

Project 1 FYS-STK3155

Lars Johan Brodtkorb October 9, 2018

Contents

1	Introduction	2
2	Methods	2
2.1	Measures of prediction accuracy	2
2.2	Ordinary least square regression	2
2.3	K-fold cross-validation algorithm	3
2.4	Bootstrap	4
2.5	Ridge and Lasso regression	4
3	Code implementation	4
3.1	Bootstrap	5
4	Analysis	5
4.1	Results	5
4.1.1	OLS	5
4.1.2	Ridge	9
4.1.3	Lasso	12
4.2	Discussion	15
4.2.1	OLS	15
4.2.2	Ridge	15
4.2.3	Lasso	16
5	Conclusion	16
6	Bibliography	16

Abstract

In this exercise I have tried to figure out which method is to be preferred of ordinary least squares, Ridge and Lasso regressions. The results show a better prediction error for the ridge regression, but that may be due to a faulty implementation of the Lasso method.

1 Introduction

The main aim of this project is to study in more detail various regression methods, including the Ordinary Least Squares (OLS) method, Ridge regression and finally Lasso regression. The methods are in turn combined with resampling techniques.

We will first study how to fit polynomials to a specific two-dimensional function called Franke's function. This is a function which has been widely used when testing various interpolation and fitting algorithms. Furthermore, after having established the model and the method, we will employ resampling techniques such as the cross-validation and/or the bootstrap methods, in order to perform a proper assessment of our models.

The Franke function, which is a weighted sum of four exponentials reads as follows

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

The function will be defined for $x, y \in [0, 1]$. Our first step will be to perform an OLS regression analysis of this function, trying out a polynomial fit with an x and y dependence of the form $[x, y, x^2, y^2, xy, \dots]$. We will also include cross-validation and bootstrap as resampling techniques. As in homeworks 1 and 2, we can use a uniform distribution to set up the arrays of values for x and y , or as in the example below just a fix values for x and y with a given step size. In this case we will have two predictors and need to fit a function (for example a polynomial) of x and y . Thereafter we will repeat much of the same procedure using the the Ridge and Lasso regression methods, introducing thus a dependence on the bias (penalty) λ .

Thereafter we are going to use (real) digital terrain data and try to reproduce these data using the same methods. We will also try to go beyond the second-order polynomials mentioned above and explore which polynomial fits the data best.

2 Methods

2.1 Measures of prediction accuracy

Find the confidence intervals of the β parameters by computing their variances, evaluate the Mean Squared error (MSE)

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

and the R^2 score function. If \tilde{y}_i is the predicted value of the i -th sample and y_i is the corresponding true value, then the score R^2 is defined as

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

where we have defined the mean value of \hat{y} as

$$\bar{y} = \frac{1}{n} \sum_{i=0}^{n-1} y_i.$$

2.2 Ordinary least square regression

We have an input vector $X^T = (X_1, X_2, X_3, \dots, X_p)$ and want to predict a real-valued output Y . The linear regression model has the form: ([2], chapter 3.2)

$$f(X) = \beta_0 + \sum_{j=1}^p X_j \beta_j \quad (1)$$

The linear model either assumes that the regression function $E(Y|X)$ is linear, or that the linear model is a reasonable approximation. The β_j s are unknown parameters or coefficients, and the variables X_j can come from many different sources. One of them is that of basis expansions such as

$X_2 = X_1^2, X_3 = X_1^3$, leading to a polynomial representation;

We can use the basis expansions as a device to achieve more flexible representations for $f(X)$. Polynomials are an example of this. Having great flexibility in approximation allows the basis expansions to work for many occurrences, but their variability may increase heavily if they are applied outside their intended scope. The coefficients to achieve a functional form in one region can cause the function to flap about madly in remote regions ([2], chapter 5.1, p. 140).

The polynomial method produces a dictionary D consisting of typically a very large number of basis functions $|D|$, far more than we can afford to fit to our data. Along with the dictionary we require a method for controlling the complexity of our model, using basis functions from the dictionary. There are three common approaches: Restriction, selection and regularization methods, which we will revisit in the Ridge and Lasso segment.

No matter the source of the X_j , the model is linear in the parameters. Typically we have a set of training data $(x_1, y_1) \dots (x_N, y_N)$ from which to estimate the parameters β . Each $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})^T$ is a vector of feature measurements for the i -th case.

Using linear least squared fittings for $X \in \mathbb{R}^2$ we can get the linear function of X that minimizes the sum of the sum of squared residuals of Y .

A unique solution for β has the form:

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

We have made minimal assumptions about the true distribution of the data. In order to pin down the sampling properties of $\hat{\beta}$, we now assume that the observations y_i are uncorrelated and have constant variance σ^2 , and that the x_i are fixed (non random). The matrix of the least squares parameter estimates is then given by:

$$Var(\hat{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2$$

2.3 K-fold cross-validation algorithm

Cross-validation is probably the simplest and most widely used method for estimating prediction error. This method directly estimates the expected extra-sample error $Err = E[L(Y, f(\hat{X}))]$, the average generalization error when the method $f(\hat{X})$ is applied to an independent test sample from the joint distribution of X and Y . As mentioned earlier, we might hope that cross-validation estimates the conditional error, with the training set T held fixed. Unfortunately, cross-validation typically estimates well only the expected prediction error ([2], chapter 7.10 - 7.12).

K-fold crossvalidation uses part of the available data to fit the model, and a different part to test it. We split the data into K roughly equal-sized parts. For the k th part, we fit the model to the other $K - 1$ parts of the data, and calculate the prediction error of the fitted model when predicting the k th part of the data. We do this for $k = 1, 2, \dots, K$ and combine the K estimates of prediction error.

Given a set of models $f(x, \alpha)$ indexed by a tuning parameter α , denote $\hat{f}^{-\kappa(x, \alpha)}$ by the α -th model fitted with the k -th part of the data removed. Then for this set of models we define:

$$CV(\hat{f}, \alpha) = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}^{-\kappa(i)}(x_i, \alpha)) \quad (2)$$

where L is the loss function, κ is the removed data.

