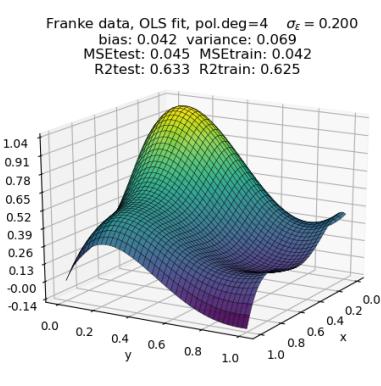
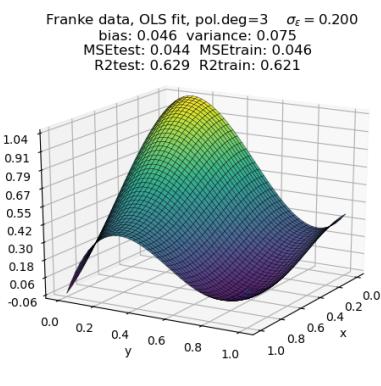
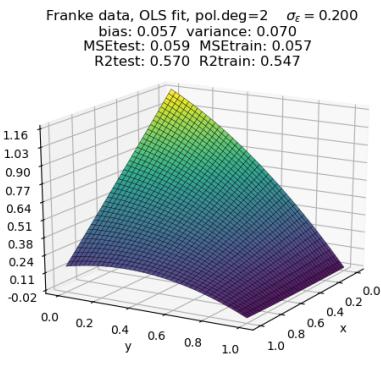
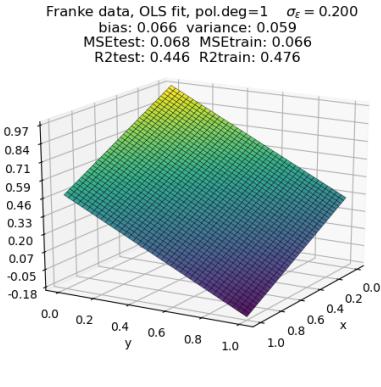


FIG. 3.

FIG. 4.

As can be seen, the first order polynomial is not complex enough to capture the characteristics of the Franke function. Our model is underfitted. As we increase the complexity, we see the bias go down and that it resembles the Franke function more and more. However, as indicated by the R2test score, $p = 5$ is still not an amazing fit. We can plot the MSEtest and MSEtrain errors against p to determine which OLS model yields the best fit. As seen in the first panel in Fig.5, for very low amplitude noise, these two values are almost the same, regardless of complexity. This is because even though the interpolating polynomial starts to fit the errors/noise, the amplitude of this noise is sufficiently small that the polynomial closely resembles the original function, and is therefore able to extrapolate well to training data. Thus, the MSPE(MSEtest) goes to zero at high complexity, while at low complexity, the model is not sufficiently complex to give us a good fit of the data, i.e. it has high bias. As the amplitude of the noise¹⁰ increases, MSEtrain keeps its shape, viz. it decreases to 0 with complexity, while MSEtest diverges from MSEtrain and starts increasing. This is because the model is sufficiently complex to start fitting the noise. Since the amplitude of this noise is not negligible, the fitted function will differ significantly from the target function. So even though it fits the training data well, it extrapolates poorly to test data. It is then in the high variance region.

Studying panel 3, which corresponds to our noise amplitude, we see that the minimum of MSEtest is given by $p = 8$. However, we will reproduce this plot with the use of resampling later in order to say with more certainty that this is indeed the optimal OLS model.

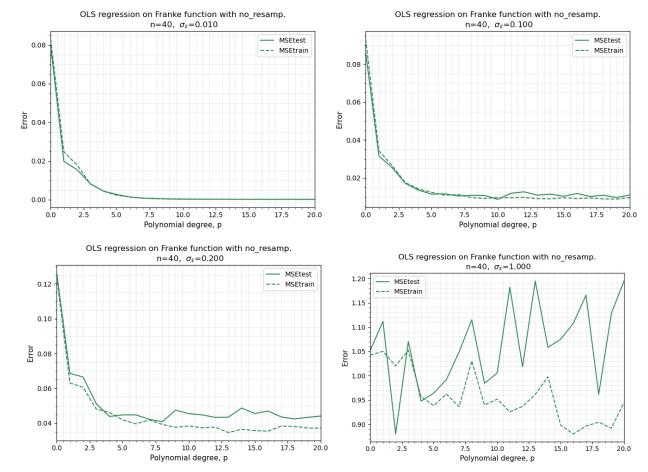


FIG. 5. Here we see the mean squared prediction error for both the training- and test data plotted against the complexity of our OLS model, i.e. the degree of the interpolating polynomial.

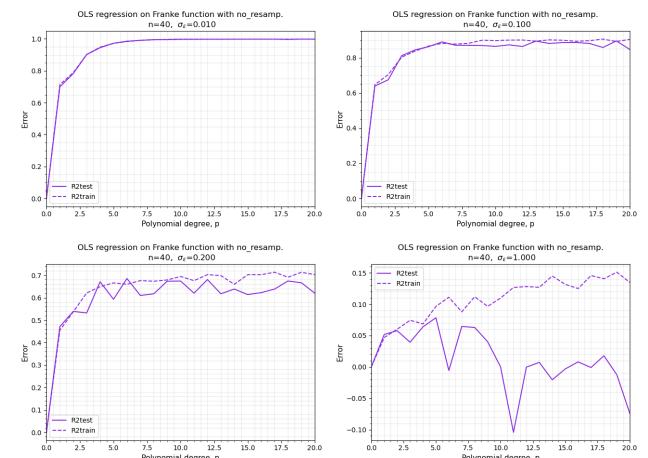


FIG. 6. R2test vs R2train scores plotted against model complexity for various noise amplitudes.

¹⁰ We use the term amplitude very loosely here. It would be more accurate to refer to the standard deviation, but more of a mouthful.

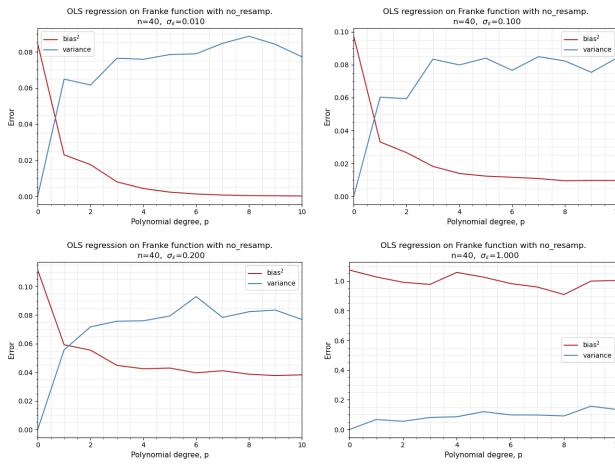


FIG. 7. Bias and variance plotted against model complexity for various noise levels. Here we only plotted up to $p = 10$ to see the intersect better.

As seen in Fig.6, for low noise levels, both R2test and R2train reach a plateau at around $p = 5$, barely inching closer to 1 as model complexity increases. In panels 2 and 3 we see that this plateau is shifted down from 1.0 to 0.9 and 0.7, respectively. This is rather peculiar since the R2 score is given by $R^2 = 1 - \frac{\sum_{i=0}^{n-1} (z_i - \bar{z}_i)^2}{\sum_{i=0}^{n-1} (z_i - \tilde{z}_i)^2}$ and as the model fits the data better and better, $z_i - \tilde{z}_i$ should go to zero for the training data, making R2train approach 1. The plateau was found to stay at 0.7 for $\sigma = 0.2$ even up to $p = 50$, so it is not simply that it is having trouble fitting the noise in the data. The reason for this behavior may be due to some mistake in our R2-score function:

```
def getR2(z,z_tilde):
    num = np.sum((z-z_tilde)**2)
    denom = np.sum((z-np.mean(z_tilde))**2)
    return 1-num/denom
#Calculate R2train for a given fit:
scoreVal = getR2(z_test, z_predict)
```

Although we have not been able to find it.

Regardless, Fig.6 shows that R2test starts to deviate downwards from the plateau as the complexity increases. This is due to overfitting, which explains why this downwards deviation is more pronounced for larger noise amplitudes, as the training data it is fitting differs more from the target function.

R2train does not exhibit this downwards deviation from the plateau, since this is the data the model was trained to fit, and does so too well in this region where R2test is decreasing.

