

Project 1

Edvard B. Rørnes* and Isak O. Rukan†
Institute of Physics, University of Oslo,
0371 Oslo, Norway
 (Dated: September 18, 2024)

Abstracting very cool

CONTENTS

1. Introduction	1
2. Theory	1
2.1. OLS	1
2.2. Ridge	2
2.3. LASSO	2
2.4. Resampling	2
2.5. Bias-Variance	2
3. Implementation	3
4. Results	4
4.1. OLS	4
4.2. Ridge	4
4.3. LASSO	4
5. Discussion	4
6. Conclusion	4
Part e)	4
References	4

The function f will then be approximated with a model $\tilde{\mathbf{y}}$ in which we will consider a polynomial expansion with coefficients β_i :

$$\tilde{y}_i = \sum_{j=0}^{p-1} \beta_j x_i^j \quad (2)$$

defining the $n \times p$ design matrix $(\mathbf{X})_{ij} = (x_i)^j$ we can rewrite this as

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

Further, each model will be defined with a different cost function $C(\boldsymbol{\beta})$ which we minimize to find the coefficients for each respective model.

2.1. OLS

OLS is a primitive method used in linear regression to estimate coefficients of a linear model. The cost function in OLS is simply defined as the residual sum of squares (RSS)

$$C_{\text{OLS}}(\boldsymbol{\beta}) = \text{RSS}(\boldsymbol{\beta}) = (\mathbf{y} - \tilde{\mathbf{y}})^2 = (y_i - X_{ij}\beta_j)^2$$

where we employ the summation notation where repeated indices are summed over. As mentioned prior, the coefficients $\boldsymbol{\beta}$ are found by minimizing the cost function, i.e. taking the derivative w.r.t. $\boldsymbol{\beta}$. This results in

$$\boldsymbol{\beta}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

which yields the model

$$\tilde{\mathbf{y}}_{\text{OLS}} = \mathbf{X}\boldsymbol{\beta}_{\text{OLS}} \quad (4)$$

Assuming our data takes the form of (1) then the expectation value \mathbf{y} is

$$\mathbb{E}(y_i) = \mathbb{E}(f(x_i)) = \mathbf{X}_{i,*}\boldsymbol{\beta}$$

since $\mathbb{E}(\varepsilon_i) = 0$ follows from its definition. The variance of \mathbf{y} is given by

$$\begin{aligned} \text{Var}(y_i) &= \mathbb{E}\{[y_i - \mathbb{E}(y_i)]^2\} = \mathbb{E}\{(\mathbf{X}_{i,*}\boldsymbol{\beta} + \varepsilon_i)^2\} - (\mathbf{X}_{i,*}\boldsymbol{\beta})^2 \\ &= (\mathbf{X}_{i,*}\boldsymbol{\beta})^2 + \mathbb{E}(\varepsilon_i^2) + 2\mathbb{E}(\varepsilon_i)\mathbf{X}_{i,*}\boldsymbol{\beta} - (\mathbf{X}_{i,*}\boldsymbol{\beta})^2 \\ &= \text{Var}(\varepsilon_i^2) = \sigma^2 \end{aligned}$$

1. INTRODUCTION

The methods used in this project are Ordinary Least Squares (OLS), Ridge regression and Least Absolute Shrinkage and Selection Operator (LASSO) regression.

2. THEORY

The general structure of all our models is that we have some data set $\{x_i, y_i\}$ where $i \in \{0, 1, \dots, n-1\}$ where x_i are independent variables whilst y_i are dependent variables. The data is assumed to be described by

$$\mathbf{y} = f(\mathbf{x}) + \boldsymbol{\varepsilon} \quad (1)$$

where f is some continuous function which takes \mathbf{x} as input and $\boldsymbol{\varepsilon}$ is a normal distributed error $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2)$.

