

## • Simple Linear Regression •

Sr No	cgpa	Package
1	6.66	3.01
2	7.1	3.5
3	4.7	1.2
4	8.9	4.2
5	8.1	3.9
6	:	:
7	:	:
200	n	20·n

① Plot the dataset

Ques, Why dataset is not linear?

→ Real world dataset it is.

- The factors which cannot be determined are called 'Stochastic Errors'.

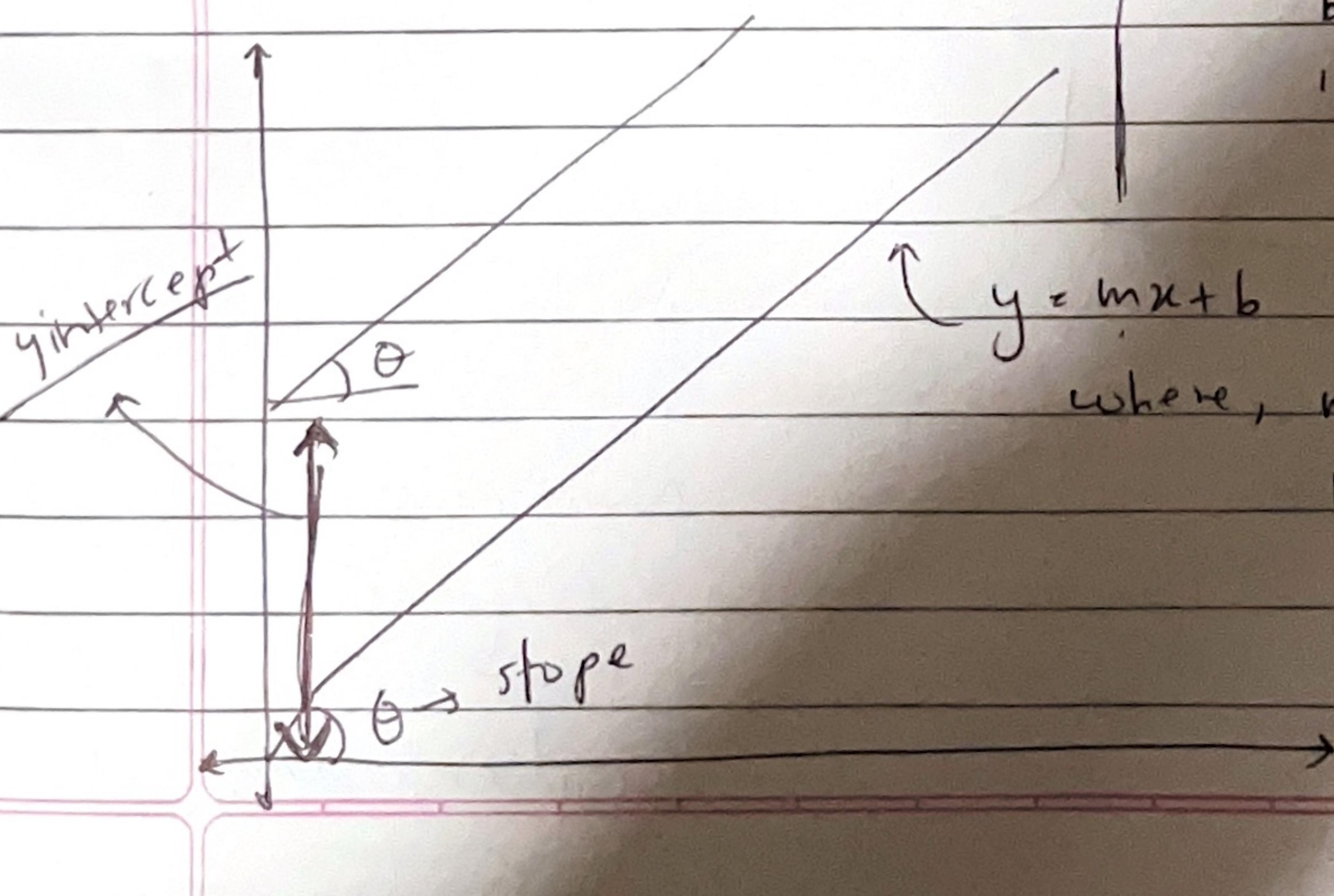
Eg India USA UK

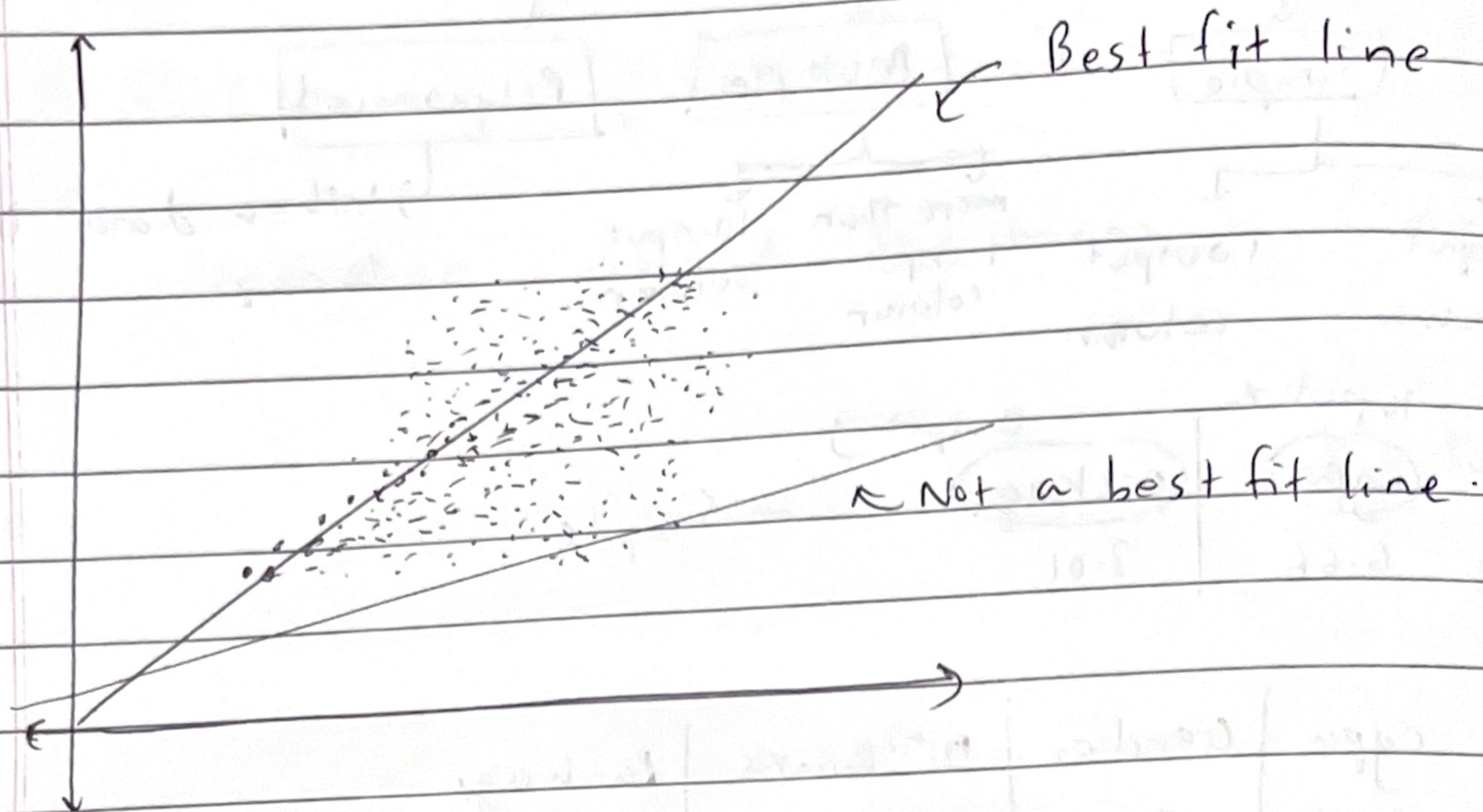
Interview?

$$y = mx + b$$

where, m = slope

b = y intercept





### Steps to work on any dataset:-

(1) Plot the dataset and analyse.