The function $CV(\hat{f}, \alpha)$ provides an estimate of the test error curve, and we find the tuning parameter $\hat{\alpha}$ that minimizes it. Our final chosen model is $f(x, \hat{\alpha})$, which we then fit to all the data.

The steps of the cross-validation method:

- Find a subset of good predictors that show fairly strong (univariate) correlation with the class labels, using all of the samples except those in fold k .
- Using just this subset of predictors, build a multivariate classifier, using all of the samples except those in fold k .
- Use the classifier to predict the class labels for the samples in fold k .

2.4 Bootstrap

The bootstrap method has been used in this project to estimate the variance and bias. As with cross-validation, the bootstrap seeks to estimate the conditional error, but typically estimates well only the expected prediction error.

Suppose we have a model fit to a set of training data. We denote the training set by $Z = (z_1, z_2, \dots, z_N)$ where $z_i = (x_i, y_i)$. The basic idea is to randomly draw datasets with replacement from the training data, each sample the same size as the original training set. This is done B times ($B = 200$ in my application), producing B bootstrap datasets. Then we refit the model to each of the bootstrap datasets, and examine the behavior of the fits over the B replications ([2] p.249)

2.5 Ridge and Lasso regression

Ridge regression

Ridge regression is a simple example of a regularization approach, while the lasso is both a regularization and selection method.

Selection methods adaptively scan the dictionary and include only those basis functions h_m that contribute significantly to the fit of the model.

Regularization methods use the entire dictionary but restrict the coefficients. ([2], chapter 5.1 p.141)

The method used follows [3] quite closely. In chapter 1.4.2 on page 8 in [3] we find an estimator for β and variance in the ridge regression:

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

$$\text{Var}(\hat{\beta}) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \sigma^2$$

Ridge regression shrinks the regression coefficients by imposing a penalty on their size. The ridge coefficients minimize a penalized residual sum of squares.

$$\hat{\beta}_{\text{ridge}} = \arg \min_{\beta} \left(\sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right) \quad (3)$$

Here $\lambda \geq 0$ is a complexity parameter that controls the amount of shrinkage. The idea of penalizing by the sum-of-squares of the parameters is also used in neural networks, where it is known as weight decay ([2] p. 63)

Lasso regression

The Lasso method is a shrinking method like Ridge. Just as in ridge regression, we can re-parametrize the constant β_0 by standardizing the predictors; the solution for $\hat{\beta}_0$ is \hat{y} , and thereafter we fit a model without an intercept. In the signal processing literature, the Lasso is also known as basis pursuit (Chen et al., 1998). We can also write the Lasso problem in the equivalent Lagrangian form:

$$\hat{\beta}_{\text{lasso}} = \arg \min_{\beta} \left(\frac{1}{2} \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j| \right) \quad (4)$$

The difference from the ridge method is that the ridge penalty $\sum_{j=1}^p \beta_j^2$ is replaced with the Lasso penalty $\sum_{j=1}^p |\beta_j|$

3 Code implementation

For most of the code implementation see github repository: <https://github.com/larsjbro/FYS-STK3155/tree/master/project1>. I will include the bootstrap method with comments in the report.

3.1 Bootstrap

The bootstrap method was implemented with the following code, which was used to generate the prediction error, bias and variance shown in the figures.

```
def bootstrap_bias_variance(x, y, model, n_bootstraps=200, test_size=0.2):

    x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=test_size)
    y_pred = np.empty((y_test.shape[0], n_bootstraps))

    for i in range(n_bootstraps):
        x_, y_ = resample(x_train, y_train)
        y_pred[:, i] = model.fit(x_, y_).predict(x_test).ravel()

    mean_y = np.mean(y_test)
    errors = np.mean((y_test - y_pred) ** 2, axis=1, keepdims=True)
    bias_squares = (y_test - np.mean(y_pred, axis=1, keepdims=True)) ** 2
    variances = np.var(y_pred, axis=1, keepdims=True)

    r2 = 1.0 - np.mean(errors)/mse(y_test, mean_y)
    error = np.mean( errors )
    bias_squared = np.mean(bias_squares)
    variance = np.mean( variances )

    return error, bias_squared, variance, r2
```

4 Analysis

4.1 Results

4.1.1 OLS

Table 1 shows the fitted coefficients for the best ordinary least squares method fitted to the Frankfunktion along with its standard deviation, z-score and 95 percent confidence interval. As we can see from the table, only a few of the coefficients are significantly different from 0. In total table 1 shows there are only 5 such coefficients. The other coefficients could have possibly been set to zero.

	coef	std	z_score	Confidence interval 0.025	Confidence interval 0.975
0	-1.648196	0.532137	-3.097317	-2.691165	-0.605227
1	-1.919886	1.848958	-1.038361	-5.543776	1.704005
2	6.699335	3.720084	1.800856	-0.591897	13.990566
3	1.577407	10.841403	0.145498	-19.671353	22.826166
4	0.248456	0.463744	0.535762	-0.660466	1.157378
5	3.048496	2.491046	1.223781	-1.833865	7.930857
6	-20.088662	12.870111	-1.560877	-45.313616	5.136292
7	-17.914271	16.822026	-1.064929	-50.884836	15.056294
8	81.253441	62.107653	1.308268	-40.475323	202.982204
9	0.006444	1.659116	0.003884	-3.245363	3.258252
10	-10.048598	13.844107	-0.725839	-37.182548	17.085353
11	122.772857	58.256501	2.107453	8.592213	236.953501
12	18.202467	89.033985	0.204444	-156.300937	192.705871
13	-464.346238	295.293981	-1.572488	-1043.111806	114.419331
14	-5.577415	2.836466	-1.966325	-11.136786	-0.018044
15	-17.666890	15.364312	-1.149865	-47.780388	12.446608
16	188.349237	79.758452	2.361496	32.025544	344.672930
17	127.373969	100.448526	1.268052	-69.501525	324.249463
18	-781.831660	379.979896	-2.057561	-1526.578570	-37.084750
19	0.140168	8.835873	0.015864	-17.177825	17.458162
20	72.613196	63.561066	1.142416	-51.964203	197.190596
21	-603.170312	289.455482	-2.083810	-1170.492631	-35.847993
22	-176.329738	401.060283	-0.439659	-962.393449	609.733972
23	2420.894666	1422.585633	1.701757	-367.321940	5209.111272