Now looking at the bias-variance plots of Fig.7., we see in the first panel that the bias rapidly goes to zero. Since the amplitude of the noise is so small, our model shapes itself to fit the training data with minimal ease. We therefore see the bias, which measures how much the model deviates from the target function, or in our case

the training data¹¹, goes to zero quite fast. As we increase the noise, we see that this decline is not as rapid. This is because the model isn't complex enough yet to fit the noise and which leaves us with a greater bias. In addition, it reaches a plateau with a higher value. We plotted the same graph up to $p = 70$ and it never seems to each zero.

The variance also seems to reach a plateau at around $p = 3$. Plotting for larger ranges of p reveals that they are not just simply approaching the same values at slower rates.

This behavior seems peculiar to us. It makes sense for the MSPE to approach a plateau of higher value at higher noise amplitudes, since it cannot go lower than the irreducible error, σ_e^2 . Why this also seems to be the case for the bias, we are not sure.

B. Bias-variance tradeoff and Bootstrap

Fig.8 shows the bias squared¹², variance and MSPE(MSEtest) again plotted vs. complexity, but this time we have used bootstrap resampling with 30 samples to improve the precision of our measurements. From Eq.(32) one would expect that neither the bias nor the variance should ever exceed the MSPE since all terms on the LHS are positive. One might think that this is because we decided on using a different expression for the bias in Fig. 8;¹³

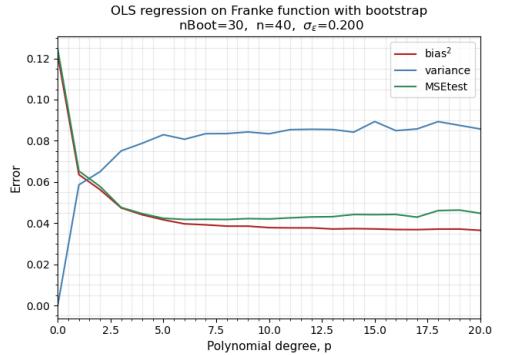


FIG. 8.

But using the original expression for bias leaves us no happier. Now the bias never approaches zero, and lies above the variance even for higher complexities where we would expect it to approach zero, as seen in Fig.9. We

¹¹ Since we usually don't have the target function, the training data is used instead.

¹² We may have used bias and bias squared somewhat interchangeably, but the legends are always correctly labeled.

¹³ We now realized that the expression we used for bias in the above figures is the MSEtrain, but it's too late to change.

also see that both the variance and bias now have a higher value than the MSPE, which doesn't make sense. We have either made some error in our code¹⁴ or some/all of our assumptions in the derivation of the expression for MSPE were wrong.

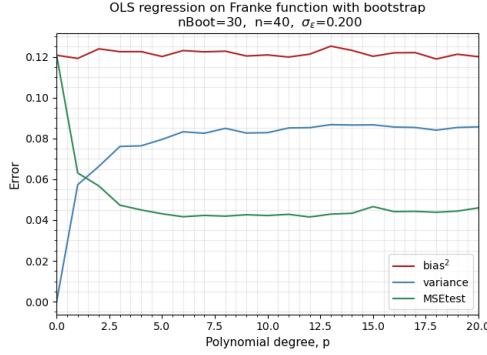


FIG. 9.

What we should expect is to see a decreasing bias and an increasing variance with complexity. The MSPE should initially be high at low complexity, since the bias is high, and reach a minimum around the point where the bias (squared) is equal to the variance. It is at this minimum that we find the optimal model complexity. Our suspicion is that the expressions we use for bias and variance are wrong, since the MSE appears to behave as expected, as can be seen in panel 3 in Fig.5, which corresponds to our chosen noise level.

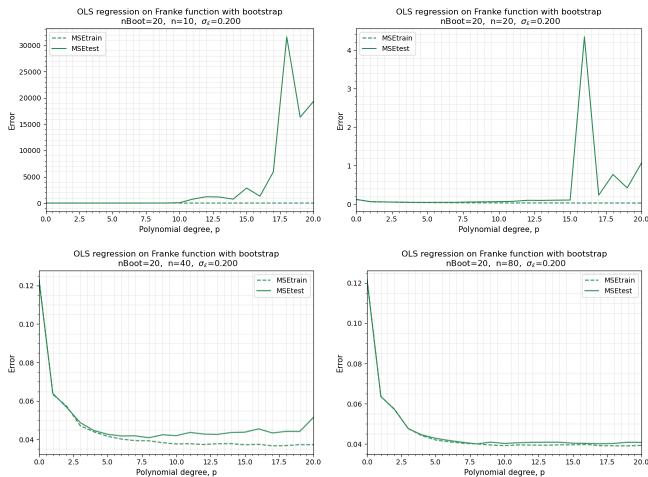


FIG. 10.

Each panel here plots MSEtest and MSEtrain against

model complexity for a specific number of gridpoints, given by n^2 where n is specified in the title.

In contrast to panel 3 in Fig.5., this is done using bootstrap resampling with $B = 30$ bootstrap samples. As can be seen, the result is a much smoother curve. This is because the variance of our estimators MSEtrain and MSEtest goes as $\sim \frac{1}{B} \sigma_{\text{MSE}}^2$ as discussed in the Theory section. To determine the optimal polynomial degree for OLS we therefore reproduce the MSE and R2 score plots, but now with a higher number of bootstrap samples. The results of this are shown in Fig.11., where we used 500 bootstrap samples. We see that the optimal polynomial degree is given at $p = 8$ in both cases, with MSEtest = 0.04179 and R2=0.6591 at these extrema. These results will come in handy when comparing with Ridge and Lasso regression.

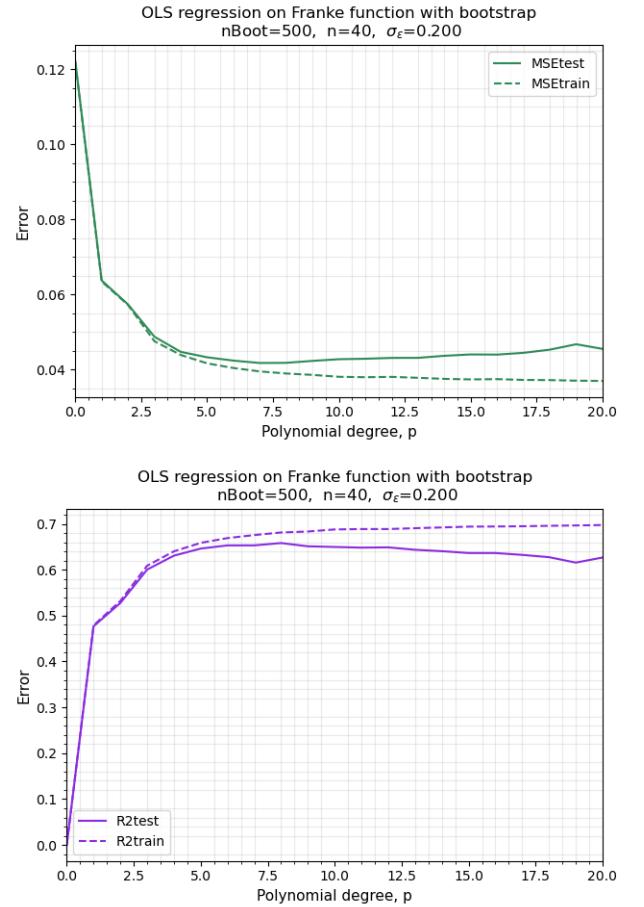


FIG. 11. MSEtest minimum: 0.04179, at $p=8$. R2test maximum: 0.6591, at $p=8$. run times: 400.01 and 432.64sec, respectively.

¹⁴ All the scores are calculated using the same six expressions given in "getTheScore()" found in functions.py, which makes it easy to

try new expressions.

