

1 The Generalized Form of Least Squares Method

The purpose of this section is to show how we solve the OLS problem in vectors and matrices form. Firstly, we define

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} 1 & X_{11} & \cdots & X_{1p} \\ 1 & X_{21} & \cdots & X_{2p} \\ \vdots & \vdots & & \vdots \\ 1 & X_{n1} & \cdots & X_{np} \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_p \end{bmatrix}, \quad \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix} \quad (1)$$

Here, we use bold-face notation for representing vectors (lower case) and matrices (upper case), and non-bold-face notation for representing scalars.

Now, consider the goal of the least squares method, that is

$$\underset{\mathbf{w}}{\operatorname{argmin}} \boldsymbol{\epsilon}^T \boldsymbol{\epsilon} \quad (2)$$

We know that $\boldsymbol{\epsilon} = \mathbf{y} - \mathbf{X}\mathbf{w}$, and we know that matrices multiplication has the property that $(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$. Thus,

$$\begin{aligned} \boldsymbol{\epsilon}^T \boldsymbol{\epsilon} &= (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \\ &= (\mathbf{y} - \mathbf{X}\mathbf{w})^T \mathbf{y} - (\mathbf{y} - \mathbf{X}\mathbf{w})^T \mathbf{X}\mathbf{w} \\ &= \mathbf{y}^T \mathbf{y} - (\mathbf{X}\mathbf{w})^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\mathbf{w} + (\mathbf{X}\mathbf{w})^T \mathbf{X}\mathbf{w} \\ &= \mathbf{y}^T \mathbf{y} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w} \end{aligned} \quad (3)$$

Both the results of $\mathbf{w}^T \mathbf{X}^T \mathbf{y}$ and $\mathbf{y}^T \mathbf{X}\mathbf{w}$ are 1×1 scalars, and we know that for a scalar a , its transpose is equal to itself ($a^T = a$). Based on these facts, we have $\mathbf{w}^T \mathbf{X}^T \mathbf{y} = (\mathbf{w}^T \mathbf{X}^T \mathbf{y})^T = \mathbf{y}^T \mathbf{X}\mathbf{w}$. Thus,

$$\begin{aligned} \boldsymbol{\epsilon}^T \boldsymbol{\epsilon} &= (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \\ &= \mathbf{y}^T \mathbf{y} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w} \\ &= \mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w} \end{aligned} \quad (4)$$

To find the minimum value of $\boldsymbol{\epsilon}^T \boldsymbol{\epsilon}$, let $\nabla \boldsymbol{\epsilon}^T \boldsymbol{\epsilon} = 0$. Therefore,

$$\begin{aligned} \nabla \boldsymbol{\epsilon}^T \boldsymbol{\epsilon} &= \nabla (\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w}) \\ &= 2\mathbf{X}^T \mathbf{X}\mathbf{w} - 2\mathbf{X}^T \mathbf{y} = 0 \end{aligned} \quad (5)$$

Then we have

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (6)$$

Notice that the $(\mathbf{X}^T \mathbf{X})^{-1}$ exists only when $\mathbf{X}^T \mathbf{X}$ is a full rank matrix, that is the column vectors (or row vector because it is a square matrix here thus it does not matter) inside the

matrix are linear independent ¹.

2 Hypothesis Testings on Coefficients

In this section, we are focusing on answering why the coefficients follow a t-distribution. From the previous section, we know that, if the following assumptions hold,

1. The linear relationship $\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}$ exists;
2. \mathbf{X} is known and it is an invertible matrix with full rank, that is, there is no multi-collinearity among the features.
3. The noise (residual) is homoscedastic and uncorrelated, that is, $\text{Var}(\epsilon_1) = \text{Var}(\epsilon_2) = \dots = \text{Var}(\epsilon_n) = \sigma^2$ and $\forall i \neq j, \text{Cov}(\epsilon_i, \epsilon_j) = 0$
4. The noise $\boldsymbol{\epsilon}$ follows gaussian distribution: $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$

then we can fit a linear model $\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{w}}$. $\hat{\mathbf{w}}$ is estimated by least square method, and the estimation is

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{w} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon} \quad (7)$$

