FYS-STK4155 - Project1

Gard, Are, David Andreas Bordvik

October 8, 2021

Motivation

In Project 1, we are tasked to study various regressions methods, such as Ordinary Least Squares, Ridge and Lasso. Our first area of study is how to fit polynomials to a specific two-dimensional function called Franke's Function. Our motivation behind fitting polynomials to Frank's function is to test the implementation of our regression algorithms, as well as studying various techniques such as bootstrapping and measurements such as the bias-variance tradeoff. Finally, we will move on to use real digital terrain data for our analysis.

The Franke Function is given on the form

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

with a 3-dimensional plot given in Figure 1

Plot of the Franke Function Variance:. 0.06034

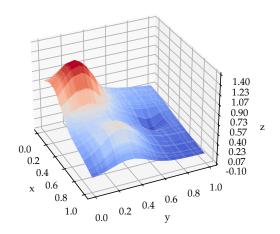


Figure 1: Plot of the Franke Function

Exercise 1

In Machine Learning, we are studying the problem of optimization, that is, finding the optimal parameter β such that $C(X,\beta) = \min_{\beta \in \mathbb{R}^p} \frac{1}{n} \left\{ (y - X\beta)^T (y - X\beta) \right\}$ our cost function is minimized. For the following exercise, where we will study Ordinary Least Squares regression, the previously stated cost function is the one we will minimize in order to fit the Franke Function, both without (as in Figure (1)) and with noise as in Figure (2).

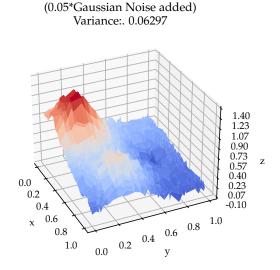


Figure 2: Plot of the Franke Function with added stochastic noise

When constructing our OLS model, we start off by minimizing the cost function with regards to β . That is, we take the derivative of $C(X, \beta)$ and set it 0. The following derivation shows how we end up with the optimal parameters $\hat{\beta}$, which in turn can be used to predict new values for the function which we are fitting.

$$\frac{\partial C(\boldsymbol{X}, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = 0$$

$$\frac{\partial C(\boldsymbol{X}, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = -\frac{2}{n} \boldsymbol{X}^{\mathrm{T}} (\boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta}) = 0$$

$$\boldsymbol{X}^{\mathrm{T}} (\boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta}) = 0$$

$$\boldsymbol{X}^{\mathrm{T}} \boldsymbol{X} \boldsymbol{\beta} = \boldsymbol{X}^{\mathrm{T}} \boldsymbol{y}$$

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

For consistency, it is noted that $\hat{\boldsymbol{\beta}} \in \mathbb{R}^p$, $\boldsymbol{y} \in \mathbb{R}^n$ and $\boldsymbol{X} \in \mathbb{R}^{n \times p}$.

With an expression for the predictors $\hat{\beta}$ derived, fitting a new value \tilde{y} is simply $\tilde{y} = X\hat{\beta}$.

Our own implementation of Ordinary Least Square regression is implemented such that its use mimics that of SciKit-learn. [6] That is, we also include separate steps for initializing our model, fitting the model and predicting using the model using the just derived mathematical expression. For further inquiry about implementation, refer to the github repository linked in the Appendix.

Confidence Intervals

Confidence intervals can be used to asses the uncertainty of a parameter. In our case, we will define confidence limits following the understanding of a Confidence Interval given in "Bootstrap Methods and their Application". That is, given a computed confidence region, any value inside the confidence region should be more likely than all values outside the confidence region. [2]

Furthermore, when computing the confidence interval for the parameters β , we first compute the variance $\sigma^2(\beta_j) = \sigma^2 \left[(XX^T)^{-1} \right]_{jj}$. Where XX^T is the Hessian matrix. Furthermore, it can be shown that the Hessian matrix can be given as the second derivative of the Cost function with respect to β . I.e.

$$\frac{\partial^2 C}{\partial \boldsymbol{\beta}^{\mathrm{T}} \partial \boldsymbol{\beta}} = \frac{2}{n} \boldsymbol{X} \boldsymbol{X}^{\mathrm{T}}$$

Mean Squared Error and R2 score

Two metrics that can be used to asses the quality of a model is its Mean Square Error (MSE) and R2 score. The MSE for any estimator is defined as

$$MSE(\boldsymbol{y}, \tilde{\boldsymbol{y}}) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2$$

It will be shown later that the MSE can be broken into two components, namely the variance and the squared bias. As can be seen from the equation, the MSE would attain the value of 0 if $y_i = \tilde{y}_i$. Moreover, by rewriting the mean squared error as $\text{MSE} = \frac{1}{m} \| \boldsymbol{y} - \tilde{\boldsymbol{y}} \|_2^2$, it can be seen that the error increases as the Euclidean distance between the prediction targets increase. [3]

The R2 score (coefficient of determination) is another metric that can be related to how the model covers its own variance. Defined as,

$$R^{2}(\boldsymbol{y}, \tilde{\boldsymbol{y}}) = 1 - \frac{\sum_{i=0}^{n-1} (y_{i} - \tilde{y}_{i})^{2}}{\sum_{i=0}^{n-1} (y_{i} - \bar{y})^{2}}$$

the closer the R2 score is to its maximum value 1, the more the variance of the model is explained by the model parameters. The R2 score gives a measure of model skill as a higher R2 generally indicates that the model is better at making new predictions. However, a perfect R2 score of 1 would result in the model covering it's entire variance, thus the model is overfitted and will not perform well in a general use case beyond the scope of it's initial training-data.

