

Linear Regression

Hanxiao Du

Contents

1	The 1-D case	1
2	Multidimensional Inputs	3
3	Multidimensional Outputs	4
4	Appendix	4
4.1	Gauss-Markov Theorem	4

1 The 1-D case

Our goal is to learn a mapping $y = f(x)$, where x and y are both real-valued scalars (i.e. $x \in \mathbb{R}$, $y \in \mathbb{R}$). Take f to be an linear function (actually an affine function) of the form:

$$y = wx + b \tag{1}$$

where w is a weight and b is a bias.

We wish to estimate w and b from the N training pairs $\{(x_i, y_i)\}_{i=1}^N$. Then, once we have to estimates of w and b , we can compute y for some new x .

We would like to find the parameters (i.e. w and b) such that minimize the residual errors (i.e. $y_i - f(x_i) = y_i - (wx_i + b)$).

The most commonly-used way to estimate the parameters is by least-square regression/estimation or [ordinary least square \(OLS\) estimation](#). OLS chooses the parameters of a linear function of a set of explanatory variables by the principle of least squares. We define an energy function (a.k.a. objective function):

$$E(w, b) = \sum_{i=1}^N (y_i - (wx_i + b))^2 \tag{2}$$

This is merely the summation of the squared residual errors. Here, the reasons of using squared residual instead of summing all residuals directly or summing all the absolute values of the residuals are:

1. It is mathematically easier to find the derivatives of the objective function so that the optimization process is much easier.
2. By [Gauss-Markov theorem](#), assuming the Gauss-Markov assumptions, the OLS estimator has the lowest sampling variance within the class of linear unbiased estimators (most efficient). The Gauss-Markov assumptions are:

- The residuals ϵ_i have mean zero, i.e. $\mathbb{E}[\epsilon_i] = 0$
- The residuals are [homoscedastic](#), i.e. they all have the same finite variance:

$$Var(\epsilon_i) = \sigma^2 < \infty$$

- Distinct error terms are uncorrelated, i.e. $Cov(\epsilon_i, \epsilon_j) = 0 \forall i \neq j$.

Proof of Gauss-Markov theorem can be found in Section 4.1

3. The OLS estimator is equivalent to the maximum likelihood estimator (MLE) under the normality assumption for the error terms.

Thus we would like to find \hat{w}_{OLS} and \hat{b}_{OLS} that minimize E . That is:

$$\hat{b}_{OLS} = \underset{b}{\operatorname{argmin}} E(w, b) \quad (3)$$

$$\hat{w}_{OLS} = \underset{w}{\operatorname{argmin}} E(w, \hat{b}_{OLS}) \quad (4)$$

Here, we can use the [derivative test](#) to find the global extrema.
Solve for \hat{b}_{OLS} :

$$\frac{\partial E}{\partial b} = \frac{\partial \sum_{i=1}^N (y_i - (wx_i + b))^2}{\partial b} \quad (5)$$

$$= \sum_{i=1}^N -2(y_i - (wx_i + b)) = 0 \quad (6)$$

$$\implies \sum_{i=1}^n y_i - w \sum_{i=1}^n x_i - Nb = 0 \quad (7)$$

$$\implies \hat{b}_{OLS} = \frac{\sum_{i=1}^n y_i - w \sum_{i=1}^n x_i}{N} \quad (8)$$

$$\implies \hat{b}_{OLS} = \bar{y} - w\bar{x} \quad (9)$$

where $\bar{y} = \frac{\sum_{i=1}^n y_i}{N}$ is the average of y_i 's, and $\bar{x} = \frac{\sum_{i=1}^n x_i}{N}$ is the average of x_i 's.

Notice that we did not check if $\frac{\partial^2 E}{\partial b^2} > 0$ to ensure that \hat{b}_{OLS} minimizes E , this is because we know that E is a quadratic function of b and we can conclude that \hat{b}_{OLS} minimizes E by just inspecting E . The same principle applied in finding \hat{w}_{OLS} .

Solve for \hat{w}_{OLS} :

$$\frac{\partial E}{\partial w} = \frac{\partial \sum_{i=1}^N (y_i - (wx_i + \hat{b}_{OLS}))^2}{\partial w} \quad (10)$$

$$= \frac{\partial \sum_{i=1}^N (y_i - (wx_i + \bar{y} - w\bar{x}))^2}{\partial w} \quad (11)$$

$$= \frac{\partial \sum_{i=1}^N (y_i - \bar{y} - w(x_i - \bar{x}))^2}{\partial w} \quad (12)$$

$$= -2 \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y} - w(x_i - \bar{x})) = 0 \quad (13)$$

