# Recap of Previous Lecture

Topic — Model

Topic — Optimisation

Topic — Training

Topic — data

Topic

Basic

# Topics to be Covered

Topic

Topic

Topic

Topic

Topic

{Gate exam}**

Linear regression

How We do optimisation

loss function

Numerical Question

Success is walking from failure to failure with no loss of enthusiasm.

WINSTON CHURCHILL

Patience

**Fill in the blanks :**

1. The target/Goal of the ML is _"To predict $y$ for new $x$"_

2. The best optimized model is that which minimize the error in _"Training data"_

3. The problem with the simple model is And not learning Pattern of data ⇒ It has lot of error.

**Fill in the blanks :**

4. The problem with highly complicated model is

noise in data is also included in Analysis/Rote learning

_____

5. The data is used to _____Train_____ the ML model

model ⇒ fxn
$y = f(x)$

6. The data is collected from ___Survey/experiment___ .

**19. The output of training process in machine learning is**

A. machine learning model

B. machine learning algorithm

C. null

D. accuracy

• we get y and x

Relation ⇒ we Call

this mL model.

**34. In simple term, machine learning is**

A. training based on historical data

B. prediction to answer a query

C. both a and b??

D. automization of complex tasks

**Problem 2 – Predict Sale of I-phone based on Age of customer**

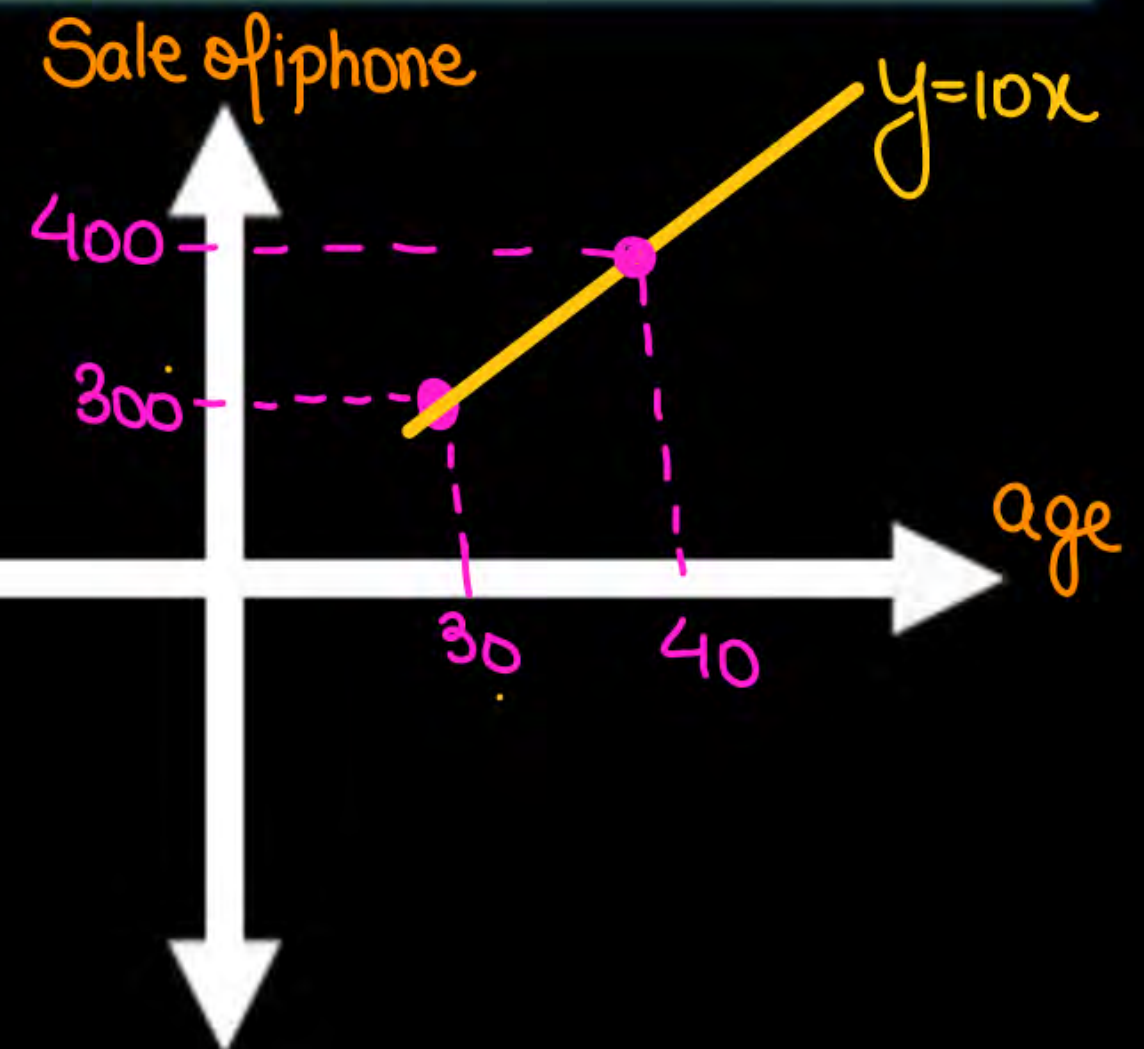## We must create a model with following data