Now, the first thing to do is to figure out the distribution of $\hat{\mathbf{w}}$ (reflecting how it varies under repeated sampling). Since $\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}$. Also the linear combination of a gaussian distribution is also a gaussian distribution. Combining these two facts with conclusion (7), we know that $\hat{\mathbf{w}}$ also follows gaussian distribution. Furthermore, we have

$$\text{E}(\hat{\mathbf{w}}) = \text{E}(\mathbf{w}) + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{E}(\boldsymbol{\epsilon}) = \mathbf{w}^2 \quad (8)$$

$$\begin{aligned} \text{Var}(\hat{\mathbf{w}}) &= \text{Var}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}) \\ &= \text{Var}(\mathbf{X}^{-1} (\mathbf{X}^T)^{-1} \mathbf{X}^T \boldsymbol{\epsilon}) \\ &= \text{Var}(\mathbf{X}^{-1} (\mathbf{X}^{-1})^T \mathbf{X}^T \boldsymbol{\epsilon}) \\ &= \text{Var}(\mathbf{X}^{-1} (\mathbf{X} \mathbf{X}^{-1})^T \boldsymbol{\epsilon}) \\ &= \text{Var}(\mathbf{X}^{-1} \mathbf{I}^T \boldsymbol{\epsilon}) \\ &= \text{Var}(\mathbf{X}^{-1} \boldsymbol{\epsilon}) \\ &= \mathbf{X}^{-1} \text{Var}(\boldsymbol{\epsilon}) (\mathbf{X}^{-1})^T \\ &= \mathbf{X}^{-1} \sigma^2 \mathbf{I} (\mathbf{X}^{-1})^T \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \quad 3 \end{aligned} \quad (9)$$

¹ The is one of the reasons that why we need the features to meet the linear-independent condition.

² \mathbf{w} is a constant vector, thus $\text{E}(\mathbf{w}) = \mathbf{w}$. $\text{E}(\boldsymbol{\epsilon}) = 0$ base on the fourth assumption.

³ The derivative process is based on following facts: $(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$; $(\mathbf{AB})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1}$; $(\mathbf{A}^T)^{-1} = (\mathbf{A}^{-1})^T$; $\mathbf{AI} = \mathbf{A}$ where \mathbf{I} is the identity matrix; $\text{Var}(aX + b) = a^2 \text{Var}(X)$; $\text{Var}(\mathbf{Ax}) = \mathbf{A} \text{Var}(\mathbf{x}) \mathbf{A}^T$.

Thus,

$$\hat{\mathbf{w}} \sim \mathcal{N}(\mathbf{w}, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1}) \quad (10)$$

For the j^{th} item \hat{w}_j in $\hat{\mathbf{w}}$, its mean is w_j , and its variance is the j^{th} item on the diagonal of matrix $\sigma^2(\mathbf{X}^T \mathbf{X})^{-1}$. Thus,

$$\hat{w}_j \sim \mathcal{N}(w_j, \sigma^2(\mathbf{X}^T \mathbf{X})_{jj}^{-1}) \quad (11)$$

Then, we have

$$\frac{\hat{w}_j - w_j}{\sqrt{\sigma^2(\mathbf{X}^T \mathbf{X})_{jj}^{-1}}} \sim \mathcal{N}(0, 1) \quad (12)$$

Lemma 2.1. *if $X \sim \mathcal{N}(0, 1)$, and $Y \sim \chi^2(n)$, then $t = \frac{X}{\sqrt{Y/n}} \sim t(n)$, which is a t -distribution with n degree of freedom.*

Lemma 2.2. *if $Z_j \sim \mathcal{N}(0, 1)$, then $\chi^2 = Z_1^2 + Z_2^2 + \dots + Z_n^2 \sim \chi_n^2$*

Now, we are focusing on getting the unbiased estimation $\hat{\sigma}^2$ of σ^2 , and we are using the sample variance for estimating.

$$\begin{aligned} \hat{\sigma}^2 &= S_\epsilon^2 = \frac{1}{n-p-1} \sum_{j=1}^n (\epsilon_j - \bar{\epsilon})^2 \\ &= \frac{\boldsymbol{\epsilon}^T \boldsymbol{\epsilon}}{n-p-1} \quad 4 \end{aligned} \quad (13)$$

where p is the number of predictors in the regression model. Thus,

$$\boldsymbol{\epsilon}^T \boldsymbol{\epsilon} = (n-p-1)\hat{\sigma}^2 \quad (14)$$

We also know that $\epsilon_j \sim \mathcal{N}(0, \sigma^2)$. Therefore, $\frac{\epsilon_j}{\sigma} \sim \mathcal{N}(0, 1)$. Based on lemma 2.2,

$$\sum_{j=1}^n \left(\frac{\epsilon_j}{\sigma}\right)^2 = \frac{1}{\sigma^2} \boldsymbol{\epsilon}^T \boldsymbol{\epsilon} \sim \chi_{n-p-1}^2 \quad (15)$$

Then we know that

$$\frac{(n-p-1)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-p-1}^2 \quad (16)$$

⁴ Based on $\epsilon_j = y_j - \hat{y}_j$ and $E(\epsilon) = 0$.

