

REGRESSION AND RESAMPLING TECHNIQUES
FYS-STK4155
ASSIGNMENT ONE

Ivar Haugaløkken Stangeby

September 12, 2019

Contents

Contents	1
1 Introduction	2

2	Regression Techniques	2
2.1	Ordinary Least Squares	4
2.2	Ridge Regression	7
2.3	Lasso Regression	7
3	Implementation and methodology	8
4	Numerical Results	9
4.1	Toy Data	9
4.2	Terrain Data	9
A	Derivation of the bias-variance-noise relationship	9
	References	10

1 Introduction

Regression analysis is the act of estimating relationships among variables. In this project, we study various regression methods in more detail. In particular, we compare the *ordinary least squares* (OLS) method with the *Ridge regression* and *Lasso regression* techniques. As we shall see, these three methods are all variations over the same theme. We start by testing the methods on noisy data sampled from a function known as *Franke's* bivariate test function (visualized on the title page), which has been in widespread use in the evaluation of interpolation and data fitting techniques. Finally, we run regression on real terrain data, comparing the aforementioned methods.

2 Regression Techniques

In general, the goal of regression analysis is to fit a *model function* $f(\mathbf{x}, \beta)$ to a set of n data points $\Omega = (\mathbf{x}_i, y_i)_{i=1}^n$. A simple example is that of a linear polynomial with two parameters:

$$f(x, \beta) = \beta_0 + \beta_1 x. \quad (1)$$

The *model parameters* β are determined in order to minimize a suitable *cost function* $C(\Omega, \beta)$ which measures to which extent the model function manages to capture trends in the data Ω . It is the choice of cost function $C(\Omega, \beta)$ which distinguishes the three regression techniques, OLS, Lasso regression, and Ridge regression.

Assumption 1. It is often assumed a priori that the data is infact generated from a noisy model, such that each y_i can be described as

$$y_i = f(\mathbf{x}_i) + \varepsilon_i \quad (2)$$

where each $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ is normally distributed with zero mean and variance σ^2 . This assumption on the error gives rise to what is known as *general linear models*.

The design matrix. We are often interested in finding the best model function in a specific function space. Assuming this space has a basis $\mathcal{B} = \{\varphi_i\}_{i=1}^M$, we may write our model function in terms of the basis functions and the model parameters as

$$f(\mathbf{x}, \beta) := \sum_{i=1}^M \beta_i \varphi_i(\mathbf{x}). \quad (3)$$

As an example, if we were to use the space Π_2^2 of bi-quadratic polynomials, our basis functions would be

$$\mathcal{B} = \{1, x, x^2, y, y^2, xy\}. \quad (4)$$

With this formulation, we can represent the approximation $\hat{\mathbf{y}}$ to \mathbf{y} as a matrix product

$$\hat{\mathbf{y}} = \mathbf{X}\beta \quad (5)$$

where \mathbf{X} is the *design matrix* defined by components $X_{ij} = \varphi_j(\mathbf{x}_i)$, and $\beta = [\beta_1, \dots, \beta_M]$ is the model parameters.

Performance metrics. In order to evaluate the how accurately the model function $f(\Omega, \beta)$ captures trends in the target data \mathbf{y} , a few standard performance metrics are used. Firstly, the *mean squared error*:

$$\text{MSE}(\mathbf{y}, \hat{\mathbf{y}}) := \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (6)$$

which simply averages the squared error over all samples and estimates. Secondly, we have the *coefficient of determination* or R^2 -score:

$$R^2(\mathbf{y}, \hat{\mathbf{y}}) := 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}. \quad (7)$$

The R^2 -score measures much of the variation in \mathbf{y} which can be attributed to a simple linear relation between \mathbf{x} and \mathbf{y} . It is a ratio, thus a value of one tells us that all the variation in the data can be attributed to an approximate linear relationship between the data \mathbf{x} and the response variable \mathbf{y} . A value of zero indicates that a non-linear model may be preferable.

2.1 Ordinary Least Squares

One common cost function is the one which involves the sum of squared residuals (or squared errors):

$$C_{\text{OLS}}(\Omega, \beta) = \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 := \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (8)$$

where $\hat{y}_i := f(\mathbf{x}_i, \beta)$. The method involving the minimization of this specific cost function is known as *least squares*. With the matrix notation from above in mind, we can also write the cost function as

$$C_{\text{OLS}}(\Omega, \beta) = \frac{1}{n} (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}). \quad (9)$$

Optimizing the parameters in OLS.