* e.b.rornes@fys.uio.no

† Insert Email

which shows that $y_i \sim \mathcal{N}(\mathbf{X}_{i,*}\boldsymbol{\beta}, \sigma^2)$. The expectation value of the optimal parameters $\hat{\boldsymbol{\beta}}$ can be found to be

$$\begin{aligned}\mathbb{E}(\hat{\boldsymbol{\beta}}_{\text{OLS}}) &= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}] \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}[\mathbf{y}] \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}.\end{aligned}$$

with the variance

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

2.2. Ridge

Ridge regression is an extension of OLS where define the cost function as a modified version of the OLS cost function with an added penalty term which is proportional to the coefficients β_i^2 :

$$C_{\text{Ridge}}(\boldsymbol{\beta}) = C_{\text{OLS}}(\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^2 \quad (5)$$

Here $\lambda \geq 0$ is a regularization parameter which controls the strength of this additional penalty. This regulator essentially drives the magnitude of these coefficients allowing for more tweaking in the parameter space. This parametrization of course includes the constraint that $\boldsymbol{\beta}^2 \leq t$ for some $t < \infty$ such that we can choose our arbitrary parameter $\lambda \geq 0$ to be sufficiently small s.t. the cost function (5) does not diverge. The optimal parameters for Ridge regressions can again be found by the same process as for OLS:

$$\begin{aligned}0 &= \frac{\partial C_{\text{Ridge}}}{\partial \boldsymbol{\beta}} = -\frac{2}{n} (\mathbf{y} - \mathbf{X} \boldsymbol{\beta})^T \mathbf{X} + 2\lambda \boldsymbol{\beta}^T \\ &= \frac{2}{n} (\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} - \mathbf{y}^T \mathbf{X}) + 2\lambda \boldsymbol{\beta}^T \\ 0 &= \boldsymbol{\beta}^T (\mathbf{X}^T \mathbf{X} + \tilde{\lambda} \mathbf{I}) - \mathbf{y}^T \mathbf{X} \\ \boldsymbol{\beta}^T &= \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X} + \tilde{\lambda} \mathbf{I})^{-1} \\ \boldsymbol{\beta} &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

where we defined $\tilde{\lambda} \equiv n\lambda$, renamed $\tilde{\lambda} \rightarrow \lambda$ and used that the matrix in the parenthesis is a symmetric matrix and thus its inverse must also be symmetric. Here we can see that the effect of adding this penalty term is essentially taking $(\mathbf{X}^T \mathbf{X})^{-1} \rightarrow (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1}$ when compared to the OLS case. In the past this was generally the starting

point for Ridge regression in the cases where the matrix $\mathbf{X}^T \mathbf{X}$ was not invertible. (If we want to keep this next part we need to mention SVD) A direct way of seeing the effect of the regulator is by considering

$$\begin{aligned}\tilde{\mathbf{y}}_{\text{Ridge}} &= \mathbf{X} \boldsymbol{\beta}_{\text{Ridge}} = \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \\ &= \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T ((\mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T)^T \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T + \lambda \mathbf{I})^{-1} (\mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T)^T \mathbf{y} \\ &= \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T (\mathbf{V} \boldsymbol{\Sigma}^T \boldsymbol{\Sigma} \mathbf{V}^T + \lambda \mathbf{I})^{-1} \mathbf{V} \boldsymbol{\Sigma}^T \mathbf{U}^T \mathbf{y} \\ &= \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T (\mathbf{V} (\boldsymbol{\Sigma}^T \boldsymbol{\Sigma} + \lambda \mathbf{I}) \mathbf{V}^T)^{-1} \mathbf{V} \boldsymbol{\Sigma}^T \mathbf{U}^T \mathbf{y} \\ &= \mathbf{U} \boldsymbol{\Sigma} (\boldsymbol{\Sigma}^T \boldsymbol{\Sigma} + \lambda \mathbf{I})^{-1} \boldsymbol{\Sigma}^T \mathbf{U}^T \mathbf{y} \\ &= \sum_{j=0}^{p-1} \mathbf{u}_j \mathbf{u}_j^T \frac{\sigma_j^2}{\sigma_j^2 + \lambda} \mathbf{y}\end{aligned}$$

where the last step is valid due to the orthogonality of \mathbf{U} and σ_j are the elements on the diagonal of $\boldsymbol{\Sigma}$. Since $\lambda \geq 0$ then this added factor compared to OLS is ≤ 1 . The larger λ is the smaller this factor becomes and is the so-called a "shrinkage" factor.

2.3. LASSO

Similarly to Ridge, LASSO also includes a penalty factor. The cost function in this case is instead defined to be

$$C_{\text{LASSO}}(\boldsymbol{\beta}) = C_{\text{OLS}}(\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1 \quad (6)$$

where

$$\|\boldsymbol{\beta}\|_k \equiv \sum_{i=0}^{n-1} |\beta_i|^k$$

is the L^k norm of $\boldsymbol{\beta}$. Taking the derivative of (6) w.r.t. $\boldsymbol{\beta}$ and requiring that this becomes zero we have

$$0 = \frac{\partial C_{\text{LASSO}}}{\partial \boldsymbol{\beta}} = -2\mathbf{X}^T (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}) + \lambda \text{sgn}(\boldsymbol{\beta}) \quad (7)$$

This has the added benefit of being able to set certain parameters to be 0 instead of suppressing them, at the cost of losing analytical expressions for $\hat{\boldsymbol{\beta}}$ in non-trivial cases.

