

FYS-STK4155 project 1

Regression analysis and resampling methods

Mikael B. Kiste, Erik F. Skaar, Sondre K. Torp

October 7, 2018

Abstract

Lorem ipsum dolor sit amet, consectetur adipisicing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

Contents

1	introduction	2
2	Theory	3
2.1	Ordinary least squares	3
2.1.1	Mean square error	4
2.2	R^2 score - The Coefficient of Variation	5
2.3	Ridge	5
2.4	Lasso	6
2.5	K-fold	7

1 introduction

A typical problem within the natural sciences is how to interpret the trends and behavior of the results and data from an experiment. As a first approximation this will often be done qualitatively, e.g. "These values appear to increase linearly with time", but a more rigorous approach through regression analysis and resampling methods is eventually more preferable. Lorem ipsum dolor sit amet, consectetur adipisicing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

2 Theory

2.1 Ordinary least squares

We want to get a specific solution to the equation

$$\hat{y} = \hat{X}\hat{\beta} + \hat{\epsilon} \quad (1)$$

Where \hat{y} is a vector of our measured values, \hat{X} is a matrix containing variables and determines how we want to fit our data, $\hat{\beta}$ is a vector of the parameters for our fit and $\hat{\epsilon}$ is a vector representing the error in our data-points (often termed the residuals, representing how far off our prediction is from the measurements). The variables \hat{y} and \hat{X} are fixed and we want to choose parameters $\hat{\beta}$ in such a way that the errors $\hat{\epsilon}$ are minimized. An example might help clarify the situation

Lets say we have conducted an experiment where we have measured the position of a ball launched straight up into the air from a cannon. Neglecting air resistance we know that the analytical solution is on the form of a second order polynomial $x(t) = x_0 + v_0t + at^2$, where x_0 and v_0 are the initial conditions for the position and velocity respectively. This analytical solution is our model, but the actual measured values could (and indeed probably would) differ from this, simply due to errors in the measurement or other effects coming into play that the model has not accounted for (like air resistance). In any case, if we measured the position n times our linear algebra problem could be stated like this:

$$\begin{bmatrix} x(t_0) \\ x(t_1) \\ \vdots \\ x(t_{n-1}) \\ x(t_n) \end{bmatrix} = \begin{bmatrix} (t_0)^0 & (t_0)^1 & (t_0)^2 \\ (t_1)^0 & (t_1)^1 & (t_1)^2 \\ \vdots & \vdots & \vdots \\ (t_{n-1})^0 & (t_{n-1})^1 & (t_{n-1})^2 \\ (t_n)^0 & (t_n)^1 & (t_n)^2 \end{bmatrix} \begin{bmatrix} x_0 \\ v_0 \\ a \end{bmatrix} + \begin{bmatrix} \epsilon_0 \\ \epsilon_1 \\ \vdots \\ \epsilon_{n-1} \\ \epsilon_n \end{bmatrix}$$

Or, again, stated through vectors and matrix notation

$$\hat{x} = \hat{T}\hat{\beta} + \hat{\epsilon}$$

With some of the variable names adjusted simply to better indicate the represented values we have in this problem. Essentially we want to determine the variables x_0 , v_0 and a so that the error terms are minimized (the equation has solutions for all $\hat{\beta}$ s, but most of those are horrible fits with huge error terms). Note that there are two dimensionalities coming into play here. n is the number of measurements and determines the length of the column vectors \hat{x} and $\hat{\epsilon}$. In addition there is a second dimension that determines the number of columns in the matrix \hat{T} and the length of the vector $\hat{\beta}$. This number, say m , indicates the complexity of our model. When doing a polynomial fit, $m - 1$ is the order of the polynomial we want to fit our data to. So in this case, where we want to fit a second degree polynomial, we have $m = 3$.

The example above gives an impression of what the variables in the linear algebra equation represents but is quite specific and we can generalize a bit. For instance, it is not necessary to fit a polynomial at all. By changing our \hat{X} matrix we could fit to any orthogonal function that we would like. The number of datapoints can be anything we want as long as $m \leq n$.