| Age | Sale of I-Phone (in a month) |
|-----|------------------------------|
| 30  | 300                          |
| 40  | 400                          |

$x$ { (Age) }

{ (Sale of I-Phone) } $y$

Predict Sale of iphone @ age=20

$\rightarrow x$

$\rightarrow y$

## Now predict the Sale of I-Phone at Age = 20

**Problem 2 – Predict Sale of I-phone based on Age of customer**

We don't have any expert now, and data has only two Points.
So _____

$$y = mx + c$$

@ $x = 30$   $y = 300$

$300 = 30m + c$ ←

@ $x = 40$   $y = 400$

$400 = 40m + c$ ←

**What is the best model now ?**

Sale of iphone

$y = 10x$

400

300

Age

30    40

$$30m + c = 300$$
$$40m + c = 400$$

$$10m = 100$$
$$m = 10 \checkmark$$

$$30 \times 10 + c = 300$$
$$c = 0 \checkmark$$

So we will try to fit a line on the data $\rightarrow$

So the Best line is that which has min gap b/w Yactual and Ypredicted values by model.

Since only 2 points are given hence $\Rightarrow$ we can draw a line that can pass through both the points $\Rightarrow$ Yactual = Ypredicted

**Problem 2 – Predict Sale of I-phone based on Age of customer**

**Now we have to find the best parameters..**

So Simple Case ⇒ 2 point in data easily we can draw a
line passing through points

