

# FYS-STK 3155 project 1

Theodor Midtbø Alstad

10th October 2020

## Abstract

The project detailed in this report is unfinished due to time constraints, and has many ways it could be improved.

This project aims to explore linear regression, with the final goal of learning and predicting the geography of an area. The project is split into 7 sections, labelled a) through g), culminating in application to geographical data in section g), and will be referred to as such in the report.

The final goal of fitting models to geographical data is not met, and the best fit gotten is a model with an  $R^2$  score of -125.

## 1 Introduction

This project implements the regression methods Ordinary Least Squares, LASSO, and Ridge, with the resampling methods bootstrap and k-fold cross-validation.

The project starts by modelling the Franke function in the for the domain  $x, y \in [0, 1]$ , before moving on to geographical data at the end.

Throughout the project, multiple bias-variance and test-train MSE analyses are done, combining the regression- and resampling methods to do so.

The goal of this project is two-fold; to model geographical data and to learn more about, and evaluate, the regression methods implemented.

### 1.1 The structure of the source code

The code is written in python3, formatted as a cascading inheriting class structure. That is, section b) inherits from the class of section a), section c) inherits from the class of section b), and so on.

All outputs are found in the outputs folder, as detailed in the .README file.

### 1.2 Franke function

The Franke function is modelled in sections a) through e), as a placeholder for geographical data in section g). It is here defined for the domain  $x, y \in [0, 1]$

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

### 1.3 Ordinary Least Squares (OLS)

The OLS method attempts to model a dataset by minimizing the square error between the dataset,  $y$ , and a model in the form of  $\hat{y} = \mathbf{X}\beta$ , where  $\mathbf{X}$  is a design matrix and  $\beta$  is a vector of parameters.

OLS successfully and analytically minimizes the cost function  $C = \mathbb{E}[(y - \hat{y})^2]$  by using parameters

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

### 1.4 Bias & Variance

Using a cost function  $C$  based on the Mean Square Error between our prediction,  $\hat{y}$ , and the recorded value,  $y = f(x) + \epsilon$ , the bias and variance can be found as

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

$$\begin{aligned} C &= \mathbb{E}[(y - \hat{y})^2] \\ &= \mathbb{E}[(y - \mathbb{E}[\hat{y}])^2 + (\mathbb{E}[\hat{y}] - \hat{y})^2 \\ &\quad + 2(y - \mathbb{E}[\hat{y}])(\mathbb{E}[\hat{y}] - \hat{y})] \\ &= \mathbb{E}[(y - \mathbb{E}[\hat{y}])^2] + \mathbb{E}[(\mathbb{E}[\hat{y}] - \hat{y})^2] \\ &\quad + \mathbb{E}[2(y - \mathbb{E}[\hat{y}])(\mathbb{E}[\hat{y}] - \hat{y})] \end{aligned}$$

$$\begin{aligned} \mathbb{E}[2(y - \mathbb{E}[\hat{y}])(\mathbb{E}[\hat{y}] - \hat{y})] &= 2(y - \mathbb{E}[\hat{y}])\mathbb{E}[(\mathbb{E}[\hat{y}] - \hat{y})] \\ &= 2(y - \mathbb{E}[\hat{y}]) \times 0 \\ &= 0 \end{aligned}$$

$$\begin{aligned} \mathbb{E}[(y - \mathbb{E}[\hat{y}])^2] &= \mathbb{E}[(f(x) + \epsilon - \mathbb{E}[\hat{y}])^2] \\ &= \mathbb{E}[(f(x) - \mathbb{E}[\hat{y}]) + \epsilon]^2 \\ &= \mathbb{E}[(f(x) - \mathbb{E}[\hat{y}])^2 + \epsilon(f(x) - \mathbb{E}[\hat{y}]) + \epsilon^2] \\ &= \mathbb{E}[(f(x) - \mathbb{E}[\hat{y}])^2 + 0 + \epsilon^2] \\ &= \mathbb{E}[(f(x) - \mathbb{E}[\hat{y}])^2] + \mathbb{E}[\epsilon^2] \\ &= \mathbb{E}[(f(x) - \mathbb{E}[\hat{y}])^2] + \sigma^2 \end{aligned}$$

$$C = \mathbb{E}[(f(x) - \mathbb{E}[\hat{y}])^2] + \mathbb{E}[(\mathbb{E}[\hat{y}] - \hat{y})^2] + \sigma^2$$

where  $\mathbb{E}[(f(x) - \mathbb{E}[\hat{y}])^2]$  is the bias,  $\mathbb{E}[(\mathbb{E}[\hat{y}] - \hat{y})^2]$  is the variance, and  $\sigma^2$  is the irreducible error.

Bias can be interpreted as the expected mean square error between a prediction and their true values, and is often used to determine how well fitted a model is.

It usually goes down with the complexity of a model, as the model is more able to bend to the training points.

Variance can be interpreted as how much a model varies over multiple iterations, and is often used to determine whether a model is overfitted to its test data. It usually goes up with the complexity of a model, as the model will twist, bend, and turn more to fit the training points.

The irreducible error is an error that cannot be overcome by improving the model, and can be attributed to things such as noise.

## 1.5 Bootstrap Resampling

The bootstrap resampling method aims to evaluate a regression method by generating multiple training sets to a single test set. It accomplishes this by initially splitting a complete data set into training and test sets, then repeatedly creating *sample* sets by picking out random points from the training set (duplicates are allowed in the sample set) until the sample set has a set size (in this report sample sets are given 80% size of the initial training set). Then the regression model is applied to the sample set, evaluated against the testing set, another sample set is drawn, and this is repeated as many times as is specified. An average calculated per evaluation method, and used as an estimate for the evaluation of the regression method.

## 1.6 K-fold cross-validation resampling

Another resampling method, and the last that will be used here, k-fold cross-validation splits a dataset into k equally sized sets, and treats one of them at a time as a testing set while treating the collection of the remaining  $k - 1$  sets collectively as a training set. Any evaluations using k-fold cross-validation are evaluated as an average over the k test sets.

## 1.7 Ridge

The ridge method is an adjustment to the OLS method, and involves introducing a term that is added to the matrix about to be inverted, in the form of a scalar constant  $\lambda$  multiplied into an identity matrix. This helps makes uninvertible matrices invertible, and assists the models accuracy.

$$\beta = \left( \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} \right)^{-1} \mathbf{X}^T y$$

The key to using Ridge is finding an optimal  $\lambda$ .

## 1.8 Lasso

The lasso method differs slightly from ridge in that it attempts to minimize  $\sum_{i=1}^n (y_i - \beta_0 - x_i^T \beta)^2$  in terms of  $\beta$ . In this report, scikit learn's method is used in place of creating one;

```
1 # Python example of lasso implementation using the sklearn module
```

```

2 # Assuming dataset has been split into X_train, X_test, y_train, and ↵
   y_test
3
4 import sklearn.linear_model as linmod
5 lassoAlpha = 0.1
6 clf = linmod.Lasso(alpha=lassoAlpha)
7 clf.fit(X_train, y_train)
8 ypredict = clf.predict(X_test)

```

## 2 Method

### 2.1 Data generation

The class is given 3 initial arguments;  $n$ ,  $p$ , and  $\text{noisefactor}$ .  $n$  refers to the square of the number of data points used,  $p$  refers to the polynomial degree of the design matrix, and  $\text{noisefactor}$  refers to the amount of noise that should be present in the data.

Raw data generation:

First,  $\frac{n}{2}$  points, rounded up, are distributed with linear distribution within the domain of Franke; between 0 and 1.

Second,  $\frac{n}{2}$  points, rounded down, are distributed with uniform random distribution within the same domain.

Third, these two sets are combined to form a 1D dataset.

Fourth, this is repeated and the two datasets,  $x$  and  $y$ , are made to form a meshgrid which is then flattened, creating an even and reliable full, 2D dataset. It is noted that this is a better dataset than one could expect normally, but a similar method is used to extract data from the geographical dataset when the required space or processing power for the full dataset is not available.

Fifth, expected data is generated from these data with the Franke function.

