



UNIVERSITY OF OSLO

FYS-STK3155

---

## Project 1

---

Eirik  
Felix  
Brage Wiseth

EIRIK@IFI.UIO.NO  
FELIX@IFI.UIO.NO  
BRAGEWI@IFI.UIO.NO

September 22, 2023

[HTTPS://GITHUB.COM/BRAGEWISETH/MACHINELEARNINGPROJECTS](https://github.com/bragewiseth/machinelearningprojects)

# Contents

1	Introduction . . . . .	3
2	Ordinary Least Squares . . . . .	4
	2.1 Discussion on Scaling . . . . .	4
	2.2 Analysis . . . . .	6
3	Ridge . . . . .	6
4	Lasso . . . . .	6
5	Paper & Pencil . . . . .	6
6	Bias & Variance . . . . .	7
	Bibliography . . . . .	10

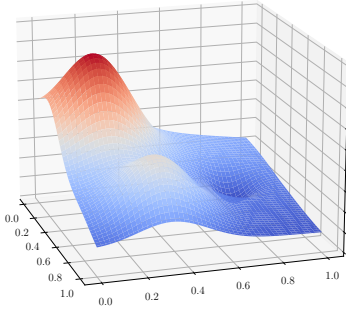
## Abstract

In this project, we explored the use of different methods for solving linear regression on topological data. We employed OLS, ridge and lasso methods for linear regression on simulated height data using *Franke's Function*. Where the goal was to fit polynomials to minimize the mean square error. As well as real topological data with the same goal.

**Keywords:** Linear Regression, Fitting Polynomials, Scaling, Bias-Variance

## 1. Introduction

Figure 1: Franke's Function



To get a smooth start on this project we begin to construct some fake topological data. This can be constructed with *Franke's Function*. This gives us a nice dataset to start linear regression. A natural start would be to use ordinary least squares to try to fit a polynomial to some degree. This is done by constructing a design matrix  $X$  with the polynomial terms and then solving the linear regression equation

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

where  $\beta$  is the coefficients of the polynomial and  $\mathbf{y}$  is the data. This is a nice start, but we can do better. We can add a penalty term to the cost function, this is called ridge regression.

$$\beta = (X^T X + \lambda I)^{-1} X^T \mathbf{y} \quad (2)$$

where  $\lambda$  is a hyperparameter that we can tune to get a better fit. This is a nice start, but we can do better. We can add a penalty term to the cost function, this is called ridge regression.

$$\beta = (X^T X + \lambda I)^{-1} X^T \mathbf{y} \quad (3)$$

where  $\lambda$  is a hyperparameter that we can tune to get a better fit. This is a nice start, but we can do better. We can add a penalty term to the cost function, this is called ridge regression.

$$\beta = (X^T X + \lambda I)^{-1} X^T \mathbf{y} \quad (4)$$

where  $\lambda$  is a hyperparameter that we can tune to get a better fit. This is a nice start, but we can do better. We can add a penalty term to the cost function, this is called ridge regression.

$$\beta = (X^T X + \lambda I)^{-1} X^T \mathbf{y} \quad (5)$$

where  $\lambda$  is a hyperparameter that we can tune to get a better fit. This is a nice start, but we can do better. We can add a penalty term to the cost function, this is called ridge regression.

$$\beta = (X^T X + \lambda I)^{-1} X^T \mathbf{y} \quad (6)$$

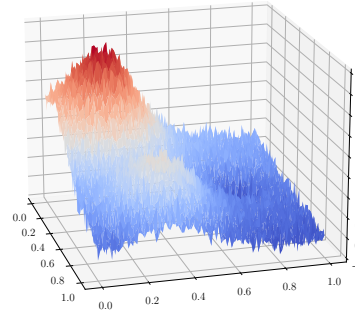


Figure 2: Franke's Function with noise

## 2. Ordinary Least Squares

### 2.1 Discussion on Scaling

Table 1: Unscaled sample design matrix fitting one-dimensional polynomial of degree 5

1.	0.	0.	0.	0.	0.
1.	0.25	0.0625	0.01562	0.00391	0.00098
1.	0.5	0.25	0.125	0.0625	0.03125
1.	0.75	0.5625	0.42188	0.31641	0.2373
1.	1.	1.	1.	1.	1.