Table 1: Values for the best fit ordinary least squares method for the data from the Frankefunction with m=300 and sigma=0.5

Below are figures from the ordinary least squares method:

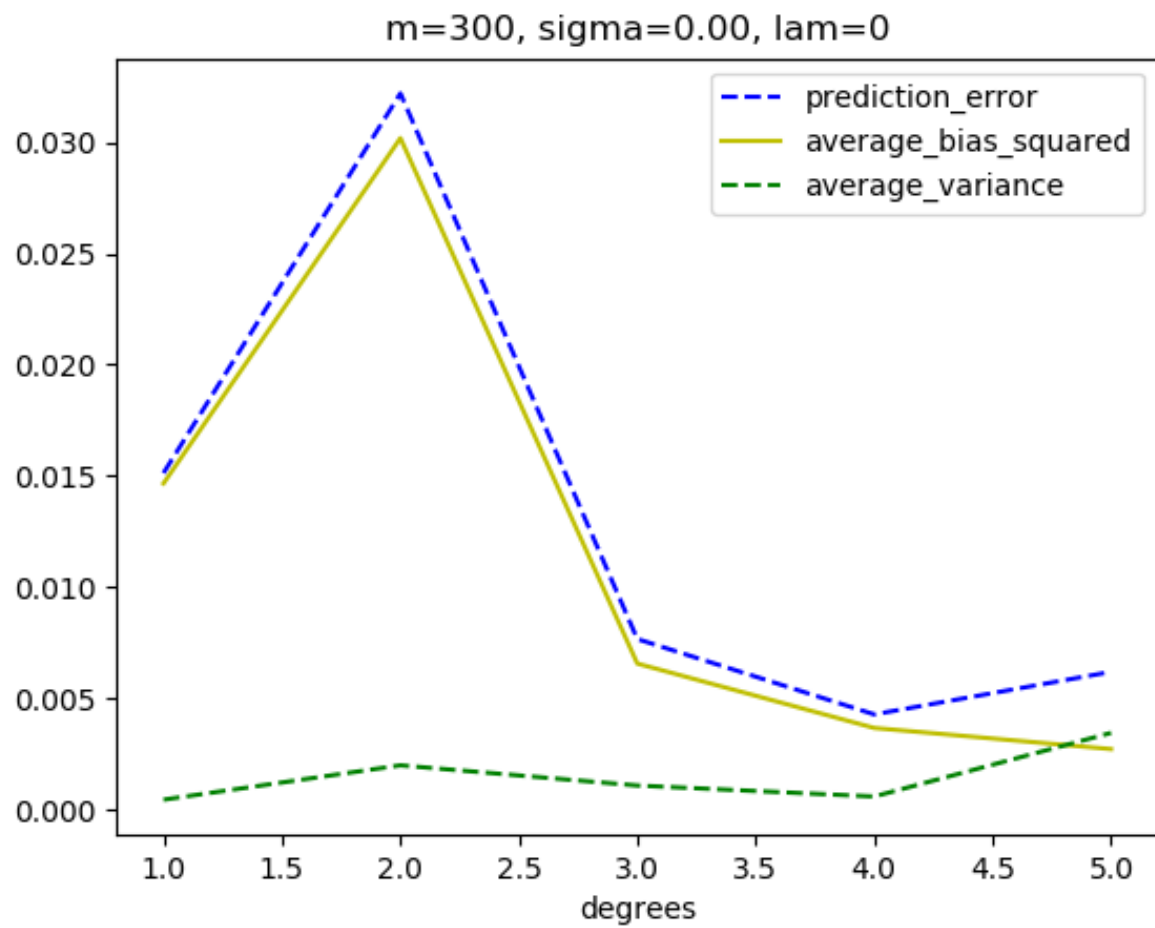


Figure 1: This shows the bias-variance tradeoff as function of the maximum polynomial degree used in the OLS method, with a sample size of 300.