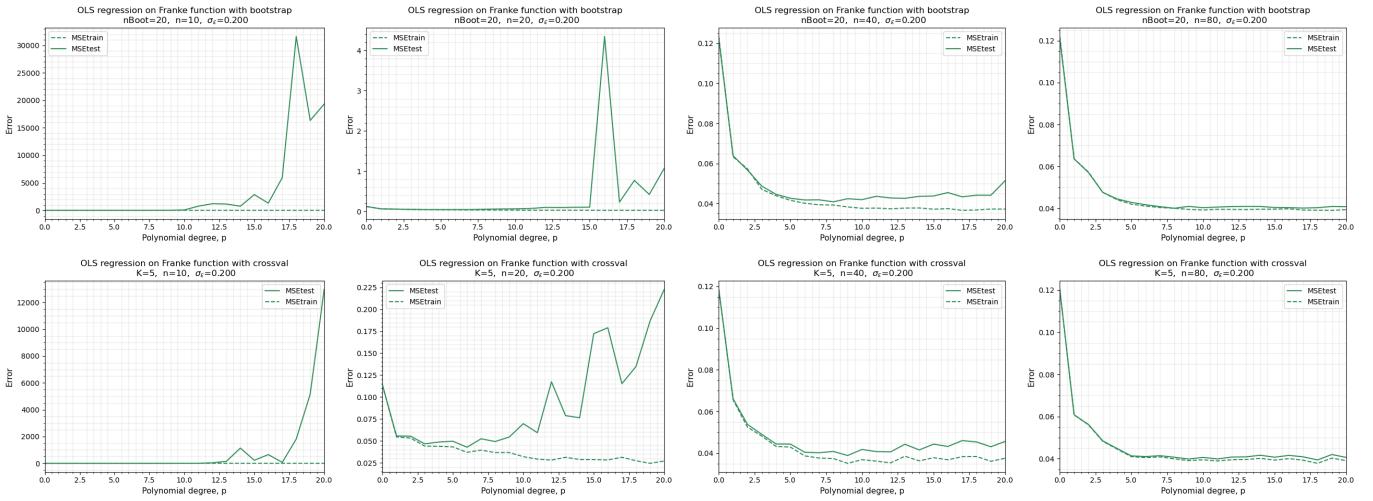


FIG. 12. In the top row, we see MSEtest vs train results using bootstrap for various number of gridpoints, while in the bottom row we used cross validation. For $n=10$ and $n=20$, we see that cross validation gives much a much lower error for high complexity. This is to be expected, as discussed in the theory section. The reason we see such high errors at low number of gridpoints is because it is harder for the model to pick up the underlying shape of the data. To see this, consider for example $n = 1$. Here it is impossible for the model to deduce the target value, as there is no way to differentiate the noise and the true value. However, as the number gridpoints increases the errors will fluctuate around zero and not add any specific shape to the "surface". This reveals the surface characteristics and allows the model to deduce the shape of the true function and fit to it, rather than the noise. This is demonstrated in Fig.13. And so for $n=40$ and $n=80$, corresponding to 1600 and 6400 gridpoints, we can evaluate the shape of the curves since it is not drowned out by overfitting as is the case for $n=10$ and $n=20$. We see that even at 20 cycles, the bootstrap method outperforms cross validation in that it produces a smoother curve and thus less variance in the measurements. We therefore opted to use bootstrap as our resampling method when determining the optimal parameter values.

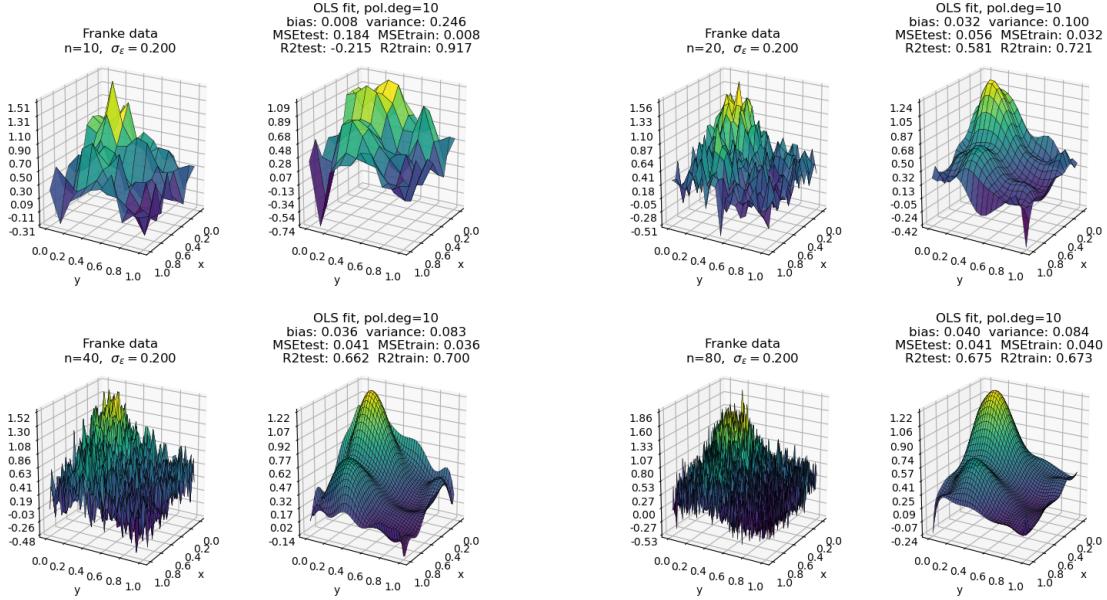


FIG. 13. To reproduce these plots, one can run "run_exercises.py" in the same directory as our other scripts and choose exercise1,task1.

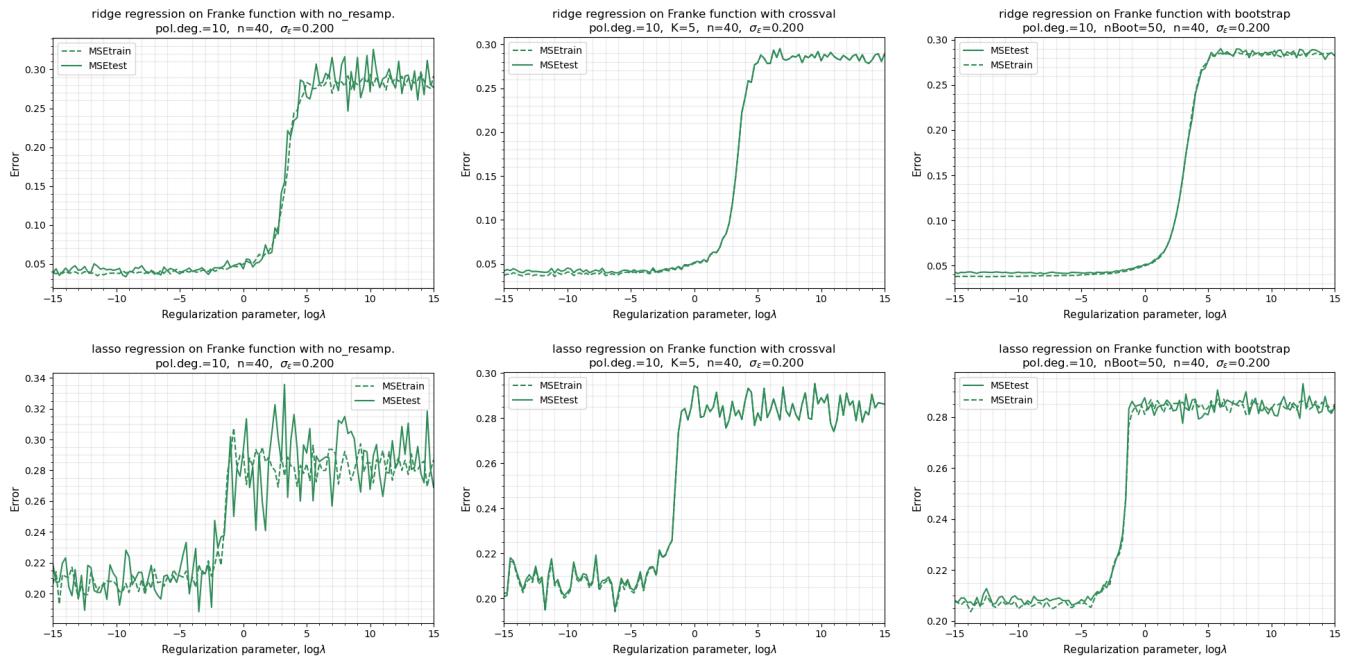


FIG. 14.