Table 2: Scaled sample design matrix fitting one-dimensional polynomial of degree 5

0.	-1.41421	-1.0171	-0.83189	-0.728	-0.66226
0.	-0.70711	-0.84758	-0.7903	-0.71772	-0.65971
0.	0.	-0.33903	-0.49913	-0.56348	-0.58075
0.	0.70711	0.50855	0.29116	0.10488	-0.0433
0.	1.41421	1.69516	1.83016	1.90431	1.94603

MSE for OLS on unscaled data:	0.010349396022903145
MSE for OLS on scaled data:	0.010349396024145656
MSE for Ridge on unscaled data:	0.02106077418650843
MSE for Ridge on scaled data:	0.01782525371566323

the code for generating this output <sup>1</sup>

First  $x_{\text{unscaled}}$  and  $x_{\text{scaled}}$  is not that different, the original data was close to zero-centered and not that spread out, which means that initially by just looking at the data scaling is not that necessary. When we add the polynomial terms we can now see that some of the entries of  $X_{\text{unscaled}}$  get really small as an example  $0.1^5 = 0.00001$  this makes the columns of  $X_{\text{unscaled}}$  live in their own order of magnitude and scaling should be considered to bring them back to the same ish order of magnitude <sup>2</sup>. The act of not scaling results in  $\beta$  spanning from  $-50$  to  $48$  while scaling gives a smaller span from  $-10$  to  $13$ . Now this alone may not justify why we should scale this dataset, as scaled and unscaled OLS yields the same MSE. However when doing ridge regression the cost function is directly dependent on the magnitude of  $\beta_i$ . Now with each  $\beta_i$  varying a lot, some are getting more penalised than others. As we can see the MSE for unscaled data is much higher than for scaled data in the ridge case. This leads us to conclude that we should scale the data, making it easier to tweak  $\lambda$  and giving us nicer numbers to work with.