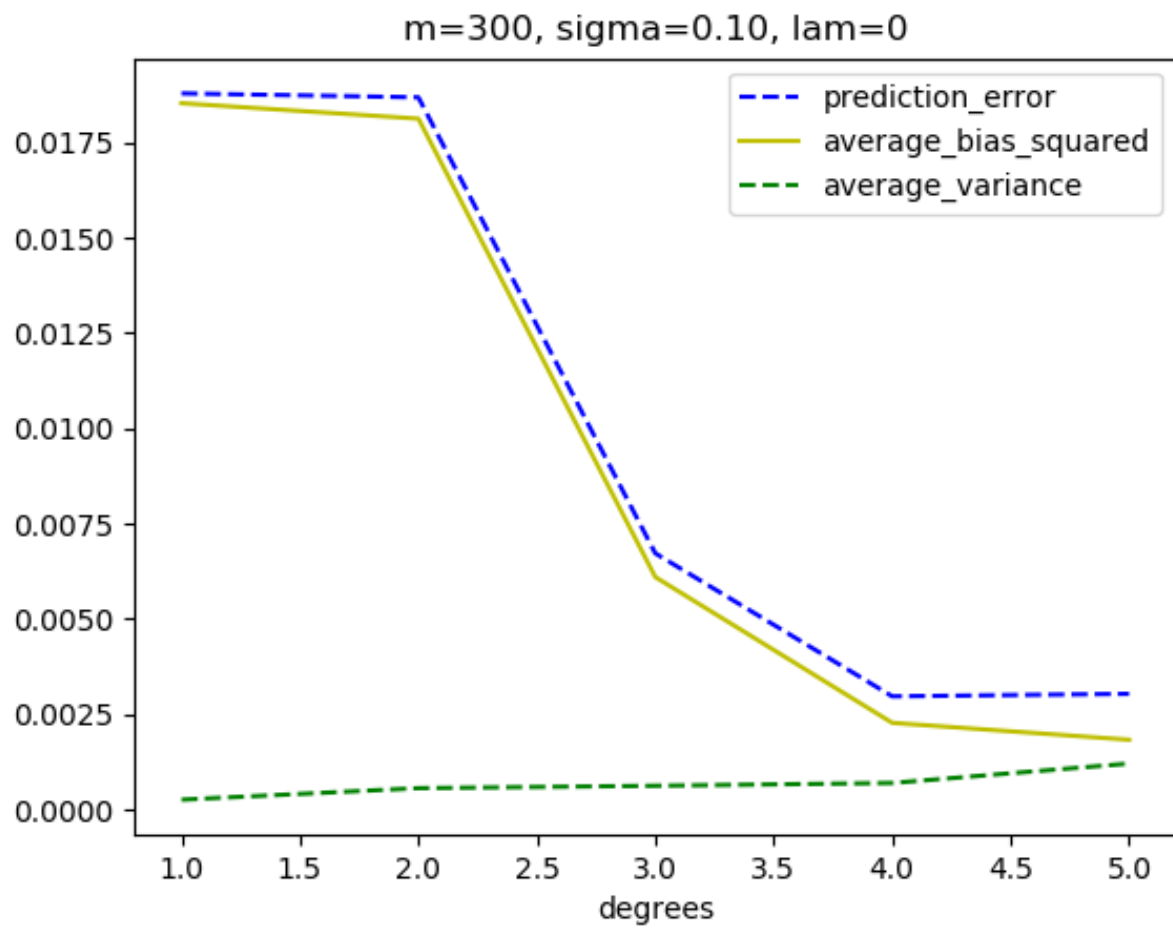


Figure 2: This shows the bias-variance tradeoff as function of the maximum polynomial degree used in the OLS method, with a sample size of 300 and a random gaussian noise with standard deviation of 0.1.

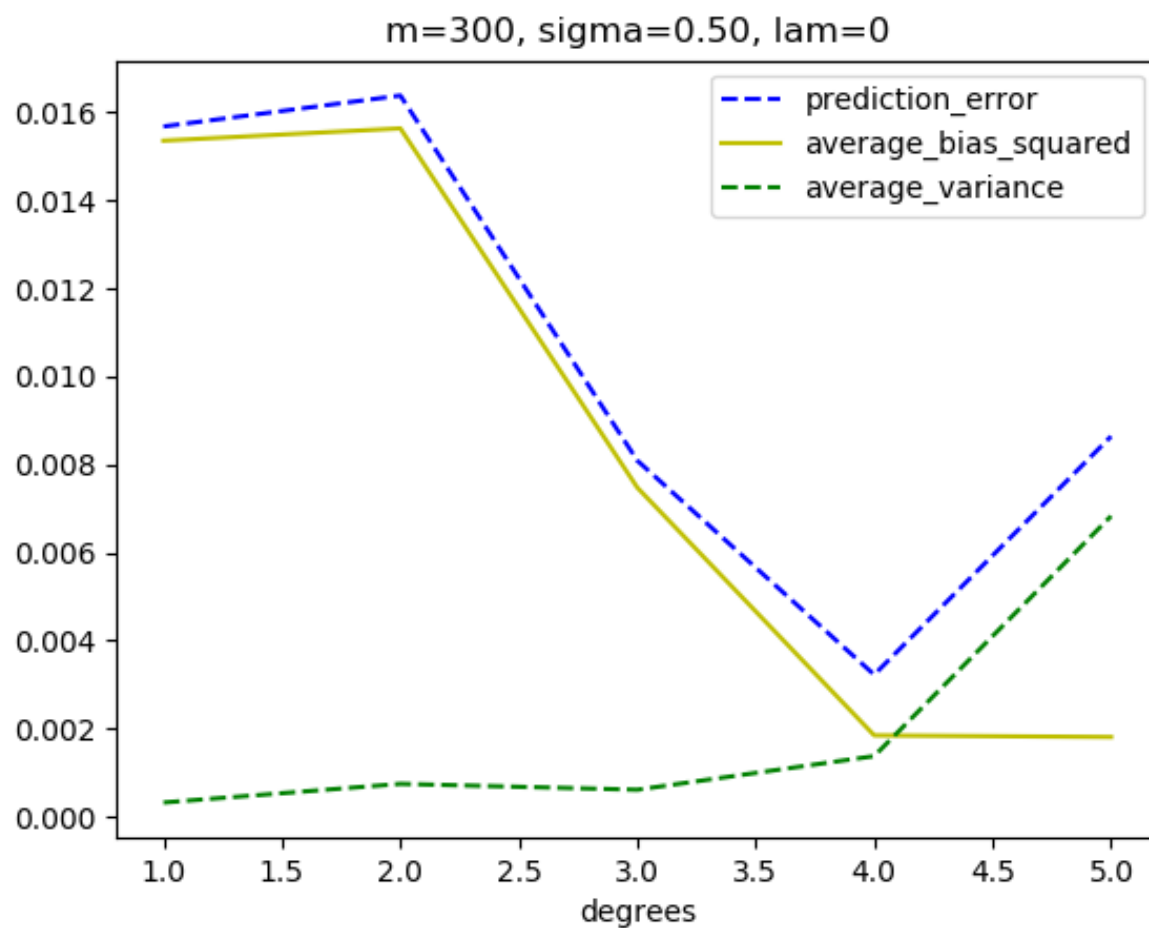


Figure 3: This shows the bias-variance tradeoff as function of the maximum polynomial degree used in the OLS method, with a sample size of 300 and a random gaussian noise with standard deviation of 0.5.

4.1.2 Ridge

Below are the figures for varying sample sizes, degree of polynomials and penalty with Ridge.

