# Linear Regression

11 February 2024　22:15

- data points
— Linear model — Linear regression

BMI — (upward trend) | Supply — downward trend | No linear relationship → curve

Weight (in kg)

Demand

(Tomato)

$$BMI = f(weight, height)$$

Polynomial (degree ≥ 2)

$x^2$: quadratic

$x^3$: cubic

$x^4$: Bi-quadratic

- A technique of finding the relationship between two or more variables
- Change in dependent variable is associated with a change in one or more independent variables.
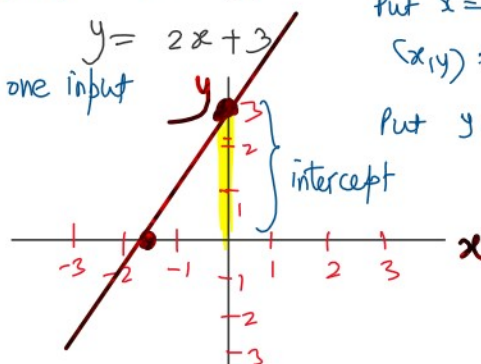
## What is Regression?

$R^2$ 0.06 | REXTHOR, THE DOG-BEARER

I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER TO GUESS THE DIRECTION OF THE CORRELATION FROM THE SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

$$y = f(x_1, x_2, x_3 \cdots x_n)$$

- a functional mapping
- to find the relationship between y and two or more x variables:

$$y = mx + c$$

→ slope
→ independent
(slope point form)

(target)

constant (intercept)

[ **simple** Linear Regression ]
↓
it has only one input

$$y = 2x + 3$$

intercept

Put $x = 0$, $y = 2 \times 0 + 3 = 3$
$(x, y) : (0, 3)$

Put $y = 0$, $x = -\frac{3}{2} = -1.5$
$(x, y) : \left(-\frac{3}{2}, 0\right)$

Slope of the line : $m = 2$
intercept of the line : $c = 3$

$$y = mx + c$$
→ slope

$$y = 2x + 3$$ ②

slope of the line : $m = 2$
intercept of the line : $c = 3$

$\dfrac{dy}{dx} = 2$

First order derivative is
slope which is $2$.

( Multiple Linear Regression )

$$y = 3 + 2x_1 + 5x_2 + 3.7\,x_3 - 0.8\,x_4$$
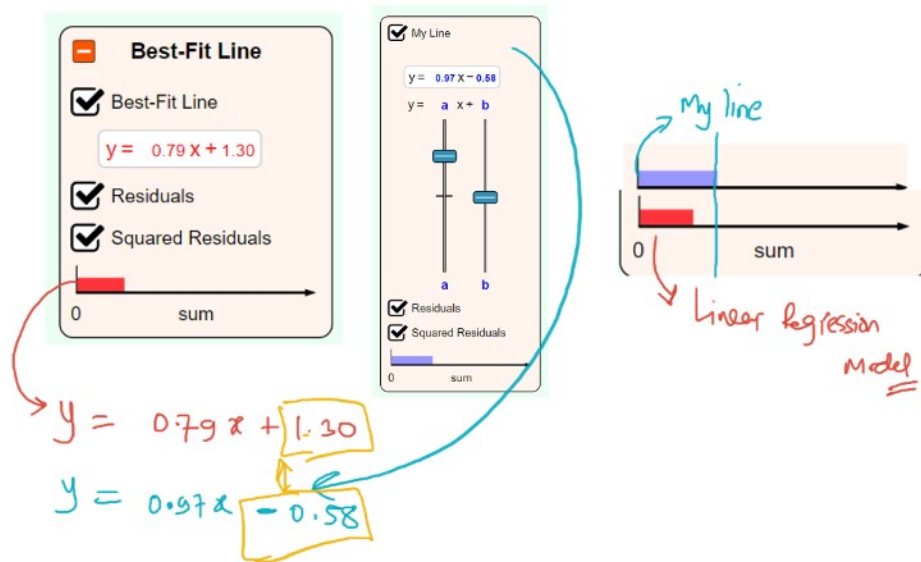
Notations