A key part of the fitting is how to minimize the so called "cost-function". In our case we want to minimize the ϵ 's. There are different ways of doing this, but perhaps the simplest one is to do an **Ordinary Least Squares** (OLS) fit. That is to say when taking the difference between our predicted values and the measured values (essentially being the residuals) we want the squared sum of the difference for each datapoint to be as low as possible. That is to say that we want to minimize the function

$$Q = \sum_{i=0}^{n-1} \epsilon_i^2 = \epsilon^T \hat{\epsilon} = (\hat{y} - \hat{X}\hat{\beta})^T (\hat{y} - \hat{X}\hat{\beta})$$

We can see that taking the squared sum of all the elements in $\hat{\epsilon}$ is the same as taking the inner product of the vector with itself, and that we can use equation 1 to develop the expression further. We are interested in finding the parameters $\hat{\beta}$ that leads to the minimization of the squared sum of the residuals. So now that we have a function of the squared sum of the residuals with $\hat{\beta}$ as a variable it is a simple matter of finding when the derivative of this function with respect to $\hat{\beta}$ is zero; as this will give us the minimum (it must, of course, be a minimum as only a quite specific $\hat{\beta}$ will give a good fit and other values, deviating from this, would only increase the residuals squared sum. The possibility of more than one minimum is not something we need to worry about at this point). We will here state the derivative without further explanation, suffice it to say that looking at the expanded indexed expression (i.e. not on vector form) one quickly comes to the conclusion that this must indeed be the correct derivative.

$$\begin{aligned} \frac{\partial Q(\hat{\beta})}{\partial \hat{\beta}} &= 0 = \hat{X}^T (\hat{y} - \hat{X}\hat{\beta}) \\ \hat{X}^T \hat{y} &= \hat{X}^T \hat{X} \hat{\beta} \\ \hat{\beta} &= (\hat{X}^T \hat{X})^{-1} \hat{X}^T \hat{y} \end{aligned} \tag{2}$$

And by the magic of linear algebra we have arrived at an analytical solution for the parameters $\hat{\beta}$ we need for an OLS regression. It is, however, important to note that we assume that the matrix \hat{X} here is invertible.

2.1.1 Mean square error

The Mean Square Error (**MSE**) can give a measure of the quality of our estimator. It is defined as

$$MSE(\epsilon) = \frac{1}{n} \sum_n^{n-1} \epsilon^2 \tag{3}$$

As such it can be thought of as the average of the square of our residuals. We see that our OLS method minimizes the MSE of our predictor variables. Of course, as easily seen from the definition, the MSE can never be negative and lower values means that we have a better prediction (at zero there is a perfect fit).

2.2 R^2 score - The Coefficient of Variation

In regression validation the R^2 is the gold standard when it comes to measuring goodness of fit. In straight terms it is the proportion of the variance in the dependent variable that is predictable from the independent variable(s) [source wikipedia - coefficient of determination].

$$R^2 = 1 - \frac{\sum (y_i - \tilde{y}_i)^2}{\sum (y_i - \bar{y})^2} \quad (4)$$

Where y_i are the indexed response variables (data we want to fit) and \tilde{y}_i is the predictor variables from our model (so $\epsilon_i = y_i - \tilde{y}_i$). The average of the response variables is denoted \bar{y} . In our case it the second term can also be considered as the ratio of **MSE** to the variance (the $1/n$ factors kill each other in a fraction). Let's interpret the formula step by step to get an impression of how it differentiates a poor from a good fit. If the residual sum of squares (SS_{res}) is low we have a good fit. However, we should compare this to the spread of our response variables. After all, if the response variables are all nicely distributed close to the mean then getting a good SS_{res} is not that impressive. We therefore do a sort of normalization in the fraction, taking the scale of our data into consideration. In the simplest polynomial fit, using a zeroth order polynomial (just a constant), we see that our model would just be a constant function of the mean. The sums would be equal, returning unity on the fraction and the total R^2 score would be zero. In the other extreme, if our model fits perfectly, then SS_{res} would be zero and the R^2 score would be one. In this sense we have a span of possible R^2 scores between zero and one, from the baseline of the simplest model at zero and a perfect fit at one. The R^2 score is useful as a measure of how good our model is at predicting future samples.

