

Simple Linear Regression - Math In-Depth

1) Imagine you have some data points:

x = number of hours you study

y = marks you get

The points don't fall on a perfect line (real life is messy), but we want to draw one line - that best shows the trend

that line is

$$\hat{y} = mx + b \quad \text{or} \quad \hat{y} = b_0 + b_1 x$$

where :

m : slope (same as b_1)

b = intercept (same as b_0)

X = input (like hours studied)

\hat{y} = predicted output (like marks)

$$\text{Error} = y - \hat{y}$$

- Slope m : tells us how much y changes if x goes up by 1
- Intercept b : tells us the value of y when $x = 0$

so it's exactly like the line equation : $y = mx + c$

The formula

$$m = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} - x \text{ and } y \text{ move together (correlation)}$$

- is used to find the slope m of the best-fit line in simple linear regression

2) What each symbol means

- x_i = the i -th input value (like one student's study hours)
- y_i = the i -th output value (like that student's marks)
- \bar{x} = the mean of all x values
- $\bar{x} = \frac{\text{sum of all } x\text{'s}}{\text{number of data points}}$
- $\bar{y} = \frac{\text{sum of all } y\text{'s}}{\text{number of data points}}$

- simple linear regression - have only one input (x) and one output (y).
- In multiple regression we can have many inputs

3) Intuition of the formula

- The top part $\sum (x_i - \bar{x})(y_i - \bar{y})$ measures how x and y move together. (This is like correlation)
- The bottom part $\sum (x_i - \bar{x})^2$ measures how spread out x is

so slope = (how much x and y move together)
 (how much x varies)

That gives us: - If x goes up by 1 unit, how much does y go up on average?

Eq

x - (Input) Independent Variable

x (hours)	y (marks)
1	2
2	4
3	5

y - (Output) Dependent Variable

Step 1 Find the means

$$\bar{x} = \frac{1+2+3}{3} = \frac{6}{3} = 2$$

$$\bar{y} = \frac{2+4+5}{3} = \frac{11}{3} = 3.6$$

Step 2 Apply the slope formula

$$m = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

x_i	y_i	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$
1	2	-1	-1.6	-1.6	1
2	4	0	0	0	0
3	5	1	1.4	1.4	1

Step 3 Add them up

$$\sum (y_i - \bar{y})(x_i - \bar{x}) = -1.6 + 0 + 1.4 = -0.2$$

$$\sum (x_i - \bar{x})^2 = 1 + 0 + 1 = 2$$

Step 4 Find slope m

$$m = \frac{-0.2}{2} = -0.1$$

so, the slope is $m = -0.1$

step 5

Find intercept b

$$\begin{aligned} b &= \bar{Y} - m\bar{X} = \bar{Y} - m(\bar{X}) \\ &= 3.6 - (1.5 \times 2) \\ &= 3.6 - 3 \\ &= 0.6 \end{aligned}$$

for eg $x = 3$ and $x = 4$
 $1.5(3) + 0.6 = 4.5 + 0.6 = 5.1$

so intercept is $b = 0.6$

Final Line Equation

$$y = 1.5x + 0.6$$

These predicted values
are close to the actual value
but not always exact. That's
a sign that the model is
capturing the overall trend and
memorizing the data precisely.