Sixth, normally distributed noise of mean 0 and variance 1, multiplied by the noise factor and the mean of the Franke data, are added to the Franke data.

```

1 def __init__(self, n, p, noisefactor = 0.1):
2     x1 = y1 = np.linspace(0, 1, mt.ceil(n/2))
3     x2 = np.random.uniform(0, 1, mt.floor(n/2))
4     y2 = np.random.uniform(0, 1, mt.floor(n/2))
5
6     x = np.concatenate((x1, x2))
7     y = np.concatenate((y1, y2))
8
9     x, y = np.meshgrid(x,y)
10    self.x = x.flatten()
11    self.y = y.flatten()
12
13    self.ytrue = self.FrankeFunction(self.x, self.y)
14    self.yData = self.ytrue + noisefactor*np.random.randn(n**2)*self.↵
        ytrue.mean()

```

Design Matrix generation: The design matrix is made as a polynomial set of  $x$  and  $y$ , where each column is a polynomial product of the two, and the rows are points of data. The non-constant columns (i.e. all columns but the first) are scaled by subtracting the average value of each respective column from itself.

```

1 def craftX(self, scaling=True):
2     self.nfeatures = int(((self.p+1)*(self.p+2))/2)
3     self.X = np.zeros((len(self.x), self.nfeatures))

```

```

4     ind = 0
5     for i in range(self.p+1):
6         for j in range(self.p+1-i):
7             self.X[:,ind] = self.x**i * self.y**j
8             ind += 1
9
10    if scaling:
11        self.X[:,1:] -= np.mean(self.X[:,1:], axis=0)
12

```

## 2.2 a) Ordinary Least Square (OLS) on the Franke function

Implemented in this section:

- Data generation
- $R^2$  score function
- Mean Squared error (MSE) evaluation function
- Design Matrix & Franke values
- OLS
- Method to find the 95% confidence intervals of the beta parameters of the OLS-method

```

1  # Main OLS method
2  beta = np.linalg.pinv(X_train.T @ X_train) @ X_train.T @ y_train

```

```

1  def MSE(self, y_data, y_model):
2      # Optimal value is 0, with higher values beign worse
3      return np.sum((y_data-y_model)**2) / np.size(y_model)
4
5  def R2(self, y_data, y_model):
6      # Optimal value is 1, with 0 implying that model performs
7      # exactly as well as predicting using the data average would.
8      # Lower values imply that predicting using the average would be ↵
9      # better
10     top = np.sum((y_data - y_model)**2)
11     bot = np.sum((y_data - np.mean(y_data))**2)
12     return 1 - top/bot

```

This section generated data as described in Section 2.1, attempted to model it with an OLS approximation, and evaluated the MSE,  $R^2$  score, and the confidence interval of the  $\beta$ -parameters.

The confidence interval of the  $\beta$ -parameters was found by running an OLS model 300 times, finding 300 examples of the  $\beta$ -parameters, calculating the standard deviation,  $\sigma$ , for the parameters, then adding and subtracting  $1.645\sigma$  to and from the mean of the parameters.

## 2.3 b) Bias-variance trade-off and resampling techniques

Implemented in this section:

- OLS plotting training and test MSE's over complexity

- OLS with bootstrap resampling technique
- Bias-variance analysis

This section aimed to evaluate OLS by plotting error in testing and training data together against a varying complexity in the model, evaluate OLS by using bootstrap resampling, and perform a bias-variance trade-off analysis using bootstrap.

Noisefactor is set to 0 in this section to find an error in the calculations for bias and variance.

```

1 def sample(self, sourceX, sourceY, Nsamples = 0.6):
2     if isinstance(Nsamples, float):
3         Nsamples = int(Nsamples*sourceX.shape[0])
4
5     sampleArrayX = np.zeros(sourceX[:Nsamples,:].shape)
6     sampleArrayY = np.zeros(Nsamples)
7     for i in range(Nsamples):
8         ind = np.random.randint(Nsamples)
9         sampleArrayX[i] = sourceX[ind]
10        sampleArrayY[i] = sourceY[ind]
11
12    return sampleArrayX, sampleArrayY

```

```

1 def bootstrap(self, sampleSize=None, sampleN=None):
2     if sampleSize is None: sampleSize = 0.8
3     if sampleN is None: sampleN = 5
4     self.craftX()
5     if isinstance(sampleSize, float): sampleSize = int(sampleSize*len(↵
        self.X))
6     r2list, SElist = np.zeros((2, sampleN))
7     X_train, X_test, y_train, y_test = train_test_split(self.X, self.↵
        yData)
8     for i in range(sampleN):
9         X_sample, y_sample = self.sample(X_train, y_train)
10        r2list[i], SElist[i] = self.OLS([X_sample, X_test, y_sample, ↵
        y_test])[:2]
11    return r2list, SElist

```

```

1 def bias(self, y_data, y_model):
2     return np.mean((y_data - np.mean(y_model, axis=1))**2)
3
4 def variance(self, y_model):
5     return np.mean(np.var(y_model, axis=1))

```

The bias-variance analysis was implemented by using bootstrap OLS over a range of polynomial complexities and plotting the bias, variance, and the MSE. It similarly plotted the training MSE and test MSE over complexity.

## 2.4 c) Cross-validation as resampling techniques, adding more complexity

Implemented in this section:

- k-fold cross-validation

This section implements k-fold cross-validation resampling and compares it to bootstrap resampling, using the same evaluation methods of  $R^2$  and MSE.

k-fold cross-validation and bootstrap are compared for equivalent number of folds and number of samples, 5 and 10.

The k-fold crossvalidation was implemented in three functions; `kfold_splitter()`, `kfold_yielder()`, and `kfold()`. `kfold_splitter()` shuffles and splits the data into several versions of training sets and testing sets, as a 3D-array for the inputs and a 2D-array for the outputs, functioning as collections of 2D-arrays and 1D-arrays. `kfold_yielder()` then functions as a generator, supplying these sets as appropriate 2- and 1D-arrays to `kfold()`, which performs OLS on the sets generated in turn.

## 2.5 d) Ridge Regression on the Franke function with resampling

Implemented in this section:

- Ridge
- Ridge with bootstrap
- Ridge with k-fold cross-validation

This section implements Ridge and applies bootstrap and k-fold cross-validation resampling techniques to it. It was also supposed to implement a bias-variance analysis, but as of this moment, it is only a copy of the function defined in `partB.py`, and does not do the job it's supposed to.

## 2.6 e) Lasso Regression on the Franke function with resampling

Implemented in this section:

- Lasso method
- Bias-variance analysis using lasso and bootstrap

This section implements scikit learn's lasso method with no resampling and inserts it into a bias-variance analysis function with bootstrap resampling, similar to what was done in subsection 2.3.

```
1 def lasso(self, alpha=None):
2     if alpha is None: alpha = 5
3     clf = linmod.Lasso(alpha=alpha)
4     self.craftX()
5     X_train, X_test, y_train, y_test = train_test_split(self.X, self.yData)
6     clf.fit(X_train, y_train)
7     ypredict = clf.predict(X_test)
8     return self.R2(y_test, ypredict), self.MSE(y_test, ypredict)
```

## 2.7 f) Introducing real data and preparing the data analysis

Implemented in this section:

- Geographical data

This section imports geographical data from [1]

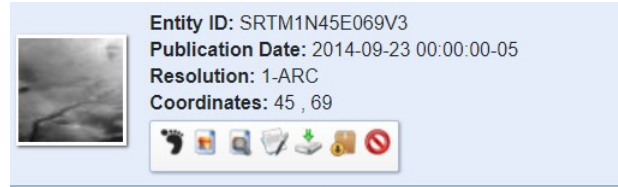


Figure 1: Metadata for the data downloaded.

## 2.8 g) OLS, Ridge and Lasso regression with resampling

Implemented in this section:

- New data generation

This section was intended to use all methods developed previously on the geographical data, compared and evaluated them against each other as viable methods. The ridge method was to be prepared as a 2D heatmap of polynomial complexity and  $\lambda$ -value. The methods were to be compared to scikit learn's methods to evaluate them as relatively simple implementations against more complex ones. Given the time constraints, the majority of this has not been implemented.

This section, instead, implements geographical data extraction and the train-test MSE over polynomial complexity from section b) over this dataset.

There also exists in the code an unfinished k-fold cross-validation version of the this train-test plot, but this was not completed in time.

## 3 Results/Discussion

Results (printouts) are described in more detail for each section in the folder outputs/.

### 3.1 a) Ordinary Least Square (OLS) on the Franke function

With 10,000 datapoints and a polynomial degree of 5, OLS reaches an  $R^2$  value of 0.95 with an MSE of 0.004. If the data generation method is changed to being  $n^2$  points linearly distributed between (0,0), (0,1), (1,0) and (1,1), the  $R^2$  rises to 0.96, indicating that the extra steps taken to vary the dataset actually hinders the quality of the model.



### 3.2 b) Bias-variance trade-off and resampling techniques

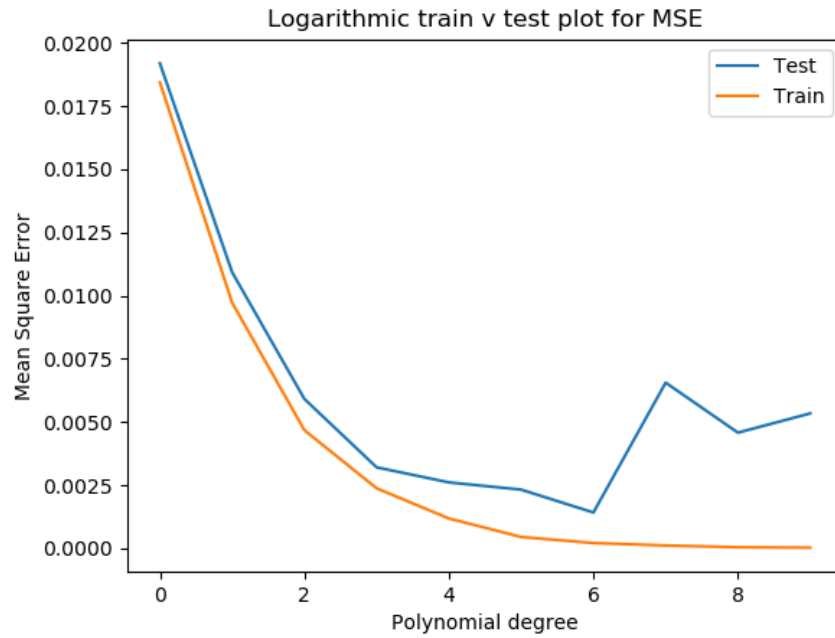


Figure 2: Testing and training errors proceeding as expected over model complexity; training error is continually falling whilst testing data falls to a point and then rises again due to overfitting.