TODO Possibly more one implementation

Scaling the data

Our motivation for scaling the data arise in the context of fitting a design matrix of predictors with different units. To avoid having to disregard some predictors in favor of others based on their unit, not necessarily their contribution to the function fit, a scaling of the data is performed. By scaling the data using one of several scaling techniques, we ensure that all predictors lie in the same reference domain, resulting in a more accurate representation of the predictors. Furthermore, scaled data generally increases model skill.

SciKit-learn includes several different Scalers, such as the StandardScaler and MinMaxScaler. [6] For this discussion, the StandardScaler will be inspected.

The idea behind scaling the data with regards to the StandardScaler method, is to subtract the mean value and then divide by the variance. By performing these two operations, we ensure that the data have a mean value equal to 0 (i.e. are standardized/zero centered) and variance equal to 1.

Before scaling the dataset, an assessment of how to deal with the intercept has to be made. For reference, the intercept is defined as the first column of the design matrix, and for a polynomial fit would represent where the function intercepts with the y-axis when all other features are set to zero. Moreover, the intercept is a constant value (for our design matrices equal to 1), thus zero centering the intercept would result in a singular matrix, rendering the optimization problem unsolvable. Throughout this assignment, we are going to readd the intercept after scaling, essentially leaving the intercept untouched as to avoid any singular matrices. Moreover, for ordinary least squares regression, as there is no regularization of the predictors during the model fit, a model fitted on scaled contra unscaled data would attain the same mean squared error.

However, other regression methods such as Ridge and Lasso, which will be discussed in greater detail in further sections, have a dependance on the intercept through the hyperparamter λ when computing their regularization term. Not including the intercept when computing the regularization term would give rise to a divergence between the mean square error for the model scaled with compared to the model scaled without the intercept. Thus, by computing the regularization term with disregard to the intercept, the first β_j i.e. that of the intercept will be skipped in the computation. This will in most cases lead to a better mean square error for both Ridge and Lasso regression. Skipping the intercept when computing the regularization term also follows the definition of Ridge regression as in. [4]

Furthermore, while on the topic of Ridge and Lasso regression, a scaling of the data should always be performed. This is due to the regularized linear models such as Ridge and Lasso being sensitive to the scale of the input features. [4]

To determine whether scaling is a propriate for the current problem, that being fitting the Franke Function using orid nary least squares, an inspection of the generated data is made in light of the just discussion. For Exercise 1, the data points $x,y \in [0,1]$. This would indicate that the data is already scaled to a unit reference system. Moreover, as we are training a model based on the ordinary least squares, there is no dependance on scale of the input features as for regularized linear models. However, for consistency with further models, the data will be scaled with respect to the training data. However, the target variables however will not be scaled.

Splitting the data

As we want our model to perform well in general cases, we split the data into a training and testing set to simulate model prediction using new data. This is achieved by the aforementioned split, since the training and test data are kept entirely separated. In practice, we fit the model using the training data, then perform a test of the model using the test data. The error rate for new cases predicted by the model using the test data, can be used to understand how the model will perform on new untrained data. [4] Moreover, by assessing the deviation between training error and test error, it can be seen whether the model is overfitting or not. It would be a case of overfitting if the training error is low,

whereas the test error is high.

Throughout this assignment, we will split the data into a train and test set. The data could be split into an additional validation set, which is normally used for hyperparamter adjustment. However, as we don't see any practical use for a validation set for this assignment, we will skip out on splitting the data into an additional validation set. Though a validation set could be used for tuning the hyperparameter to select an optimal model, it could also lead to suboptimal model selection caused by imprecise prediction from a too small validation set. [4] Though our data consists of potentially unlimited data points, due to computational time constraints, we will resort to a relatively sparse dataset. Hence, we split into a training and test set to avoid sacrificing the training data in favor of a sufficient validation set.

Had we not omitted the validation set for hyperparameter tuning, the process of studying regularized models would have deviated somewhat from the study of the ordinary least squares regression. The process would then be that we split the data into a train-test split, as before. Moreover, the training data would have been split once more into a train and validation set, with an approximate ratio of 80 - 20 percent respectively. Furthermore, the hyperparameters are tuned with the MSE obtained from predicting using the validation set. However, as the validation set typically reports a lower error than the test set, the generalization error of the optimized model is studied using the test set. [3]

Comparing our OLS implementation to the one delivered by SciKit-learn

With our Ordinary Least Squares model implemented as described above, we start off by benchmarking our implementation to the LeastSquares method found in SciKit-learn. [6] We start off by setting up a uniform 2-dimensional grid and initialize a Franke Function with some added stochastic noise.

```
np.random.seed(4155)

n = 100 # The number of points in direction for the Franke Function

x = np.sort(np.random.uniform(0, 1, n))

y = np.sort(np.random.uniform(0, 1, n))

x, y = np.meshgrid(x,y)

z = FrankeFunction(x, y) + noise_factor(n, factor=0.05)
```

For this initial run, we are interested in studying the least squares fit of the Franke Function up to the fifth order.

By inspecting Figure (3), we can see that there are no visual differences between our implementation of OLS compared to the SciKit-learn implementation. Moreover, there is a reduction in MSE as model complexity increases. Thus it is clear that a higher order fit of the model results in a lower error when predicting new values using the test data.

Comparing Figure (3) and (4), where the first figure features half the amount of added noise compared to the second figure, there are two things of note. Firstly, Figure (3) attains an overall lower MSE, even as the order of the fitted function increases. Secondly, the initial difference between the train and test MSE is greater for Figure (4) than for Figure (3), though the difference converges to zero for higher order fits in both figures. As such, we make a preliminary conclusion that fitting functions with a higher amount of added Gaussian noise is more difficult than for smooth "predictable" functions.

