

An Unknown Signal Report

George Herbert
cj19328@bristol.ac.uk

April 6, 2021

Abstract

This report demonstrates my understanding of the methods I have used, the results I have obtained and my understanding of issues such as overfitting for the ‘An Unknown Signal’ coursework.

1 Equations for linear regression

For a set of points that lie along a line with Gaussian noise $\mathbf{y} = \mathbf{X}\mathbf{w} + \epsilon$ where $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$, the maximum likelihood estimation of \mathbf{w} is equivalent to the least square error estimation and is given by the equation:

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

This equation is implemented in my code as the following method:

```
def regressionNormalEquation(self, X, y):  
    return np.linalg.inv(X.T @ X) @ X.T @ y
```

\mathbf{X} can take one of the following three forms:

$$\mathbf{X} = \begin{bmatrix} x_1 & 1 \\ \vdots & \vdots \\ x_{20} & 1 \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^n & x_1^{n-1} & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ x_{20}^n & x_{20}^{n-1} & \dots & 1 \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} f(x_1) & 1 \\ \vdots & \vdots \\ f(x_{20}) & 1 \end{bmatrix}$$

depending on whether the line is linear, polynomial of degree n , or the unknown function f , respectively.

2 Choice of polynomial degree

Having created a program, `display.py`, to visualise the graphs, I drew up a list of line segments that appeared to be non-linear. Then, I created a program, `order.py`, that calculated the cross-validation error for each of these line segments when trained using a model from a polynomial of degree 2, to a polynomial of degree 10. A small section of the output from this program is shown in Table 1.

Having analysed this output, it was clear that the polynomial degree that most consistently gave a low cross-validation error was a polynomial of degree 3 (i.e. a cubic).

Table 1: Section of the output from `order.py`

Filename	Line segment	Polynomial order	Cross-validation error
basic_3.csv	0	2	7.3947610358752875
basic_3.csv	0	3	1.2989585613760917e-23
\vdots	\vdots	\vdots	\vdots
adv_3.csv	5	9	318.8443359827487
adv_3.csv	5	10	279.2750683133305

3 Choice of unknown function

Initially, I used an ‘eyeball estimation’ to draw up a list of candidates for the unknown function.

4 Overfitting

Overfitting occurs when a machine learning algorithm produces a model that has learnt the noise in the data as if it represents the structure of the underlying model. [1]

In the case of linear regression, overfitting is most likely to occur by producing a model with too complex a function type, such that it would fail to predict future observations.

5 Model selection

To prevent overfitting, I have used leave-one-out cross-validation when producing a model for each 20-point line segment. Leave-one-out cross-validation is an extreme case of k -fold cross validation such that $k = n$, where n is the number of data-points (in this case 20). Despite being computationally expensive, I believe that leave-one-out cross-validation is an appropriate technique to prevent overfitting in this case, owing to the limited sample size of each line segment.

Leave-one-out cross-validation involves using each of the 20 data-points exactly once as validation data for a model trained using the other 19 data-points. The cross-validation error for each function type is calculated as follows [2]:

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}^{(-i)})^2$$

where n is the number of datapoints in a line segment (i.e. 20); y_i is the actual y -value for the i -th datapoint; and $\hat{y}^{(-i)}$ is the predicted y -value for the i -th datapoint, when trained without using the i -th sample.

The function type with the lowest cross-validation error is then selected.

6 Optimisations and improvements

To begin with, computing the matrix inverse using the `np.linalg.inv` method is computationally expensive and unnecessary. Instead, given \mathbf{X} and \mathbf{y} , the maximum likelihood

estimation could be computed directly as follows: `np.linalg.solve(X.T @ X, @ X.T @ y)`. This would be faster, as `np.linalg.inv` computes the inverse of a matrix \mathbf{A} by solving for \mathbf{A}^{-1} in $\mathbf{A}\mathbf{A}^{-1} = \mathbf{I}$ [3]. Thus, there would be a performance benefit by solving for $\hat{\mathbf{w}}$ in $\mathbf{X}^T\mathbf{X}\hat{\mathbf{w}} = \mathbf{X}^T\mathbf{y}$ directly.

Another computationally expensive operation in my algorithm is that used to calculate the cross-validation error using leave-one-out cross-validation. This is because it involves fitting the model and calculating the sum squared error n times. Instead, there exists a faster method I could have adopted that involves calculating the leverage. Despite this, I opted not to include this method; owing to the fact that my program as it currently stands can be easily adapted to use k -fold cross-validation for any value of k that is a factor of 20 by changing the constant `K` in the code.

1. Not shuffling the data before leave-one-out cross-validation
2. Not finding the average cross-validation error

7 Testing

I created a file `test.py` to test each of the functions in my program.

References

- [1] Burnham, K. P. and Anderson, D. R. (2002) *Model Selection and Multimodel Inference*. 2nd ed. Springer-Verlag.
- [2] Taylor, J. (2020) *Leave one out cross-validation (LOOCV) — STATS 202* <https://web.stanford.edu/class/stats202/notes/Resampling/LOOCV.html>
- [3] Muldal, A. (2017) *Why does numpy.linalg.solve() offer more precise matrix inversions than numpy.linalg.inv()?* <https://stackoverflow.com/a/31257909/8540479>