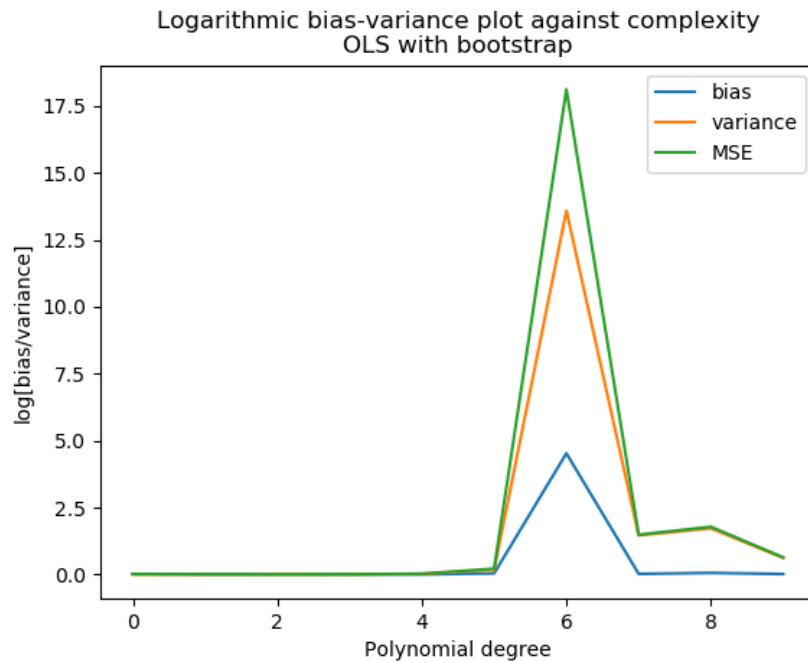


Figure 3: Bias and variance behaving unexpectedly, having a massive spike 17 orders of magnitude above the rest in both bias and variance.

Figure 3 suggests that something has gone wrong during the calculations of bias and variance. Possible reason include:

- A poor model: In addition to the well-behaved Figure 2, my model has been compared to sklearn's model (see lines 101-104 of partB.py) and both models exhibit this behaviour.
- Incorrect equations for bias, variance or MSE: All three values are calculated independently, and the relationship  $MSE = Bias + Variance$  is found consistently (see last line of printout in outputs/partB.txt or line 114 in partB.py), implying that the calculations are correct.

I have been unable to find any other source of error in the time I have for this project.

### 3.3 c) Cross-validation as resampling techniques, adding more complexity

```

1 Folds/number of samples: 5
2 k-fold [R2, MSE] : [0.95563823 0.00368468]
3 Bootstrap [R2, MSE] : [0.9545114 0.00377724]
4 kfold/bootstrap [R2, MSE] : [1.00118054 0.9754959 ]
5
6 Folds/number of samples: 10
7 k-fold [R2, MSE] : [0.95312732 0.00396122]
8 Bootstrap [R2, MSE] : [0.9526719 0.00386939]
9 kfold/bootstrap [R2, MSE] : [1.00047805 1.02373248]

```

Comparing 5-fold cross-validation against 5-sample bootstrap and 10-fold cross-validation against 10-sample bootstrap shows that the  $R^2$  and MSE values are very similar, suggesting that the resampling methods have no impact on the actual accuracy of the model.

### 3.4 d) Ridge Regression on the Franke function with resampling

```

1 Normal OLS; average values:
2 OLS, multiple times [R2, MSE] : [0.9533851 0.0038293]
3 Bootstrap           [R2, MSE] : [0.9545558 0.0038718]
4 kfold               [R2, MSE] : [0.9540727 0.0037702]
5
6 Ridge application; optimal values:
7 Bootstrap           [R2, MSE] : [0.9524293 0.0039346]
8 k-fold              [R2, MSE] : [0.9540655 0.0037696]

```

The optimal value for  $\lambda$  is  $\lambda = 0$ , meaning that for this application, ridge provides no benefit, but it also does not detract from the accuracy in any substantial way.

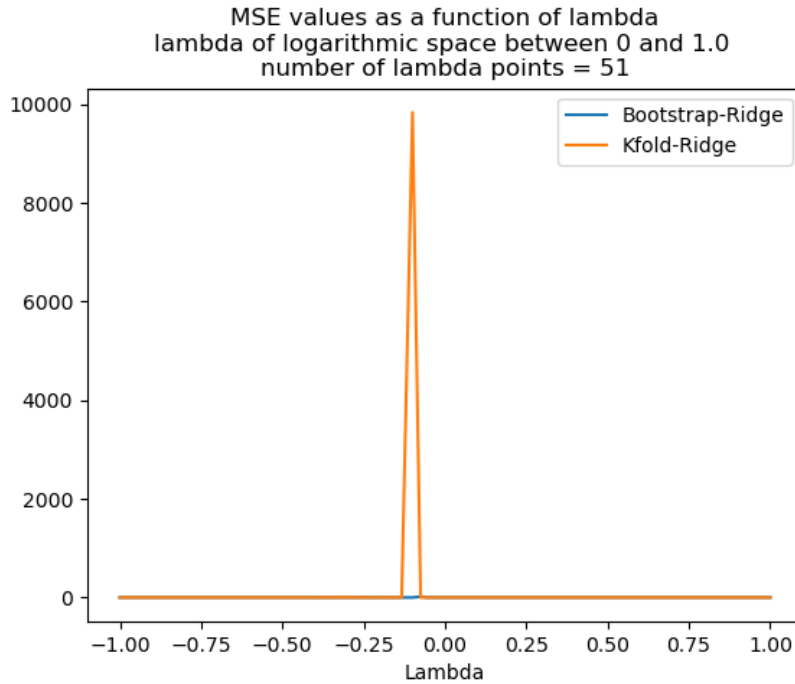


Figure 4: Ridge implemented a variable for MSE, using both bootstrap and k-fold cross-validation resampling methods. There is a big spike on the negative side of 0, possibly implying that the matrix is not invertible in that area, or that the issue seen in Figure 3 may extend to this.

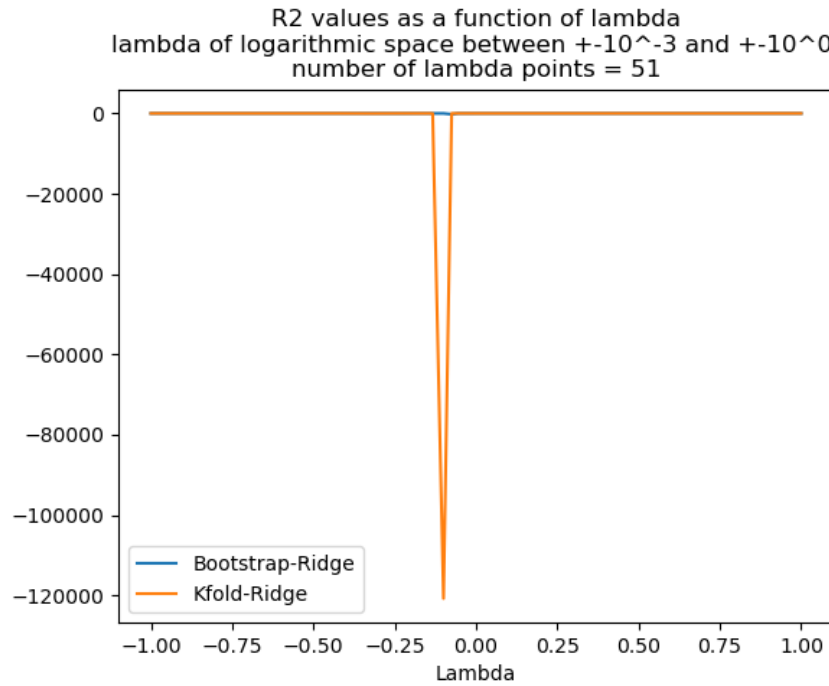


Figure 5: Similar to Figure 4, this figure presents the  $R^2$ -evaluation, seeing a big dip on the negative side of 0.

### 3.5 e) Lasso Regression on the Franke function with resampling

```

1 Lasso evaluation [R2, MSE] at alpha = 0.01:
2   (-0.00040418848759649073, 0.09415771032493848)
3 Lasso evaluation [R2, MSE] at alpha = 0.1:
4   (-0.0008014043714434926, 0.09901736017684995)
5 Lasso evaluation [R2, MSE] at alpha = 0.3:
6   (-6.013665053505868e-05, 0.09889216375005189)
7 Lasso evaluation [R2, MSE] at alpha = 1:
8   (-0.00011187832184433866, 0.09252839549202323)
9 Lasso evaluation [R2, MSE] at alpha = 3:
10  (-0.00013499053518417625, 0.09445379258455494)
11 MSE/(bias+variance) = [1. 1. 1. 1. 1. 1. 1. 1. 1.]