(2) Separate x and y columns.

(3) Divide data in train and test

x-train | x-test

y-train | y-test

29x2

(4) Import linear regression model.

(5) Use fit to train data

(6) Now test the model.

Ques

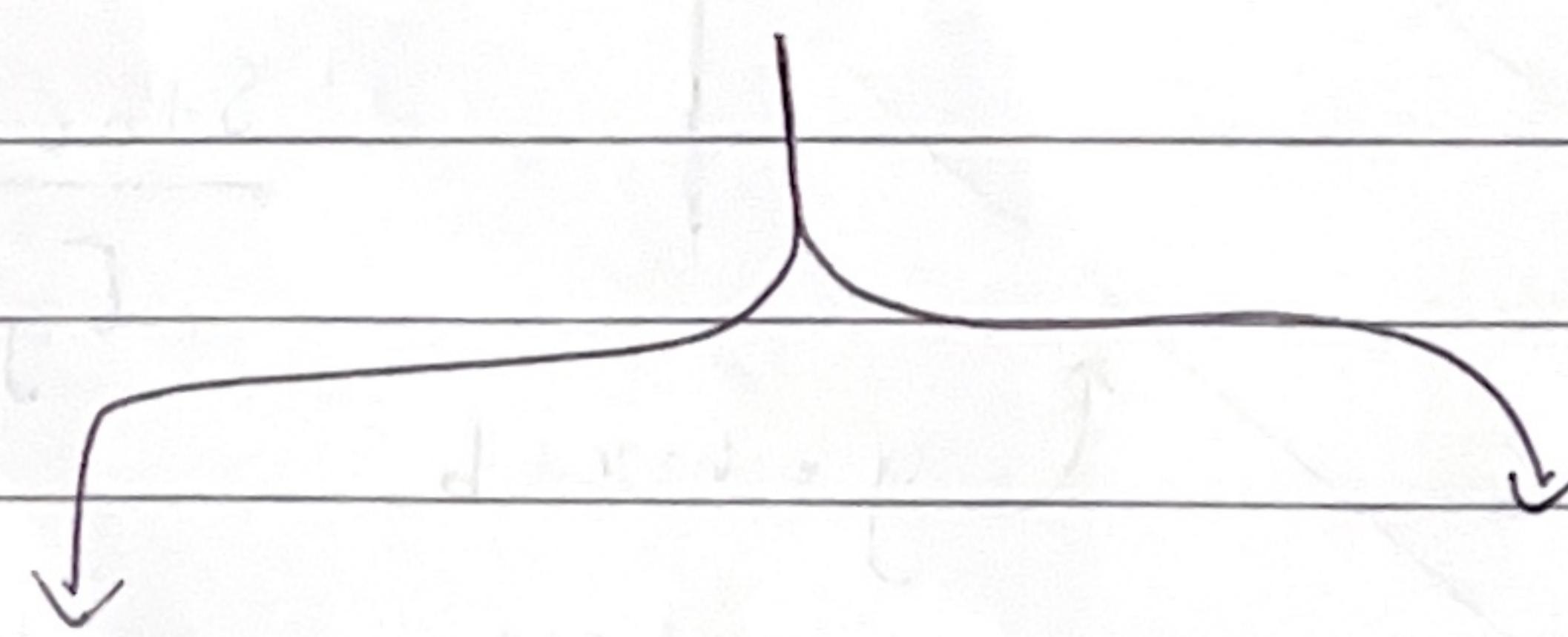
What is  
input shape?

↳ Size of  
dataset.

Eg. 64x6y

or

To find  $(m, b)$  there are two  
methods:



Closed form

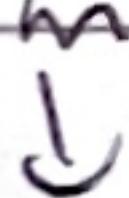
sol<sup>r</sup>



OLS

Non-closed

form sol<sup>n</sup>



Gradient Descent

(Direct formula)

(closed form solution  $\rightarrow +, -, \times, \div$  use Karke  
jis bhi chiz ka formula ban sakta hai but no  
differentiation and integration.)