Finally, based on conclusion (11), (12), conclusion (16), and lemma 2.1, we can prove that

$$\frac{\frac{\hat{w}_j - w_j}{\sqrt{\sigma^2(\mathbf{X}^T \mathbf{X})_{jj}^{-1}}}}{\sqrt{\frac{(n-p-1)\hat{\sigma}^2}{\sigma^2}} / (n-p-1)} = \frac{\hat{w}_j - w_j}{\hat{\sigma}\sqrt{(\mathbf{X}^T \mathbf{X})_{jj}^{-1}}} = \frac{\hat{w}_j - w_j}{\sqrt{\hat{\sigma}^2(\mathbf{X}^T \mathbf{X})_{jj}^{-1}}} \sim t_{n-p-1} \quad (17)$$

Since $\sqrt{\hat{\sigma}^2(\mathbf{X}^T \mathbf{X})_{jj}^{-1}}$ in (17) is the unbiased estimation of $\sqrt{\sigma^2(\mathbf{X}^T \mathbf{X})_{jj}^{-1}}$, which is the standard error of \hat{w}_j . Notice that $\text{SE}(\hat{w}_j)$ is also the unbiased estimation of the standard error of \hat{w}_j . Therefore,

$$\frac{\hat{w}_j - w_j}{\text{SE}(\hat{w}_j)} \sim t_{n-p-1} \quad (18)$$

3 Multilinearity

In the previous section, we have already discussed why multilinearity will break our assumptions during solving OLS. In this section, we will use another perspective, model interpretation, for exploring why we need to eliminate multilinearity.

Example 3.1. Assume

$$\begin{aligned} X_1 &: \text{TV} + 2 \times \text{Newspaper} + \text{Radio} \\ X_2 &: \text{TV} \\ X_3 &: \text{Radio} \\ X_4 &: \text{Newspaper} \\ Y &: \text{Sales} \end{aligned}$$

and we are building the regression model $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \hat{\beta}_3 X_3 + \hat{\beta}_4 X_4$. We will show that due to the existence of multilinearity, the advertisement effect from newspaper can have a positive influence on sales even if $\hat{\beta}_4 < 0$.

$$\begin{aligned} \hat{Y} &= \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \hat{\beta}_3 X_3 + \hat{\beta}_4 X_4 \\ &= \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \hat{\beta}_3 X_3 + \hat{\beta}_4 X_4 + X_2 - X_2 + 2X_4 - 2X_4 + X_3 - X_3 \\ &= \hat{\beta}_0 + (\hat{\beta}_1 - 1)X_1 + (\hat{\beta}_2 - 1)X_2 + (\hat{\beta}_3 - 1)X_3 + (\hat{\beta}_4 - 2)X_4 \end{aligned} \quad (19)$$

Let $\hat{\beta}_0' = \hat{\beta}_0$, $\hat{\beta}_1' = \hat{\beta}_1 - 1$, $\hat{\beta}_2' = \hat{\beta}_2 - 1$, $\hat{\beta}_3' = \hat{\beta}_3 - 1$, $\hat{\beta}_4' = \hat{\beta}_4 - 2$. Thus, if the final model is

$$\hat{Y} = \hat{\beta}_0' + \hat{\beta}_1' X_1 + \hat{\beta}_2' X_2 + \hat{\beta}_3' X_3 + \hat{\beta}_4' X_4 \quad (20)$$

and $\hat{\beta}_4' < 0$, we cannot say the level of the advertisement effect from newspaper increases will make sales decrease based on $\hat{\beta}_4' < 0$ since $\hat{\beta}_4$ still have a chance to be a value that is greater than zero (e.g. $\hat{\beta}_4 = 1$).

4 References

- [1] J.W, Li, Regression Based on Least Square: The Matrices Prospective, May. 13, 2018. Accessed on: <https://zhuanlan.zhihu.com/p/33899560>
- [2] J.W, Li, The Hypothesis Testing for Linear Regression: The T test and the F test, May. 14, 2018. Accessed on: <https://zhuanlan.zhihu.com/p/36782834>