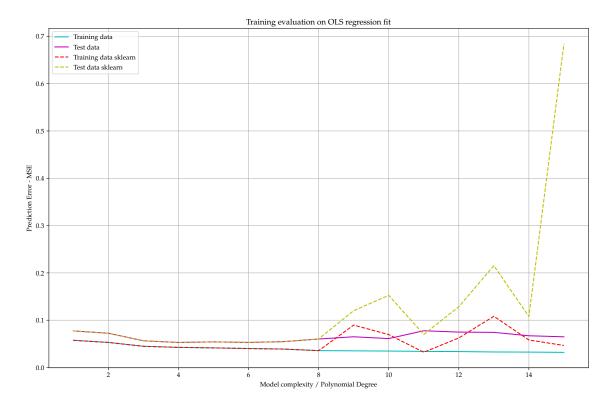


Figure 3: Benchmark run of the Franke Function fit with added Gaussian noise (factor of 0.05) with degree up to the fifth order of our OLS implementation compared to the similar LinearRegression() from SciKit-learn

The Confidence Interval for the predictors β_j is constructed for the fit using up to fifth order polynomials. As the predictors are based not only of the x and y parameters in isolation, but also their interaction terms, a fifth order fitted model includes 21 different predictors. The result of computing the Confidence Interval for the 21 different predictors can be seen in Figure 5. Note that the predictors From Figure 5, it can be seen that the lowermost and higher order terms have the smallest confidence intervals. Thus some of the predictors related to x and y of the third and fourth order pose the highest uncertainty.

Following the comparison of Figure (3) and (4), comparing the two confidence intervals reveals that as the added Gaussian noise is increased in size, the confidence intervals for the predictors increase as well. This can be seen especially for the predictors with some variability for the 0.05 noise factor case. Thus, as we increase the Gaussian noise added to the Franke Function, our predictors become less uncertain, which in turn could result in less precise model predictions.

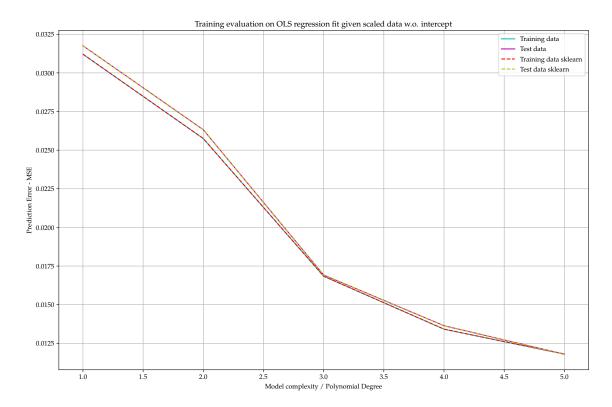


Figure 4: Benchmark run with doubled the amount of added Gaussian noise (factor 0f 0.1) of our OLS implementation to the similar LinearRegression() from SciKit-learn

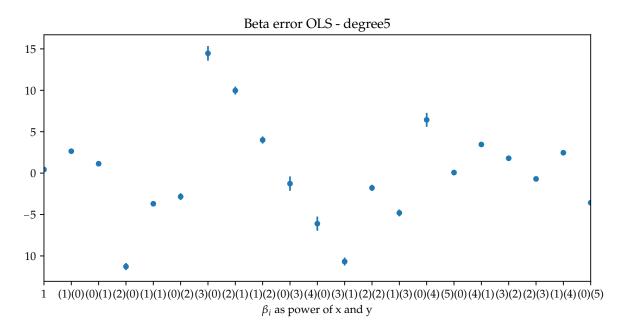


Figure 5: 95% Confidence intervals for the predictors of an OLS model with polynomials up to the fifth order and added Gaussian noise multiplied with a factor of 0.05.

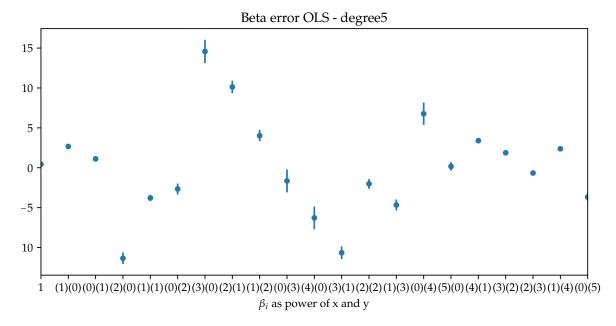


Figure 6: 95% Confidence intervals for the predictors of an OLS model with polynomials up to the fifth order and added Gaussian noise multiplied with a factor of 0.1.

Exercise 2

Some words on bootstrapping

As described in the preceding task, we split our data into several subsets used for different purposes. The overall objective is to train a model on a given subset and then measure how well it performs on a new unseen data set. *Figure 3* illustrates the results of this method, where the average prediction error is plotted against the model complexity. However, a dilemma that may emerge from this approach is that our data only represents a sample from an unknown distribution. We run the risk of measuring a predicted error that deviates from the actual probability by iterating through our data points only once.