```

The lasso implementation reaches an  $R^2$  value of -0.001, which is a very poor value, but it also reaches an MSE value down to 0.08, implying that the model might not be bad, but that something else could have gone wrong. This may be a similar issue to what occurs in Figure 3

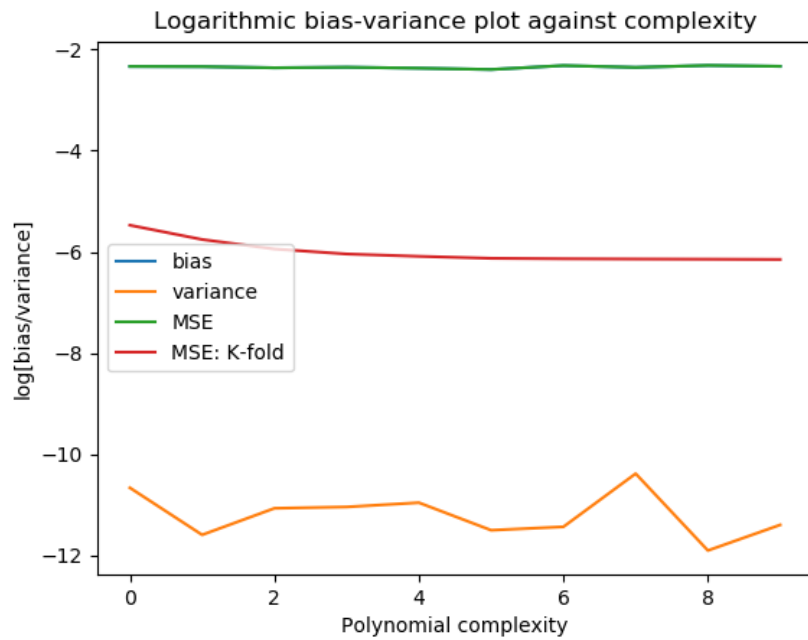


Figure 6: Bias is hidden by the MSE. MSE:K-fold is the MSE calculated during the caluculations, rather than after-the-fact, and the two MSE's disagree, implying that something has gone wrong.

### 3.6 f) Introducing real data and preparing the data analysis

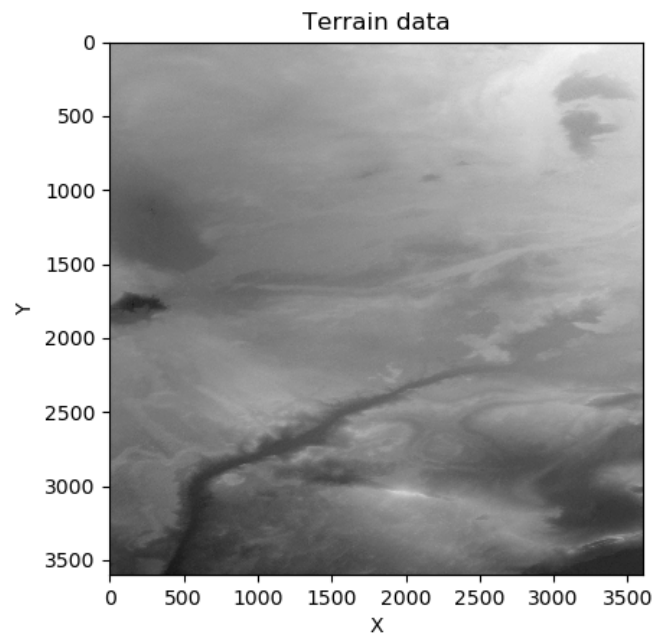


Figure 7: Geographical data, represented on a gray scale.

### 3.7 g) OLS, Ridge and Lasso regression with resampling

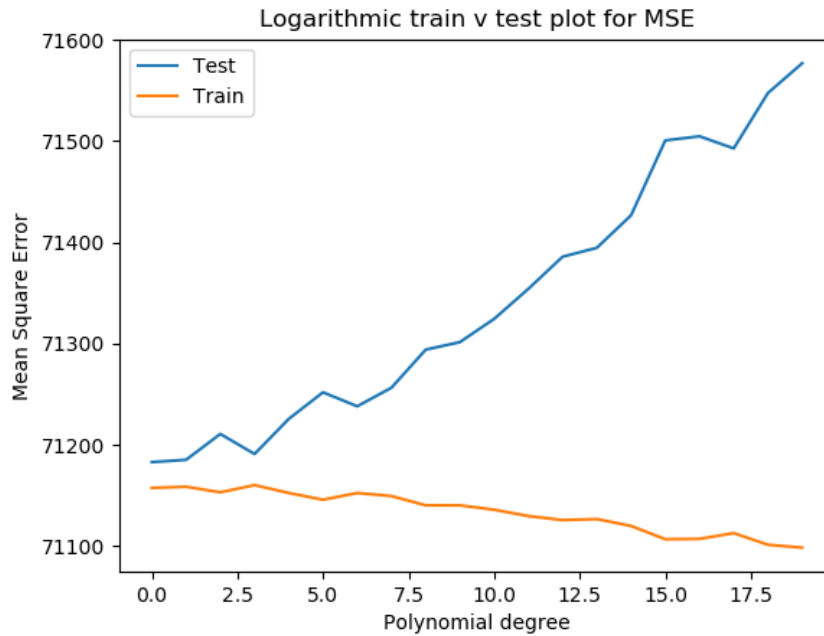


Figure 8

```
1 R2: ( avg = -126.007750214505 )
2     ( max = -125.28996659328048 )
3     ( min = -126.63064184509612 )
4     [-125.28996659 -126.00075666 -125.99061752
5      -126.63064185 -126.33000825 -125.80451042]
6
7 MSE: ( avg = 71177.81149110844 )
8     ( max = 71218.32821713865 )
9     ( min = 71113.55185997034 )
10    [71175.41435112 71218.32821714 71167.27618466
11     71178.49262646 71213.80570729 71113.55185997]
```

Given the lack of different evaluations, it is difficult to see anything in the data; the poor values could be due to OLS being a poor fit for this data, or it could be a poor implementation. Without additional approaches or different real datasets to compare to, it is difficult to say.

## 4 Conclusion

This project is riddled with symptoms of being unfinished, and there are few conclusions to be drawn. There are, however, many things to improve on. The initial train-test plot in section b) seems promising, indicating that the base Franke data generation, the basic OLS implementation, and the base evaluation methods are all implemented properly, and fit well with the data. I don't

currently have enough information to be able to tell what is wrong with the bias-variance plot in the same section, but given how early it is implemented (although chronologically it was implemented near the end of the project) it should be one of the first places that are looked to for errors and mistakes.

The large spike in the ridge plots in section d) implies that there might be an issue regarding the invertibility of  $\mathbf{X}^T \mathbf{X}$ , which might lead to very unstable solutions near  $\lambda = 0$ , which includes the general OLS-solution. This can serve as part of a one-of-many-problems explanation, but is unlikely to be the entire issue, given that the train-test plot and bias-variance plots in section b) behave so differently.

Section c) shows that the resampling methods, bootstrap and k-fold cross-validation, seem to behave well, indicating that they are well implemented and function well for the Franke dataset. Although their aptness for the dataset is derived purely from the base model they use in this instance, OLS, and from the resampling methods themselves, so their fitting the Franke dataset is no surprise given that OLS already fit.

The Lasso implementation, as mentioned in subsection 3.5, has a curious case of a consistently poor  $R^2$  and a consistently good MSE. This is made even more curious by the fact that this is scikit learn's own method implemented, and that it predicts with similar accuracy to a straight average ( $R^2$  score of near 0). Unfortunately, there has not been enough time to look into this curiosity.

The only conclusion that could be drawn about the final section, section g), is that the model used to model it is a poor fit at all complexities used, unfortunately. This section could benefit most from more time, even though I believe that time would have to start at the beginning and root out errors that may propagate throughout the entire project.

Finally, the entire source code needs to have bugs rooted out, functions commented, a greater generalization, a refactoring, and to be finished. Unfortunately, I cannot do this given that I don't have the time to.