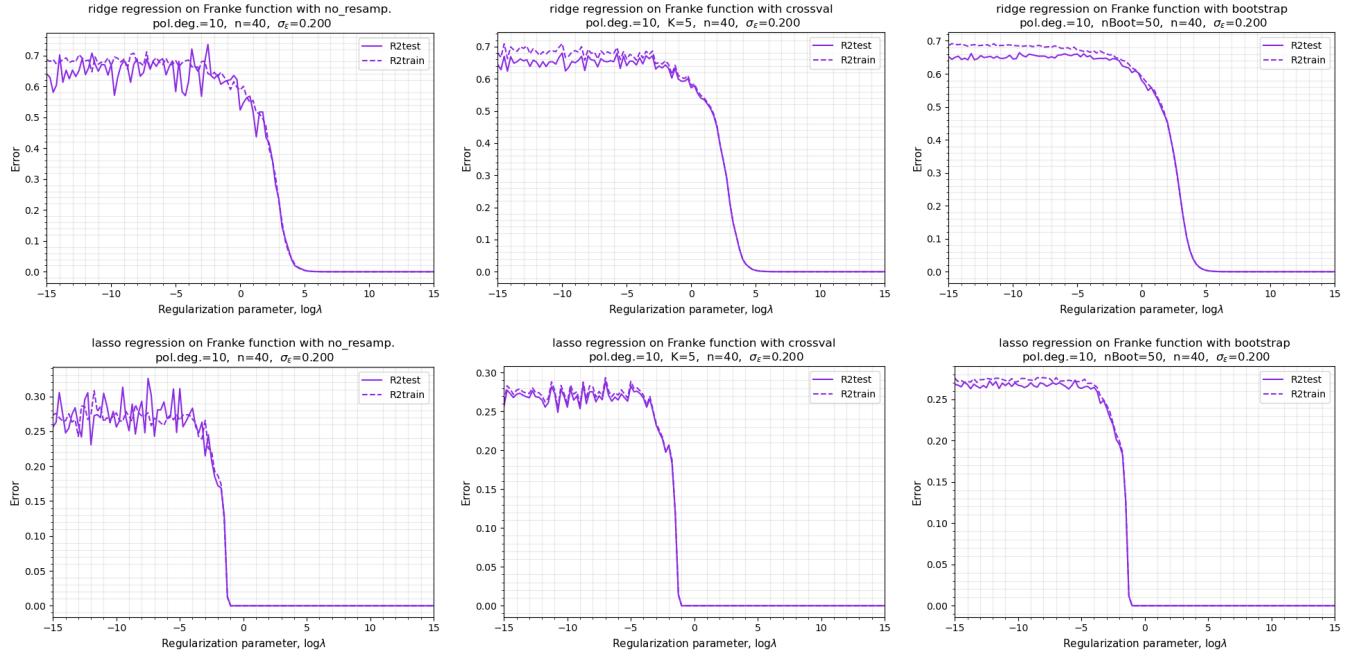


FIG. 15.

Figures 14 and 15 show the MSPE and the R²-score, respectively. In both figures we compare the results from Ridge(1st row) and Lasso(2nd row) regression with no resampling, cross validation($k=5$) and lastly bootstrap($B=50$). We see that bootstrap produces the smoothest curve in all cases, meaning it gives the lowest variance in our MSPE and R²-score estimates. We also see that Lasso regression has a higher variance for the MSPE than Ridge, while this difference doesn't seem as great for the R²-score.

Furthermore, we see both models performing well up to some threshold value of λ , which seems to be reached

at around $\log_{10}(\lambda) = 4$ for Ridge and $\log_{10}(\lambda) = -1$ for Lasso. Beyond this point, both methods have a quickly decreasing performance with λ before reaching a plateau, although this happens faster for Lasso. As the plateau is reached, the shrinkage parameter λ has probably shrunk the regression coefficients close to zero. The fact that Lasso regression can make them go to zero instead of approaching it might explain why we observe a steeper curve towards the plateau.

From these figures we can also compare their performance in their "functional" region of λ . We in Fig.14 see that the MSPE of Lasso regression is about four times that of Ridge in this region. In Fig.15 we see that the R2-score for Lasso settles at a value of somewhere between 0.25 and 0.3 while Ridge does so at 0.6-0.7 as we decrease λ .

For this polynomial degree, viz. $p = 10$, Ridge certainly outperforms Lasso on our Franke data with $\sigma_\varepsilon = 0.2$. In order to evaluate their performance as a function of both λ and polynomial degree, we produce heat maps in their functional range, as seen in Fig.16 and Fig.17.

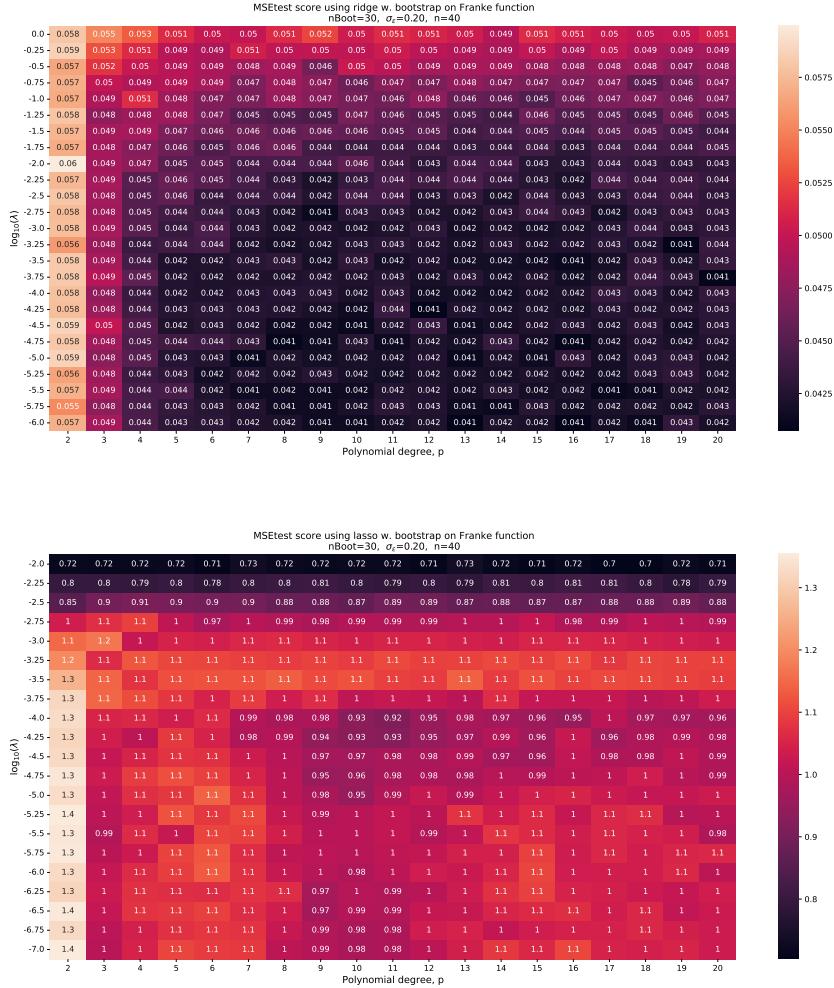


FIG. 16.

As suspected, Lasso performs worse for a wide range of values of λ and p . We also realize that our Ridge MSEtest measurement lacks the number of decimals required to determine whether its minimum is greater or smaller than that of OLS. However, looking at the table below, we see that Ridge gave a higher R2-score. That being said, OLS have unbiased estimators for the regression coefficients, while this is not the case for Ridge. It is therefore probably preferable to simply use OLS with a polynomial of degree 8.

Further improvements would include determining the minimum of MSEtest for Ridge more precisely.

TABLE I. MSEtest minima and R2test maxima of the three different regression methods, obtained using bootstrap with B=500 for OLS and B=30 for Ridge and Lasso due to time considerations.

Reg. Meth	min(MSE _{test})	max(R _{test} ²)
OLS	0.04179	0.6591
Ridge	0.041	0.67
Lasso	0.7	0.079

