

### V-3-4. Optimization and pseudo-inverse matrix

In multiple regression analysis, we estimate optimal coefficient of variables to explain explained variables from data sets. When we express explained variable as  $\mathbf{y}$  and explanatory variables as  $\mathbf{x}$ , example of data set is as follow.

$$\begin{aligned} & y_1, x_{11}, x_{21} \cdots x_{p1} \\ & y_2, x_{12}, x_{22} \cdots x_{p2} \\ & y_3, x_{13}, x_{23} \cdots x_{p3} \\ & \vdots \\ & y_m, x_{1m}, x_{2m} \cdots x_{pm} \end{aligned}$$

We are expecting to write the relation as follow.

$$\begin{aligned} y_1 &= a_1 x_{11} + a_2 x_{21} + \cdots + a_p x_{p1} \\ y_2 &= a_1 x_{12} + a_2 x_{22} + \cdots + a_p x_{p2} \\ y_3 &= a_1 x_{13} + a_2 x_{23} + \cdots + a_p x_{p3} \\ & \vdots \\ y_m &= a_1 x_{1m} + a_2 x_{2m} + \cdots + a_p x_{pm} \end{aligned}$$

However, it is generally impossible because explained variable data includes impacts of unknown variables, so we have to add error term as follow.

$$\begin{aligned} y_1 &= a_1 x_{11} + a_2 x_{21} + \cdots + a_p x_{p1} + e_1 \\ y_2 &= a_1 x_{12} + a_2 x_{22} + \cdots + a_p x_{p2} + e_2 \\ y_3 &= a_1 x_{13} + a_2 x_{23} + \cdots + a_p x_{p3} + e_3 \\ & \vdots \\ y_m &= a_1 x_{1m} + a_2 x_{2m} + \cdots + a_p x_{pm} + e_m \end{aligned}$$

Multiple regression analysis is optimization of error term. There several ideas of the optimization, Least square method is one of them. The idea of least square method is minimizing of error term.

When we express the relation using matrix.

$$\mathbf{E} = \mathbf{Y} - \mathbf{XA}$$

$$\mathbf{E} = \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_m \end{pmatrix}, \mathbf{Y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}, \mathbf{X} = \begin{pmatrix} x_{11} & x_{21} & \cdots & x_{p1} \\ x_{12} & x_{22} & \cdots & x_{p2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1m} & x_{2m} & \cdots & x_{pm} \end{pmatrix}, \quad \mathbf{A} = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_p \end{pmatrix}$$

we consider minimizing of magnitude of  $\mathbf{E}$ . In this case, magnitude of  $\mathbf{E}$  is not the sum of error term  $\sum_{i=1}^m e_i$ . The magnitude of vector and matrix is norm.

The definition of norm is as follow

$$\mathbf{X} = (x_1, x_2, \cdots, x_n), \quad 1 \leq q \leq \infty$$

$\sqrt[q]{|x_1|^q + |x_2|^q + \dots + |x_n|^q}$  is  $L^q$  norm and is expressed as  $\|\mathbf{X}\|_q$

$$\begin{aligned}\|\mathbf{X}\|_1 &= |x_1| + |x_2| + \dots + |x_n| \\ \|\mathbf{X}\|_2 &= \sqrt{|x_1|^2 + |x_2|^2 + \dots + |x_n|^2}\end{aligned}$$

$L^2$  norm is Euclid distance. Least square method is minimization of  $L^2$  norm.

$$\begin{aligned}\|\mathbf{E}\|_2 &= \sqrt{|e_1|^2 + |e_2|^2 + \dots + |e_m|^2} \\ &= \sqrt{e_1^2 + e_2^2 + \dots + e_m^2} \\ \|\mathbf{E}\|_2^2 &= e_1^2 + e_2^2 + \dots + e_m^2 \\ &= (e_1, e_2, \dots, e_m) \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_m \end{pmatrix} \\ \mathbf{E}^T &= (e_1, e_2, \dots, e_m)\end{aligned}$$

o

$$\begin{aligned}\mathbf{E} &= \mathbf{Y} - \mathbf{XA} \\ \mathbf{E}^T &= (\mathbf{Y} - \mathbf{XA})^T \\ \|\mathbf{E}\|_2^2 &= (\mathbf{Y} - \mathbf{XA})^T(\mathbf{Y} - \mathbf{XA}) \\ &= \mathbf{Y}^T\mathbf{Y} - \mathbf{A}^T\mathbf{X}^T\mathbf{Y} - \mathbf{Y}^T\mathbf{XA} + \mathbf{X}^T\mathbf{A}^T\mathbf{XA} \\ &= \mathbf{Y}^T\mathbf{Y} - 2\mathbf{Y}^T\mathbf{XA} + \mathbf{X}^T\mathbf{A}^T\mathbf{XA}\end{aligned}$$

$\|\mathbf{E}\|_2^2$  is positive definite and we can obtain minimum value by following differential.

