

Linear Regression Mathematics

Charlie Clark

January 6, 2022

1 Simple Linear Regression

Model:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i,$$

where $1 \leq i \leq N$ and $N \geq 1$ is the sample size.

Ordinary Least Squares Cost Function:

$$J(\hat{\beta}_0, \hat{\beta}_1) = \sum_{i=1}^N (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

Solve for Coefficients:

$$\min_{(\hat{\beta}_0, \hat{\beta}_1)} J(\hat{\beta}_0, \hat{\beta}_1) = \min_{(\hat{\beta}_0, \hat{\beta}_1)} \sum_{i=1}^N (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \quad \because \text{Minimize cost function.}$$

$$\Rightarrow \frac{\partial J}{\partial \hat{\beta}_0} = \frac{\partial}{\partial \hat{\beta}_0} \sum_{i=1}^N (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \quad \because \text{Partial with respect to } \hat{\beta}_0.$$

$$= -2 \sum_{i=1}^N (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \quad \because \text{Chain rule.}$$

$$= -2 \left(\sum_{i=1}^N y_i - N\hat{\beta}_0 - \hat{\beta}_1 \sum_{i=1}^N x_i \right) \quad \because \text{Distribute summation.}$$

$$= 0 \quad \because \frac{\partial J}{\partial \hat{\beta}_0} = 0 \text{ at minimum.}$$

$$\Leftrightarrow 0 = -2 \left(\sum_{i=1}^N y_i - N\hat{\beta}_0 - \hat{\beta}_1 \sum_{i=1}^N x_i \right) \quad \because \text{Rewrite.}$$

$$= \sum_{i=1}^N y_i - N\hat{\beta}_0 - \hat{\beta}_1 \sum_{i=1}^N x_i \quad \because \text{Divide by } -2.$$

$$\Leftrightarrow N\hat{\beta}_0 = \sum_{i=1}^N y_i - \hat{\beta}_1 \sum_{i=1}^N x_i \quad \because \text{Add } N\hat{\beta}_0.$$

$$\Leftrightarrow \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \because \text{Divide by } N, \text{ substitute averages.}$$

$$\begin{aligned}
\Rightarrow \frac{\partial J}{\partial \hat{\beta}_1} &= \frac{\partial}{\partial \hat{\beta}_1} \sum_{i=1}^N \left(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \right)^2 && \because \text{Partial with respect to } \hat{\beta}_1. \\
&= -2 \sum_{i=1}^N x_i \left(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \right) && \because \text{Chain rule.} \\
&= -2 \sum_{i=1}^N \left(y_i x_i - \hat{\beta}_0 x_i - \hat{\beta}_1 x_i^2 \right) && \because \text{Distribute } x_i. \\
&= -2 \left(\sum_{i=1}^N y_i x_i - \hat{\beta}_0 \sum_{i=1}^N x_i - \hat{\beta}_1 \sum_{i=1}^N x_i^2 \right) && \because \text{Distribute summation.} \\
&= 0 && \because \frac{\partial J}{\partial \hat{\beta}_1} = 0 \text{ at minimum.} \\
\Leftrightarrow 0 &= -2 \left(\sum_{i=1}^N y_i x_i - \hat{\beta}_0 \sum_{i=1}^N x_i - \hat{\beta}_1 \sum_{i=1}^N x_i^2 \right) && \because \text{Rewrite.} \\
&= \sum_{i=1}^N y_i x_i - \hat{\beta}_0 \sum_{i=1}^N x_i - \hat{\beta}_1 \sum_{i=1}^N x_i^2 && \because \text{Divide by } -2. \\
&= \overline{yx} - \hat{\beta}_0 \bar{x} - \hat{\beta}_1 \overline{x^2} && \because \text{Divide by } N. \\
&= \overline{yx} - \left(\bar{y} - \hat{\beta}_1 \bar{x} \right) \bar{x} - \hat{\beta}_1 \overline{x^2} && \because \text{Substitute } \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}. \\
&= \overline{yx} - \bar{y} \cdot \bar{x} - \hat{\beta}_1 \overline{x^2} + \hat{\beta}_1 \bar{x}^2 \\
&= \overline{yx} - \bar{y} \cdot \bar{x} - \hat{\beta}_1 \left(\overline{x^2} - \bar{x}^2 \right) && \because \text{Factor } -\hat{\beta}_1. \\
\Leftrightarrow \hat{\beta}_1 \left(\overline{x^2} - \bar{x}^2 \right) &= \overline{yx} - \bar{y} \cdot \bar{x} && \because \text{Add } \hat{\beta}_1 \left(\overline{x^2} - \bar{x}^2 \right). \\
\Leftrightarrow \hat{\beta}_1 &= \frac{\overline{yx} - \bar{y} \cdot \bar{x}}{\overline{x^2} - \bar{x}^2} && \because \text{Divide by } \overline{x^2} - \bar{x}^2.
\end{aligned}$$