$y_{act} = y_{pred}$

**Problem 2 – Predict Sale of I-phone based on Age of customer**

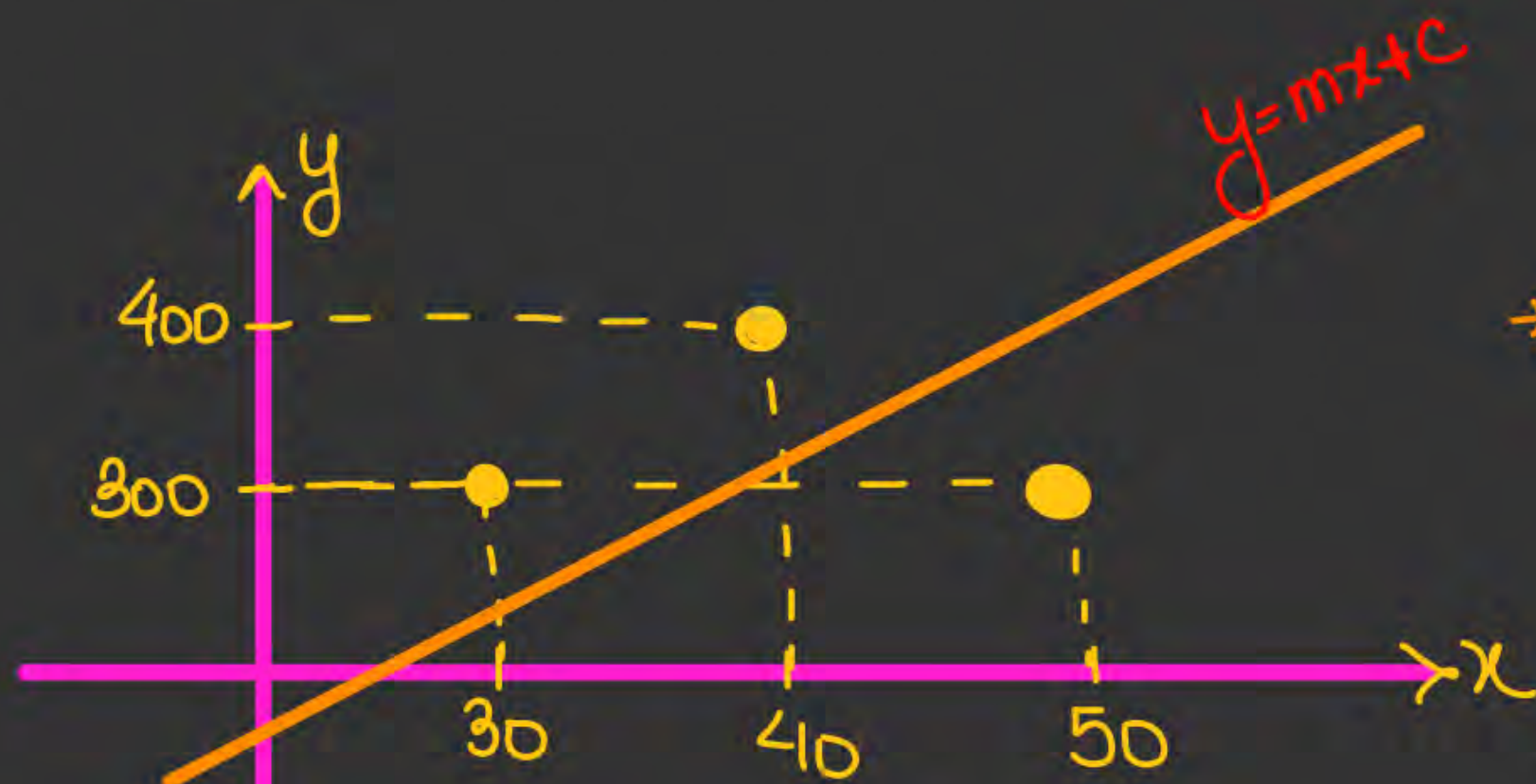Now we have to find the best parameters..

(done)

**Problem 3 – Predict Sale of I-phone based on Age of customer**

## We must create a model with following data

| Age | Sale of I-Phone (in a month) |
|-----|------------------------------|
| 30  | 300                          |
| 40  | 400                          |
| 50  | 300                          |