Assume now that our model is of the form given in Equation (3) and that our cost function is the mean squared error defined in Equation (8). We are interested in finding the parameters β that minimize the cost function $C(\Omega, \beta)$. Since $C(\Omega, \beta)$ is convex, it suffices to differentiate with respect to β and equating to zero.

We have that

$$\frac{\partial C_{\text{OLS}}(\Omega, \beta)}{\partial \beta} = \mathbf{X}^T (\mathbf{y} - \mathbf{X}\beta), \quad (10)$$

and equating this to zero yields the familiar *normal equations*

$$\mathbf{X}^T \mathbf{y} = \mathbf{X}^T \mathbf{X} \beta. \quad (11)$$

If the matrix $X^T X$ is invertible, we may obtain the solution by direct numerical matrix inversion. In this case, the optimal model parameters are found directly by

$$\beta_{\text{OLS}} = (X^T X)^{-1} X^T y. \quad (12)$$

However, these matrices may be ill-conditioned when the number of equations are very large, and it is therefore common to apply approximate solvers for the inverse, for instance using low-rank approximation based on SVD-decomposition, which we will briefly turn to in the following.

Singular value decomposition.

Recall that any matrix $A \in \mathbb{C}^{n,m}$ can be decomposed as

$$A = U \Sigma V^T, \quad (13)$$

where U and V are comprised of the eigenvectors of AA^T and $A^T A$ respectively. As these eigenvectors are orthonormal, it follows that both U and V are unitary matrices. Furthermore, Σ is a matrix

$$\Sigma = \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{bmatrix}, \quad (14)$$

where Σ_1 is a square diagonal matrix of size $r \times r$ with the non-zero singular values of A . The integer r is the *rank* of A . As the matrix Σ is mostly containing zeros, the information stored in A is attributed to only some parts of U and V . We can remove the redundant parts, and more compactly express A as $A = U_1 \Sigma_1 V_1^T$ without loss of accuracy in the decomposition. Here, U_1 is $m \times r$ and V_1 is $n \times r$.

Applying the singular value decomposition to the matrix $X = U \Sigma V^T$ we can analyze the expression for the prediction \hat{y} in terms of the matrix $X^T X$. First of all, we have that

$$\begin{aligned} X^T X &= V \Sigma^T U^T U \Sigma V^T \\ &= V D V^T, \end{aligned} \quad (15)$$

where we have used that $U U^T = I$ and defined $D = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$. Thus, plugging this into Equations (5) and (12), we obtain the expression

$$\hat{y} = X \beta = X (V D V^T)^{-1} X^T y. \quad (16)$$

Finally, substituting $X = U\Sigma V^T$ and noting that diagonal matrices always commute, we end up with

$$\hat{y} = UU^T y. \quad (17)$$

The bias-variance tradeoff

Since we are attempting to estimate a true distribution based on finite samples, there are certain statistical properties it is worthwhile discussing. We follow the derivation from [Meh+19]. Recall Assumption 1 which states that our true data is generated from a noisy model

$$y = f(\mathbf{x}) + \varepsilon \quad (18)$$

with $\varepsilon \sim \mathcal{N}(0, \sigma^2)$. We select a finite sample from the true data, which constitutes our finite data-set Ω . Furthermore, we have a model function $f(\Omega, \beta)$ which depends on the selected data-set and the model parameters β . Thus, our cost-function C will depend on the specific sample taken from the true distribution. It is of interest to examine the expected out-of-sample error $\mathbb{E}_{\Omega, \varepsilon}[C(\Omega, \beta)]$ of the model function over the sampled data set Ω and the noise ε . This can be decomposed into three distinct parts, *variance*, *bias*, and *noise*:

$$\mathbb{E}_{\Omega, \varepsilon}[C(\Omega, \beta)] = \frac{1}{n} \sum_{i=1}^n \left(f(\mathbf{x}_i) - \mathbb{E}_{\Omega} [\hat{y}_i] \right)^2 + \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\Omega} \left[\left(\hat{y}_i - \mathbb{E}_{\Omega} [\hat{y}_i] \right)^2 \right] + \sigma^2. \quad (19)$$

See Appendix A for the full derivation. The bias-term measures the deviation of the expectation of the model function $f(\Omega, \beta)$ from the true value f , while the variance-term measures the fluctuations arising due to using different samples of the true data.