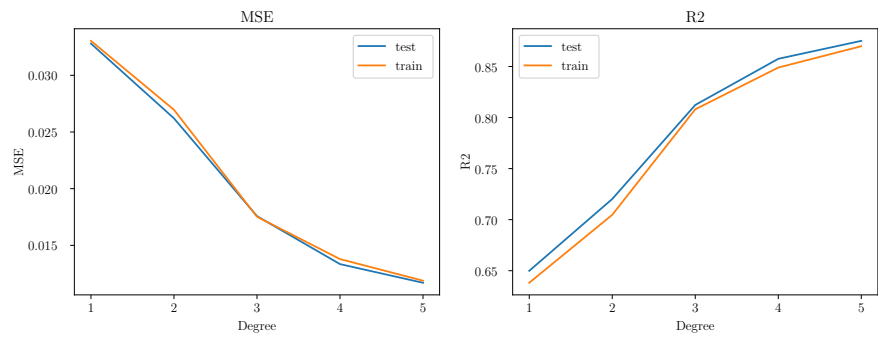


Figure 3: Franke's Function

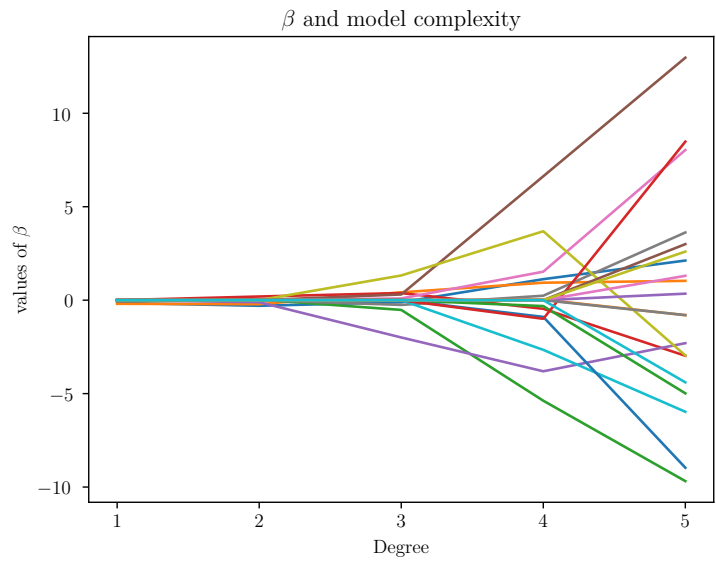


Figure 4: Franke's Function

## 2.2 Analysis

## 3. Ridge

## 4. Lasso

## 5. Paper & Pencil

The assumption we have made is that there exists a continuous function  $f(\mathbf{x})$  and a normal distributed error  $\epsilon \sim N(0, \sigma^2)$  which describes our data  $y = f(\mathbf{x}) + \epsilon$ . We then approximate this function  $f(\mathbf{x})$  with our model  $\hat{y}$  from the solution of the linear regression equations (ordinary least squares OLS), that is our function  $f$  is approximated by  $\hat{y}$  where we minimized  $(y - \hat{y})^2$ , with  $\hat{y} = \mathbf{X}\beta$ . The matrix  $\mathbf{X}$  is the so-called design or feature matrix. Show that the expectation value of  $y$  for a given element  $i$   $E(y_i) = \sum_j x_{ij} \beta_j = \mathbf{X}_{i,*} \beta$ , and that its variance is  $\text{Var}(y_i) = \sigma^2$ . Hence,  $y_i \sim N(\mathbf{X}_{i,*} \beta, \sigma^2)$ , that is  $y$  follows a normal distribution with mean value  $\mathbf{X}\beta$  and variance  $\sigma^2$ . With the OLS expressions for the optimal parameters  $\hat{\beta}$  show that  $E(\hat{\beta}) = \beta$ . Show finally that the variance of  $\hat{\beta}$  is  $\text{Var}(\hat{\beta}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$ . We can use the last expression when we define a so-called confidence interval for the parameters  $\beta$ . A given parameter  $\beta_j$  is given by the diagonal matrix element of the above matrix.

We can assume that  $y$  follows some function  $f$  with some noise  $\epsilon$

$$\mathbf{y} = f(\mathbf{x}) + \epsilon$$

$$E[\mathbf{y}] = E[f(\mathbf{x}) + \epsilon] = E[f(\mathbf{x})] + E[\epsilon]$$

The expected value of  $\epsilon_i$  is 0,  $f(x)$  is a non-stochastic variable and is approximated by  $\mathbf{X}\beta$

$$E[y_i] = \mathbf{X}_{i,*} \beta$$

The variance is defined as

$$\text{Var}(y_i) = E[(y_i - E[y_i])^2] = E[y_i^2 - 2y_i E[y_i] + E[y_i]^2] = E[y_i^2] - E[y_i]^2$$

$$\text{Var}(y_i) = E[(\mathbf{X}_{i,*} \beta)^2 + 2\epsilon \mathbf{X}_{i,*} \beta + \epsilon^2] - (\mathbf{X}_{i,*} \beta)^2$$

$$\text{Var}(y_i) = (\mathbf{X}_{i,*} \beta)^2 + 2E[\epsilon] \mathbf{X}_{i,*} \beta + E[\epsilon^2] - (\mathbf{X}_{i,*} \beta)^2$$

$$\text{Var}(y_i) = E[\epsilon^2] = \sigma^2$$

for the expected value of  $\beta$  we can insert the definition of  $\beta$  from earlier

$$E[\hat{\beta}] = E[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}] = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E[\mathbf{Y}] = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta = \beta.$$

$$\begin{aligned} \text{Var}(\hat{\beta}) &= E\{[\beta - E(\beta)][\beta - E(\beta)]^T\} \\ &= E\{[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} - \beta][(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} - \beta]^T\} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E\{\mathbf{y} \mathbf{y}^T\} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} - \beta \beta^T \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \{\mathbf{X} \beta \beta^T \mathbf{X}^T + \sigma^2 \mathbf{I}\} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} - \beta \beta^T \\ &= \beta \beta^T + \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} - \beta \beta^T = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}, \end{aligned}$$

1. <https://github.com/bragewiseth/MachineLearningProjects/blob/main/project1/src/linearRegression.py>[5]

2. We would imagine that keeping everything in the same 0 order of magnitude is something that the computer likes, perhaps improving performance

## 6. Bias & Variance

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

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

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

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

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

The term *variance* refers to how spread out our observations are. We can for our case think of the observations  $\tilde{y}_i$  as either individual predictions, or the mean of predictions from one test set.

In a sense the bias resembles the mathematical definition for variance, We can see that the bias term is a distance measure of  $\mathbf{y}$  from the expectation value of  $\tilde{\mathbf{y}}$ . If we take the individual observations approach this can be tough of as A high bias will then result in a large band around  $\tilde{\mathbf{y}}$  in which  $y_i$  will exist. The variance term is a measure of how much the  $\tilde{\mathbf{y}}$  varies around its own mean. In other terms a high variance model will experience large variance in the predicted values *'boldsymbol{\tilde{y}}* from one test set to another. Or as a result of

this a large variance in score if you will. Now these two terms are in a sense opposites, it is likely that a model with high bias will have low variance and vice versa. This is something we as designers of the model get to tune, if we value a lower bias more than we value a low variance or vice versa, we can choose a model accordingly

## **Acknowledgments**



## Appendix A.

# Bibliography

- [1] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer, 1 edition, 2007.
- [2] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. Book in preparation for MIT Press.
- [3] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. Springer Series in Statistics. Springer New York Inc., New York, NY, USA, 2001.
- [4] Wessel N. van Wieringen. Lecture notes on ridge regression, 2023.
- [5] Brage Wiseth, Eirik ?, and Felix ? MachineLearningProjects. <https://github.com/bragewiseth/MachineLearningProjects>, September 2023.