Therefore, we impose resampling methods on our data to produce a more realistic estimate of the error. The following algorithm will demonstrate how *bootstrap* is performed to achieve this[?]:

We have some data, which consists of training data:

$$X = [x_1, x_2, ..., x_n]$$

and their respective target labels:

$$y = [y_1, y_2, ..., y_n]$$

(I) A new sample is then drawn. This dataset will have the same dimension as D, but as it is drawn with replacements, they will differ from eachother as the latter has a high probability of drawing the same sample mulitple times:

$$X^* = x_1^*, x_2^*, ..., x_n^*$$

$$\mathbf{y}^* = y_1^*, y_2^*, ..., y_n^*$$

(II) This sample is then used to create new set of predictors:

$$\hat{oldsymbol{eta}}^* = \left(oldsymbol{X}^{* ext{T}}oldsymbol{X}^*
ight)^{-1}oldsymbol{X}^{* ext{T}}oldsymbol{y}^*$$

(III)which in turn is used to predict new values on our test data. It is of great importance to emphasize that that X_{test} is completly separated from the rest of data, and only used to predict the new values for each sample.

$$\tilde{y}^* = X_{test} \hat{\beta}^*$$

•

(IV)We then compute the error in our predictions:

$$MSE(\boldsymbol{y}, \tilde{\boldsymbol{y}}^*) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i^*)^2$$

Step I-IV is then repeated B times, while keeping track of each sample MSE. After the completion of B bootstraps, an average of all MSE's is calculated. Ideally, B is less than the number of datapoints.

The bias-variance trade-off analysis

Moving on to the bias-variance trade-off analysis, we start off by showing that

$$C(\boldsymbol{X}, \boldsymbol{\beta}) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2 = \mathbb{E}\left[(\boldsymbol{y} - \tilde{\boldsymbol{y}})^2 \right]$$

Can be rewritten as

$$\mathbb{E}\left[(\boldsymbol{y}-\tilde{\boldsymbol{y}})^2\right] = \frac{1}{n}\sum_{i}(f_i - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right])^2 + \frac{1}{n}\sum_{i}(\tilde{y}_i - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right])^2 + \sigma^2.$$

where the terms are respectively $(bias)^2$, variance and noise. For simplicity, assume that we have a dataset where the data is generated from a noisy model

$$y = f(x) + \epsilon$$

Furthermore, we will assume that the residuals ϵ are independent and normally distributed $\epsilon \sim \mathcal{N}(0, \sigma^2)$. Finally, $\tilde{\boldsymbol{y}} = \boldsymbol{X}\boldsymbol{\beta}$ is our approximation to the functions f. Start off by adding and subtracting $\mathbb{E}\left[\tilde{\boldsymbol{y}}\right]$ inside the expectation value.

$$\mathbb{E}\left[(\boldsymbol{y} - \tilde{\boldsymbol{y}})^2 \right] = \mathbb{E}\left[(\boldsymbol{y} - \tilde{\boldsymbol{y}} + \mathbb{E}\left[\tilde{\boldsymbol{y}}\right] - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right])^2 \right]$$
$$= \mathbb{E}\left[((\boldsymbol{y} - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right]) - (\tilde{\boldsymbol{y}} - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right]))^2 \right]$$

By using the fact that $y = f(x) + \epsilon$ we can rewrite this as

$$= \mathbb{E}\left[\left(\left(f(\boldsymbol{x}) - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right]\right) - \left(\tilde{\boldsymbol{y}} - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right]\right) + \epsilon\right)^{2}\right]$$

Computing the square inside the expectation value gives us

$$= \mathbb{E}\left[(f(\boldsymbol{x}) - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right])^2 + (\tilde{\boldsymbol{y}} - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right])^2 + \epsilon^2 \right]$$

$$+ 2\left(\mathbb{E}\left[\epsilon(f(\boldsymbol{x}) - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right]) \right] - \mathbb{E}\left[\epsilon(\tilde{\boldsymbol{y}} - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right]) \right] - \mathbb{E}\left[(f(\boldsymbol{x}) - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right])(\tilde{\boldsymbol{y}} - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right]) \right]$$

Moreover, as ϵ are independent variables, the expectation value involving them as a product can be written as a product of expectation values. Knowing that $\mathbb{E}\left[\epsilon\right]=0$, the third and second to last term is equal to zero. Also, knowing that $\mathbb{E}\left[\tilde{y}\right]=\tilde{y}$, the last t

$$= \mathbb{E}\left[(f(\boldsymbol{x}) - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right])^{2} \right] + \mathbb{E}\left[(\tilde{\boldsymbol{y}} - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right])^{2} \right] + \mathbb{E}\left[\epsilon^{2} \right]$$
(1)

The first term in (1) can be discretized as

$$\mathbb{E}\left[(f(\boldsymbol{x}) - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right])^2\right] = \frac{1}{n} \sum_{i} (f_i - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right])^2$$

Which is the bias squared as we were to show.

The second term in (1) is also discretized, yielding

$$\mathbb{E}\left[(\tilde{\boldsymbol{y}} - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right])^2\right] = \frac{1}{n} \sum_{i} (\tilde{\boldsymbol{y}} - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right])^2$$

Which takes form of the variance, as was set out to show.

Finally, it can be shown that $var(y) = var(f + \epsilon) = \mathbb{E}\left[(f + \epsilon)^2\right] - (\mathbb{E}\left[(f + \epsilon)\right])^2 = \mathbb{E}\left[\epsilon^2\right]$ As such, we can use that

$$var(y) = \sigma^2 = \mathbb{E}\left[\epsilon^2\right]$$

to see that the final term in (1) is equal to the noise. Thus we have shown that

$$\mathbb{E}\left[(\boldsymbol{y} - \tilde{\boldsymbol{y}})^2 \right] = \mathbb{E}\left[(f(\boldsymbol{x}) - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right])^2 \right] + \mathbb{E}\left[(\tilde{\boldsymbol{y}} - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right])^2 \right] + \mathbb{E}\left[\epsilon^2 \right]$$
$$= \frac{1}{n} \sum_{i} (f_i - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right])^2 + \frac{1}{n} \sum_{i} (\tilde{y}_i - \mathbb{E}\left[\tilde{\boldsymbol{y}}\right])^2 + \sigma^2.$$

The final expression consists of three terms. The first term is a sum of the squared bias, the second term is the variance and the final term is a constant noise term. The sum of squared bias and variance make up the mean square error of our model. [5]

The bias component of the mean square error measures the difference from the true mean to the desired regression model. As function of model complexity, the bias will decrease as complexity increases. The second term, the variance, gives a measurement of the variation of the model values around their average value. The variance will increase with model complexity. The constant noise term σ^2 is an irreducible error which can only be reduced by cleaning up the data beforehand. [4]. As the irreducible error is out of our control, it does not contribute when analyzing the bias-variance tradeoff.

