

Multivariable Linear Regression

X : Input matrix of shape (n, d) where

→ n = number of samples

→ d = number of features (columns)

→ θ = weight vector (parameter) of shape $(d, 1)$
these are the coefficient need to learn.

→ y = True output vector (target values) of shape $(n, 1)$

Linear Model Prediction:

$$\hat{y} = X\theta \quad \hat{y} = \text{predicted output vector of shape } (n, 1)$$

For each sample i , $\hat{y}_i = \sum_{j=1}^d X_{ij} \cdot \theta_j$

Define the Loss function (MSE)

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$= \frac{1}{n} e^T e$$

$$= \frac{1}{n} (y - X\theta)^T (y - X\theta)$$

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - x_i \theta)^2 = \boxed{\frac{1}{n} \sum_{i=1}^n e_i^2}$$

Error vector

$$e = y - \hat{y} = y - X\theta$$

↑
shape $(n, 1)$

This is just summing up squares of error for each sample.

what does $e^T e$ means

$$e = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} (n \times 1)$$

$$e^T = [e_1, e_2, \dots, e_n] (1 \times n)$$

$$e^T \cdot e = [e_1 \ e_2 \ \dots \ e_n] \cdot \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$$

$$= e_1^2 + e_2^2 + \dots + e_n^2$$

$$e^T e = \sum_{i=1}^n e_i^2$$

$$e^T e = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$e = y - x\theta$$

vector for loss:

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{n} e^T e = \frac{1}{n} (y - x\theta)^T (y - x\theta)$$

Our goal is to minimize $L(\theta)$ with respect to θ

→ Find the gradient of $L(\theta)$ with respect to θ

$$\nabla_{\theta} L = \frac{\partial L}{\partial \theta}$$

To find the derivative $L(\theta)$ we have to simplify $L(\theta)$

$$L(\theta) = \frac{1}{n} (y - x\theta)^T (y - x\theta)$$

$$= \frac{1}{n} (y^T - \theta^T x^T) (y - x\theta) \quad \begin{bmatrix} (a-b)^T = a^T - b^T \\ (ab)^T = b^T a^T \end{bmatrix}$$

$$= \frac{1}{n} (y^T y - y^T x\theta - \theta^T x^T y + \theta^T x^T x\theta)$$

$$L(\theta) = \frac{1}{n} (y^T y - 2\theta^T x^T y + \theta^T x^T x\theta)$$

$$\begin{bmatrix} y^T x\theta = \theta^T x^T y \\ \text{we will see it later} \end{bmatrix}$$

Final simplified loss function

Take the derivative of cost function:

① $\frac{\partial}{\partial \theta} (y^T y) = 0 \rightarrow$ doesn't depend on θ . so constant $\left[\frac{dc}{d\theta} = 0 \right]$

② derivative of $-2\theta^T X^T y$

$$\frac{\partial}{\partial \theta} (-2\theta^T X^T y) = -2X^T y \quad \left[\begin{array}{l} \text{vector calculus} \\ \frac{d}{dA} A^T = 1 \end{array} \right]$$

③ derivative of $\theta^T X^T X \theta$

$$\frac{\partial}{\partial \theta} (\theta^T X^T X \theta) = \frac{\partial}{\partial \theta} (X^T X \theta^T \theta) = 2X^T X \theta \quad \left[\frac{d}{dA} (A^T A) = 2A \right]$$

∴ Finally after combining all the derivatives

$$\nabla_{\theta} J = \frac{\partial J}{\partial \theta} = (-2X^T y + 2X^T X \theta) \times \frac{1}{n}$$

$$\Rightarrow \frac{\partial J}{\partial \theta} = \frac{-2}{n} (X^T y - X^T X \theta) \quad \left\{ \begin{array}{l} \text{both are okay} \end{array} \right.$$
$$= \frac{2}{n} (X^T X \theta - X^T y)$$

Express the gradient in terms of vector

$$\frac{\partial J}{\partial \theta} = \frac{-2}{n} (X^T y - X^T X \theta)$$

$$= \frac{-2}{n} \cdot X^T (y - X \theta)$$

$$= \frac{-2}{n} \cdot X^T \cdot e \quad [e = y - X \theta]$$

Final gradient $\boxed{\nabla_{\theta} (J(\theta)) = \frac{-2}{n} X^T \cdot e}$

Proof $y^T x \theta = \theta^T x^T y$ basically we have to show that they are symmetric.

y shape $(n \times 1)$
 y^T shape $(1 \times n)$
 x " $(n \times d)$
 x^T " $(d \times n)$
 θ " $(d \times 1)$
 $\theta^T \Rightarrow (1 \times d)$

$$\text{LHS} = y^T \cdot x \cdot \theta$$

$$= (1 \times n) (n \times d) (d \times 1)$$

$$= (1 \times d) (d \times 1)$$

$$= (1 \times 1) \rightarrow \text{scalar (a single value)}$$

$$\text{RHS} = \theta^T x^T y$$

$$= (1 \times d) (d \times n) (n \times 1)$$

$$= (1 \times n) (n \times 1)$$

$$= (1 \times 1) \rightarrow \text{again a single value}$$

$$\boxed{\text{RHS} = \text{LHS (proved)}}$$

gradient algo steps:

(i) take any values of theta

(ii) Do iterations n times

(iii) At each iteration Find the gradient and update theta

$$\theta := \theta - \alpha \nabla_{\theta} J(\theta)$$

learning rate

a very small number to make change of θ stable.

Till now we calculated the gradient descent without the bias term to make things simple.

now we will just add the bias term and we will find the gradient with respect to theta and bias. But guess we already calculated the gradient with respect to theta.

now prediction, $\hat{y} = x\theta + b$
 \uparrow bias term (intercept)

$$\begin{aligned}\text{Loss function, } L(\theta, b) &= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (y_i - (x_i\theta + b))^2 \\ &= \frac{1}{n} e^T \cdot e \quad [e = y - x\theta - b] \\ &= \frac{1}{n} (y - x\theta - b)^T (y - x\theta - b)\end{aligned}$$

gradient w.r.t. θ

$$\nabla_{\theta} L = \frac{\partial L}{\partial \theta} = -\frac{2}{n} x^T e$$

$$\boxed{\frac{\partial L}{\partial \theta} = -\frac{2}{n} x^T (y - x\theta - b)}$$

gradient with respect to b

$$\begin{aligned}\frac{\partial L}{\partial b} &= \frac{\partial}{\partial b} \left[\frac{1}{n} \sum_{i=1}^n (y_i - x_i\theta - b)^2 \right] \\ &= \frac{2}{n} \sum_{i=1}^n (y_i - x_i\theta - b) \cdot (-1) \\ &= -\frac{2}{n} \sum_{i=1}^n e_i = -\frac{2}{n} \sum_{i=1}^n (y - x\theta - b)\end{aligned}$$