## 5 References

- 1 <https://earthexplorer.usgs.gov/>
- 2 Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. Springer, 2017



## 6 Appendix

### 6.1 Source code

#### 6.1.1 a) Ordinary Least Square (OLS) on the Franke function

```
1 import numpy as np
2 from random import random, seed
3 from sklearn.model_selection import train_test_split
4 from matplotlib.pyplot import plot, show, legend
5 import math as mt
6
7
8 # Uses a two-dimensional polynomial of degree p to model Franke's ↵
9   Function
10  class learningA:
11      def __init__(self, n, p, noisefactor=None):
12          if noisefactor is None: noisefactor = 0.1
13          x1 = y1 = np.linspace(0, 1, mt.ceil(n/2))
14          x2 = np.random.uniform(0, 1, mt.floor(n/2))
15          y2 = np.random.uniform(0, 1, mt.floor(n/2))
16
17          x = np.concatenate((x1, x2))
18          y = np.concatenate((y1, y2))
19
20          x, y = np.meshgrid(x, y)
21          self.x = x.flatten()
22          self.y = y.flatten()
23
24          self.ytrue = self.FrankeFunction(self.x, self.y)
25          self.yData = self.ytrue + noisefactor*np.random.randn(n**2)*↵
26              self.ytrue.mean()
27
28          self.n = n
29          self.p = p
30          self.noisefactor = noisefactor
31          self.nfeatures = int(((self.p+1)*(self.p+2))/2)
32
33      def MSE(self, y_data, y_model):
34          # Optimal value is 0, with higher values being worse
35          return np.sum((y_data - y_model)**2) / np.size(y_model)
36
37      def R2(self, y_data, y_model):
38          # Optimal value is 1, with 0 implying that model performs
39          # exactly as well as predicting using the data average would.
40          # Lower values imply that predicting using the average would ↵
41          # be better
42          top = np.sum((y_data - y_model)**2)
43          bot = np.sum((y_data - np.mean(y_data))**2)
44          return 1 - top/bot
45
46      def FrankeFunction(self, x, y):
47          term1 = 0.75*np.exp(-(9*x-2)**2/4.00 - 0.25*((9*y-2)**2))
48          term2 = 0.75*np.exp(-(9*x+1)**2/49.0 - 0.10*(9*y+1))
49          term3 = 0.50*np.exp(-(9*x-7)**2/4.00 - 0.25*((9*y-3)**2))
50          term4 = -0.20*np.exp(-(9*x-4)**2 - (9*y-7)**2)
51          return term1 + term2 + term3 + term4
52
53      def generateNewDataset(self, n=None, p=None):
54          if n is None: n = self.n
55          else: self.n = n
56          if p is None: p = self.p
57          else: self.p = p
58
59          x1 = y1 = np.linspace(0, 1, mt.ceil(n/2), dtype=int)
60          x2 = np.random.uniform(0, 1, mt.floor(n/2))
61          y2 = np.random.uniform(0, 1, mt.floor(n/2))
62
63          x = np.concatenate((x1, x2))
64          y = np.concatenate((y1, y2))
```

```

62         x, y = np.meshgrid(x,y)
63         self.x = x.flatten()
64         self.y = y.flatten()
65
66
67         self.ytrue = self.FrankeFunction(self.x, self.y)
68         self.yData = self.ytrue + self.noisefactor*np.random.randn(n←
        **2)*self.ytrue.mean()
69
70     def craftX(self, scaling=True):
71         self.nfeatures = int(((self.p+1)*(self.p+2))/2)
72         self.X = np.zeros((len(self.x), self.nfeatures))
73
74         ind = 0
75         for i in range(self.p+1):
76             for j in range(self.p+1-i):
77                 self.X[:,ind] = self.x**i * self.y**j
78                 ind += 1
79
80         if scaling:
81             self.X[:,1:] -= np.mean(self.X[:,1:], axis=0)
82
83
84     def OLS_core(self, dataset):
85         X_train, X_test, y_train = dataset[:3]
86         beta = np.linalg.pinv(X_train.T @ X_train) @ X_train.T @ ←
            y_train
87         ypredict = X_test @ beta
88         return ypredict, beta
89
90     # Ordinary Least Square solution
91     def OLS(self, dataset=None):
92         if dataset is None:
93             self.craftX()
94             dataset = train_test_split(self.X, self.yData)
95             ypredict, beta = self.OLS_core((dataset))
96         else:
97             ypredict, beta = self.OLS_core(list(dataset))
98         return self.R2(dataset[3], ypredict), self.MSE(dataset[3], ←
            ypredict), beta
99
100     # Generates R2- and MSE-values from the OLS function
101     def OLSeval(self, reps):
102         r2list, SElist = np.zeros((2, reps))
103         for i in range(reps):
104             r2list[i], SElist[i] = self.OLS()[2:]
105
106         return r2list, SElist
107
108     # Calculates confidence interval of the parameters for OLS
109     def confidenceIntervalOLS(self, reps=None):
110         if reps is None: reps = 300
111         betas = np.zeros((reps, self.nfeatures))
112         for rep in range(reps):
113             betas[rep,:] = self.OLS()[2]
114
115         # Calculating the standard deviation
116         betasAvg = np.sum(betas, axis = 0)/reps
117         betasDiff = betas - betasAvg
118
119         STD = np.sqrt((1/(reps-self.p-1))*np.sum(betasDiff**2, axis = ←
            0))
120
121         magicNumberFor95PercentConfidenceInterval = 1.645
122         ShortHand = magicNumberFor95PercentConfidenceInterval
123
124         confidenceInterval = [betasAvg-ShortHand*STD, betasAvg+←
            ShortHand, STD**2]
125
126         # returns lower and upper boundaries of confidence interval ←
            with the variance
127         return confidenceInterval
128
129

```

```

130
131 if __name__ == "__main__":
132     print("""Task as interpreted:
133     (x) Generate a dataset as a bipolynomial to a general order
134     (x) Implement and execute OLS on bipolynomial and Franke
135     (x) Find confidence interval of parameters beta
136     (x) Evaluate R2 and mean square error
137     (x) Scale and split the data
138     """)
139
140
141
142     Npoints = 100
143     polydegree = 5
144     OLSruns = 18
145
146     OLSlearner = learningA(Npoints, polydegree)
147
148     print("Npoints = ", Npoints, ", giving a number of data points = "↵
149           , Npoints**2, sep="")
150     print("Polynomial degree = ", polydegree, ", giving a number of ↵
151           features = ",int(((polydegree+1)*(polydegree+2))/2), sep="")
152     print("OLS run", OLSruns, "times")
153
154     for name, result in zip(["R2", "MSE"], OLSlearner.OLSseval(OLSruns)↵
155           ):
156         print("\n" + name +
157               ": ( avg =", sum(result)/len(result),
158               ", ( max = ", max(result),
159               ", ( min = ", min(result), ")\n",
160               result)
161
162     lower, upper, STD = OLSlearner.confidenceIntervalOLS()
163     print("""\n95 Confidence Interval for beta:
164     Lower: {}↵
165     Upper: {}↵
166     STD : {}""".format(lower, upper, STD))
167
168     # Changing the data generation method
169     x = np.linspace(0,1,OLSlearner.n)
170     y = np.linspace(0,1,OLSlearner.n)
171     x, y = np.meshgrid(x,y)
172     OLSlearner.x = x.flatten()
173     OLSlearner.y = y.flatten()
174     OLSlearner.ytrue = OLSlearner.FrankeFunction(OLSlearner.x, ↵
175           OLSlearner.y)
176     OLSlearner.yData = OLSlearner.ytrue + OLSlearner.noisefactor*np.↵
177           random.randn(OLSlearner.n**2)*OLSlearner.ytrue.mean()
178
179     print("\n\nOLS with data determined just by linspace")
180     for name, result in zip(["R2", "MSE"], OLSlearner.OLSseval(OLSruns)↵
181           ):
182         print("\n" + name +
183               ": ( avg =", sum(result)/len(result),
184               ", ( max = ", max(result),
185               ", ( min = ", min(result), ")\n",
186               result)