The bias-variance tradeoff can be used as a method for model selection. As has been alluded to in the previous paragraph, the variance is inverse proportional to the bias. As

such, there is a trade-off between bias and variance. As the bias-variance is directly related to the mean square error for both the training data but also test and production data used for making new predictions. When selecting a model we wish to balance and minimize both quantities, as that leads to the model with the best predictive capabilities. [1]

Furthermore, the problem of overfitting can also be discussed in light of the bias variance tradeoff. Overfitting is proportional to the model variance, as such a high variance leads to an overfitted model. An overfitted model is one that has learned the noise of the training data, resulting in a perfect fit for the training data but a high mean square error when predicting using new data. Hence, a higher bias can be considered more useful to circumvent overfitting.

Lastly, it should be noted that the bias-variance tradeoff measurement is performed for a single dataset of limited size. As such, if there were more datapoints to train the model with, the model would attain a better overall fit. Moreover, a greater amount of training data would reduce the level of overfitting for a given model complexity. [1] Though there are some limitations to the measurement, the bias-variance tradeoff can be used to estimate where the trained model is general enough to both avoid under -and overfitting (i.e. too high bias or too high variance).

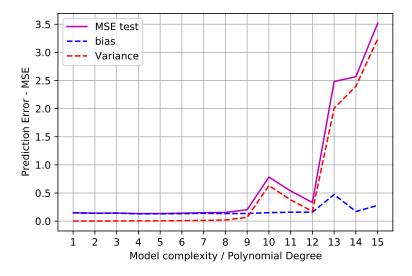


Figure 7: bias and variance trade-off as function of model complexity

Figure 8 visualizes where this occurs. As the complexity increases beyond degree 8, one can see an apparent increase in the variance. At this point, one can assume that the model has started overfitting its predictors to the training data and lost some of its ability to generalize on unseen data.

Exercise 3

In this project, an infinite number of data points could easily be created by simply adjusting the number of *n* decimals created in the span [0,1] for x and y. Setting a large *n* could then, in turn, yield a test set of substantial size, which ables the prediction of estimated average test error with statistical certainty. However, for real-life purposes, this is rarely the case. More often than not, the challenges related to the validation of models in question

lies in the limitations of available data. A widely used approach for combating these issues is implementing the cross-validation method.

This algorithm is implemented by dividing the dataset into K parts of roughly the same size. In this project, we try K values ranging from 5 to 10, which also coincides with best-practice values [5]. We have implemented the split into K-folds by using the method K-fold from Scikit learn model-selection [?]. This method provides us with arrays of approximately equal size, consisting of indices for choosing test- and train data. Each fold is then used once as a test set, while the remaining K-1 folds make up the training set. The overall purpose of this method is to produce an estimate of the prediction error, given by:

$$CV(\hat{f}) = \frac{1}{K} \sum_{i=1}^{K} L(y_i, \hat{f}^{-k(i)}(x_i))$$

where $\hat{f}^{-k}(x)$ denotes the fitted function, in this case the linear regression model, with the kth part of the data removed.[5]

The following figures provide a visualizes the results of 5-10 fold cross-validation on a dataset consisting of a total of 22^2 points[?]. The MSE for polynomial degrees ranging from 3 to 14 is then calculated. Here we present the results from degrees 4 and 12, while the remaining is available in

The MSE estimated by bootstrap for the corresponding polynomial degree is also plotted. As figure? displays, at degree 4, the size of K seems to be of little importance, as the mean MSE of every fold is moderately low at approximately 0.10, with a standard deviation ranging from 0.8 to .012.

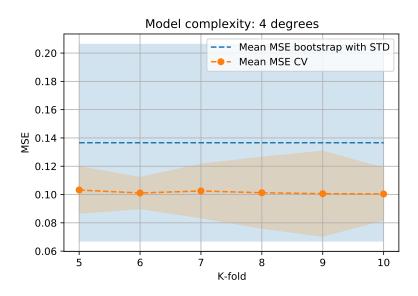


Figure 8: Mean MSE from CV and bootstrap at ploynomial degree = 4.

However, when evaluating the results provided, it is of great importance to also assess the total size of the dataset as this could bias the estimate of error:

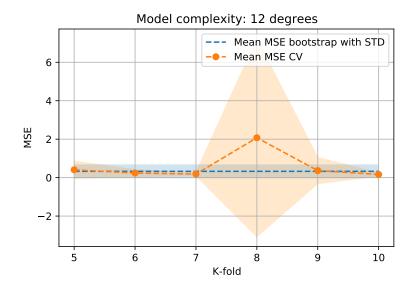


Figure 9: ean MSE from CV and bootstrap at ploynomial degree = 12.

Comparing methods

The scores from Scikit-learn needed to be multiplied with -1 before beeing compared to our own results. This is due to the feact that this method expects a "greater is better" function rather than a cost function, resulting in an output opposite of the MSE.[4] As illustrated by figure 8, we can conclude that our own implemented function produce the same results as the one fra Scikit-lean. As shuffeling is set to False as default in Scikit-learn, the same is done for our implemented method.