$$\Rightarrow \sum_{i=1}^N x_i y_i - \bar{y} \sum_{i=1}^N x_i - \bar{x} \sum_{i=1}^N y_i + \bar{x} \bar{y} - w \sum_{i=1}^N (x_i - \bar{x})^2 = 0 \quad (14)$$

$$\Rightarrow \hat{w}_{OLS} = \frac{\sum_{i=1}^N (x_i y_i - \bar{y} x_i - \bar{x} y_i + \bar{x} \bar{y})}{\sum_{i=1}^N (x_i - \bar{x})^2} \quad (15)$$

$$\Rightarrow \hat{w}_{OLS} = \frac{\sum_{i=1}^N (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^N (x_i - \bar{x})^2} \quad (16)$$

Notice that for this equation, the b in E is dependent to w , thus we need to express b in terms of w in order to find the derivative of w properly.

2 Multidimensional Inputs

For a D -dimensional linear regression, we can express the linear model in terms of matrices:

$$f(X) = X\beta \quad (17)$$

where β is a $(D + 1) \times 1$ matrix (or a column vector) of parameters in the form

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_D \end{bmatrix} \quad (18)$$

β_0 is the intercept term and $\beta_i \forall i \in [D]$ is the coefficient of the i -th attribute/feature. X is a $N \times (D + 1)$ design matrix of the form:

$$X = \begin{bmatrix} 1 & x_1^{(1)} & x_1^{(2)} & \cdots & x_1^{(D)} \\ 1 & x_2^{(1)} & x_2^{(2)} & \cdots & x_2^{(D)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_N^{(1)} & x_N^{(2)} & \cdots & x_N^{(D)} \end{bmatrix} \quad (19)$$

where $x_i^{(j)} \forall i \in [N] \forall j \in [D]$ denotes the j -th attribute of the i -th instance. Notice that all the values in the first column of X are all 1's, this is because when we do the matrix multiplication $X\beta$, the intercept term β_0 is added to each column of the product. i.e. this is a dummy column for the intercept term.

The least-squares objective function is then $E(\beta) = \|y - X\beta\|_2^2$ if we treat y and $X\beta$ as column vectors:

$$X\beta = \begin{bmatrix} \sum_{j=1}^D \beta_j x_1^{(j)} + \beta_0 \\ \sum_{j=1}^D \beta_j x_2^{(j)} + \beta_0 \\ \vdots \\ \sum_{j=1}^D \beta_j x_N^{(j)} + \beta_0 \end{bmatrix}, y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \quad (20)$$

Notice that

$$y^T X\beta = (y^T X\beta)^T = (X\beta)^T y = \beta^T X^T y \quad (21)$$

since $y^T X\beta$ is just a 1×1 matrix, thus symmetric. Therefore, the squared Euclidean norm

$$E(\beta) = \|y - X\beta\|_2^2 \quad (22)$$

$$= (y - X\beta)^T (y - X\beta) \quad (23)$$

$$= y^T y - y^T X\beta - (X\beta)^T y + (X\beta)^T X\beta \quad (24)$$

$$= \beta^T X^T X\beta - 2\beta^T X^T y + y^T y \text{ by (21)} \quad (25)$$

The derivatives of matrices are calculated by [scalar-by-matrix derivative](#). This is a derivative formula sheet:

[Matrix_derivatives_cribsheet.pdf](#)

Now, solve for $\hat{\beta}_{OLS}$ by using the [derivative test](#):

$$\frac{\partial E}{\partial \beta} = 2X^T X\beta - 2X^T y = 0 \quad (26)$$

$$\implies \hat{\beta}_{OLS} = (X^T X)^{-1} X^T y \quad (27)$$

Or equivalently, $\hat{\beta}_{OLS} = X^+ y$, where $X^+ = (X^T X)^{-1} X^T$ is called the pseudo-inverse of X . This is also a projection matrix.

3 Multidimensional Outputs

In the most general case, both the inputs and outputs may be multidimensional.

4 Appendix

4.1 Gauss-Markov Theorem

Theorem 1 (Gauss-Markov theorem). *Suppose we have a linear model $y = X\beta + \epsilon$, where ϵ is a column vector containing all residuals ϵ_i 's.*

Under the Gauss-Markov assumptions,

1. *The residuals ϵ_i have mean zero, i.e. $\mathbb{E}[\epsilon_i] = 0$ or $\mathbb{E}[\epsilon] = \vec{0}$.*

2. They are *homoscedastic*, i.e. they all have the same finite variance:

$$\text{Var}(\epsilon_i) = \sigma^2 < \infty$$

3. Distinct error terms are uncorrelated, i.e. $\text{Cov}(\epsilon_i, \epsilon_j) = 0 \forall i \neq j$.

or equivalently to assumption 2 and 3, $\text{Var}(\epsilon) = \sigma^2 I$ where $\sigma^2 < \infty$, we have that the ordinary least squares (OLS) of the parameter:

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

is the BLUE (Best Linear Unbiased Estimator), i.e. it is the most efficient estimator (has lowest variance) among all linear unbiased estimators.