```

### 6.1.2 b) Bias-variance trade-off and resampling techniques

```
1 from partA import learningA
2 import numpy as np
3 from random import random, seed
4 from sklearn.model_selection import train_test_split
5 from matplotlib.pyplot import plot, show, legend, title, ylabel, ←
   xlabel, savefig
6
7 class learningB(learningA):
8     def __init__(self, n, p, noisefactor=None):
9         super().__init__(n, p, noisefactor)
10        self.imageFilePath = "../outputs/images/partB/"
11
12    def sample(self, sourceX, sourceY, Nsamples=None):
13        if Nsamples is None: Nsamples = 0.8
14        if isinstance(Nsamples, float):
15            Nsamples = int(Nsamples*sourceX.shape[0])
16
17        sampleArrayX = np.zeros(sourceX[:Nsamples,:].shape)
18        sampleArrayY = np.zeros(Nsamples)
19        for i in range(Nsamples):
20            ind = np.random.randint(Nsamples)
21            sampleArrayX[i] = sourceX[ind]
22            sampleArrayY[i] = sourceY[ind]
23
24        return sampleArrayX, sampleArrayY
25
26    def plotOLS_trainvtest(self, reps, p_range):
27        # Made to be similar to Fig. 2.11 of Hastie, Tibshirani, and ←
           Friedman
28        p_old = self.p
29        p_rangeObject = range(p_range[0], p_range[1])
30        SElist_train, SElist_test = np.zeros((2,p_range[1] - p_range←
           [0]))
31
32        for i in p_rangeObject:
33            self.p = i + p_old
34            self.craftX()
35
36            for _ in range(reps):
37                X_train, X_test, y_train, y_test = train_test_split(←
                    self.X, self.yData)
38                ypredict_test, beta = self.OLS_core([X_train, X_test, ←
                    y_train, y_test])
39                ypredict_train = X_train @ beta
40
41                SElist_test[i-p_range[0]] += self.MSE(y_test, ←
                    ypredict_test)/reps
42                SElist_train[i-p_range[0]] += self.MSE(y_train, ←
                    ypredict_train)/reps
43
44        self.p = p_old
45
46        imageFileName = "trainTestMSE"
47        plot(p_rangeObject, (SElist_test), label="Test")
48        plot(p_rangeObject, (SElist_train), label="Train")
49        legend()
50        title("Logarithmic train v test plot for MSE")
51        xlabel("Polynomial degree")
52        ylabel("Mean Square Error")
53        savefig(self.imageFilePath + imageFileName + ".png")
54        show()
55
56    def bootstrap(self, sampleSize=None, sampleN=None):
57        if sampleSize is None: sampleSize = 0.8
58        if sampleN is None: sampleN = 5
59
60        self.craftX()
61        X_train, X_test, y_train, y_test = train_test_split(self.X, ←
            self.yData)
62
```

```

63         if isinstance(sampleSize, float): sampleSize = int(sampleSize*←
        X_test.shape[0])
64
65         r2list, SElist = np.zeros((2, sampleN))
66         for i in range(sampleN):
67             X_sample, y_sample = self.sample(X_train, y_train, ←
        sampleSize)
68             r2list[i], SElist[i] = self.OLS([X_sample, X_test, ←
        y_sample, y_test])[:2]
69
70         return r2list, SElist
71
72     def bias(self, y_data, y_model):
73         return np.mean((y_data - np.mean(y_model, axis=1))**2)
74
75     def variance(self, y_model):
76         return np.mean(np.var(y_model, axis=1))
77
78     def biasVarianceAnalysis_bootstrap(self, p_range, sampleSize=None, ←
        sampleN=None):
79         if sampleSize is None: sampleSize=0.8
80         if isinstance(sampleSize, float): sampleSize = int(sampleSize*←
        len(self.x))
81         if sampleN is None: sampleN=5
82
83         p_old = self.p
84         p_rangeObject = range(p_range[0], p_range[1])
85
86         bias = np.zeros(len(p_rangeObject))
87         variance = bias.copy()
88         SElist = bias.copy()
89
90         # Imported to compare to my own model, see line 102
91         from sklearn.linear_model import LinearRegression
92
93         for i in p_rangeObject:
94             self.p = i + p_old
95             self.craftX()
96
97             X_train, X_test, y_train, y_test = train_test_split(self.X←
        , self.yData)
98
99             ypredict_models = np.zeros((y_test.shape[0], sampleN))
100
101             for j in range(sampleN):
102                 X_sample, y_sample = self.sample(X_train, y_train)
103
104
105                 ypredict = self.OLS_core([X_sample, X_test, y_sample, ←
        y_test])[0]
106                 # My OLS method was compared to ScikitLearns method, ←
        and was found not to be the issue
107                 # reg = LinearRegression().fit(X_sample, y_sample)
108                 # ypredict = reg.predict(X_test)
109                 # print("HERE WE GO", np.mean(test_var/ypredict))
110                 ypredict_models[:,j] = ypredict
111
112                 variance[i-p_range[0]] = self.variance(ypredict_models)
113                 # bias[i-p_range[0]] = np.mean((y_test - np.mean(←
        ypredict_models, axis=1))**2)
114                 bias[i-p_range[0]] = self.bias(y_test, ypredict_models)
115
116                 SElist[i-p_range[0]] = np.mean( np.mean((y_test.reshape←
        (-1,1) - ypredict_models)**2, axis=1))
117
118         print("\nMSE divided by (bias+variance):", SElist/(bias+←
        variance))
119
120         self.p = p_old
121
122         imageFileName = "biasVariance"
123         filetype = ".png"
124
125         plot(p_rangeObject, (bias), label="bias")

```

```

126     plot(p_rangeObject, (variance), label="variance")
127     plot(p_rangeObject, (SElist), label="MSE")
128
129     legend()
130     title("Logarithmic bias-variance plot against complexity\nOLS ←
        with bootstrap")
131     xlabel("Polynomial degree")
132     ylabel("log[bias/variance]")
133     savefig(self.imageFilePath + imageFileName + filetype)
134     show()
135
136
137
138
139
140
141 if __name__ == "__main__":
142     print("""Task as interpreted:
143     (?) "general aim: study bias-variance trade-off by implementing ←
        bootstrap resampling"
144     (x) Implement OLS with resampling techniques (bootstrap)
145     (x) Generate a figure similar to Fig. 2.11, Hastie (showing test&←
        training MSE's)
146     (x) perform a bias-variance analysis by: Comparing MSE to ←
        complexity (MSE v polydegree), make a graph
147     (x) Do some equations (mentioned in the introduction part of the ←
        report)
148     (x) Describe bias and variance as part of those equations;
149     (x) Discuss said bias and variance trade-off as a function of ←
        complexity, # of data points, and possibly also bootstrap
150     """)
151
152
153
154     Npoints = 15
155     polydegree = 2
156     sampleN = 100
157     polydegree_range = [0,10]
158     trainvtest_reps = 18*10
159     noisefactor = 0
160
161     Bootstraplearner = learningB(Npoints, polydegree, noisefactor)
162
163     for name, result in zip(["R2", "MSE"], Bootstraplearner.bootstrap(←
        sampleN = sampleN)):
164         print("\n" + name +
165             ": ( avg = ", sum(result)/len(result),
166             ", ( max = ", max(result),
167             "), ( min = ", min(result), " )\n",
168             result)
169
170     sampleN *= 1
171
172     Bootstraplearner.plotOLS_trainvtest(trainvtest_reps, ←
        polydegree_range)
173     Bootstraplearner.biasVarianceAnalysis_bootstrap(polydegree_range, ←
        sampleN=sampleN)

```

### 6.1.3 c) Cross-validation as resampling techniques, adding more complexity

```
1 from partB import learningB
2 from sklearn.model_selection import train_test_split
3 import numpy as np
4
5 class learningC(learningB):
6
7     # I stole the shit out of this from StackOverflow
8     # https://stackoverflow.com/questions/4601373/better-way-to-shuffle-two-numpy-arrays-in-unison
9     def unison_shuffled_copies(self, a, b):
10         assert len(a) == len(b)
11         p = np.random.permutation(len(a))
12         return a[p], b[p]
13
14     def kfold_splitter(self, k):
15         X_dimx, X_dimy = self.X.shape
16         X, y = self.unison_shuffled_copies(self.X, self.yData)
17         splitsX = np.zeros((k, int(X_dimx/k), X_dimy))
18         splitsY = np.zeros((k, int(X_dimx/k)))
19         frac = np.floor(X_dimx/k)
20
21         for i in range(k):
22             index_0 = int(i*frac)
23             index_1 = int((i+1)*frac)
24             try:
25                 splitsX[i, :, :] = X[index_0:index_1, :]
26                 splitsY[i, :] = y[index_0:index_1]
27             except ValueError:
28                 raise Exception("kfold function does not split correctly")
29
30         return splitsX, splitsY
31
32     def kfold_yielder(self, k):
33         splitsX, splitsY = self.kfold_splitter(k)
34
35         pointsPerSplit = splitsX.shape[1]
36
37         for i in range(k):
38             X_test = np.zeros((pointsPerSplit, self.nfeatures))
39             Y_test = np.zeros(pointsPerSplit)
40             X_train = np.zeros(((k-1) * pointsPerSplit, self.nfeatures))
41             Y_train = np.zeros((k-1) * pointsPerSplit)
42
43             X_test[:, :] = splitsX[i, :, :]
44             Y_test[:] = splitsY[i, :]
45
46             flattenedX = splitsX.reshape(-1, splitsX.shape[-1])
47             flattenedY = splitsY.reshape(-1)
48
49             X_train[:i*pointsPerSplit, :] = flattenedX[:i*pointsPerSplit, :]
50             X_train[i*pointsPerSplit:, :] = flattenedX[(i+1)*pointsPerSplit:, :]
51
52             Y_train[:i*pointsPerSplit] = flattenedY[:i*pointsPerSplit]
53             Y_train[i*pointsPerSplit:] = flattenedY[(i+1)*pointsPerSplit:]
54
55             yield X_train, X_test, Y_train, Y_test
56
57     def kfold(self, k, solver=None):
58         self.craftX()
59         dataset = self.kfold_yielder(k)
60         r2list, SElist = np.zeros((2, k))
61         for i in range(k):
62             data = next(dataset)
```

```

63         R2, MSE = self.OLS(data)[:2]
64         r2list[i] = R2
65         SElist[i] = MSE
66         return r2list, SElist
67
68
69
70
71
72
73
74 if __name__ == "__main__":
75     print("""Task as interpreted:
76     (X) Implement k-fold cross-validation, evaluate MSE from this
77     (X) Compare MSE from k-fold and bootstrap
78     (X) try 5-10 folds
79     """)
80
81     Npoints = 100
82     polydegree = 5
83     # folds = 10
84     for folds in [5,10]:
85         sampleN = folds
86         kfoldlearner = learningC(Npoints, polydegree)
87         kfoldResults = np.asarray(kfoldlearner.kfold(folds)).sum(axis=
88                                 =1)/folds
89         print("\nFolds/number of samples: ", folds)
90         print("k-fold          [R2, MSE] :", kfoldResults)
91
92         bootstrapResults = np.asarray(kfoldlearner.bootstrap(sampleN =
93                                 sampleN)).sum(axis=1)/sampleN
94         print("Bootstrap          [R2, MSE] :", bootstrapResults)
95
96         print("kfold/bootstrap [R2, MSE] :", kfoldResults/
97             bootstrapResults)