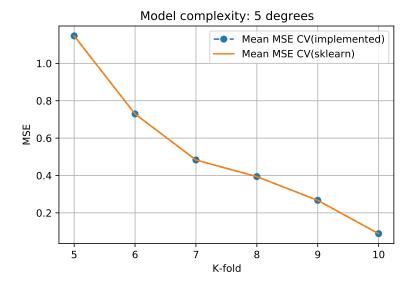


Figure 10: Comparison of results from sklearn vs implemented method for Cross-Validation

Exercise 4

Throughout Exercise 1 - 3, we have studied function fitting using the Ordinary Least Squares regression. When comparing models, we have plotted the test error as a function of the model complexity and chosen a model based on the bias-variance tradeoff. The bias-variance tradeoff has also been studied in light of resampling techniques such as the bootstrap and k-fold cross-validation. As we have seen when interpreting the result of the resampling, an increased degree of freedom for the model will make the model more prone to overfit.

Instead of reducing the degree of freedom of a model by limiting it's order, we will in this section look at a different approach for model selection through the introduction of a regularization parameter λ . With the regularization parameter, individual predictors in the model can be restrained without resorting to reducing the overall order of the model. The regularization term will punish the predictors such that they stray from reaching large values, adding another method to control overfitting. [1]

Ridge regression is a regression model where the weights are constrained. [4] This is done by adding a regularization term to the cost function, such that the function we are to optimize becomes

$$C\left(\boldsymbol{X},\boldsymbol{\beta}\right) = \left\{ \left(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\right)^{\mathrm{T}} \left(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\right) \right\} + \lambda \boldsymbol{\beta}^{\mathrm{T}}\boldsymbol{\beta}$$

Taking the derivatives of this function in terms of beta returns the closed form solution for the optimal $\boldsymbol{\beta}$

$$\hat{oldsymbol{eta}} = \left(oldsymbol{X}^{\mathrm{T}} oldsymbol{X} + oldsymbol{\lambda} oldsymbol{I}
ight)^{-1} oldsymbol{X}^{\mathrm{T}} oldsymbol{y}$$

As for Ordinary Least Squares regression, by using the closed form solution for the predictors we can predict new values \tilde{y} as $\tilde{y} = X\hat{\beta}$.

The No Free Lunch Theorem

As further motivation as to why we want to apply Ridge (and further on also Lasso) regression for fitting the Franke Function is due to the implications of the No Free Lunch Theorem. The theorem states that for any two learning models A and B, there are just as many situations where algorithm is A is superior to algorithm B as vice versa. [7] The theorem implies that there are no sets of model that are guaranteed to fit a problem better than others. As a result of the theorem, regardless of the data being fitted, different models have to be individually evaluated for the given task at hand.

As there are no universally best learning algorithm for fitting the problem at hand, we expand our analysis of Fitting the Franke Function in light the No Free Lunch Theorem to also include an assessment of Ridge and Lasso Regression. Both models handles feature selection through the introduction of the regularization parameter lambda, thus potentially achieving a lower MSE than for Ordinary Least Squares. Introducing several models also follows the implications of the Free Lunch Theorem, notably that we have to fine tune and assess our models to perform well on the current problem. [3]

Weight Decay

The Cost function for Ridge regression can be seen as the Cost function for Ordinary Least Squares with an added weight decay term $\lambda \beta^T \beta$. The parameter λ is known as a hyperparameter, i.e. the parameter can be controlled before fitting the model and changes how the model behaves during fitting and predicting. [3] Moreover, $\lambda \geq 0$, such that λ close to 0 implies Ordinary Least Squares and λ large implies that the predictors are close to 0. Hence weight decay, as the hyperparameter can be tuned to enforce some predictors (weights) to be excluded from the model fit.

For this Exercise, we will analyse the effect of the weight decay when fitting the Franke Function using Ridge regression. We will analyse the polynomial of optimal degree that we determined in the previous exercises through assessing the bias-variance tradeoff for an increased model complexity. Firstly with the bootstrap resampling technique, then compared the bootstrap MSE with what was obtained using the k-fold cross validation. With the model complexity

Implementing Ridge regression

Our implementation of Ridge regression follow closely that of Ordinary Least Squares, but with an input parameter λ determining the strength of the regularization term. Furthermore, the implementation differs as the normal equation for computing the optimal predictors have changed to $\hat{\boldsymbol{\beta}} = (\boldsymbol{X}^T \boldsymbol{X} + \lambda \boldsymbol{I})^{-1} \boldsymbol{X}^T \boldsymbol{y}$. When implementing Ridge regression, we follow the convention that the intercept is skipped when fitting the function, hence we start computing the regularization term starting at β_1 . [4] if the intercept were included, the regularizer would be dependent on the choice of origin for the target, which generally causes the mean squared error to increase as touched upon during the first Exercise.

Exercise 5

The second regularized linear model which we will study is Least Absolute Shrinkage and Selection Operator (LASSO) Regression. Lasso adds a regularization term to the Cost function similarly to Ridge regression, though Lasso utilizes the L1 norm of the predictors

instead of the L2 as in Ridge. [4] Mathematically, adding the L1 regularizer gives us the following Cost function

$$C\left(\boldsymbol{X},\boldsymbol{\beta}\right) = \left\{ \left(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\right)^{\mathrm{T}} \left(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\right) \right\} + \lambda \|\boldsymbol{\beta}\|_{1}$$

We note that the L1-norm is the same as the sum of the absolute value for each β_i . Moreover, the derivative of the absolute value is the sign-function. The resulting normal equation for Lasso regression is thus not a nice analytical expression, and is omitted for this Exercise. As a result, only the Sci-Kit implementation of Lasso regression will be studied.