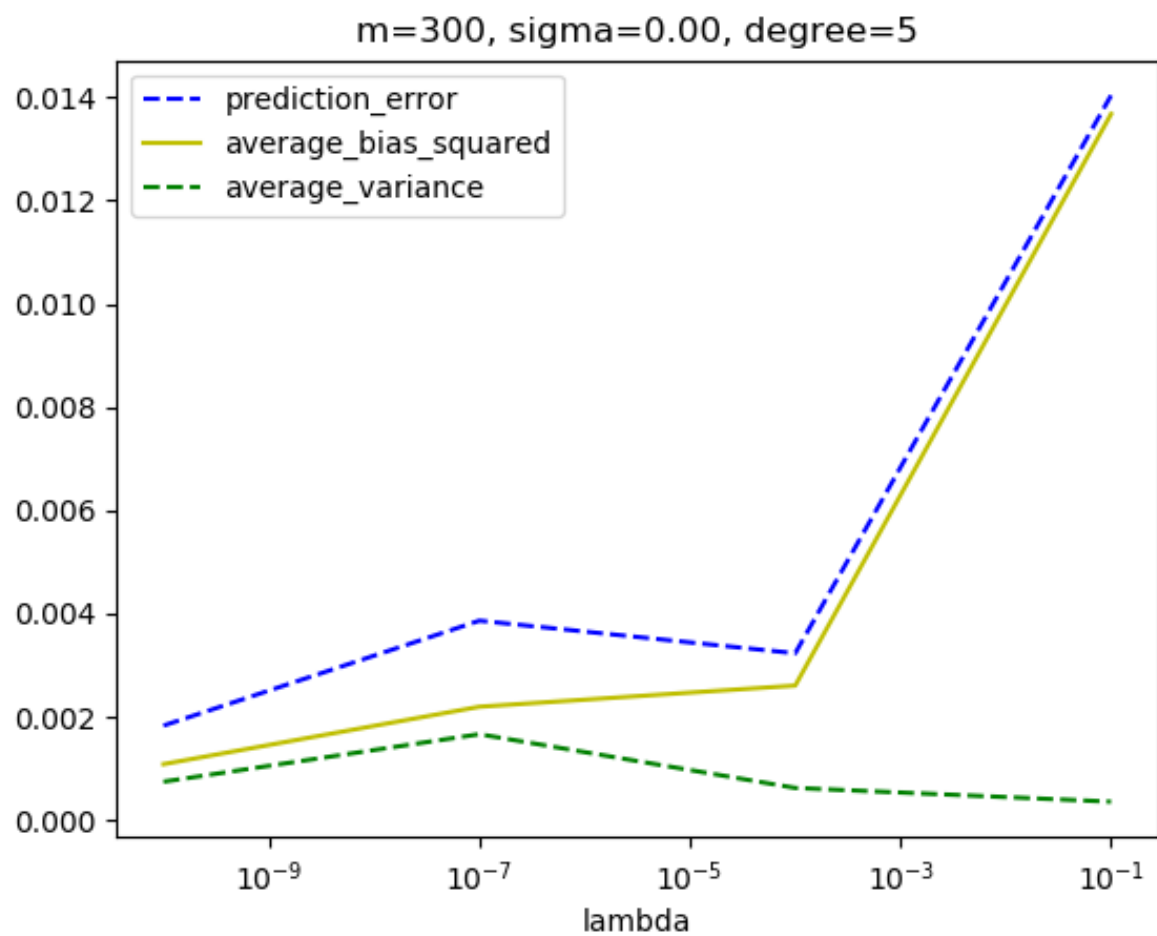


Figure 4: This shows the bias-variance tradeoff as function of the regularization parameter, λ , for Ridge regression with a sample size of 300 and a polynomial degree of 5

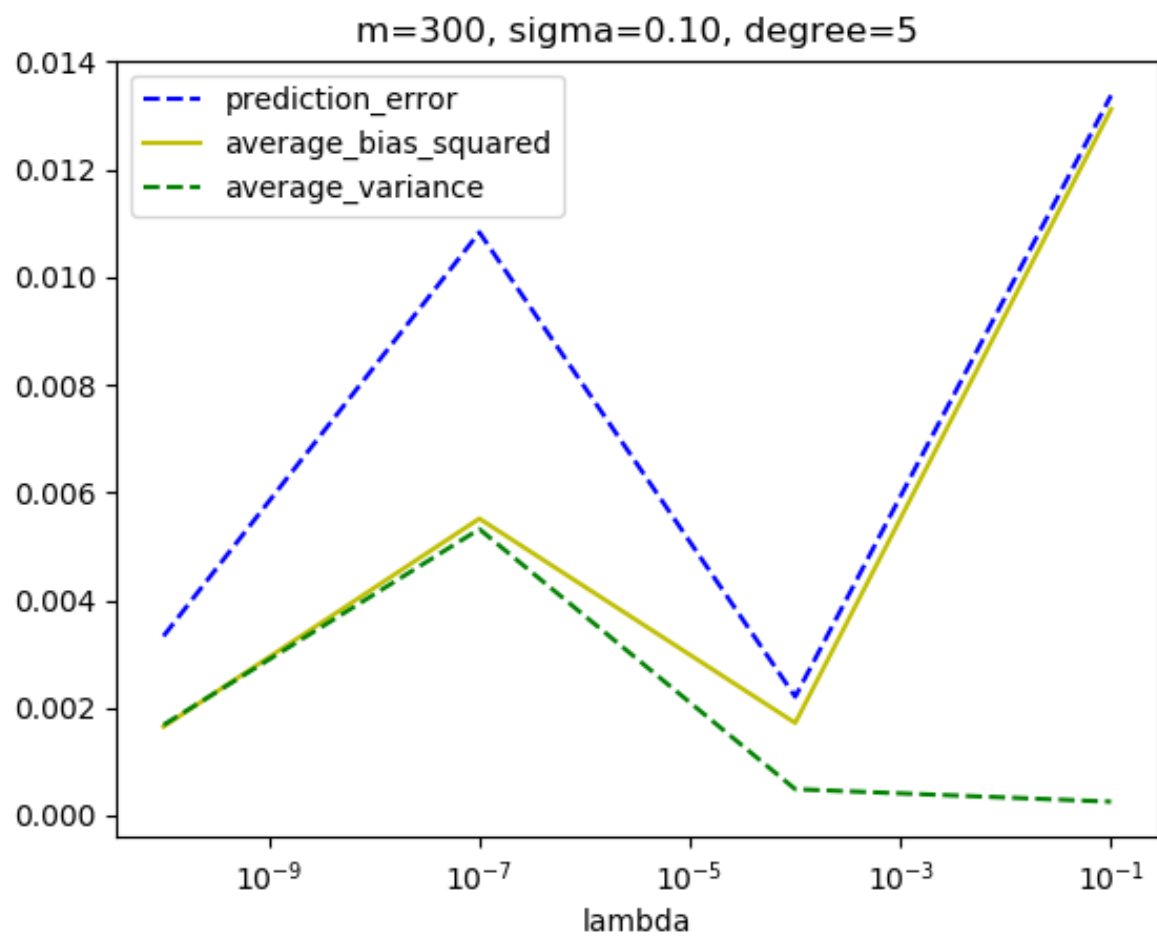


Figure 5: This shows the bias-variance tradeoff as function of the regularization parameter, λ , for Ridge regression with a sample size of 300, polynomial degree of 5 and standard deviation of 0.1

