

Normal Equation Solution

Closed-form solution for linear regression

Matrix Form

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

\mathbf{X} is the design matrix



Normal Equation

$$\mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^T \mathbf{Y}$$



General Solution

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

When $\mathbf{X}^T \mathbf{X}$ is invertible

Simple Regression

$$\hat{\beta}_1 = \text{Cov}(\mathbf{X}, \mathbf{Y}) / \text{Var}(\mathbf{X})$$

$$\hat{\boldsymbol{\beta}} = \bar{\mathbf{Y}} - \hat{\beta}_1 \mathbf{X}^T$$

Complexity

Computational cost

$$O(n p^2 + p^3)$$

n: samples

p: features

Alternative

For large-scale:

Gradient Descent

(iterative method)

Advantage

Closed-form

Direct computation

No iterations needed

Gradient Descent Process

Parameter update with MSE Loss - Step by Step

Mean Squared Error (MSE) Loss

$$L(\beta) = (1/n) \sum (y_i - \hat{y}_i)^2 = (1/n) \sum (y_i - \beta_0 - \beta_1 x_i)^2$$



Compute Partial Derivatives

1

$$\frac{\partial L}{\partial \beta_0} = -(2/n) \sum (y_i - \beta_0 - \beta_1 x_i)$$

$$\frac{\partial L}{\partial \beta_1} = -(2/n) \sum x_i (y_i - \beta_0 - \beta_1 x_i)$$

Calculate gradient for each parameter

Evaluate at Current Parameters

2

$$g_0 = \frac{\partial L}{\partial \beta_0} |_{\beta^{(t)}}, \quad g_1 = \frac{\partial L}{\partial \beta_1} |_{\beta^{(t)}}$$

Substitute current $\beta_0^{(t)}$ and $\beta_1^{(t)}$ values

Update Parameters

3

$$\beta_0^{(t+1)} = \beta_0^{(t)} - \alpha \cdot g_0$$

$$\beta_1^{(t+1)} = \beta_1^{(t)} - \alpha \cdot g_1$$

α is the learning rate (step size)

Check Convergence

4

$$|L^{(t+1)} - L^{(t)}| < \varepsilon \text{ or } ||\nabla L|| < \varepsilon$$

Repeat steps 1-3 until convergence criterion is met

Learning Rate α

Controls step size

Too large: diverge

Too small: slow

Matrix Form

$$\beta^{(t+1)} = \beta^{(t)} - \alpha \cdot \nabla L$$

$$\nabla L = -(2/n) X^T (Y - X\beta)$$

Convergence

Iterative approach

Works for large-scale

$O(np)$ per iteration