```



### 6.1.4 d) Ridge Regression on the Franke function with resampling

```

1 from partC import learningC
2 from sklearn.model_selection import train_test_split
3 import numpy as np
4 from matplotlib.pyplot import plot, show, legend, title, ylabel, ←
   xlabel, figure, savefig
5
6
7 class learningD(learningC):
8     def __init__(self, n, p, noisefactor=None):
9         super().__init__(n, p, noisefactor)
10        self.imageFilePath = "../outputs/images/partD/"
11
12    def ridge_core(self, dataset):
13        X_train, X_test, y_train, _, lamb = dataset
14        beta = np.linalg.pinv(X_train.T @ X_train + lamb*np.identity(←
            self.nfeatures)) @ X_train.T @ y_train
15        ypredict = X_test @ beta
16        return ypredict, beta
17
18    def ridge(self, dataset=None):
19        if dataset is None:
20            self.craftX()
21            dataset = train_test_split(self.X, self.yData)
22            ypredict, beta = self.ridge_core((dataset))
23        else:
24            ypredict, beta = self.ridge_core(list(dataset))
25        return self.R2(dataset[3], ypredict), self.MSE(dataset[3], ←
            ypredict), beta
26
27    def bootstrapRidge(self, lambSpace, sampleSize=0.8, sampleN=5):
28        lambN = len(lambSpace)
29        self.craftX()
30        if isinstance(sampleSize, float):
31            sampleSize = int(sampleSize*len(self.x))
32        r2list, SElist = np.zeros((2, lambN))
33
34        X_train, X_test, y_train, y_test = train_test_split(self.X, ←
            self.yData)
35
36        for lamb in range(lambN):
37            for _ in range(sampleN):
38                X_sample, y_sample = self.sample(X_train, y_train)
39
40                r2, SE = self.ridge([X_sample, X_test, y_sample, ←
                    y_test, lambSpace[lamb]])[:2]
41                r2list[lamb] += r2/sampleN
42                SElist[lamb] += SE/sampleN
43
44        return r2list, SElist
45
46    def kfoldRidge(self, lambSpace, k, solver=None):
47        self.craftX()
48        lambN = len(lambSpace)
49        R2avg, MSEavg = np.zeros((2, lambN))
50        for lamb in range(lambN):
51            dataset = self.kfold_yielder(k)
52            for data in dataset:
53                R2, MSE = self.ridge(list(data) + [lambSpace[lamb]]) ←
                    [:2]
54                R2avg[lamb] += R2/k
55                MSEavg[lamb] += MSE/k
56        return R2avg, MSEavg
57
58    def biasVarianceAnalysis_ridge(self, p_range, sampleSize=None, ←
        sampleN=None):
59        #TODO: CURRENTLY JUST A COPY OF THE ORIGINAL FUNCTION FROM b), ←
        STILL NEEDS TO BE ADAPTED FOR RIDGE
60        p_old = self.p
61        p_rangeObject = range(p_range[0], p_range[1])
62        if sampleSize is None: sampleSize=0.8

```

```

63         if isinstance(sampleSize, float): sampleSize = int(sampleSize*len(self.x))
64         if sampleN is None: sampleN=5
65
66         bias = np.zeros(len(p_rangeObject))
67         variance = bias.copy()
68         SElist = bias.copy()
69
70         for i in p_rangeObject:
71             self.p = i + p_old
72             self.craftX()
73
74             X_train, X_test, y_train, y_test = train_test_split(self.X, self.yData)
75
76             ypredict_models = np.zeros((y_test.shape[0], sampleN))
77
78             for j in range(sampleN):
79                 X_sample, y_sample = self.sample(X_train, y_train)
80
81                 ypredict = self.OLS_core([X_sample, X_test, y_sample, y_test])[0]
82
83                 ypredict_models[:,j] = ypredict
84
85                 test_var = np.var(ypredict_models, axis=1)
86
87                 print(test_var, len(test_var), max(test_var), np.mean(test_var), sep="\n", end="\n\n")
88                 variance[i-p_range[0]] = self.variance(ypredict_models)
89                 bias[i-p_range[0]] = np.mean((y_test - np.mean(ypredict_models, axis=1))**2)
90
91                 SElist[i-p_range[0]] = np.mean(np.mean((y_test.reshape(-1,1) - ypredict_models)**2, axis=1))
92
93         self.p = p_old
94
95         plot(p_rangeObject, np.log(bias), label="bias")
96         plot(p_rangeObject, np.log(variance), label="variance")
97         plot(p_rangeObject, np.log(SElist), label="MSE")
98
99         legend()
100        title("Logarithmic bias-variance plot against complexity")
101        xlabel("Polynomial complexity")
102        ylabel("log[bias/variance]")
103        show()
104
105    if __name__ == "__main__":
106        print("***Task as interpreted")
107        (x) Implement ridge method with kfold and bootstrap
108        (-) Change bootstrap and kfold to take method as argument
109        Too much fixing required
110        (x) Create bootstrapRidge & kfoldRidge
111        (x) Compare ridge-Bootstrap, ridge-kfold, Bootstrap, and kfold
112        (-) Study bias-variance trade-off of various values of lambda using bootstrap.
113        (-) Comment on the above
114        """
115
116        Npoints = 100
117        polydegree = 5
118        lambN = 50
119        logLow = -3
120        logHigh = 0
121
122        ridgelearner = learningD(Npoints, polydegree)
123        lambSpace = np.logspace(logLow, logHigh, num=int(lambN/2))
124        lambSpace = np.concatenate((np.flip(-lambSpace), [0], lambSpace))
125
126        folds = sampleN = OLSruns = 5
127
128        print(" " * 35 + "*** NORMAL OLS: ***")

```

```

129     for name, result in zip(["OLS, multiple times", "bootstrap", "↵
130         kfold"],
131                             [ridgelearner.OLSeval(OLSruns),
132                             ridgelearner.bootstrap(sampleN=sampleN),
133                             ridgelearner.kfold(folds)]):
134         print("Method used: ", name)
135         for name, result in zip(["R2", "MSE"], result):
136             print("\n" + name +
137                   ": ( avg =", sum(result)/len(result),
138                   "), ( max = ", max(result),
139                   "), ( min = ", min(result), " )\n",
140                   result)
141             print("-----"*4+"\n")
142     print(" "35 + "*** RIDGE: ***")
143     for name, result in zip(["BRidge", "KRidge"],
144                             [ridgelearner.bootstrapRidge(lambSpace, ↵
145                             sampleN=sampleN),
146                             ridgelearner.kfoldRidge(lambSpace, folds)↵
147                             ]):
148         print("Method used: ", name)
149         for name, result in zip(["R2", "MSE"], result):
150             print("\n" + name +
151                   ": ( avg =", sum(result)/len(result),
152                   "), ( max = ", max(result),
153                   "), ( min = ", min(result), " )\n",
154                   result)
155             print("-----"*4+"\n")
156     for name, result in zip(["Bootstrap-Ridge", "Kfold-Ridge"],
157                             [ridgelearner.bootstrapRidge(lambSpace, ↵
158                             sampleN=sampleN),
159                             ridgelearner.kfoldRidge(lambSpace, folds)↵
160                             ]):
161         # print("Method used: ", name)
162         # for i in range(len(lambSpace)):
163             # print("lambda =", lambSpace[i], "R2: ", result[0][i], "↵
164             # MSE: ", result[1][i])
165         # print("-----"*4)
166         figure(1)
167         plot(lambSpace, result[0], label=name)
168         xlabel("Lambda")
169         title("R2 values as a function of lambda\nlambda of ↵
170             logarithmic space between +-10^" + str(logLow) +
171             " and +-10^" + str(logHigh) + " \nnumber of lambda ↵
172             points = " + str(len(lambSpace)))
173         legend()
174         figure(2)
175         plot(lambSpace, result[1], label=name)
176         title("MSE values as a function of lambda\nlambda of ↵
177             logarithmic space between " + "0" +
178             " and " + str(lambSpace[-1]) + " \nnumber of lambda ↵
179             points = " + str(len(lambSpace)))
180         xlabel("Lambda")
181         legend()
182         # show(block=False)
183         imageFileNameMSE = "MSERidge"
184         imageFileNameR2 = "R2Ridge"
185         filetype = ".png"
186         figure(1)
187         savefig(ridgelearner.imageFilePath + imageFileNameR2 + filetype)
188         figure(2)
189         savefig(ridgelearner.imageFilePath + imageFileNameMSE + filetype)
190         show()
191     # ridgelearner.biasVarianceAnalysis_ridge([0,10])