Difference between L1 and L2-norm regularization

The difference between the L1 and L2-norm is how it can be related to our understanding of distance. The L1-norm, mathematically defined as $\|\boldsymbol{\beta}\|_1 = \sum_{i=0}^{n-1} |\beta_i|$ is the "Manhattan Distance" of $\boldsymbol{\beta}$, i.e. the distance from two points a to b using only straight and normal lines. On the other hand, the L2-norm is mathematically defined as $\|\boldsymbol{\beta}\|_2 = \sqrt{\sum_{i=0}^{n-1} \beta_i^2}$. The L2-norm is more commonly referred to as the Euclidean distance, or the distance of a straight line drawn between a and b.

Exercise 6

The main goal of exercise 6 is to analyze the performance of our models on real digital terrain data using the approach conducted in previous exercises 1-5. Before commencing model analysis in exercise 6, we had to understand the terrain data's nature and how it differs from the constructed franke function dataset. To recap on the franke function; one must decide on the number of points for x and y at a certain range, construct a meshgrid of x and y, then z is created based on the franke function. For the franke function we randomly and uniformly chose x, y such that $\forall x, y \in [0.0, 1.0]$. Comparing the franke function to the terrain data, we get the z variable from the terrain image itself. The values of the z variable in the terrain data represent height or depth similarly as for the franke function. However, the height or depth in the terrain data image represents the real height in the actual topographical terrain having an entirely different scale comparing the constructed range between 0-1 of the franke function. Looking at the image data, the x and y variables must be constructed as a meshgrid spanning the entire range of all x and y pixel positions within the terrain image. Since the range of x and y is determined by the span of pixel positions corresponding to the image width and height, we must scale both the x and y variables from the terrain data. The terrain input data utilized is of dimension 3601x1801, inducing 6485401 data points for the entire dataset. The terrain image is therefore not quadratic meaning the range of x and y is not equal. This motivates scaling of x and y to prevent one feature, by its magnitude, from dominating over the other. One must avoid the model to favor features by their magnitude since the contribution from features of smaller magnitude can be equally strong or important. Hence, the x-direction in real terrain is equally important as the y-direction even though the higher end of x is at 3600 and the higher end of y is at 1800. Another approach to level out the importance of the x and y variable (features) is to resize the image into a quadratic form where image height equals image width. However, resizing the entire image into a quadratic form will transform and distort the nature of the topographical terrain we try to fit. We solved this by first resizing while keeping the scale factor between width and height to make the number of data points feasible to work with. While keeping the terrain structures intact, we rescaled the image

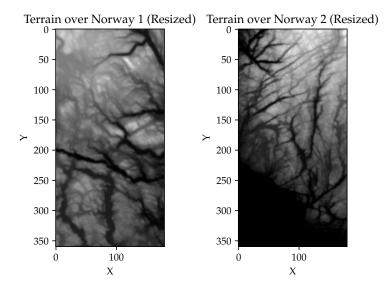


Figure 11: Terrain data resized

down to 20% of its original size, figure 11. Moreover, we created 8 patches of image data from the resized images where each patch being quadratic of dimension 180x180 (32400 data points), figure 12. This preprocessing serves the purpose of feeding the model with only one patch at a time, making the problem more digestible for obtaining a good fit concering subproblems instead. Another difference separating the terrain data from the synthetically created franke function data is that noise is incorporated within the image itself. For the franke function, we had to add noise in order for us to analyze the robustness of the models and to which degree the models were generalized. However, having control over the influence and magnitude of the noise is a luxury situation in real data situations, and for the real terrain image data, the noise stems from sensors and the assembly or construction of the image.

After investagating the different patches we narrowed down to two patches, one from each terrain. Two different patches were considered for our main analysis, figure 13. We concluded on mainy focusing on patch1 from terrain1 since it was easier for our model to fit. For the terrain data we scale both target and input data using standard scaler for reasons described above. The franke function is relatively smooth comparing the real terrain. The frank function has three dominating structures in its landscape, one small dip together with two Gaussian peaks. Looking at the terrain landscape, it is obviouse that its structures and shape are very complex. We also see that the terrain patch we chose has a higher normalized variance than what the frank function has using a noise factor of 0.05 Due to the complexity in the terrain data, we expected models of relatively high complexity before running into overfitting.

Exercise 6.1 (Exercise 1 approach)

In this task we conducted a OLS fit on the the terrain1 patch using similar appraoch as done in exercise1. We fitted the OLS model between degree 1-24. Looking at figure 14 we see that the model is overfitting to the training data at degree 19. We have a relatively large dataset of 32400 data points (180x180). We dicovered that the need for a slightly larger test sample fraction to more accuratly determain which complexity the model starts

