

Machine Learning

Bhaswar Chakma

2021-07-06

Contents

	5
1 OLS	7
1.1 SLR	7

Source: Google ML Comic

Chapter 1

OLS

1.1 SLR

We define our column vector of parameters \mathbf{w} (β in econometrics):

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

A single observation \mathbf{x}_n is:

$$\mathbf{x}_n = \begin{bmatrix} 1 \\ x_n \end{bmatrix}$$

And thus the value predicted by the model can be written as:

$$f(x_n; w_0, w_1) = w_0 + w_1 x_n = [w_0 \quad w_1] \begin{bmatrix} 1 \\ x_1 \end{bmatrix} = \mathbf{w}^T \mathbf{x}_n$$

The matrix representing all the input data \mathbf{X} and the column vector of outcomes \mathbf{t} (\mathbf{y} in econometrics) are:

$$\mathbf{X} = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix}$$

$$\mathbf{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix}$$

The loss function \mathcal{L} becomes

$$\mathcal{L} = \frac{1}{N}(\mathbf{t} - \mathbf{X}\mathbf{w})^T(\mathbf{t} - \mathbf{X}\mathbf{w})$$

The column vector of values predicted by the function $f(x_n; w_0, w_1)$ are:

$$\mathbf{X}\mathbf{w} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix} \times \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} w_0 & w_1 x_1 \\ w_0 & w_1 x_2 \\ \vdots & \vdots \\ w_0 & w_1 x_N \end{bmatrix}$$

$$\begin{matrix} \mathbf{X} \\ (N \times 2) \end{matrix} \times \begin{matrix} \mathbf{w} \\ (2 \times 1) \end{matrix} = \begin{matrix} \mathbf{X}\mathbf{w} \\ (N \times 1) \end{matrix}$$

The column vector of residual is

$$\mathbf{t} - \mathbf{X}\mathbf{w} = \begin{bmatrix} t_1 - w_0 + w_1 x_1 \\ t_2 - w_0 + w_1 x_2 \\ \vdots \\ t_N - w_0 + w_1 x_N \end{bmatrix}$$

The sum of squared errors can be expressed as:

$$\begin{aligned} & \sum_{n=1}^N (t_n - (w_0 + w_1 x_1))^2 \\ &= [t_1 - (w_0 + w_1 x_1) \quad t_2 - (w_0 + w_1 x_2) \quad \dots \quad t_N - (w_0 + w_1 x_N)] \times \begin{bmatrix} t_1 - (w_0 + w_1 x_1) \\ t_2 - (w_0 + w_1 x_2) \\ \vdots \\ t_N - (w_0 + w_1 x_N) \end{bmatrix} = (\mathbf{t} - \mathbf{X}\mathbf{w})^T \end{aligned}$$

We can now write our loss function compactly with matrix notation, which produces a scalar (1×1) value:

$$\frac{1}{N} \begin{matrix} (1 \times 1) \end{matrix} (\mathbf{t} - \mathbf{X}\mathbf{w})^T \begin{matrix} (1 \times N) \end{matrix} (\mathbf{t} - \mathbf{X}\mathbf{w}) \begin{matrix} (N \times 1) \end{matrix} = \mathcal{L} \begin{matrix} (1 \times 1) \end{matrix}$$

We can expand the sum of squared errors as follows:

$$\begin{aligned} & (\mathbf{t} - \mathbf{X}\mathbf{w})^T (\mathbf{t} - \mathbf{X}\mathbf{w}) \\ &= (\mathbf{t}^T - (\mathbf{X}\mathbf{w})^T) (\mathbf{t} - \mathbf{X}\mathbf{w}) \end{aligned}$$

$$\begin{aligned}
&= (\mathbf{t}^T - (\mathbf{w}^T \mathbf{X}^T))(\mathbf{t} - \mathbf{X}\mathbf{w}) \\
&= \underset{(1 \times 1)}{\mathbf{t}^T \mathbf{t}} - \underset{(1 \times N)(N \times 2)(2 \times 1)}{\mathbf{t}^T \mathbf{X} \mathbf{w}} - \underset{(1 \times 2)(2 \times N)(N \times 1)}{\mathbf{w}^T \mathbf{X}^T \mathbf{t}} + \underset{(1 \times 2)(2 \times N)(N \times 2)(2 \times 1)}{\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}} \quad (1.1)
\end{aligned}$$

Note the following:

- each of the term becomes 1×1
- the terms $\mathbf{t}^T \mathbf{X} \mathbf{w}$ and $\mathbf{w}^T \mathbf{X}^T \mathbf{t}$ are transposes of each other and are scalars. So they are the same value and can be combined.

So, sum of squared residuals can be written as

$$\mathbf{t}^T \mathbf{t} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{t} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}$$

Hence, the loss function becomes:

$$\mathcal{L} = \frac{1}{N} \mathbf{t}^T \mathbf{t} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{t} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}$$