2.3 Ridge

A problem with the **OLS** method through linear regression is that the matrix is not necessarily invertible. In this case it is basically impossible to model the data using linear regression [Morten regression notat]. When there are many columns in the \hat{X} matrix it is less likely that the columns are all linearly independent, which is a requirement to get the inverse $(\hat{X}^T \hat{X})^{-1}$. Since the number of columns increase when the complexity of our fit increases (as mentioned earlier, if the polynomial we want to fit is of order n then the number of columns is $m = n + 1$) we understand that a more complex model decreases the likelihood of invertibility. Since there is no longer a unique solution to our problem, it is an example of an ill-posed problem that is either overdetermined with more equations than unknowns (oftentimes leading to no solution) or underdetermined with more unknowns than equations (resulting in many

potential solutions). These situations correspond to an over-fitting or an under-fitting respectively. [Wikipedia Tikhonob regularization] A simple solution to this linear algebra problem is to add a diagonal matrix, $\lambda \hat{I}$, term to the matrix that is to be inverted. This shrinks the regression coefficients so that $\hat{\beta}$'s with smaller norm are preferred. In this way we superimpose a quality control variable on our solution space allowing us to arrive at a unique solution with (relatively) small variance. [hastings et. al. elementsstat.pdf] In the end we only need to adjust our linear algebra equation slightly

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

There is a caveat here when considering the coefficient determining the y-intercept. If this is included in the regularization the result would depend on the origin chosen for y [hastings et. al. elementsstat.pdf]. This is a problem as we expect that adding a constant to the response variables should simply shift the predictor variables by the same amount, not change the shape of our fit. As an illustration, if we return to our previous cannon and ball experiment, we don't want the time at which we start the clock, be it at the moment of launch or a minute before, to lead to a completely different polynomial fit. To avoid this problem it is possible to center the inputs by simply estimating β_0 by the average \bar{y} and doing a ridge regression by equation 5 with the remaining coefficients but replacing the elements in \hat{X} with $x_{ij} - \bar{x}_j$.

2.4 Lasso

The Lasso (least absolute shrinkage and selection operator) method works in much the same way as Ridge in that it incorporates variable selection and regularization in the regression. But the constraints makes the solutions nonlinear and there is no closed form expression as in the ridge regression [hastings p.68]. We therefore need to specify the workings of the method more explicitly

$$\hat{\beta} = \operatorname{argmin} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \quad (6)$$

Under the restriction

$$\sum_{j=1}^p |\beta_j| \leq t \quad (7)$$

At first glance this expression may seem foreign when given on this form but it is essentially the same as for Ridge with the only difference being the restriction in Ridge

$$\sum_{j=1}^p \beta_j^2 \leq t \quad (8)$$

We recognize equation 6 as a function that, when we exclude β_0 due to centering, finds the minimum of SS_{res} with an additional constraint that discards some solutions. In the case of Lasso the constraint makes sure that the L_1 norm is less than t while Ridge considers the L_2 norm (i.e. euclidian norm).

2.5 K-fold

K-folding is a cross validation technique that allows us to generalize the trends in our data set to an independent data set. In this way we can circumvent typical problems like over-fitting and selection bias[wikipedia cross validation]. The approach for the technique is pretty straight forward. Instead of doing a regression on the entire data set, we first segment it into k number of subsets of equal size (making sure to pick out the variables randomly before distributing them to the subsets). Now we choose one subset to be the 'control' or 'validation' set while the rest of the subsets are the training sets. We then perform the regression we want to use on the training set, arriving at some data fitting that is our prediction. From here it is a straight forward process to compare how well our predicted variables compare to the validation variables, for example through the R^2 score function. However, even though our subsets are picked randomly, the validation subset we used could potentially not be a representative selection of the entire set. Therefore we make sure to repeat the process k times, each time using a new subset as the validation subset. After all this is done we can simply calculate the average of the scores to get the predictive power of our model. As an added benefit, since we are doing the calculations anyway, we can use the average of our predictions as our final fit. Cross validation techniques are extremely useful when the gathering of new data is difficult or, sometimes, even impossible, as we are using the extra computational power at our disposal to squeeze the most amount of relevant information out of our precious data.

Combining cross validation with the penalty parameter t used in Ridge or Lasso (eq: 7 and 8) can make for an even better fit. Instead of trying to pick a good t from our intuition we can actually get a numerical measure of the predictive power of our model as a function of t by performing the cross validation method for a range of ts , choosing the penalty parameter that optimizes the predictive power of the model.