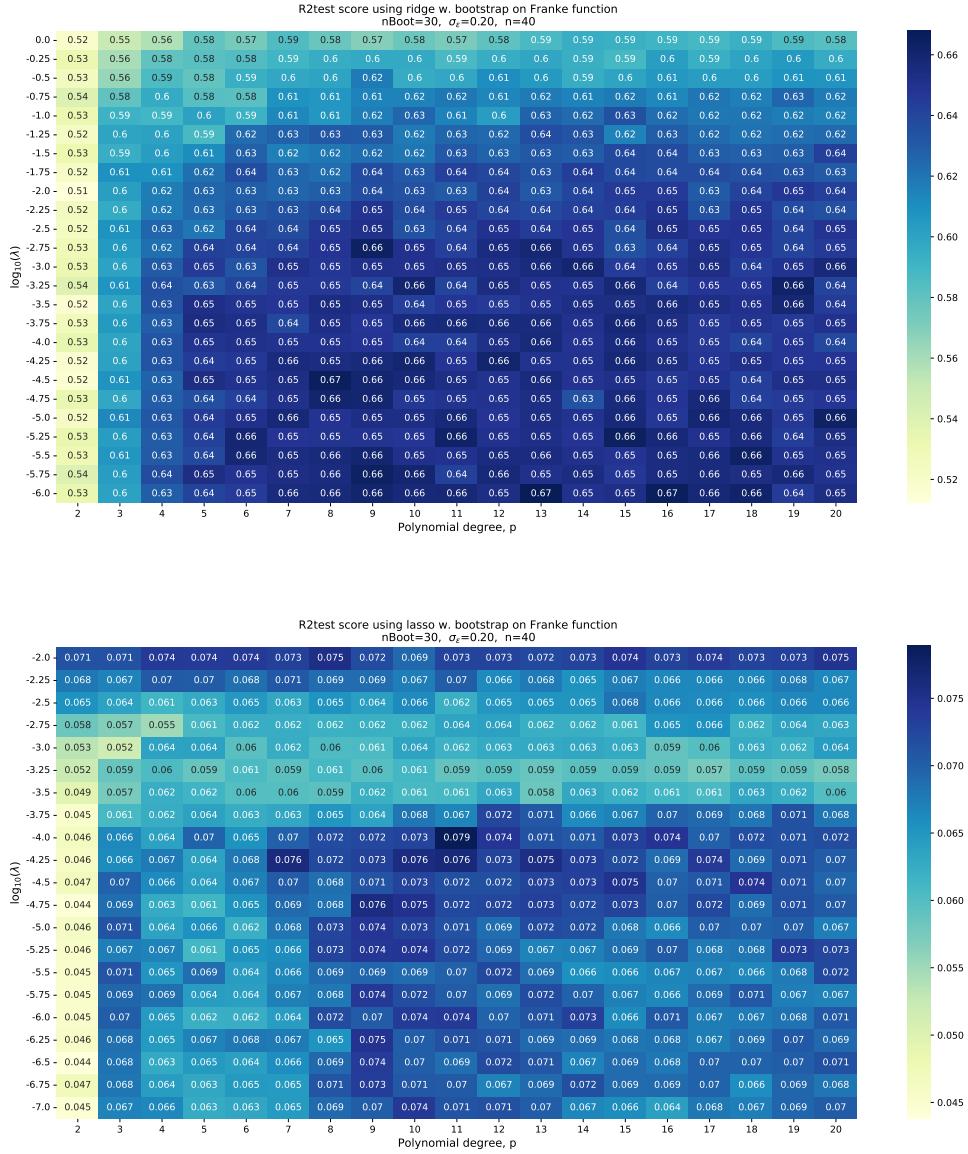


FIG. 17. Here we see the R2test scores for Ridge and Lasso regression in their "functional range", using 30 Bootstrap samples.

C. Terrain data

In this section, we assess our models on the terrain 1 data provided in [3]. While performing model assessment, we limit ourselves to using a data sample of 2500 points, corresponding to $n = 50$, since the heat maps take a while to produce.

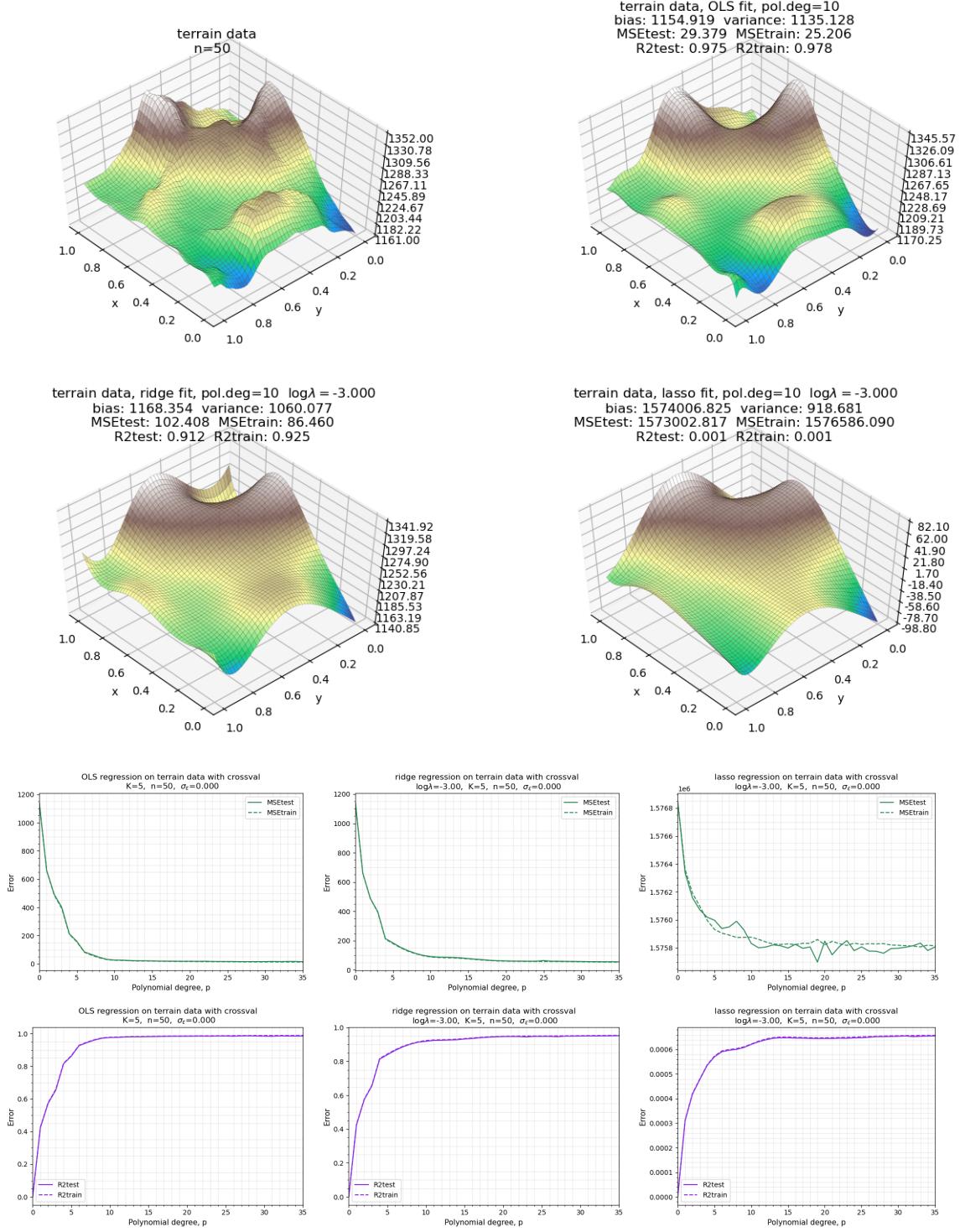


FIG. 18.

In Fig.18, we see the original terrain data along with an OLS, Ridge and Lasso fit on the same data. For this single sample, we see that OLS obtained the best fit, followed by Ridge and Lasso. We have also plotted MSEtest and R2test for these three regression methods as a function of model complexity, which shows Lasso to perform the worst, while Ridge comes in at a close second to OLS, which gives a lower MSPE at low complexities than Ridge does. For higher complexity it is not so obvious which of the two performs best. We therefore investigate this further by use of heat maps.