Eg. quadratic eqns.

↳ OLS - Ordinary Least Squares

Non - closed form sol<sup>n</sup>  $\rightarrow$  Ek approximation method  
use Karke (with using differentiation & integration)  
(if needed) such dhundhna.

↳ Gradient Descent.

Ques. If we have a fixed formula to find  
answer than why use GD?

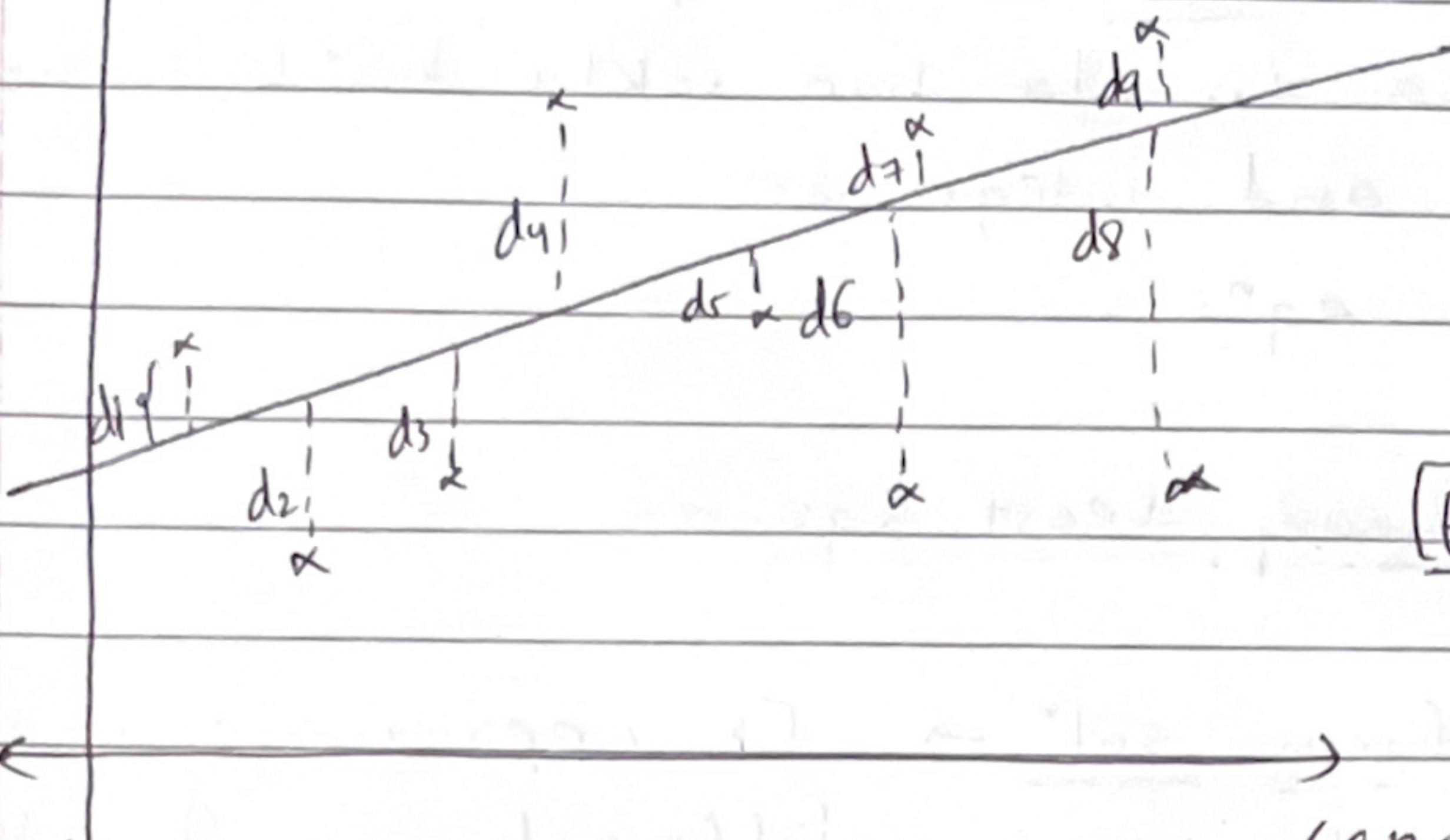
$\rightarrow$  In higher dimensions, using a formula  
is not enough to find the solution.