Predict Sale of iphone $\Rightarrow y$

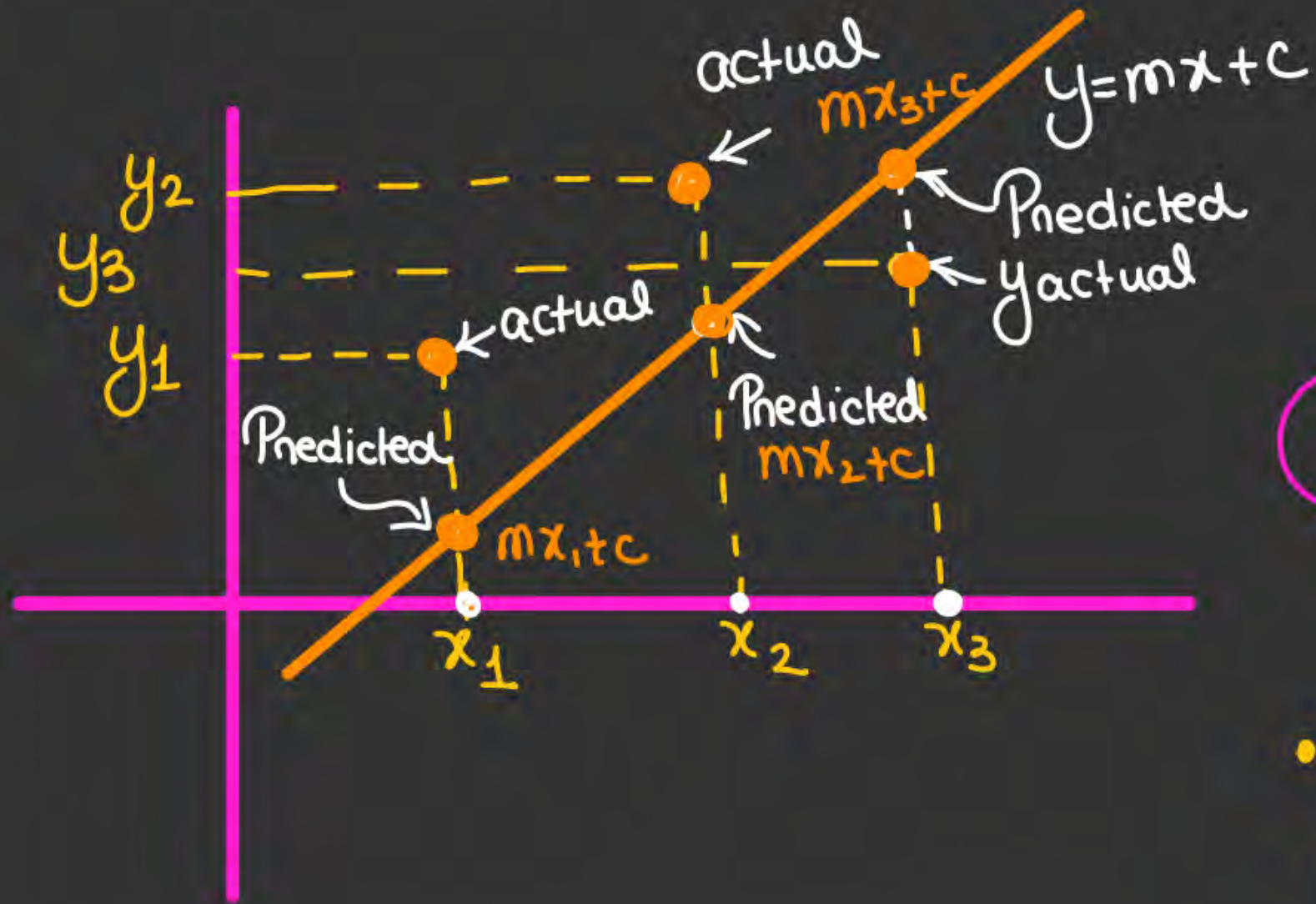age $\Rightarrow x$

## Now predict the Sale of I-Phone at Age = 20

$y = mx + c$

400

300

30    40    50

$\rightarrow x$

* So now data has 3
  Points
  $(x_1, y_1)$ $(x_2, y_2)$ $(x_3, y_3)$

* So we can see that no
  St. line can pass through
  all 3 points.

* Now we have to find best
  Straight line.

• Model's goodness is measured
  by $\rightarrow$ the gap b/w $y_{actual}$ and
  $y_{predicted}$.

actual

$mx_3+c$

$y=mx+c$

$y_2$

$y_3$

$y_1$

Predicted

$y_{actual}$

actual

Predicted

$mx_2+c$

Predicted

$mx_1+c$

$x_1$   $x_2$   $x_3$

So Shall we min $[y_{actual} - y_{predicted}]$

- Here we have a problem that is this value can be +ve/ -ve
⇒ If we add all error then -ve error will Reduce the total error, may lead to Confusion

So 2 options

$$\min \sum_{i=1}^{3} |Y_{i}^{act} - Y_{i}^{pred}|$$

$\Downarrow$ absolute error

OR

$$\min \sum_{i=1}^{3} (Y_{i}^{act} - Y_{i}^{pred})^{2}$$

- Residual Sum of Square.