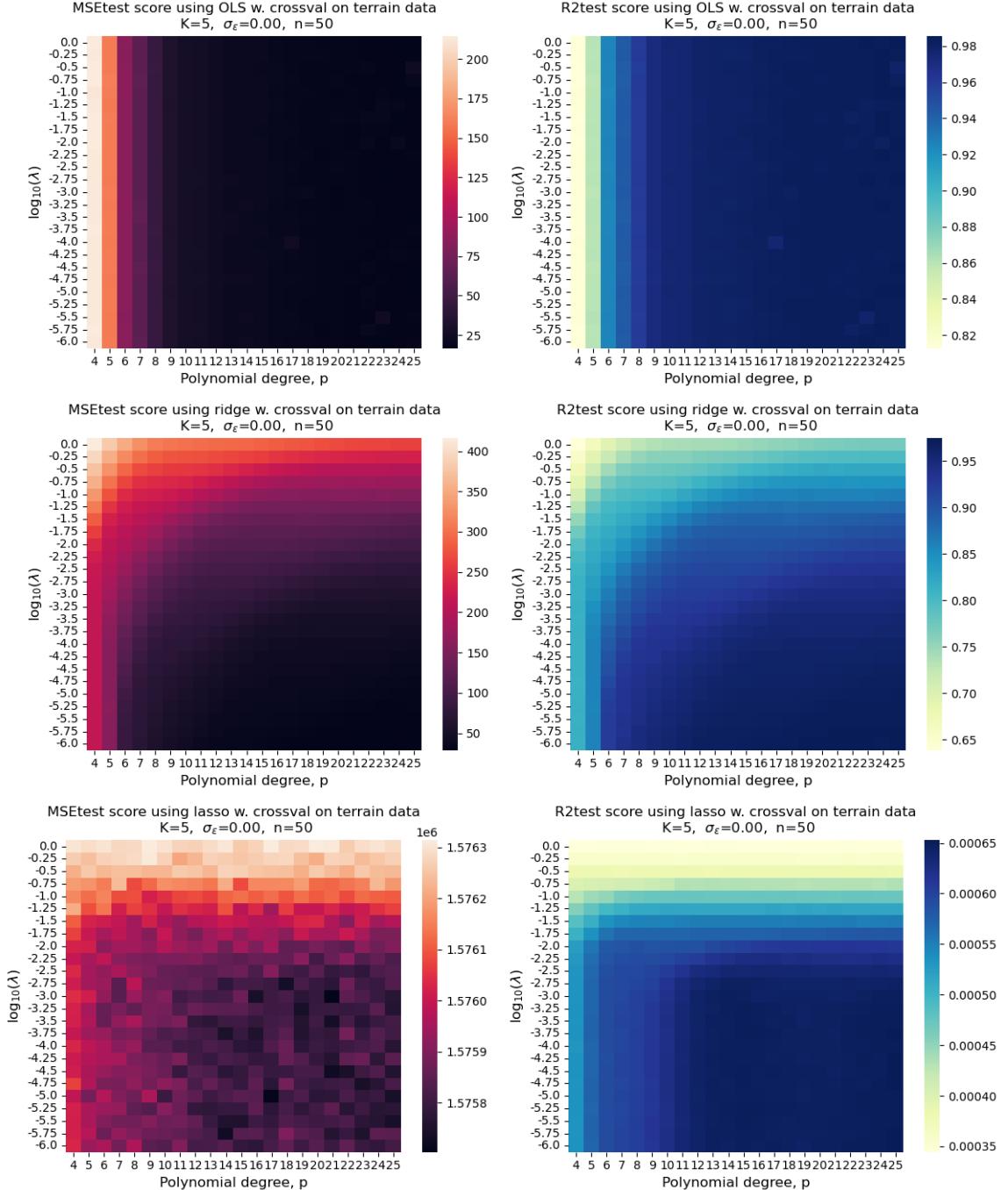


FIG. 19. From these heatmaps we see that that both Ridge and OLS outperforms Lasso regression in the selected parameter space. OLS also seems to outperform Ridge regression, although it is somewhat hard to tell for values of $\lambda < -2.5$ and $p > 10$. We therefore produce new heat maps with those values as cut-offs

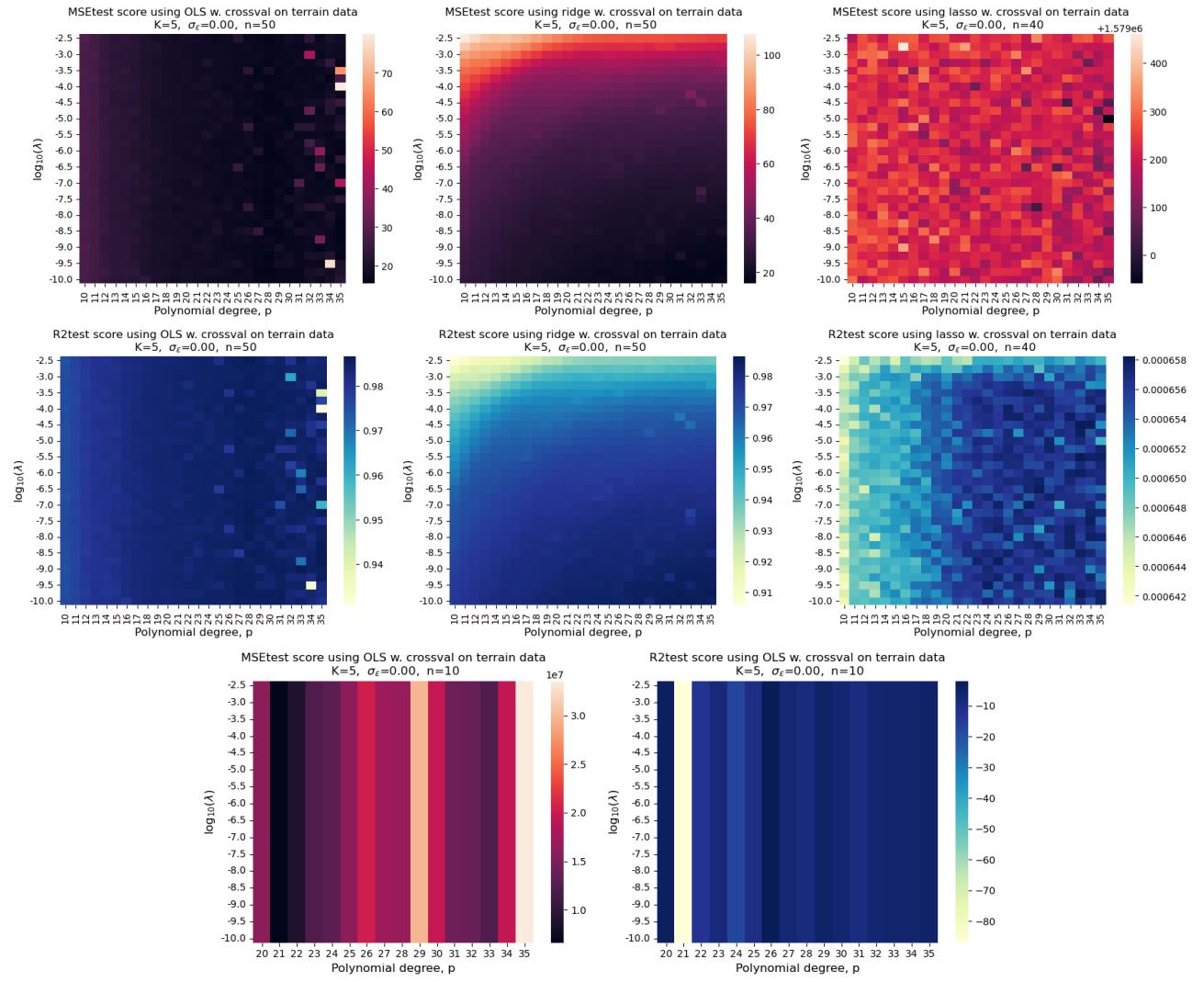


FIG. 20. We see here that for very high degree polynomials, we see start seeing variations in the MSPE and R2-score along the λ -axis of the OLS heat map. The reason for this behavior is overfitting. As our model becomes increasingly complex, it starts also fitting the noise in the data, inheriting characteristics that aren't present in the target function. If the OLS regression algorithm saw the same data set each time, the variations with λ wouldn't make sense, since it would be fitting the same noise each time, and although the fit might be terribly off, it should yield the same errors each time for a given polynomial degree. However, since the cross validation algorithm randomly shuffles our data points for each iteration of p , the model fits a new data set for each pixel in the heat map. Turning shuffling off, should therefore result in our errors being invariant along the λ axis, as is demonstrated by the two bottom heat maps. Furthermore see that, in the range $25 < p < 35$, Ridge is much more resistant to overfitting than OLS.