```

### 6.1.5 e) Lasso Regression on the Franke function with resampling

```
1 from partD import learningD
2 from sklearn.model_selection import train_test_split
3 import sklearn.linear_model as linmod
4 import numpy as np
5 from matplotlib.pyplot import plot, show, legend, title, ylabel, ←
   xlabel, figure, savefig
6
7 class learningE(learningD):
8     def __init__(self, n, p, noisefactor=None):
9         super().__init__(n, p, noisefactor)
10        self.imageFilePath = "../outputs/images/partE/"
11
12
13    def lasso(self, alpha=None):
14        if alpha is None: alpha = 5
15        clf = linmod.Lasso(alpha=alpha)
16        self.craftX()
17        X_train, X_test, y_train, y_test = train_test_split(self.X, ←
            self.yData)
18        clf.fit(X_train, y_train)
19        ypredict = clf.predict(X_test)
20        return self.R2(y_test, ypredict), self.MSE(y_test, ypredict)
21
22    def biasVarianceAnalysisLassoBootstrap(self, p_range, sampleSize←
        =None, sampleN=None):
23        p_old = self.p
24        p_rangeObject = range(p_range[0], p_range[1])
25        if sampleSize is None: sampleSize=0.8
26        if isinstance(sampleSize, float): sampleSize = int(sampleSize*←
            len(self.x))
27        if sampleN is None: sampleN=5
28
29        bias = np.zeros(len(p_rangeObject))
30        variance = bias.copy()
31        SElist = bias.copy()
32        SEkfold = bias.copy()
33        clf = linmod.Lasso(alpha=0.1)
34
35        for i in p_rangeObject:
36            self.p = i + p_old
37            self.craftX()
38
39            X_train, X_test, y_train, y_test = train_test_split(self.X←
                , self.yData)
40
41            ypredict_models = np.zeros((y_test.shape[0], sampleN))
42
43            for j in range(sampleN):
44                X_sample, y_sample = self.sample(X_train, y_train)
45
46                clf.fit(X_sample, y_sample)
47                ypredict = clf.predict(X_test)
48
49                ypredict_models[:,j] = ypredict
50
51            variance[i-p_range[0]] = self.variance(ypredict_models)
52            bias[i-p_range[0]] = self.bias(y_test, ypredict_models)
53            SElist[i-p_range[0]] = np.mean( np.mean((y_test.reshape(←
                (-1,1) - ypredict_models)**2, axis=1))
54            SEkfold[i-p_range[0]] = np.mean(self.kfold(sampleN)[1])
55
56        self.p = p_old
57        filetype = ".png"
58        imageFileName = "lassoBiasVariance"
59
60        print("MSE/(bias+variance) =", SElist/(bias+variance))
61
62        plot(p_rangeObject, np.log(bias), label="bias")
63        plot(p_rangeObject, np.log(variance), label="variance")
64        plot(p_rangeObject, np.log(SElist), label="MSE")
```

```

65     plot(p_rangeObject, np.log(SEkfold), label="MSE: K-fold")
66
67     legend()
68     title("Logarithmic bias-variance plot against complexity")
69     xlabel("Polynomial complexity")
70     ylabel("log[bias/variance]")
71     savefig(self.imageFilePath + imageFileName + filetype)
72     show()
73
74
75
76
77
78
79
80
81
82
83 if __name__ == "__main__":
84     print("""Task as interpreted:
85     (x) Implement lasso from scikit learn
86     ( ) Discuss the three methods (assume Lasso, Ridge, and Bootstrap ↵
87         (or OLS?)), and which fits the best
88     ( ) Bias-variance trade-off analysis using bootstrap and kfold
89     """)
90
91     Npoints = 100
92     polydegree = 5
93     polydegree_range = [0,10]
94     sampleN = 18
95     alphas = [0.01,0.1,0.3,1,3]
96
97     Lassolearner = learningE(Npoints, polydegree)
98     for alph in alphas:
99         print("Lasso evaluation [R2, MSE] at alpha = {:4}:".format(↵
100             alph), Lassolearner.lasso())
101     Lassolearner.biasVarianceAnalysis_lasso_bootstrap(polydegree_range=↵
102         , sampleN=sampleN)

```

### 6.1.6 f) Introducing real data and preparing the data analysis

```
1 import numpy as np
2 from imageio import imread
3 import matplotlib.pyplot as plt
4 from matplotlib.pyplot import figure, title, imshow, xlabel, ylabel, ←
   show, savefig
5 from mpl_toolkits.mplot3d import Axes3D
6 from matplotlib import cm
7 from partE import learningE
8
9 class learningF(learningE):
10     def __init__(self, n, p, noisefactor=None):
11         super().__init__(n, p, noisefactor)
12         self.imageFilePath = "geodata/"
13
14
15     def loadTerrainData(self, filename=None):
16         if filename is None: filename = "geodata/n45_e069_larc_v3.tif"
17         self.terrainData = imread(filename)
18
19 if __name__ == "__main__":
20     print("""Task as interpreted:
21     (x) download data and manage to import it
22     """)
23
24     Npoints = 2
25     polydegree = 0
26     # Datapoints moot
27
28     terrainlearner = learningF(Npoints, polydegree)
29     terrainlearner.loadTerrainData()
30
31     # Show the terrain
32     imageFileName = "geodata_area"
33     filetype = ".png"
34     figure()
35     title("Terrain data")
36     imshow(terrainlearner.terrainData, cmap="gray")
37     xlabel("X")
38     ylabel("Y")
39     savefig(terrainlearner.imageFilePath + imageFileName + filetype)
40     show()
```

### 6.1.7 g) OLS, Ridge and Lasso regression with resampling

```

1  from partF import learningF
2  from sklearn.model_selection import train_test_split
3  import numpy as np
4  import math as mt
5  from matplotlib.pyplot import plot, figure, title, imshow, xlabel, ↵
   ylabel, show, savefig, legend
6
7
8  class learningG(learningF):
9      def __init__(self, nFrac, p):
10         self.loadTerrainData()
11         nMax = self.terrainData.shape[0]
12         n = int((nMax)*nFrac)
13
14         self.p = p
15         self.nfeatures = int(((self.p+1)*(self.p+2))/2)
16
17         self.imageFilePath = "../outputs/images/partG/"
18
19         if nFrac < 0.5:
20             # Includes both random and regular data points
21             x1 = y1 = np.linspace(0, nMax-1, mt.ceil(n/2), dtype=int)
22             x2 = np.random.randint(0, nMax, mt.floor(n/2))
23             y2 = np.random.randint(0, nMax, mt.floor(n/2))
24
25             x = np.concatenate((x1, x2))
26             y = np.concatenate((y1, y2))
27
28         else:
29             x = y = np.linspace(0, nMax-1, n, dtype=int)
30
31         x, y = np.meshgrid(x,y)
32         self.x = x.flatten()
33         self.y = y.flatten()
34         self.yData = self.terrainData[self.x,self.y]
35
36     def plotkfold_trainvtest(self, k, p_range):
37         p_old = self.p
38         p_rangeObject = range(p_range[0], p_range[1])
39         SElist_test = np.zeros((2,p_range[1] - p_range↵
   [0]))
40
41         for i in p_rangeObject:
42             self.p = i + p_old
43             self.craftX()
44             dataset = self.kfold_yielder(k)
45             for i in range(k):
46                 data = next(dataset)
47                 ypredict = self.OLS_core(data)[0]
48
49
50             # X_train, X_test, y_train, y_test = train_test_split(self↵
   .X, self.yData)
51             # ypredict_test, beta = self.OLS_core([X_train, X_test, ↵
   y_train, y_test])
52             # ypredict_train = X_train @ beta
53
54             SElist_test[i-p_range[0]] += self.MSE(y_test, ↵
   ypredict_test)/k
55             SElist_train[i-p_range[0]] += self.MSE(y_train, ↵
   ypredict_train)/k
56
57         self.p = p_old
58
59         plot(p_rangeObject, np.log(SElist_test), label="Test")
60         plot(p_rangeObject, np.log(SElist_train), label="Train")
61         legend()
62         title("Logarithmic train v test plot for MSE")
63         xlabel("Polynomial complexity")
64         ylabel("Mean Square Error")

```

```

65         show()
66
67
68 if __name__ == "__main__":
69     print("""Tasks as interpreted:
70     (x) Parametrize terrain data
71     ( ) Apply all three models of [see below] to geographical data
72         ( ) OLS (k-fold)
73         ( ) Ridge (k-fold)
74         ( ) Lasso (k-fold)
75     ( ) Critically evaluate results and discuss "the applicabilty of ←
        these regression methods to the type of data presented here"
76     """)
77
78     polydegree = 5
79     datafrac = 0.1
80     OLSruns = 6
81
82     finallearner = learningG(datafrac, polydegree)
83     # finallearner.OLSeval(OLSruns)
84
85     for name, result in zip(["R2", "MSE"], finallearner.OLSeval(←
        OLSruns)):
86         print("\n" + name +
87             ": ( avg =", sum(result)/len(result),
88               ", ( max =", max(result),
89               ", ( min =", min(result), ")\n",
90               result)
91
92     finallearner.plotOLS_trainvtest(OLSruns, [0,20])

```