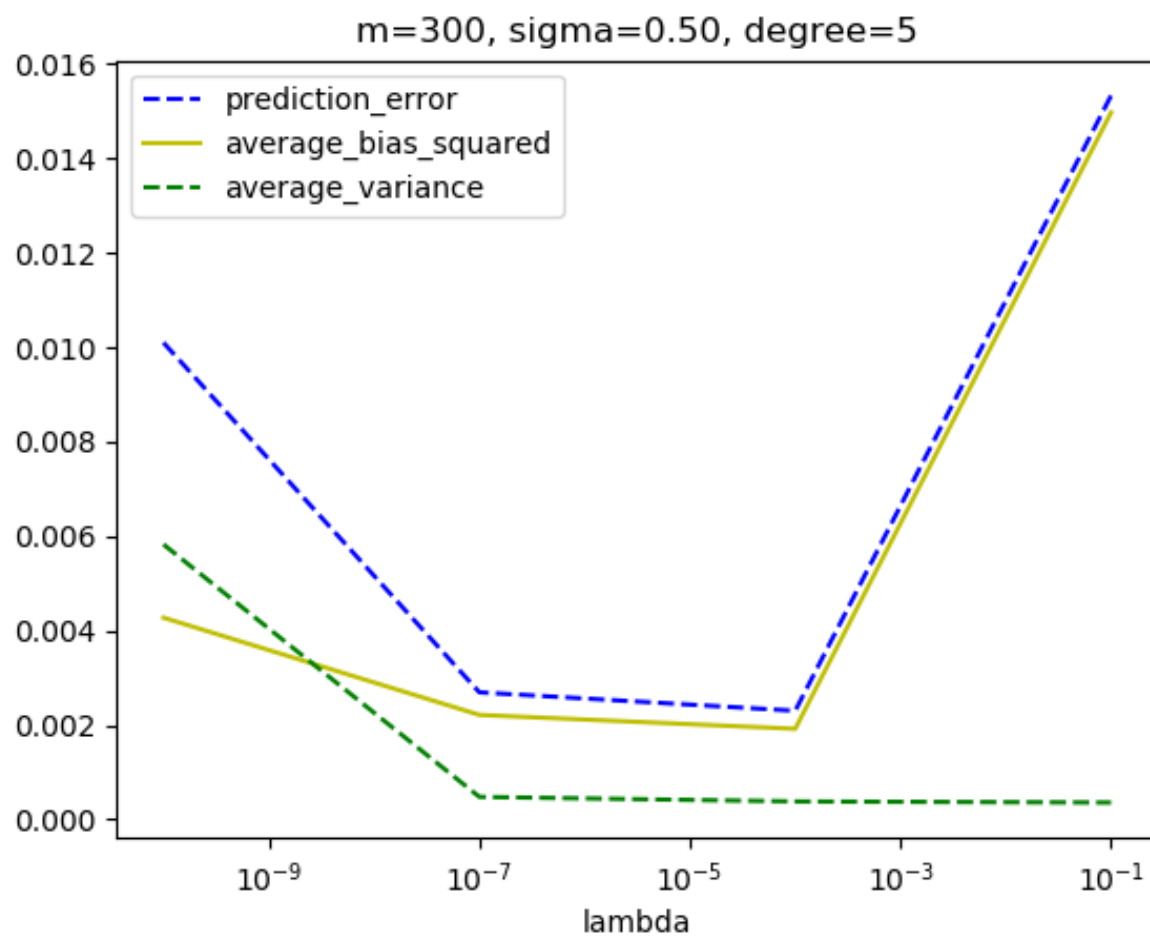


Figure 6: This shows the bias-variance tradeoff as function of the regularization parameter, λ , for Ridge regression with a sample size of 300, polynomial degree of 5 and standard deviation of 0.5

4.1.3 Lasso

Below are the figures for varying sample sizes, polynomials and penalty with Lasso.