OLS (Ordinary least squares)

$$b = \bar{y} - m\bar{x}$$

$$\text{Slope } \rightarrow (m) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

where  $\bar{x}$  and  $\bar{y}$  are means



Now to find total distance / error,

$$E = d_1 + d_2 + \dots + d_n$$

What is q.  
penalizing.

→ A technique i.e used for reducing shrinking a large amount of features to a manageable set and for making good predictions in a variety of large data sets.

OR

Giving some weightage for deviating from the objective.

Therefore, our Error function looks like this,

$$E = \sum_{i=0}^n (d_i)^2 \quad \leftarrow \text{Loss function}$$

NOTE: Also represented as 'J'

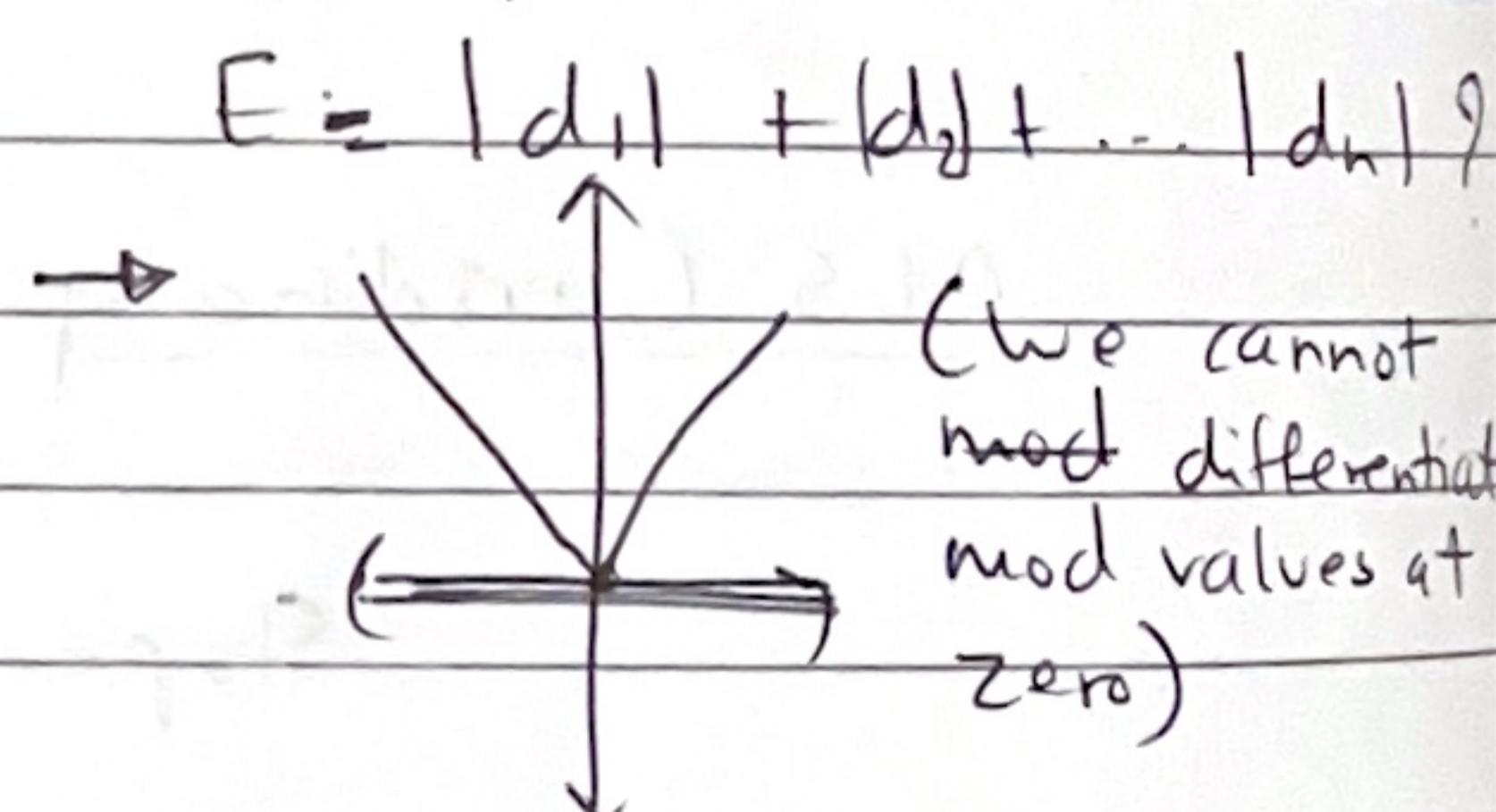
This predicted line is represented by  $\hat{y}$  (HAT)