Therefore,

$$\hat{\beta}_1 = \frac{\overline{yx} - \bar{y} \cdot \bar{x}}{\overline{x^2} - \bar{x}^2} \text{ and } \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x},$$

where

$$\begin{aligned}
\bar{y} &= \frac{1}{N} \sum_{i=1}^N y_i, \\
\bar{x} &= \frac{1}{N} \sum_{i=1}^N x_i, \\
\overline{yx} &= \frac{1}{N} \sum_{i=1}^N y_i x_i, \\
\text{and } \overline{x^2} &= \frac{1}{N} \sum_{i=1}^N x_i^2.
\end{aligned}$$

2 Multiple Linear Regression

Model:

$$y_i = \beta_0 + \beta_1 x_{i,1} + \cdots + \beta_K x_{i,K} + \varepsilon_i,$$

where $1 \leq i \leq N$, $N \geq 1$ is the sample size, and $K \geq 2$ is the number of coefficients -1 .

Model in Matrix Form:

$$\vec{y} = \mathbf{X}\vec{\beta} + \vec{\varepsilon},$$

where

$$\vec{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 1 & x_{1,1} & \cdots & x_{1,K} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N,1} & \cdots & x_{N,K} \end{bmatrix}, \vec{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_K \end{bmatrix}, \text{ and } \vec{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_N \end{bmatrix}.$$

Ordinary Least Squares Cost Function:

$$J(\vec{\beta}) = (\vec{y} - \mathbf{X}\vec{\beta})^T (\vec{y} - \mathbf{X}\vec{\beta}) = \vec{y}^T \vec{y} - \vec{y}^T \mathbf{X} \vec{\beta} - \vec{\beta}^T \mathbf{X}^T \vec{y} + \vec{\beta}^T \mathbf{X}^T \mathbf{X} \vec{\beta}.$$

Solve for Coefficients:

$$\min_{(\vec{\beta})} J(\vec{\beta}) \Rightarrow \frac{\partial J}{\partial \vec{\beta}} = 0 \quad \because \text{Minimum when slop is 0.}$$

$$\Leftrightarrow \frac{\partial}{\partial \vec{\beta}} (\vec{y}^T \vec{y} - \vec{y}^T \mathbf{X} \vec{\beta} - \vec{\beta}^T \mathbf{X}^T \vec{y} + \vec{\beta}^T \mathbf{X}^T \mathbf{X} \vec{\beta}) = 0 \quad \because \text{Definition of } J(\vec{\beta})$$

$$\Leftrightarrow 2\mathbf{X}^T \mathbf{X} \vec{\beta} - 2\mathbf{X}^T \vec{y} = 0 \quad \because \text{Simplify.}$$

$$\Leftrightarrow 2\mathbf{X}^T \mathbf{X} \vec{\beta} = 2\mathbf{X}^T \vec{y} \quad \because \text{Add } 2\mathbf{X}^T \vec{y}.$$

$$\Leftrightarrow \mathbf{X}^T \mathbf{X} \vec{\beta} = \mathbf{X}^T \vec{y} \quad \because \text{Divide by 2.}$$

$$\Leftrightarrow (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}) \vec{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{y} \quad \because \text{Multiply by } (\mathbf{X}^T \mathbf{X})^{-1}.$$

$$\Leftrightarrow \vec{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{y} \quad \because (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}) = \mathbf{I}, \mathbf{I} \vec{\beta} = \vec{\beta}.$$

Therefore,

$$\vec{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{y},$$

where

$$\vec{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_K \end{bmatrix}.$$