$x$ : input variable | features | independent | predictor

$y$ : output variable | target | dependent | response

## Intuition behind Linear Regression:

**Best-Fit Line**

☑ Best-Fit Line

$y = 0.79\,X + 1.30$

☑ Residuals

☑ Squared Residuals

0 ———— sum

➕ Correlation Coefficient

Custom ▼

20

my line

15

'model line'

10

5

0
0    5    10    15    20

☑ My Line

$y = 0.97\,X - 0.58$

$y = a\,X + b$

a    b

☑ Residuals

☑ Squared Residuals

0 ———— sum

**Best-Fit Line**

Best-Fit Line
$y = 0.79 X + 1.30$
Residuals
Squared Residuals

0    sum

My Line
$y = 0.97 X - 0.58$
$y = a x + b$

a    b

Residuals
Squared Residuals

0    sum

My line

0    sum

Linear regression model

$y = 0.79x + 1.30$

$y = 0.97x - 0.58$

## Least-Squares Regression



Stats Modelling

Machine learning

Linear Regression

$y = x$

$\hat{y} = 2.5$

$(-0.5)^2$

3

$\hat{y} = 2.5$

$y_{actual} = 3$

delta $= (3 - 2.5) = 0.5$

square delta $= (0.5)^2 = 0.25$

$+$    $0.25$

square delta $= (0.5)^2 = 0.25$  $+$  $(0.25)$

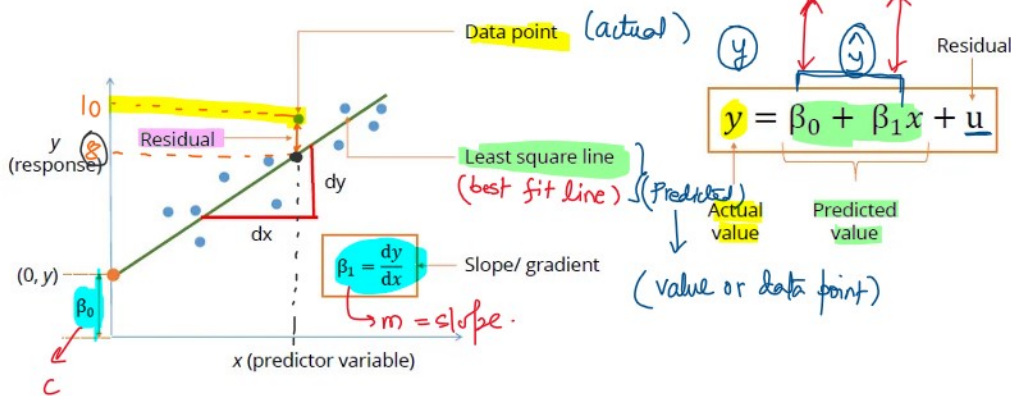sum of square residues

$y = mx + c$  (linear)

$m = -2$

$2x + y = 18$

$y = -2x + 18$

# Terminologies in Linear Regression

$c = \beta_0$  $\quad m = \text{slope} = \beta_1$

● : data points
── Linear Regression line

$y = c + mx +$



$y = \beta_0 + \beta_1 x + \underline{u}$

Data point (actual)  $\textcircled{y}$  $\textcircled{\hat{y}}$   Residual

Actual value    Predicted value

(value or data point)

Least square line
(best fit line)  (Predicted)

$\beta_1 = \dfrac{dy}{dx}$  — Slope/ gradient

→ m = slope.