Now, use square the distances and then add

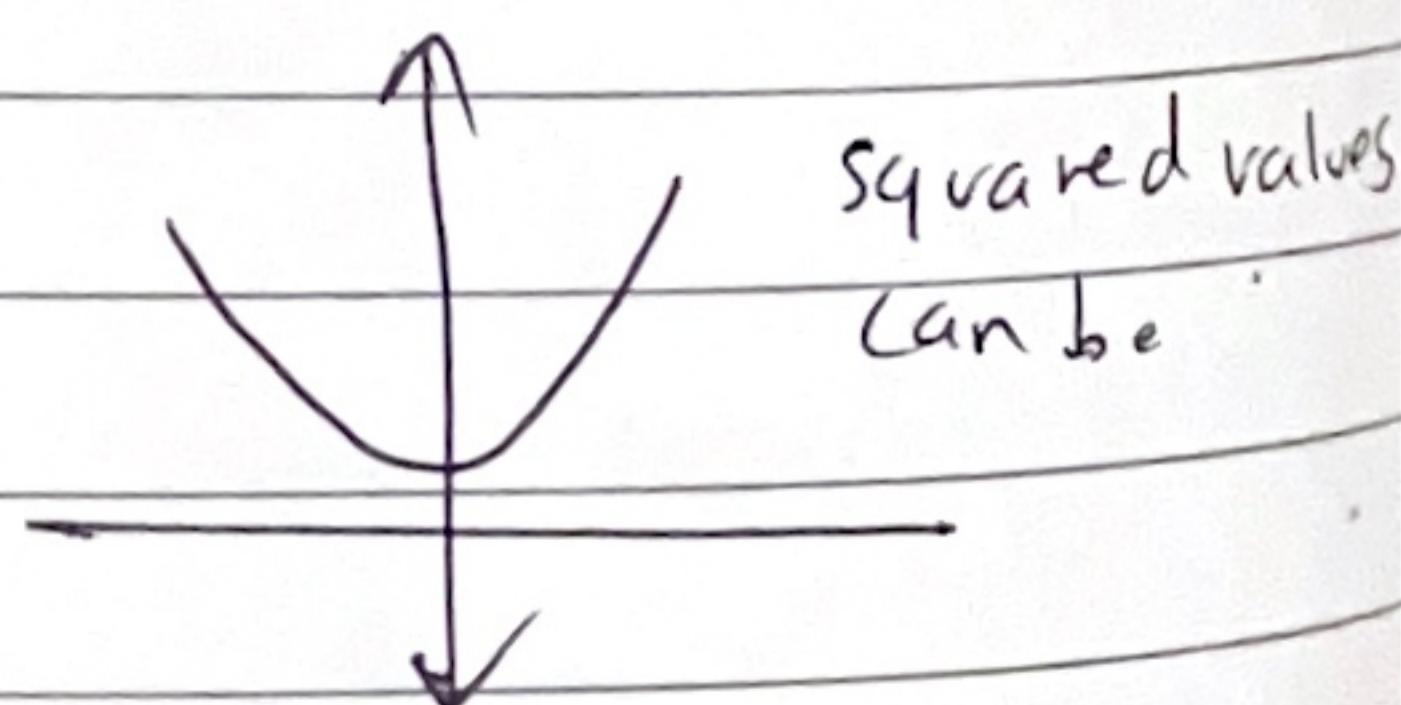
$$E = (d_1)^2 + (d_2)^2 + \dots + (d_n)^2$$