Now, note that the equation above suggests that we should minimize both the variance and the bias of the model, in order to obtain the smallest generalization error of the model function. The variance term will be larger for complex models, as performing a fit on one sampled data-set will not necessarily translate well to another sampled data-set. The bias term measures errors introduced by making certain assumptions about the relations between the covariates X and the response y . In our case, we are assuming a linear relationship (recall *linear regression*), and thus if the true underlying relation is non-linear, our bias will be high. More complex models will in general have lower bias.

This is known as the *bias-variance tradeoff* as having a highly complex model, we reduce the bias, yet increase the variance. Having a simple model (for instance a linear model), we increase the bias, and reduce the variance.

2.2 Ridge Regression

When the dimensionality of the matrix X is high, the problem of near-collinear columns often arise, as mentioned earlier. Thus, the model parameters β_{OLS} of the ordinary least squares regression cannot be determined. The proposed fix to this problem is to slightly perturb the matrix $X^T X$ by adding small positive values along the diagonal to ensure its invertability. This is known as *ridge regression* and yields a model parameter estimator

$$\beta_{\text{ridge}}(\lambda) = (X^T X + \lambda I)^{-1} X^T y \quad (20)$$

where $\lambda \geq 0$ is a hyper parameter. To see why the resulting modified matrix is invertible, one has to look at the singular values. It turns out that the model parameters $\beta_{\text{ridge}}(\lambda)$ directly minimize the *ridge loss function* defined as

$$C_{\text{ridge}}(\Omega, \beta; \lambda) := \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^p \beta_i^2 = C(\Omega, \beta) + \lambda \beta^T \beta, \quad (21)$$

which we recognize as standard penalized least squares.

Ridge regression carries several benefits. One of the more important is its resilience to overfitting. In a high-dimensional setting, if many of the data-features can be attributed to a few covariates, then the remaining covariates do nothing more than fit to the noise. Large values of β_i is then a good indication that your model might be overfitting as some covariates are given more weight than others. The Ridge regression solves this by directly penalizing large model parameters in the cost function.

2.3 Lasso Regression

One drawback of the OLS and Ridge regression approaches is that irrelevant features may still be present in the model. Lasso regression fixes this by — in addition to penalizing large model parameters — also setting sufficiently small parameters to zero.

The associated cost function reads

$$C_{\text{lasso}}(\Omega, \beta; \lambda) := \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^n |\beta_i|. \quad (22)$$

Again, $\lambda \geq 0$ is a model hyper-parameter. Note that the penalization-term now involves absolute values of the model parameters. As the Lasso estimator may set model parameters to zero, it is effectively producing a sparser model.

Attempting to differentiate the Lasso cost function with respect to the model parameters, we see that

$$\frac{\partial C_{\text{lasso}}(\Omega, \beta; \lambda)}{\partial \beta_i} = \frac{\partial C_{\text{OLS}}(\Omega, \beta)}{\partial \beta_i} + \lambda \text{sign}(\beta_i). \quad (23)$$

The analysis of the optimal Lasso parameters is made difficult due to the discontinuity of the sign function, and no closed form solution as in OLS and Ridge-regression can be found in the general case. Thus, numerical optimization techniques are often used, for instance gradient ascent/descent.

For further information on Ridge and Lasso regression, the excellent article “Lecture Notes on Ridge Regression” by van Wieringen [vWie15] is highly recommended.