$y$ (response)   10   8

Residual   dy   dx

(0, y)   $\beta_0$   c   x (predictor variable)

actual value = 10        (Predicted − Actual)

Predicted value = 8     $= (8 - 10) = -2$

$y = \beta_0 + \beta_1 x$   → $y_{\text{actual}}$   $= \hat{y}$ predicted

$y = c + mx$   $\textcircled{y}$   $\hat{\beta}_0 + \hat{\beta}_1 x = \hat{y}$

$y$ : actual value   10

$\hat{y}$ : Predicted value  8

residue $= (\hat{y} - y) = 8 - 10 = -2$

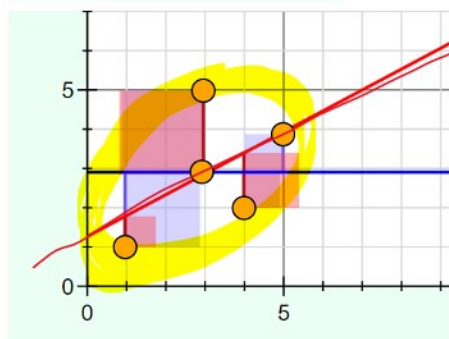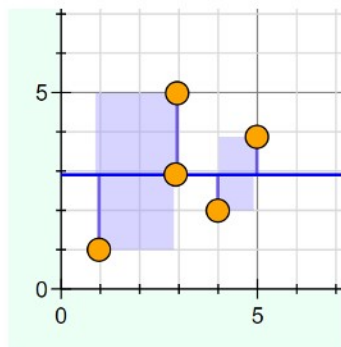# ✳ Issue with just doing — sum of residues



# (Predicted — Actual)

Residues

1st Point: (P-A) = 3-1 = 2

2nd Point: = 3-3 = 0

3rd Point = 3-5 = -2

4th Point = 3-2 = 1

5th Point = 3-4 = -1

sum of residues = $\cancel{2} + 0 - \cancel{2} + \cancel{1} - \cancel{1}$
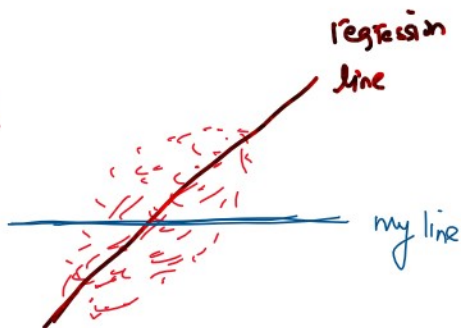
= 0

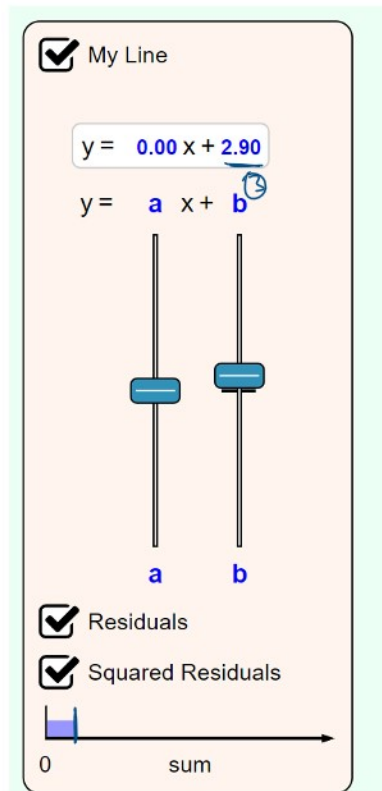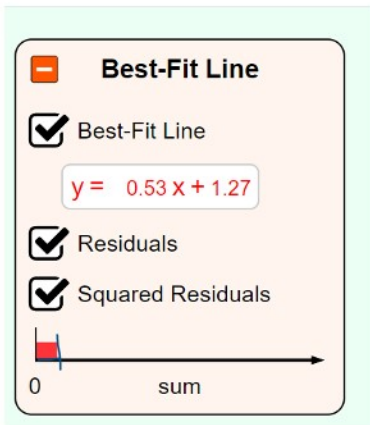① sum of residues = $\cancel{2} + 0 - \cancel{2} + 1 - 1 = 0$

(zero sum of residues)