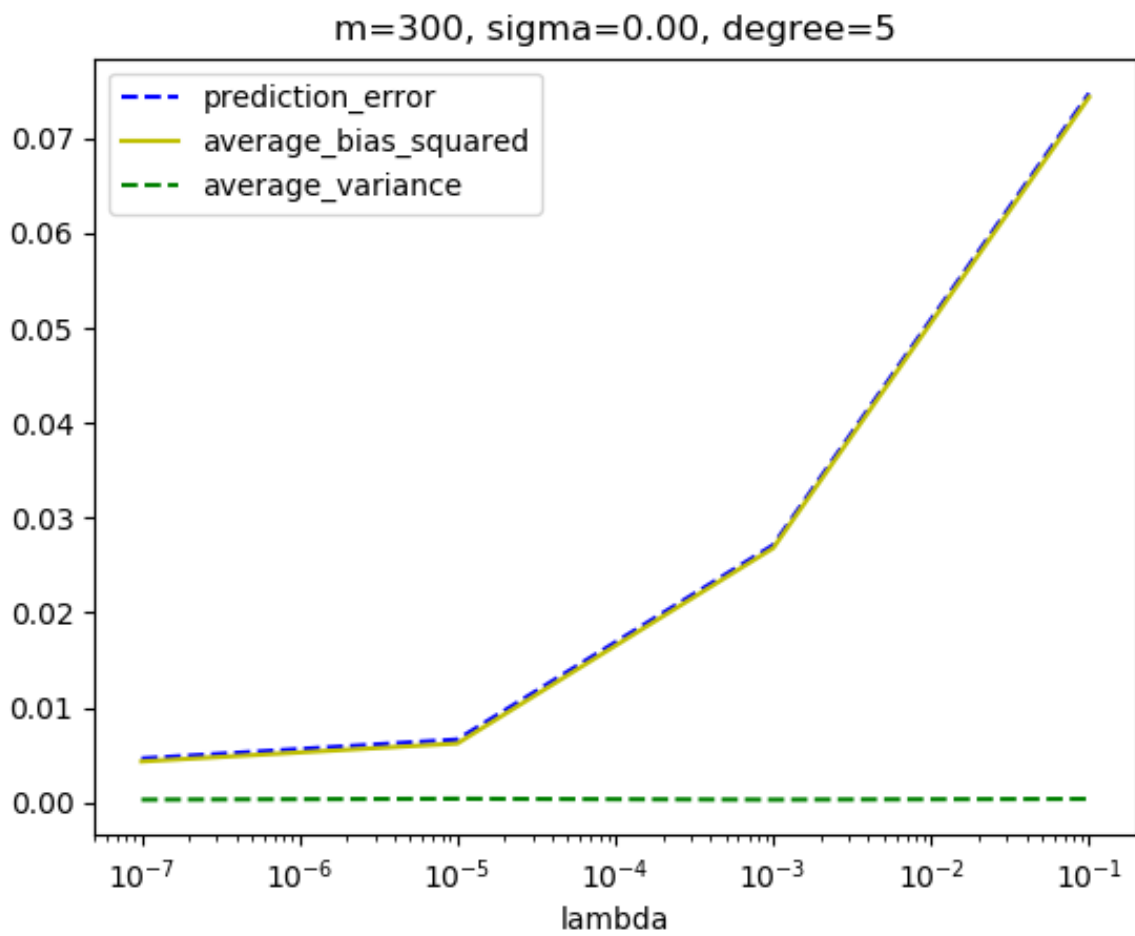


Figure 7: This shows the bias-variance tradeoff as function of the regularization parameter, λ , for Lasso regression method with a sample size of 300 and a polynomial degree of 5

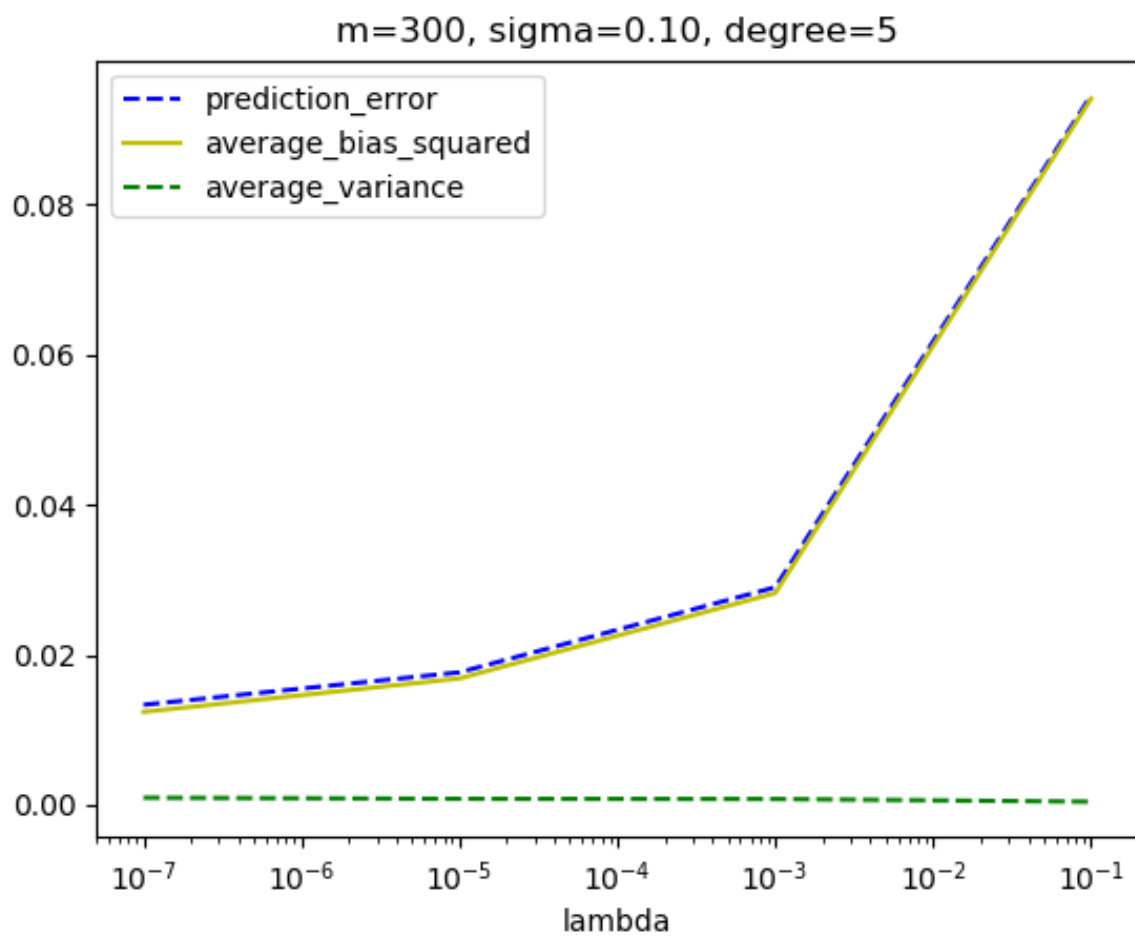


Figure 8: This shows the bias-variance tradeoff as function of the regularization parameter, λ , for Lasso regression method with a sample size of 300, polynomial degree of 5 and standard deviation of 0.1.