$$\begin{aligned}\frac{\partial \|\mathbf{E}\|_2^2}{\partial \mathbf{A}} &= 0 \\ \frac{\partial \|\mathbf{E}\|_2^2}{\partial \mathbf{A}} &= \frac{\partial (\mathbf{Y}^T\mathbf{Y} - 2\mathbf{Y}^T\mathbf{XA} + \mathbf{X}^T\mathbf{A}^T\mathbf{XA})}{\partial \mathbf{A}} \\ &= \frac{\partial (\mathbf{Y}^T\mathbf{Y})}{\partial \mathbf{A}} - 2 \frac{\partial (\mathbf{Y}^T\mathbf{XA})}{\partial \mathbf{A}} + \frac{\partial (\mathbf{X}^T\mathbf{A}^T\mathbf{XA})}{\partial \mathbf{A}} = \mathbf{0} \\ \frac{\partial (\mathbf{Y}^T\mathbf{Y})}{\partial \mathbf{A}} &= \mathbf{0} \\ \frac{\partial (\mathbf{Y}^T\mathbf{XA})}{\partial \mathbf{A}} &= \mathbf{Y}^T\mathbf{X} \\ \frac{\partial (\mathbf{X}^T\mathbf{A}^T\mathbf{XA})}{\partial \mathbf{A}} &= 2\mathbf{X}^T\mathbf{XA} \\ \frac{\partial \|\mathbf{E}\|_2^2}{\partial \mathbf{A}} &= -2\mathbf{Y}^T\mathbf{X} + 2\mathbf{X}^T\mathbf{XA} = \mathbf{0} \\ \mathbf{Y}^T\mathbf{X} &= \mathbf{X}^T\mathbf{XA} \\ \mathbf{A} &= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{Y}^T\mathbf{X} \\ \mathbf{A} &= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}\end{aligned}$$

We consider  $(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$  is a matrix

$$(X^T X)^{-1} X^T = X^\#$$

In case of  $\mathbf{Y} = \mathbf{XA}$ , We can obtain  $\mathbf{A}$  by multiplying  $\mathbf{X}^{-1}$  from left

$$\mathbf{X}^{-1} \mathbf{Y} = \mathbf{X}^{-1} \mathbf{XA} = \mathbf{A}$$

Similarly, we can obtain  $\mathbf{A}$  by multiplying  $\mathbf{X}^\#$

$$\begin{aligned} \mathbf{X}^\# \mathbf{Y} &= \mathbf{A} \\ \mathbf{X}^\# &= (X^T X)^{-1} X^T \end{aligned}$$

Formula 71

$$\mathbf{X}^\# \mathbf{X} = (X^T X)^{-1} (X^T X) = \mathbf{I}$$

The function of  $\mathbf{X}^\#$  is analogous to inverse matrix, so we call this matrix pseudo-inverse matrix. Generalized inverse matrix is the other name of this matrix. However,  $\mathbf{X}^\#$  gives least squares solution, when  $\mathbf{X}$  is non-regular matrix. When  $\mathbf{X}$  is regular matrix,  $\mathbf{X}^\# = \mathbf{X}^{-1}$ , and  $\mathbf{X}^\#$  gives solution of simultaneous equation.

For an exercise, we us this method for simple linear regression

$$\begin{aligned} y &= ax + b \\ \mathbf{Y} &= \mathbf{XA} \\ \mathbf{Y} &= \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \mathbf{X} = \begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{pmatrix}, \quad \mathbf{A} = \begin{pmatrix} a \\ b \end{pmatrix} \\ \mathbf{X}^T &= \begin{pmatrix} x_1 & x_2 & \cdots & x_n \\ 1 & 1 & \cdots & 1 \end{pmatrix} \\ \mathbf{X}^T \mathbf{X} &= \begin{pmatrix} x_1 & x_2 & \cdots & x_n \\ 1 & 1 & \cdots & 1 \end{pmatrix} \begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^n x_i^2 & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & n \end{pmatrix} \\ (\mathbf{X}^T \mathbf{X})^{-1} &= \begin{pmatrix} \sum_{i=1}^n x_i^2 & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & n \end{pmatrix}^{-1} = \frac{1}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \begin{pmatrix} n & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{pmatrix} \end{aligned}$$

$$\begin{aligned}
X^\# &= (X^T X)^{-1} X^T = \frac{1}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \begin{pmatrix} n & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{pmatrix} \begin{pmatrix} x_1 & x_2 & \cdots & x_n \\ 1 & 1 & \cdots & 1 \end{pmatrix} \\
&= \frac{1}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \begin{pmatrix} nx_1 - \sum_{i=1}^n x_i & nx_2 - \sum_{i=1}^n x_i & \cdots & nx_n - \sum_{i=1}^n x_i \\ -x_1 \sum_{i=1}^n x_i + \sum_{i=1}^n x_i^2 & -x_2 \sum_{i=1}^n x_i + \sum_{i=1}^n x_i^2 & \cdots & x_n \sum_{i=1}^n x_i + \sum_{i=1}^n x_i^2 \end{pmatrix} \\
X^\# Y &= \frac{1}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \begin{pmatrix} nx_1 - \sum_{i=1}^n x_i & nx_2 - \sum_{i=1}^n x_i & \cdots & nx_n - \sum_{i=1}^n x_i \\ -x_1 \sum_{i=1}^n x_i + \sum_{i=1}^n x_i^2 & -x_2 \sum_{i=1}^n x_i + \sum_{i=1}^n x_i^2 & \cdots & x_n \sum_{i=1}^n x_i + \sum_{i=1}^n x_i^2 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \\
&= \frac{1}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \begin{pmatrix} n \sum_{i=1}^n x_i y_i + \sum_{i=1}^n x_i \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \sum_{i=1}^n x_i + \sum_{i=1}^n y_i \sum_{i=1}^n x_i^2 \end{pmatrix} = \begin{pmatrix} a \\ b \end{pmatrix} \\
a &= \frac{n \sum_{i=1}^n x_i y_i + \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \\
b &= \frac{\sum_{i=1}^n x_i y_i \sum_{i=1}^n x_i + \sum_{i=1}^n y_i \sum_{i=1}^n x_i^2}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}
\end{aligned}$$

$a$  is regression coefficient