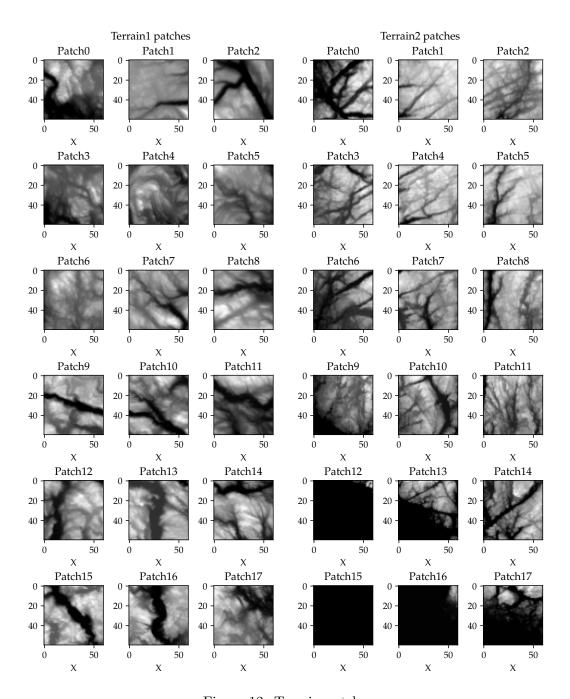
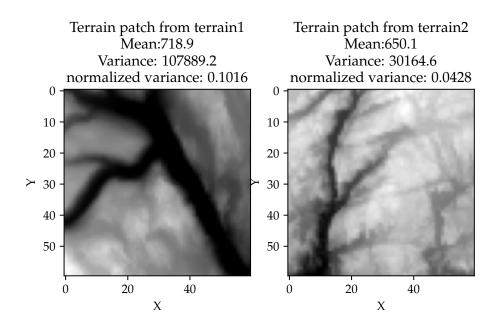


Figure 12: Terrain patches



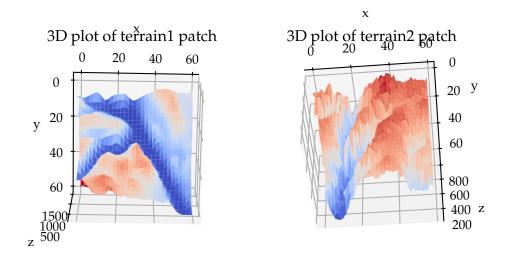


Figure 13: Chosen terrain patches

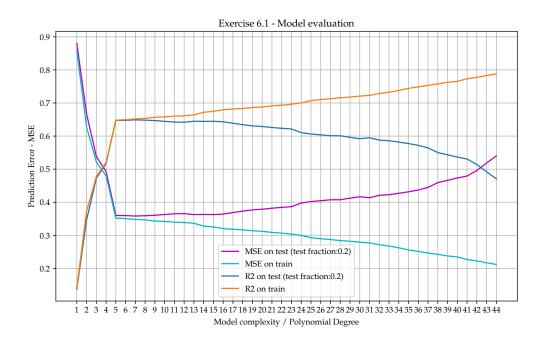


Figure 14: Model complexity vs prediction error using OLS

to overfit.

	degree	coeff name	coeff value	std error	CI lower	CI upper
0	5.0	$\beta 0$	0.0775	0.5562	-1.0126	1.1677
1	5.0	$\beta 1$	6.5442	0.5506	5.4651	7.6233
2	5.0	$\beta 2$	-4.6944	2.8491	-10.2786	0.8899
3	5.0	$\beta 3$	-14.5972	1.6039	-17.7408	-11.4535
4	5.0	$\beta 4$	-21.9678	2.8124	-27.4801	-16.4555
5	5.0	$\beta 5$	10.872	6.1673	-1.216	22.96
6	5.0	$\beta 6$	13.3363	3.1302	7.2011	19.4716
7	5.0	eta 7	46.7835	3.1077	40.6924	52.8746
8	5.0	$\beta 8$	20.4608	6.1127	8.4799	32.4416
9	5.0	$\beta 9$	1.7151	6.1094	-10.2593	13.6895
10	5.0	$\beta 10$	-28.2881	3.1003	-34.3646	-22.2116
11	5.0	$\beta 11$	-13.2902	2.6475	-18.4793	-8.1012
12	5.0	$\beta 12$	-41.1398	3.0833	-47.1829	-35.0966
13	5.0	$\beta 13$	-2.0759	6.0772	-13.9871	9.8353
14	5.0	$\beta 14$	-7.4504	2.272	-11.9035	-2.9973
15	5.0	$\beta 15$	24.3469	1.2791	21.8397	26.854
16	5.0	$\beta 16$	-15.9894	1.1068	-18.1588	-13.82
17	5.0	$\beta 17$	23.4262	1.1079	21.2547	25.5977
18	5.0	β 18	4.3273	1.2748	1.8287	6.8259
19	5.0	β 19	-2.0179	2.264	-6.4554	2.4195

Table 1: Coefficient summary of OLS fit

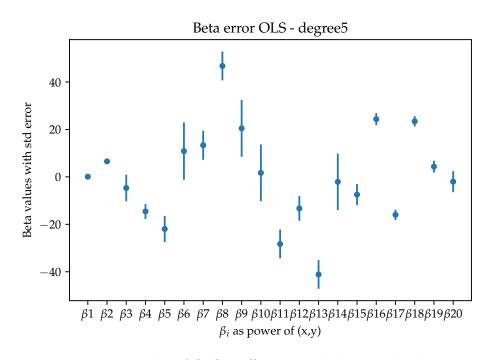


Figure 15: Plot of OLS coefficients including standard error

Appendix

 ${\bf TODO}$ her skal det ligge link til et github repo.

References

- [1] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer New York, August 2016.
- [2] A. C. Davison and D. V. Hinkley. *Bootstrap Methods and their Application*. Cambridge University Press, oct 1997.
- [3] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. http://www.deeplearningbook.org.
- [4] Aurélien Géron. Hands-on Machine Learning with Scikit-Learn, Keras, and TensorFlow. O'Reilly UK Ltd., October 2019.
- [5] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition. SPRINGER NATURE, February 2009.
- [6] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 12:2825–2830, 2011.
- [7] David H. Wolpert. The lack of a priori distinctions between learning algorithms. *Neural Computation*, 8(7):1341–1390, oct 1996.