2.4. Resampling

2.5. Bias-Variance

A key part of... is the so-called Bias-Variance Trade-Off. For ease of notation we write $f(\mathbf{x}) = f$ and simply ignore vector notation since everything is a scalar in the

end. Then we have

$$\begin{aligned}
\mathbb{E}[(\mathbf{y} - \tilde{\mathbf{y}})^2] &= \mathbb{E}[(f + \varepsilon - \tilde{\mathbf{y}})^2] \\
&= \mathbb{E}[(f - \tilde{\mathbf{y}})^2] + 2 \underbrace{\mathbb{E}[(f - \tilde{\mathbf{y}})\varepsilon]}_{=0} + \underbrace{\mathbb{E}[\varepsilon^2]}_{=\sigma^2} \\
&= \mathbb{E}[(f - \mathbb{E}[\tilde{\mathbf{y}}]) - (\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])]^2 + \sigma^2 \\
&= \mathbb{E}[(f - \mathbb{E}[\tilde{\mathbf{y}}])^2] + \mathbb{E}[(\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])^2] \\
&\quad - 2 \mathbb{E}[(f - \mathbb{E}[\tilde{\mathbf{y}}])(\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])] + \sigma^2 \\
&= \text{Bias}[\tilde{\mathbf{y}}] + \text{Var}[\tilde{\mathbf{y}}] + \sigma^2 \\
&\quad - 2 \mathbb{E}[(f - \mathbb{E}[\tilde{\mathbf{y}}])(\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])]
\end{aligned}$$

where $\mathbb{E}[(f - \tilde{\mathbf{y}})\varepsilon] = 0$ is justified by ε being independent and we note that the wrong definition of the Bias is given in the problem text (with that definition σ^2 gets put into the ‘Bias’). All that remains is to show that the last term is 0. Since $\mathbb{E}[f] = f$ and $\mathbb{E}[f \mathbb{E}[\tilde{\mathbf{y}}]] = f \mathbb{E}[\mathbb{E}[\tilde{\mathbf{y}}]] = f \mathbb{E}[\tilde{\mathbf{y}}]$ then

$$\begin{aligned}
\mathbb{E}[(f - \mathbb{E}[\tilde{\mathbf{y}}])(\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}])] &= \mathbb{E}[f\tilde{\mathbf{y}} - f\mathbb{E}[\tilde{\mathbf{y}}] - \tilde{\mathbf{y}}\mathbb{E}[\tilde{\mathbf{y}}] + \mathbb{E}^2[\tilde{\mathbf{y}}]] \\
&= f\mathbb{E}[\tilde{\mathbf{y}}] - f\mathbb{E}[\tilde{\mathbf{y}}] - \mathbb{E}^2[\tilde{\mathbf{y}}] + \mathbb{E}^2[\tilde{\mathbf{y}}] = 0
\end{aligned}$$

which proves the claim.

The LHS of (??) is the expected value of the MSE which tells us how well the model’s predictions match the true data on average. The equation shows that we can decompose this expected MSE into 3 different components.

- **Bias:** This quantity measures how much the model’s average prediction differs from its true value. A high bias implies that the model is underfitting the data or is simply too simplistic.
- **Var:** The variance measures how much the model’s predictions vary when trained on different datasets. It captures the sensitivity of the model to small changes in the training data. A high variance suggests overfitting, meaning it performs well on the training data but may be capturing noise or false patterns.
- σ^2 : This is the irreducible error or noise in the data itself which cannot be explained by the model.

The idea is to minimize the LHS of (??), so clearly we want to minimize both the bias and the variance at the same time. However these are correlated to one another, so lowering the e.g. the bias will in general increase the variance. So Bias-Variance Tradeoff is essentially trying to optimize the complexity of the model such that we neither overfit nor underfit the model such that it can be generalized to other cases. These quantities can then be used as means to fine tune a model.