3 Implementation and methodology

In order to verify the implementation, we test on sampled data from the Franke function [Fra79]. This bivariate function $F(x, y)$ is defined for $x, y \in [0, 1]$ as follows:

$$\begin{aligned} F(x, y) := & \frac{3}{4} \exp \left(-\frac{(9x-2)^2}{4} - \frac{(9y-2)^2}{4} \right) \\ & + \frac{3}{4} \exp \left(-\frac{(9x+1)^2}{49} - \frac{(9y+1)^2}{10} \right) \\ & + \frac{1}{2} \exp \left(-\frac{(9x-7)^2}{4} - \frac{(9y-3)^2}{4} \right) \\ & - \frac{1}{5} \exp \left(-(9x-4)^2 - (9y-7)^2 \right). \end{aligned} \quad (24)$$

We are interested in fitting a bivariate polynomial surface to this data, by means of the regression techniques discussed in Chapter 2. That is, our model function $f(\Omega, \beta)$ lies

in the polynomial space Π_2^d of bi-degree (d, d) polynomials. We are mainly interested in bi-quintic polynomial surfaces, and will therefore restrict our attention to $d = 5$ in the numerical experiments.

We sample the function F both with and without noise at $N = 5000$ points in each parameter direction, yielding a matrix of function values

$$\begin{aligned} Z &= [F(x_i, y_j)]_{ij}, \\ Z^* &= [F(x_i, y_j) + \varepsilon_{ij}]_{ij}, \end{aligned} \tag{25}$$

where each $\varepsilon_{ij} \sim \mathcal{N}(0, \sigma^2)$. The design matrix X is computed using the POLYNOMIALFEATURES-method from the SKLEARN package. The dimension of the bi-quintic polynomial space is 21, yielding a design matrix of size $N^2 \times 21$.

For the real data set, we use the supplied TIFF data of an area close to Stavanger in Norway, as displayed in Figure 1. The data-set consists of approximately 1800×3600 data-points.

The ordinary least squares and the Ridge regression is performed by solving the corresponding normal equations. We use the Lasso regression available in SKLEARN package.

4 Numerical Results

4.1 Toy Data

4.2 Terrain Data

A Derivation of the bias-variance-noise relationship

Combining the two previous equations, we obtain

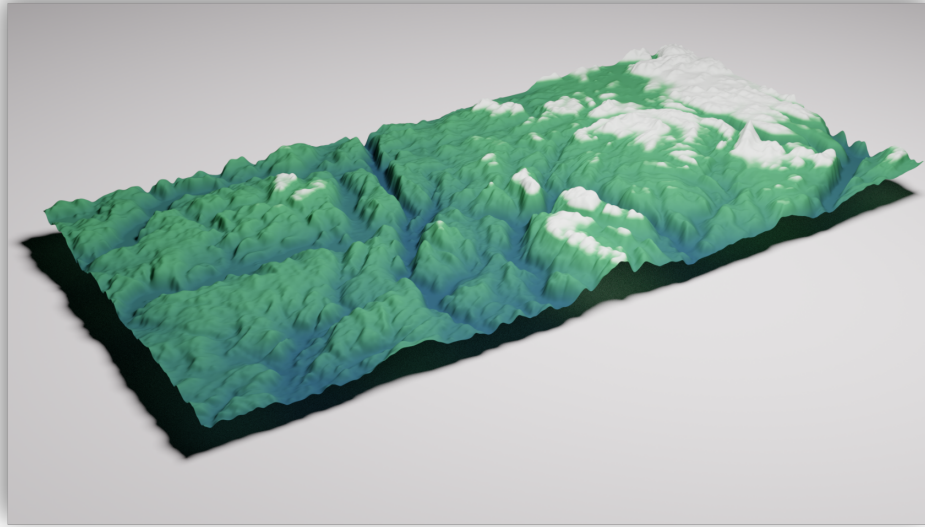


Figure 1: Real terrain data from a region close to Stavanger, Norway. For the regression tasks, we select data from small rectangular regions to keep the data-sizes manageable. The data is supplied in the file `SRTM_data_1_Norway.tif`.

References

- [Fra79] Richard Franke. *A Critical Comparison of Some Methods for Interpolation of Scattered Data*. NPS53-79-003. NAVAL POSTGRADUATE SCHOOL MONTEREY CA, Dec. 1979.
- [Meh+19] Pankaj Mehta et al. “A High-Bias, Low-Variance Introduction to Machine Learning for Physicists”. In: *Physics Reports* 810 (May 2019), pp. 1–124. ISSN: 03701573. arXiv: 1803.08823.
- [vWie15] Wessel N. van Wieringen. “Lecture Notes on Ridge Regression”. In: (Sept. 30, 2015). arXiv: 1509.09169 [stat].