Proof.

$$\hat{\beta}_{OLS} = (X^T X)^{-1} X^T y \quad (28)$$

$$= (X^T X)^{-1} X^T (X\beta + \epsilon) \quad (29)$$

$$= (X^T X)^{-1} X^T X\beta + (X^T X)^{-1} X^T \epsilon \quad (30)$$

$$= \beta + (X^T X)^{-1} X^T \epsilon \quad (31)$$

Thus, the unbiasedness of $\hat{\beta}_{OLS}$ gives:

$$\mathbb{E}[\hat{\beta}_{OLS}] = \mathbb{E}[\beta + (X^T X)^{-1} X^T \epsilon] \quad (32)$$

$$= \beta + (X^T X)^{-1} X^T \mathbb{E}[\epsilon] \quad (33)$$

$$= \beta, \text{ since by assumption 1, } \mathbb{E}[\epsilon] = \vec{0}. \quad (34)$$

Also by assumption 1,

$$\text{Var}(\epsilon) = \mathbb{E}[(\epsilon - \mathbb{E}[\epsilon])(\epsilon - \mathbb{E}[\epsilon])^T] = \mathbb{E}[\epsilon\epsilon^T] = \sigma^2 I \quad (35)$$

Thus,

$$\text{Var}(y) = \text{Var}(X\beta + \epsilon) = \text{Var}(\epsilon) = \sigma^2 I \quad (36)$$

Now, we calculate the variance of the $\hat{\beta}_{OLS}$:

$$\text{Var}(\hat{\beta}_{OLS}) = \mathbb{E}[(\hat{\beta}_{OLS} - \beta)(\hat{\beta}_{OLS} - \beta)^T] \quad (37)$$

$$= \mathbb{E}[(X^T X)^{-1} X^T \epsilon ((X^T X)^{-1} X^T \epsilon)^T] \quad (38)$$

$$= \mathbb{E}[(X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1}] \quad (39)$$

$$= (X^T X)^{-1} X^T \mathbb{E}[\epsilon \epsilon^T] X (X^T X)^{-1} \quad (40)$$

$$= (X^T X)^{-1} X^T \sigma^2 I X (X^T X)^{-1} \text{ by (35)} \quad (41)$$

$$= \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} \quad (42)$$

$$= \sigma^2 (X^T X)^{-1} \quad (43)$$

Now, for any linear unbiased estimator $\hat{\beta}$ for β , it must can be expressed in linear form:

$$\hat{\beta} = Cy$$

for some constant matrix C .

Also, let D be the matrix such that $C = (X^T X)^{-1} X^T + D$.

Additionally, by the unbiasedness assumption of the estimator $\hat{\beta}$:

$$\mathbb{E}[\hat{\beta}] = \mathbb{E}[Cy] \quad (44)$$

$$= \mathbb{E}[(X^T X)^{-1} X^T + D](X\beta + \epsilon) \quad (45)$$

$$= \beta + DX\beta = \beta \quad (46)$$

$$\implies DX\beta = \underset{\sim}{0} \implies DX = \underset{\sim}{0} \quad (47)$$

where $\underset{\sim}{0}$ is the zero matrix.

Another fact from linear algebra:

Finally, calculate the variance of $\hat{\beta}$:

$$Var(\hat{\beta}) = Var[Cy] \quad (48)$$

$$= CVar(y)C^T \quad (49)$$

$$= \sigma^2 CC^T \text{ by (36)} \quad (50)$$

$$= \sigma^2[(X^T X)^{-1} X^T + D][X(X^T X)^{-1} + D^T] \quad (51)$$

$$= \sigma^2[(X^T X)^{-1} X^T X(X^T X)^{-1} + (X^T X)^{-1}(DX)^t] \quad (52)$$

$$+ DX(X^T X)^{-1} + DD^T] \quad (53)$$

$$= \sigma^2(X^T X)^{-1} + \sigma^2 DD^T \text{ by (47)} \quad (54)$$

$$= Var(\hat{\beta}_{OLS}) + \sigma^2 DD^T \quad (55)$$

notice that $\sigma^2 > 0$ and DD^T is [definite symmetric](#).

Therefore, $\hat{\beta}_{OLS}$ has the lowest variance among all linear unbiased estimator so that it is the BLUE (Best Linear Unbiased Estimator). ■