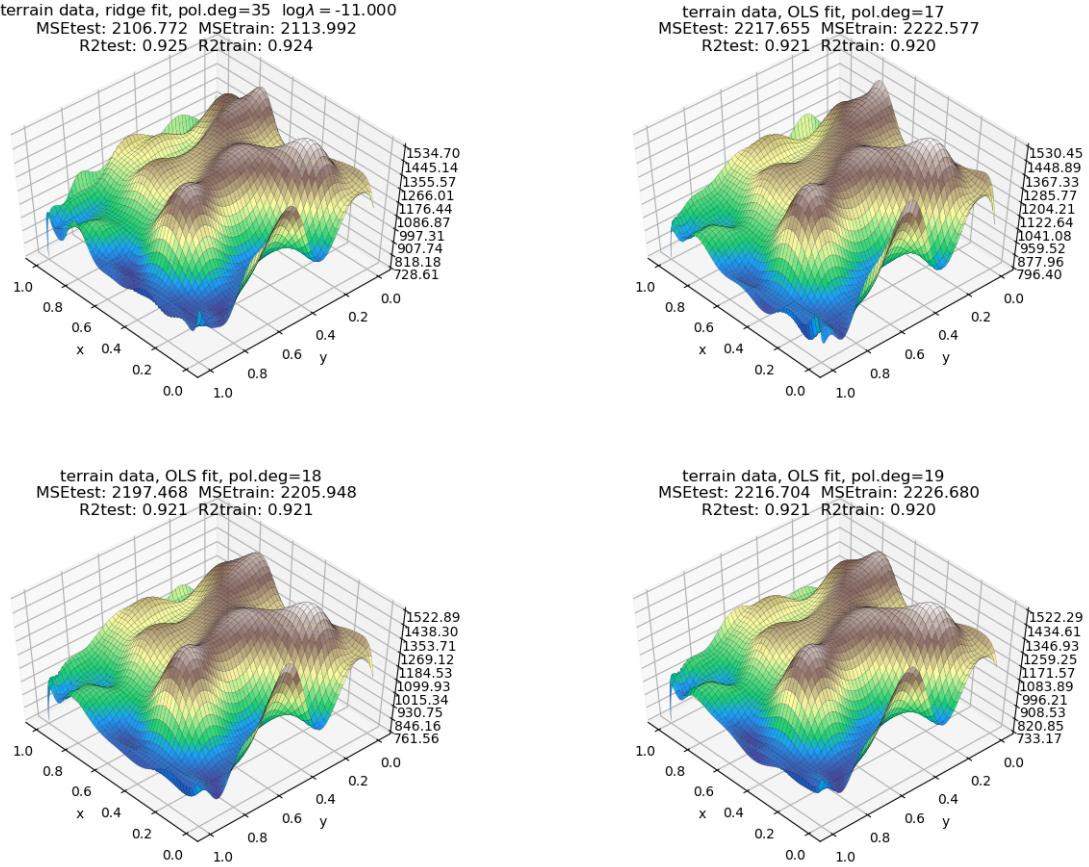


FIG. 21. The fits seen here were done using $n=500$. We see a Ridge fit of order 35 as well as OLS fits of orders 17, 18 and 19. Their run times were 16.74, 3.75, 3.94 and 4.44 seconds respectively. So even though it seems we can achieve a slightly better fit with Ridge regression, it comes at about 4 times the computational cost. Furthermore, we see that among the OLS models, the 18th degree polynomial yields the lowest MSPE for this particular sample. Fig.20. shows this (with more precision) to be the case for $n = 50$ as well, so we can suspect that $p = 18$ will remain optimal for larger data sets still. We therefore choose the OLS with polynomial degree 18 as our optimal model for the terrain data.

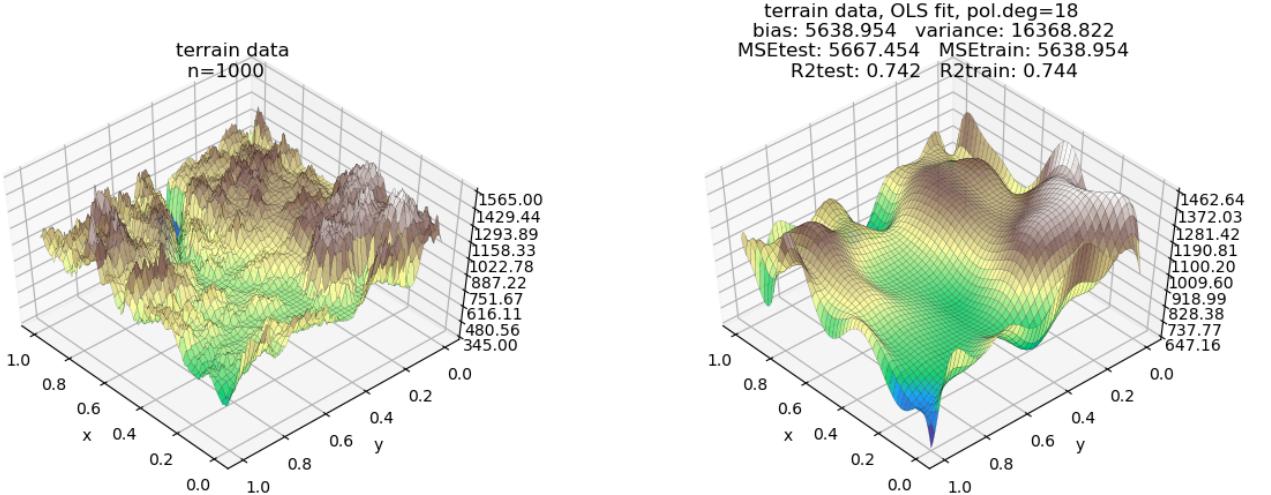


FIG. 22. Here we see our chosen model (OLS using $p = 18$) fitted on a terrain data consisting of a million datapoints. We tried using Ridge regression with the optimal parameters for comparison, but our "program" gets killed.

V. Conclusion

In this project, we assessed the performance of OLS, Ridge and Lasso regression in fitting a two dimensional polynomial to self-generated data in the form of a noisy Franke function as well as real terrain data. This was done by calculating and comparing their MSPE and R2-scores. To reduce the uncertainty in these measurements, we performed resampling. The Bootstrap and Cross Validation resampling techniques were compared, and Bootstrap was found to yield better results for larger data sets, including ours. We therefore used this as our resampling technique in performing model selection for fitting the Franke function. Cross validation was used for the terrain data due to time considerations.

Of the regression models, Lasso performed the worst by a large margin. At greater levels of noise, this might not have been the case, as it is the most resistant to overfitting.

For the Franke data, the OLS and Ridge were found to have very similar performance. Their MSPE minimas in the parameter space of our search were found to be 0.0418 and 0.041, respectively, while their R2-score maxima were 0.659 and 0.67, respectively.

Determining the MSPE measurement of Ridge to a higher decimal place and with more bootstrap samples, in addition to comparing its run time with that of OLS is required to say with more certainty which of these models performs best on our Franke data.

For the Terrain data, we found Ridge regression to yield a better fit at around $\lambda = -11$ and $25 < p < 35$, although this would likely be the case for higher polynomial degrees as well. Considering that the OLS using $p = 18$ had very similar performance at only one fourth the computational cost, this was chosen as our optimal model for the terrain data.

References

- [1] M. Hjorth-Jensen. *Applied Data Analysis and Machine Learning, FYS-STK3155/4155 at the University of Oslo, Norway*. URL: https://compphysics.github.io/MachineLearning/doc/LectureNotes/_build/html/intro.html.
- [2] M. Hjorth-Jensen. *Overview of course material: Data Analysis and Machine Learning*. URL: <https://compphysics.github.io/MachineLearning/doc/pub/week34/html/week34.html>.
- [3] *SRTM_data_Norway_1.tif*. Sept. 2021. URL: <https://github.com/CompPhysics/MachineLearning/tree/master/doc/Projects/2021/Project1/DataFiles>.

Appendices

A. The Franke function

The Franke function is a bivariate function, and so outputs a matrix $\mathbf{Z} = \mathcal{F}(xx, yy)$ given by

$$\mathbf{Z} = \begin{bmatrix} \mathcal{F}(x_0, y_0) & \mathcal{F}(x_1, y_0) & \dots & \mathcal{F}(x_{n-1}, y_0) \\ \mathcal{F}(x_0, y_1) & \mathcal{F}(x_1, y_1) & \dots & \mathcal{F}(x_{n-1}, y_1) \\ \vdots & \vdots & \vdots & \vdots \\ \mathcal{F}(x_0, y_{n-1}) & \mathcal{F}(x_1, y_{n-1}) & \dots & \mathcal{F}(x_{n-1}, y_{n-1}) \end{bmatrix}$$

We would like to use the OLS algorithm from II to produce a model of the form

$$\tilde{\mathbf{Z}} = \begin{bmatrix} \tilde{\mathcal{Z}}(x_0, y_0) & \tilde{\mathcal{Z}}(x_1, y_0) & \dots & \tilde{\mathcal{Z}}(x_{n-1}, y_0) \\ \tilde{\mathcal{Z}}(x_0, y_1) & \tilde{\mathcal{Z}}(x_1, y_1) & \dots & \tilde{\mathcal{Z}}(x_{n-1}, y_1) \\ \vdots & \vdots & \vdots & \vdots \\ \tilde{\mathcal{Z}}(x_0, y_{n-1}) & \tilde{\mathcal{Z}}(x_1, y_{n-1}) & \dots & \tilde{\mathcal{Z}}(x_{n-1}, y_{n-1}) \end{bmatrix}$$