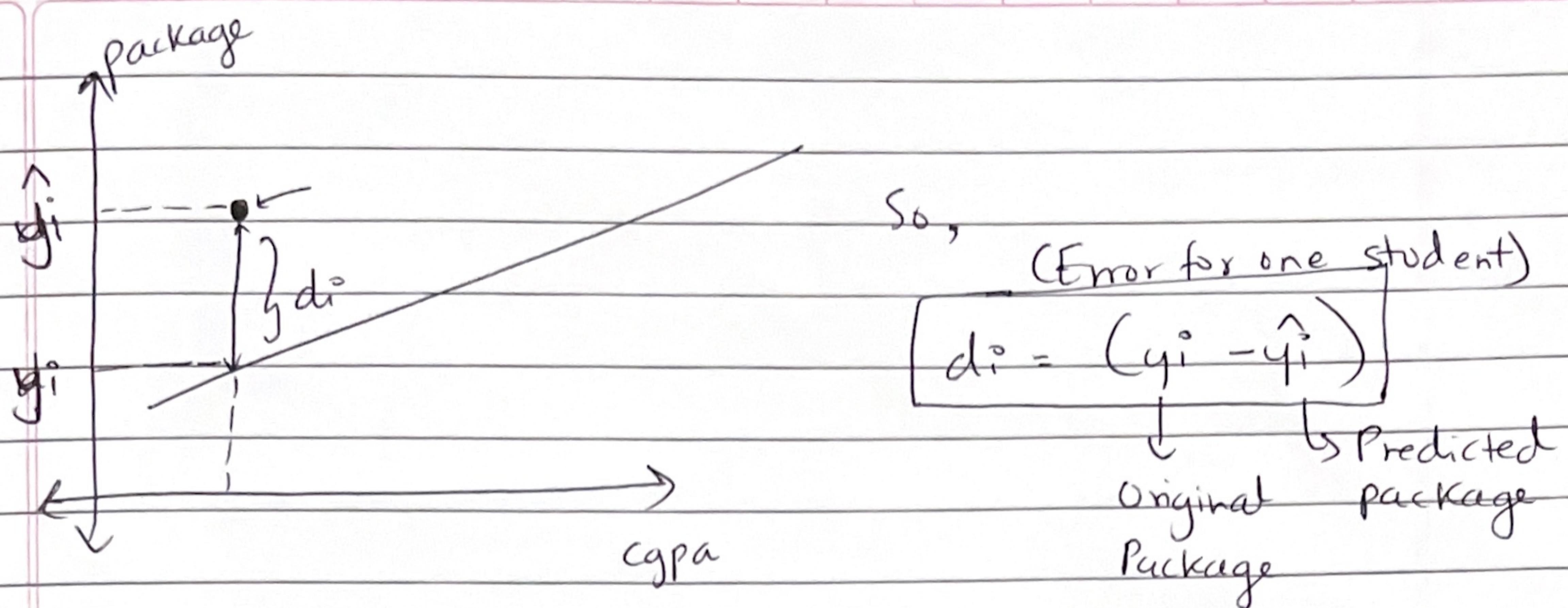
REASON: To cancel out positive and negative error

Ques. Why not mod the error?



whereas,





Hence,

(Loss / Error Function)

$$E/J = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

i.e

$$E/J = \sum_{i=1}^n (d_i)^2$$

Also,

$$E/J \text{ or } \Sigma(m, b)$$

$$= \sum_{i=1}^n (y_i - m x_i - b)^2$$