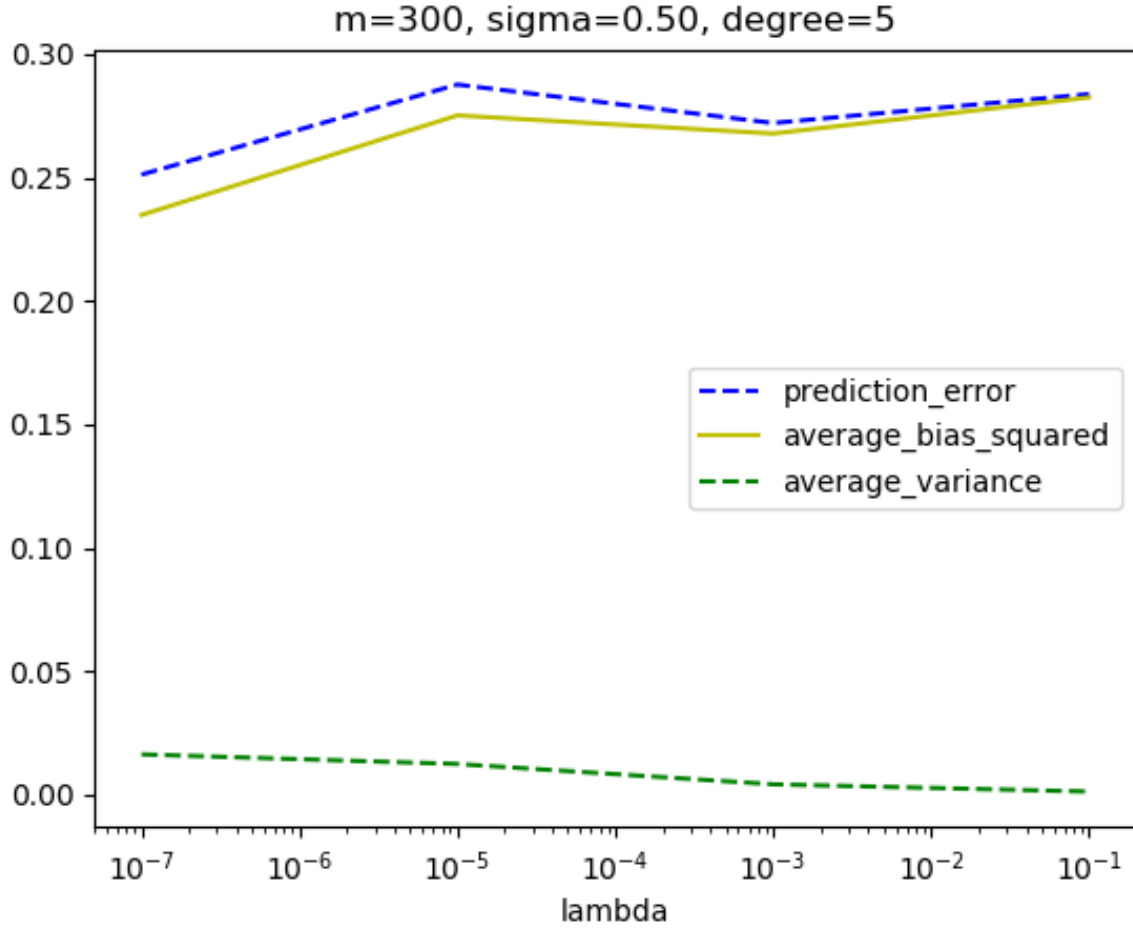


Figure 9: This shows the bias-variance tradeoff as function of the regularization parameter, lambda, for Lasso regression method with a sample size of 300, a polynomial degree of 5 and standard deviation of 0.5

4.2 Discussion

From all the figures I will now consider prediction error, bias and variance for the various methods.

4.2.1 OLS

For the OLS I can see that the minimum prediction error is for a polynomial of degree 4 for all standard deviations 0, 0.1 and 0.5. IN all cases the value of the prediction error is about 0.004. The bias is smallest at 5 degrees in all cases. The variance is smallest for a degree 4 polynomial for standard deviations of 0 and 0.1. For standard deviation of 0.5 the smallest variance is at 3. For a degree 1 polynomial, we do have a low variance, but the bias is really high.

Table 1 shows that there were only 5 significant coefficients, the other ones could possibly have been set to 0, but that was not implemented. In principle I could have found a best subset model here.

4.2.2 Ridge

For the ridge regression I can see a low prediction error for lambda being 10^{-4} and 10^{-10} for the 0 and 0.1 standard deviation cases, with a value of the prediction error of about 0.0035 and 0.002. The lowest value being for 0.1 standard deviation and lambda 10^{-4} . For the 0.5 standard deviation case, the prediction error is the smallest at about 0.25 for lambda 10^{-4} .

4.2.3 Lasso

For the lasso regression at standard deviation of 0 I can see a minimum prediction error of about 0.005 for lambda 10^{-7} . For standard deviation 0.1 the prediction error is up to 0.015 and for standard deviation 0.5 it is all the way up to 0.25. This tells me that something might have gone wrong with the implementation of lasso.

5 Conclusion

The results show that the Ridge regression gives a better prediction error than the other two methods. That does not agree with the expected, since the Lasso method should have been better. I tried to figure out the reason why that happened, and I came up with the following explanation:

The formulation of Lasso is the same as for Ridge, but for the extra β term in Ridge. It seems clear from the results that the Ridge method was to be preferred, but on further inspection, it seemed that the Lasso method may have been worse because the standardization method was not done correctly for Lasso. I suspect this was the reason for Lasso being worse than the Ridge regression.

I then implemented this standardization, and still found that the Lasso regression was inferior to the ridge regression, which I can not fully explain.

Improvements

I would have liked to complete the data analysis, and also figured out why Ridge regression showed superior results compared to the Lasso Regression. For the OLS method I could have found a best subset model after having found the significant coefficients.

6 Bibliography

References

- [1] <https://github.com/CompPhysics/MachineLearning/blob/master/doc/Projects/2018/Project1/pdf/Project1.pdf>
- [2] <https://github.com/CompPhysics/MachineLearning/blob/master/doc/Textbooks/elementsstat.pdf>
- [3] <https://arxiv.org/pdf/1509.09169.pdf>
- [4] <http://www.dtic.mil/dtic/tr/fulltext/u2/a081688.pdf>