Where we denote $\tilde{\mathbf{Z}}_{i,j}$ by $\tilde{\mathcal{Z}}(x_i, y_j)$. However, \mathbf{Z} and $\tilde{\mathbf{Z}}$ are matrices and the OLS algorithm assumes \mathbf{y} and $\tilde{\mathbf{y}}$ to

be column vectors. We can solve this problem by using the NumPy function `ravel` on \mathbf{Z} and $\tilde{\mathbf{Z}}$ to produce \mathbf{z} and $\tilde{\mathbf{z}}$, respectively:

$$\mathbf{z} = \begin{bmatrix} \mathcal{F}(x_0, y_0) \\ \mathcal{F}(x_1, y_0) \\ \vdots \\ \mathcal{F}(x_{n-1}, y_0) \\ \mathcal{F}(x_0, y_1) \\ \mathcal{F}(x_1, y_1) \\ \vdots \\ \mathcal{F}(x_{n-1}, y_1) \\ \mathcal{F}(x_0, y_{n-1}) \\ \mathcal{F}(x_1, y_{n-1}) \\ \vdots \\ \mathcal{F}(x_{n-1}, y_{n-1}) \end{bmatrix}, \quad \tilde{\mathbf{z}} = \begin{bmatrix} \tilde{\mathcal{Z}}(x_0, y_0) \\ \tilde{\mathcal{Z}}(x_1, y_0) \\ \vdots \\ \tilde{\mathcal{Z}}(x_{n-1}, y_0) \\ \tilde{\mathcal{Z}}(x_0, y_1) \\ \tilde{\mathcal{Z}}(x_1, y_1) \\ \vdots \\ \tilde{\mathcal{Z}}(x_{n-1}, y_1) \\ \tilde{\mathcal{Z}}(x_0, y_{n-1}) \\ \tilde{\mathcal{Z}}(x_1, y_{n-1}) \\ \vdots \\ \tilde{\mathcal{Z}}(x_{n-1}, y_{n-1}) \end{bmatrix} \quad (\text{A.1})$$

Since we are using polynomials of \mathbf{x} and \mathbf{y} to interpolate \mathbf{z} , we have

$$\begin{aligned} \tilde{\mathcal{Z}}(x_i, y_j) &= \beta_0 + \beta_1 x_i + \beta_2 y_i + \beta_3 x_i^2 + \beta_4 x_i y_j + \beta_5 y_j^2 \\ &\quad + \dots + \beta_{\frac{(p+2)(p+1)}{2}} y_j^p \end{aligned} \quad (\text{A.2})$$

and so we get:

$$\begin{bmatrix} \tilde{\mathcal{Z}}(x_0, y_0) \\ \tilde{\mathcal{Z}}(x_1, y_0) \\ \vdots \\ \tilde{\mathcal{Z}}(x_{n-1}, y_0) \\ \tilde{\mathcal{Z}}(x_0, y_1) \\ \tilde{\mathcal{Z}}(x_1, y_1) \\ \vdots \\ \tilde{\mathcal{Z}}(x_{n-1}, y_1) \\ \vdots \\ \tilde{\mathcal{Z}}(x_0, y_{n-1}) \\ \tilde{\mathcal{Z}}(x_1, y_{n-1}) \\ \vdots \\ \tilde{\mathcal{Z}}(x_{n-1}, y_{n-1}) \end{bmatrix} = \begin{bmatrix} 1 & x_0 & y_0 & x_0^2 & x_0 y_0 & y_0^2 & \dots & y_0^p \\ 1 & x_1 & y_0 & x_1^2 & x_1 y_0 & y_0^2 & \dots & y_0^p \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{n-1} & y_0 & x_{n-1}^2 & x_{n-1} y_0 & y_0^2 & \dots & y_0^p \\ 1 & x_0 & y_1 & x_0^2 & x_0 y_1 & y_1^2 & \dots & y_1^p \\ 1 & x_1 & y_1 & x_1^2 & x_1 y_1 & y_1^2 & \dots & y_1^p \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{n-1} & y_1 & x_{n-1}^2 & x_{n-1} y_1 & y_1^2 & \dots & y_1^p \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_0 & y_{n-1} & x_0^2 & x_0 y_{n-1} & y_{n-1}^2 & \dots & y_{n-1}^p \\ 1 & x_1 & y_{n-1} & x_1^2 & x_1 y_{n-1} & y_1^2 & \dots & y_{n-1}^p \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{n-1} & y_{n-1} & x_{n-1}^2 & x_{n-1} y_{n-1} & y_{n-1}^2 & \dots & y_{n-1}^p \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \\ \vdots \\ \beta_{\frac{(p+2)(p+1)}{2}} \end{bmatrix} \quad (\text{A.3})$$

Which we write as $\tilde{\mathbf{z}} = \mathbf{X}\beta$. We now see that the columns of the design matrix can be written as various powers of \mathbf{x}_r and \mathbf{y}_r :

$$\mathbf{X} = [\mathbf{1}, \mathbf{x}_r, \mathbf{y}_r, \mathbf{x}_r^2, \mathbf{x}_r \mathbf{y}_r, \mathbf{y}_r^2, \dots, \mathbf{y}_r^p] \quad (\text{A.4})$$

Here \mathbf{x}_r and \mathbf{y}_r are given by

$$\mathbf{x}_r, \mathbf{y}_r = \mathbf{xx}.ravel(), \mathbf{yy}.ravel()$$

Where xx and yy are meshgrids, in turn produced from \mathbf{x} and \mathbf{y} by use of numpy's meshgrid function:

$$\mathbf{xx}, \mathbf{yy} = \text{np.meshgrid}(\mathbf{x}, \mathbf{y})$$

$$\mathbf{x}_r = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_{n-1} \end{bmatrix}, \quad \mathbf{y}_r = \begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_{n-1} \end{bmatrix} \quad (\text{A.5})$$

We now see that the k -th column of \mathbf{X} corresponds to the same combination of \mathbf{x}_r and \mathbf{y}_r as that of x and y in the k -th term of the polynomial given in (A.2). Ignoring the regression coefficients, and dropping the indices, this polynomial is given by

$$P(x, y) = \sum_{l=0}^p \sum_{k=0}^l x^{l-k} y^k \quad (\text{A.6})$$

Furthermore, the total number of terms in the polynomial is given by $N = \frac{(p+2)(p+1)}{2}$, which can be found by constructing a Pascal's triangle with each row containing only the terms of a given power, and realizing that the number of terms in the rows form a geometric series. Knowing this, we can now construct our design matrix

```
#=====
#===== CONTROL PANEL =====#
ONbutton = True #Set to False to run main.py without performing regression with the parameters specified below
tinkerBool = False #If True, a folder named 'tinkerings' will be created if savePlot is True, in which the plots will be saved b
exercise = 4 #If tinkerBool=False, a folder named 'exercise_results' will be created if savePlot is True. This int determini

terrainBool = False #Use terrain data or Franke function?
n_t = 100 #How many points on the x and y axes if using terrain data
n_f = 20 #How many points on the x and y axes if using FrankeFunction
n = pick_n(n_f,n_t,terrainBool)
scaling = False #Scale the design matrix, X?
skOLS = False #Use sklearn in OLS (rather than pseudoinverse)?
skCV = False #Use sklearn's cross validation contra own code?

nBoot = 100 #Number of bootstrap samples
K = 5 #Number of folds in cross validation alg.
shuffle = True #Shuffle the data before performing crossval folds?

#-----
#S=0.1 #Use this for fixed standard deviation as a function of n
S = 40/(n**2) #Is equal to 0.1 when n=20. Ensures sigma scales as 1/n^2. Was not used in the report.
sigma_v = [0.5*S,S,2*S,5*S] #Make a separate plot for each of these sigmas
sigma_s = [10*S] #Default standard deviation of epsilon
#-----

minOrder,maxOrder = 1, 20 #Will make order vector from minOrder to maxOrder
order_s = [10] #Default pol.degree if we don't plot vs. degrees
#-----

minLoglmd, maxLoglmd = -15, 15 #Will make log(lambda) vector from minLoglmd to maxLoglmd
lambda_s = [1e-4] #Default log(lambda) value. Must be set
#-----

sigmasBool = False #If true, will produce a plot for each std. in sigma_v

# [ordersBool lamdasBool] = plotTypeInt
plotTypeInt = 0
# [ False 0 False ] Generate surf plots and print results
# [ True 1 False ] Plots error vs. pol.deg.
# [ False 2 True ] Plots error vs. lambda
# [ True 3 True ] Produces heatmap
savePlot = True #Save plots?
plotBool = True #Make error/score plots?

resampInt = 0 #0=no_resamp., 1=bootstrap, 2=crossval
regInt = 1 #0=OLS, 1=ridge, 2=lasso

dummyList = [['bias'], ['variance'], ['MSEtest'], ['MSEtrain'], ['R2test'], ['R2train']] #Which regression scores one can plot
# 0 1 2 3 4 5
scoresNames = getScores(0,1)
#Sets which scores to calculate by passing in the corresponding index from allScores.
#Function def before main().
#=====
```

from \mathbf{x}_r and \mathbf{y}_r for a given polynomial degree, p :

```
#Create the design matrix
def desMat(xr,yr,p):
    N = len(xr)
    #Number of elements in beta:
    numEl = int((p+2)*(p+1)/2)
    X = np.ones((N,numEl))
    colInd=0#Column index
    for l in range(1,p+1):
        for k in range(0,l+1):
            X[:,colInd+1] = (xr**(l-k))*yr**k
            colInd = colInd+1
    return X
```