3. IMPLEMENTATION

For this project the surface we will consider is given by the Franke function

$$\begin{aligned}
f(x, y) &= \frac{3}{4} \exp\left(-\frac{(9x-2)^2}{4} - \frac{(9y-2)^2}{4}\right) \\
&\quad + \frac{3}{4} \exp\left(-\frac{(9x+1)^2}{49} - \frac{(9y+1)^2}{10}\right) \\
&\quad + \frac{1}{2} \exp\left(-\frac{(9x-7)^2}{4} - \frac{(9y-3)^2}{4}\right) \\
&\quad - \frac{1}{5} \exp(-(9x-4)^2 - (9y-7)^2) \quad (8)
\end{aligned}$$

This function maps a surface defined on the interval $x, y \in [0, 1]$. To perform an analysis on this function we consider a polynomial fit up to degree n where

$$\begin{aligned}
\tilde{z} &= \frac{1}{n+1} \sum_{i=0}^n \left(\beta_{00} \right. \\
&\quad + \beta_{10}x_i + \beta_{11}y_i \\
&\quad + \beta_{20}x_i^2 + \beta_{21}x_iy_i + \beta_{22}y_i^2 \\
&\quad + \dots \\
&\quad \left. + \beta_{n0}x_i^n + \beta_{n1}x_i^{n-1}y_i + \dots + \beta_{n(n-1)}x_iy_i^{n-1} + \beta_{nn}y_i^n \right) \\
&= \frac{1}{n+1} \sum_{i,j=0}^n \sum_{k=0}^i \beta_{jk}x_i^{j-k}y_i^k \equiv \mathbf{X}\boldsymbol{\beta}
\end{aligned} \quad (9)$$

where the components x_i and y_i are entries in the input vectors $\mathbf{x}^T = [x_0 \dots x_n]$ and $\mathbf{y}^T = [y_0 \dots y_n]$ respectively which are our independent variables. Each β_{ij} is a $\frac{(n+1)(n+2)}{2}$ component vector with a single non-zero entry with magnitude β_{ij} and the design matrix \mathbf{X} is then an $(n+1) \times \frac{(n+1)(n+2)}{2}$ matrix of the form:

$$\mathbf{X} = \frac{1}{n+1} \begin{bmatrix} 1 & x_0 & y_0 & x_0^2 & x_0y_0 & y_0^2 & \dots & x_0^n & \dots & y_0^n \\ 1 & x_1 & y_1 & x_1^2 & x_1y_1 & y_1^2 & \dots & x_1^n & \dots & y_1^n \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 1 & x_n & y_n & x_n^2 & x_ny_n & y_n^2 & \dots & x_n^n & \dots & y_n^n \end{bmatrix}$$

and the $\boldsymbol{\beta}$ vector contains the $\frac{(n+1)(n+2)}{2}$ components

$$\boldsymbol{\beta}^T = [\beta_{00} \ \beta_{10} \ \beta_{11} \ \beta_{20} \ \beta_{21} \ \beta_{22} \ \dots \ \beta_{n(n-1)} \ \beta_{nn}]$$

It should now be clear which unit vectors correspond to each term in (9).

We then generated data with the Franke function and used an $x(x-1)$ train-test split. This was chosen because... Next we

4. RESULTS

as

4.1. OLS

$$\mathbb{E}[(\boldsymbol{y} - \tilde{\boldsymbol{y}})^2] = \text{Bias}[\tilde{\boldsymbol{y}}] + \text{Var}[\tilde{\boldsymbol{y}}] + \sigma^2$$

4.2. Ridge

4.3. LASSO

where

5. DISCUSSION

$$\text{Bias}[\tilde{\boldsymbol{y}}] = \mathbb{E}[(\boldsymbol{y} - \mathbb{E}[\tilde{\boldsymbol{y}}])^2]$$

6. CONCLUSION

and

PART E)

Show that you can rewrite

$$\text{Var}[\tilde{\boldsymbol{y}}] = \mathbb{E}[(\tilde{\boldsymbol{y}} - \mathbb{E}[\tilde{\boldsymbol{y}}])^2] = \frac{1}{n} \sum_i (\tilde{y}_i - \mathbb{E}[\tilde{\boldsymbol{y}}])^2$$

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

Test bib [1]

[1] Planck: N. Aghanim *et. al.*, *Planck 2018 results. VI. Cosmological parameters*, *Astron. Astrophys.* **641** (2020)

A6, [[arXiv:1807.06209](#)]. [Erratum: *Astron.Astrophys.* 652, C4 (2021)].