

Simple Linear Regression

It is a supervised machine learning algorithm used to predict a continuous target variable based on the values of independent variable. The main goal is to find the best fit line that represents the relationship between the input and the output data.

Best fit line

- ↳ Best captures the overall trend of the data.
- ↳ Minimizes the error between actual and predicted value.
- ↳ helps to predict future outcomes for unseen value of x .

The general form of the line is :

$$\hat{y} = w_0 + w_1 x$$

where,

\hat{y} : Predicted value

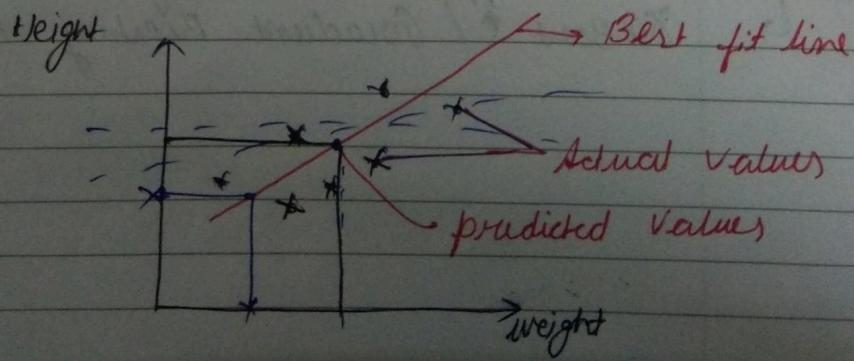
w_0 : Intercept (bias)

w_1 : slope (weight)

x : independent variable

Actual value (y): The real output from the dataset.

Predicted value (\hat{y}): The value predicted by the regression model.



As we know the equation of some simpler linear regression model is

$$\hat{y} = w_0 + w_1 x$$

we have to calculate w_0 and w_1 .

Eg:- Suppose we have a dataset and we want to predict a person height (in cm) based on their weight (in kg).

Person	Weight x	Height y
A	50	160
B	60	163
C	70	175
D	80	170
E	90	180

Suppose after training the model, we get

$$\begin{aligned}\hat{y} &= 145 + 0.4x \\ w_0 &= 145 \text{ (Intercept)} \\ w_1 &= 0.4 \text{ (slope)}\end{aligned}$$

Weight(x)	Actual Height(y)	Predicted Height(\hat{y})	Residual($y - \hat{y}$)
50	160	165.0	-5.0
60	163	169.0	-6.0
70	175	179.0	-4.0
80	170	177.0	-7.0
90	180	181.0	-1.0

Negative Residual \rightarrow Prediction too high
Positive Residual \rightarrow Prediction too low.

* Mean Squared Error (MSE)

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = 23.0$$

so, the average squared error made by the model on the dataset is 23.0 cm^2

Now, there are two ways to find the value of w_0 and w_1 .

↪ y_i

Now, we want the best fit-line $\hat{y} = w_0 + w_1 x$ that minimizes the prediction error b/w actual value y_i and predicted value \hat{y}_i .

∴ we use Mean Squared Error (MSE) as the cost function:

$$J(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

cost function

Now, there are two ways to minimize the cost function.

- ↪ ↪ Gradient Descent (Iterative)
- ↪ ↪ Normal equation (closed form)

1) Gradient descent (convergence algorithm)

Gradient Descent is an optimization algorithm used to minimize the cost function in machine learning and deep learning models.

Intuition behind Gradient descent

Imagine you are on a mountain and it's dark.

You want to reach the lowest point (valley) using just your sense of slope under your feet.

→ Each step you take is in the opp. direction of the slope (gradient)

→ The steepness of the slope tells you how far to move

→ you repeat until you reach a flat point (minimum)

Mathematical Formulation

Let's say we are trying to minimize a cost function $J(\theta)$, where θ are the parameters (like weights in linear regression).

We update the parameters using:

$$\theta := \theta - \alpha \nabla J(\theta) \rightarrow \text{update formula}$$

where;

→ θ : Parameters (like weights w_0, w_1, \dots)

→ α : Learning rate (step size)

→ $\nabla J(\theta)$: Gradient (partial derivative of cost fn w.r.t θ)

Proof of the formula using example

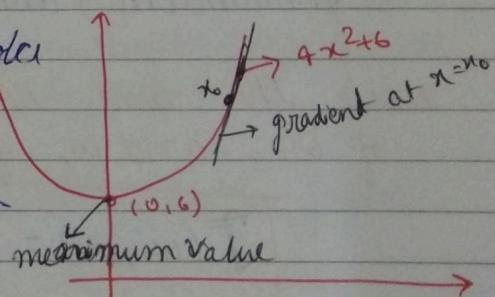
We are going to use gradient descent to minimize the function $f(x) = 4x^2 + b$.