_Observation: Residues nullify each other and may not be truly representing the residues.
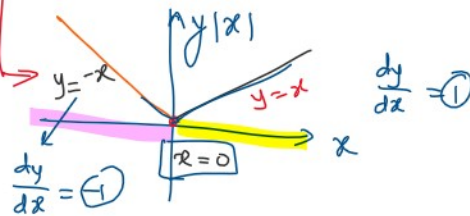




regression line

my line

regression line

my line

## Best-Fit Line

☑ Best-Fit Line

$y = 0.53 x + 1.27$

☑ Residuals

☑ Squared Residuals

0      sum

---

☑ My Line

$y = 0.00 x + 2.90$

$y = a x + b$

a     b

☑ Residuals

☑ Squared Residuals

0      sum

---

(Modulus / Absolute Function)

$$|x| = \begin{cases} x & x \geq 0 \\ -x & x < 0 \end{cases}$$

sum of residues $= 2 + 0 - 2 + 1 + 1$

$= 0$

$2 + 0 + 2 + 1 + 1 = 6$



$\frac{dy}{dx} = 1$

$\frac{dy}{dx} = -1$

$|x|$ is not differentiable at $x = 0$

↳ LINK:

minimize the **modulus** of residues

⇓

since modulus function is non-differentiable and maximal minima concept can't be used

Hence, let us introduce the concept of [OLS — ordinary least square]

# SSE, SSR and SST

i: data points



(adjust)
$Y_1$ = observed

Regression Line

SSE

SSE

SST $\Leftarrow \oplus$

$\hat{Y}_i$ = Predicted

SSR

SSR

$\bar{Y}$ = Mean

$\bar{Y}$

(avg. of Y values)



(3,5)

(5,4)

my line

(3,2)

(4,2)

(1,1)

$$\bar{Y} = \frac{1+3+5+2+4}{5}$$

$$\bar{Y} = \frac{15}{5} = 3$$

$$\bar{X} = \frac{1+3+3+5+4}{5}$$

$$\bar{X} = \frac{16}{5} = 3.33$$

SSE: Sum of squares error

SSR : sum of squares regression.

SST : sum of squares total

# sum of squares error (residues) : SSE

-it is the sum of the squared differences between the observed value (actual) and predicted value $\hat{Y}_i$ $Y_i$

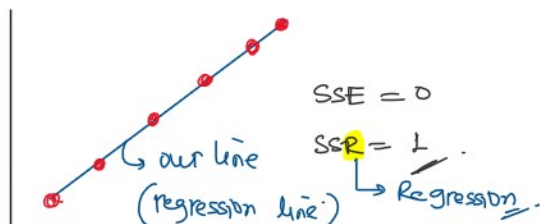- SSE shows the unexplained variance by regression.

$$SSE = \sum_{i=1}^{n} \left( \hat{Y}_i - Y_i \right)^2$$

sum of squared error

# sum of squares regression (SSR)

- it is the sum of squared differences between predicted value $(\hat{Y}_i)$ and the mean of the dependent value $(\bar{Y})$

- SSR shows the explained variance by regression
- it is a measure that describes how well our line fits the data.



SSE = 0
SSR = 1 .
→ Regression.

↳ our line
(regression line)

$$SSR = \sum_{i=1}^{n} \left( \hat{Y}_i - \bar{Y} \right)^2$$
↳ (regression)

# Sum of squares total (SST)

- it is the squared differences between **observed**
  **dependent variable** and its mean    (Yactual)
                                                       $(\bar{y})$

$$\text{SST}/\text{TSS} = \sum_{i=1}^{n}(y_i - \bar{y})^2$$

total sum of squares

- it is a measure of **total variability** of the dataset

$$\boxed{SST = SSR + SSE}$$

Total variability      =    Variability     +    Unexplained
of the data set              explained by the          variability
(SST)                         regression line              (SSE)
                                  (SSR)

$$100\% = (80\%) + (20\%)$$
                  (accuracy)       (error)

Total variance = 1000

Total error = 100

Total regressed value = 900       Accuracy = $\dfrac{900}{1000} \times 100$

                                                    = 90%