Since shape \Rightarrow U-shaped parabola

minimum \Rightarrow at $x=0$

value at minimum $\Rightarrow b$

derivative slope: $f'(x) \Rightarrow 8x$



Now, we use the formula

$$x_{\text{new}} = x - \alpha f'(x)$$

Iteration 1:

Let's take $x=3$

$$\therefore \text{slope} = 8 \times 3 = 24$$

$$\therefore x_1 = 3 - 0.1 \times 24 = 0.6$$

$$\therefore \text{New value} = f(0.6) = 7.44$$

Iteration 2

$$\rightarrow x = 0.6$$

$$\rightarrow \text{slope} = 0.6 \times 8 = 4.8$$

$$\rightarrow x_2 = 0.6 - 0.1 \times 4.8 = 0.12$$

$$f(x_2) = f(0.12) = 6.0576$$

Iteration 3

$$\rightarrow x = 0.12$$

$$\rightarrow \text{slope} = 0.12 \times 8 = 0.96$$

$$\rightarrow x_3 = 0.12 - 0.1 \times 0.96 \approx 6.002304$$

NOTE

1) α is too small \Rightarrow slow convergence (takes forever)

2) α is too large \Rightarrow overshoot minimum \rightarrow diverge

3) α is just right \Rightarrow smooth convergence \Rightarrow find minimum

Result

After step 3

$\rightarrow n$ changed: 3 \rightarrow 0.6 \rightarrow 0.12 \rightarrow 0.024

$\rightarrow f(n)$ changed: 42 \rightarrow 7.44 \rightarrow 6.0576 \rightarrow 6.002304

hence as $x \rightarrow 0$ $f(n) \rightarrow 6$

\therefore the cost function $f(n)$ will be 6 which is minimum and it occurs at $x=0$.

Now, the same we will find the minimum value of w_0 , and w_1 for the cost J

$$J(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

where $\hat{y}_i = w_0 + w_1 x_i$

Now, Gradient w.r.t w_0

$$\frac{\partial J}{\partial w_0} = -\frac{2}{n} \sum (y_i - \hat{y}_i)$$

Gradient w.r.t w_1

$$\frac{\partial J}{\partial w_1} = -\frac{2}{n} \sum (y_i - \hat{y}_i) \cdot x_i$$

Instead of keeping the $\frac{1}{2}$ in the derivative, we absorb it into the learning rate α . This amplifies the gradients to:

$$\frac{\partial J}{\partial w_0} = -\frac{1}{n} \sum (y_i - \hat{y}_i)$$

$$\frac{\partial J}{\partial w_i} = -\frac{1}{n} \sum (y_i - \hat{y}_i) \cdot n_i$$

* Update rule

$$w_0 = w_0 + \alpha \cdot \frac{1}{n} \sum (y_i - \hat{y}_i)$$

$$w_i = w_i + \alpha \cdot \frac{1}{n} \sum (y_i - \hat{y}_i) \cdot n_i$$

Repeat until convergence like we did in example.

#2 > Normal Equation (closed form)

A closed form solution means a direct mathematical formula to compute the model parameters (weights) in one-step - no need for loop or iteration.

* Setup of linear regression in vector/matrix form.

Suppose we are given

- Input matrix $X \in R^{m \times n}$

Each row = one sample

Each column = one feature

(usually, a column of 1s is added to include the bias / intercept)

- Output vector $y \in R^m$

- weights vector $w \in R^n$

* Model Hypothesis

$$\hat{y} = \mathbf{x}\mathbf{w}$$

The cost function will be :

$$J(\mathbf{w}) = \frac{1}{2m} \| \mathbf{y} - \mathbf{x}\mathbf{w} \|^2$$

As we know

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

Now, let say $\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} \Rightarrow \| \mathbf{a} \|^2 = \mathbf{a}^\top \mathbf{a} = a_1^2 + a_2^2 + a_3^2$

So, when

$$\mathbf{y} - \mathbf{x}\mathbf{w} = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_m \end{bmatrix} \Rightarrow \| \mathbf{y} - \mathbf{x}\mathbf{w} \|^2 = \sum_{i=1}^m e_i^2$$

This exactly mean

$$(\mathbf{y} - \mathbf{x}\mathbf{w})^\top (\mathbf{y} - \mathbf{x}\mathbf{w}) = e_1^2 + e_2^2 + e_3^2 + \dots + e_m^2$$

$\therefore J(\mathbf{w})$ becomes $\frac{1}{2m} (\mathbf{y} - \mathbf{x}\mathbf{w})^\top (\mathbf{y} - \mathbf{x}\mathbf{w})$

$$\Rightarrow J(\mathbf{w}) = \frac{1}{2m} [\mathbf{y}^\top \mathbf{y} - 2\mathbf{w}^\top \mathbf{x}^\top \mathbf{y} + \mathbf{w}^\top \mathbf{x}^\top \mathbf{x}\mathbf{w}]$$

taking partial derivative of $J(\mathbf{w})$ (∇) w.r.t \mathbf{w}

$$\begin{aligned} \nabla_{\mathbf{w}} J(\mathbf{w}) \text{ or } \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} &= \frac{1}{2m} [-2\mathbf{x}^\top \mathbf{y} + 2\mathbf{x}^\top \mathbf{x}\mathbf{w}] \\ &= \frac{1}{m} [\mathbf{x}^\top \mathbf{x}\mathbf{w} - \mathbf{x}^\top \mathbf{y}] \end{aligned}$$

Set derivative to zero

$$\frac{1}{m} (X^T X \omega - X^T y) = 0 \Rightarrow$$

$$\Rightarrow X^T X \omega = X^T y$$

$$\therefore \boxed{\omega = (X^T X)^{-1} X^T y}$$

This gives optimal weight (ω) in one shot (provided $X^T X$ is invertible)