- RSS

To find the min location we differentiate the fxn $\longrightarrow$

So Residual Sum of Square $\Rightarrow$ The loss function

So data 3point $\Rightarrow$ $(x_1, y_1)$ , $(x_2, y_2)$ , $(x_3, y_3)$

The St. line model $\Rightarrow$ $(y = mx + c)$

$x_1 \xrightarrow{\text{Predicted}} \hat{y}_1 = mx_1 + c$

$x_2 \xrightarrow{\text{Predicted}} \hat{y}_2 = mx_2 + c$

$x_3 \xrightarrow{\text{Predicted}} \hat{y}_3 = mx_3 + c$

The $y_1, y_2, y_3 \Rightarrow$ These are actual 'y' Values of data

Now Best model $\Rightarrow$ that minimize the loss function

$\Downarrow$

By this maths we will get m, c

$\Downarrow$

The Best model $y = mx + c$

The loss function $\Rightarrow$

$$\alpha = \sum_{i=1}^{3} \left( y_{i}^{actual} - y_{i}^{pred.} \right)^2$$

Rss $\Leftarrow \alpha = \sum_{i=1}^{3} \left( y_i - (mx_i + c) \right)^2$

So $x_i$, $y_i$ are given
But $m$, $c$ are unknown

Now we have to minimize '$\alpha$' loss fxn $\Rightarrow$

$\min \alpha \Rightarrow \min \sum_{i=1}^{3} \left( y_i - (mx_i + c) \right)^2 \Rightarrow$ for minimizing $\Rightarrow \frac{\partial L}{\partial m} = 0 , \frac{\partial L}{\partial c} = 0$

$$L = \sum_{i=1}^{3} \left( y_i^0 - (mx_i^0 + c) \right)^2$$

$$\frac{\partial L}{\partial m} = 2 \sum_{i=1}^{3} \left( y_i^0 - (mx_i^0 + c) \right) (-x_i^0) = 0$$

$$\Rightarrow \sum_{i=1}^{3} x_i^0 y_i^0 - m x_i^{0^2} - c x_i^0 = 0 \quad \text{—①}$$

$$\frac{\partial L}{\partial c} = 2 \sum_{i=1}^{3} \left( y_i^0 - (mx_i^0 + c) \right) (-1) = 0$$

$$\sum_{i=1}^{3} \left( y_i^0 - m x_i^0 - c \right) = 0 \quad \text{—②}$$

$$\frac{\partial}{\partial m} \left( f(m) \right)^2$$
$$\downarrow$$
$$2 f(m) \cdot f'(m)$$

$$2 f(m) = 0$$
$$\Rightarrow f(m) = 0$$

2 eq, 2 variable we can find $m, c$

Q $(30,300)$, $(40,400)$, $(50,300)$

$(x_i, y_i)$

op
find out best line for the data $\Rightarrow$

$y = mx + c$

$$L = \sum_{i=1}^{3} \left( y_i - y_{pred_i} \right)^2$$

$$\min L = \min \sum_{i=1}^{3} \left( y_i - (mx_i + c) \right)^2$$

$$\frac{\partial L}{\partial m} = 0 \Rightarrow \sum_{i=1}^{3} \left( y_i - (mx_i + c) \right) x_i = 0$$

$$\Rightarrow \left\{ \sum_{i=1}^{3} \left( y_i x_i - mx_i^2 - cx_i \right) = 0 \right.$$

$$\frac{\partial L}{\partial c} = 0 \Rightarrow \left. \sum_{i=1}^{3} \left( y_i - mx_i - c \right) = 0 \right\}$$

$$\sum_{i=1}^{3} (y_i x_i) - m\sum_{i=1}^{3} x_i^2 - C\sum_{i=1}^{3} x_i = 0$$

$30 \times 300 + 40 \times 400 + 50 \times 300$

$30^2 + 40^2 + 50^2$

$30 + 40 + 50$

$40,000 - 5000m - 120C = 0$

$$\sum_{i=1}^{3} y_i^0 - m\sum_{i=1}^{3} x_i^0 - C\sum_{i=1}^{3} 1 = 0$$

$300 + 400 + 300$

$30 + 40 + 50$

$1000 - 120m - 3C = 0$

$5000m + 120C = 40,000$

$120m + 3C = 1000$

$m = 0, \quad C = 1000/3$

$$Q \quad \sum_{i=1}^{3} (3i^0 + c) \Rightarrow (3(1)+c) + (3(2)+c) + (3(3)+c)$$

$$\left( 3\sum_{i=1}^{3} i^0 + c\sum_{i=1}^{3} 1 \right)$$

$$\Rightarrow 3\sum_{i=1}^{3} i^0 + 3c$$

If we had N points $\Rightarrow$ $\alpha = \min \left( \sum_{i=1}^{N} \left( y_i - (mx_i + c) \right)^2 \right)$

$\dfrac{\partial L}{\partial m} = 0 \checkmark$

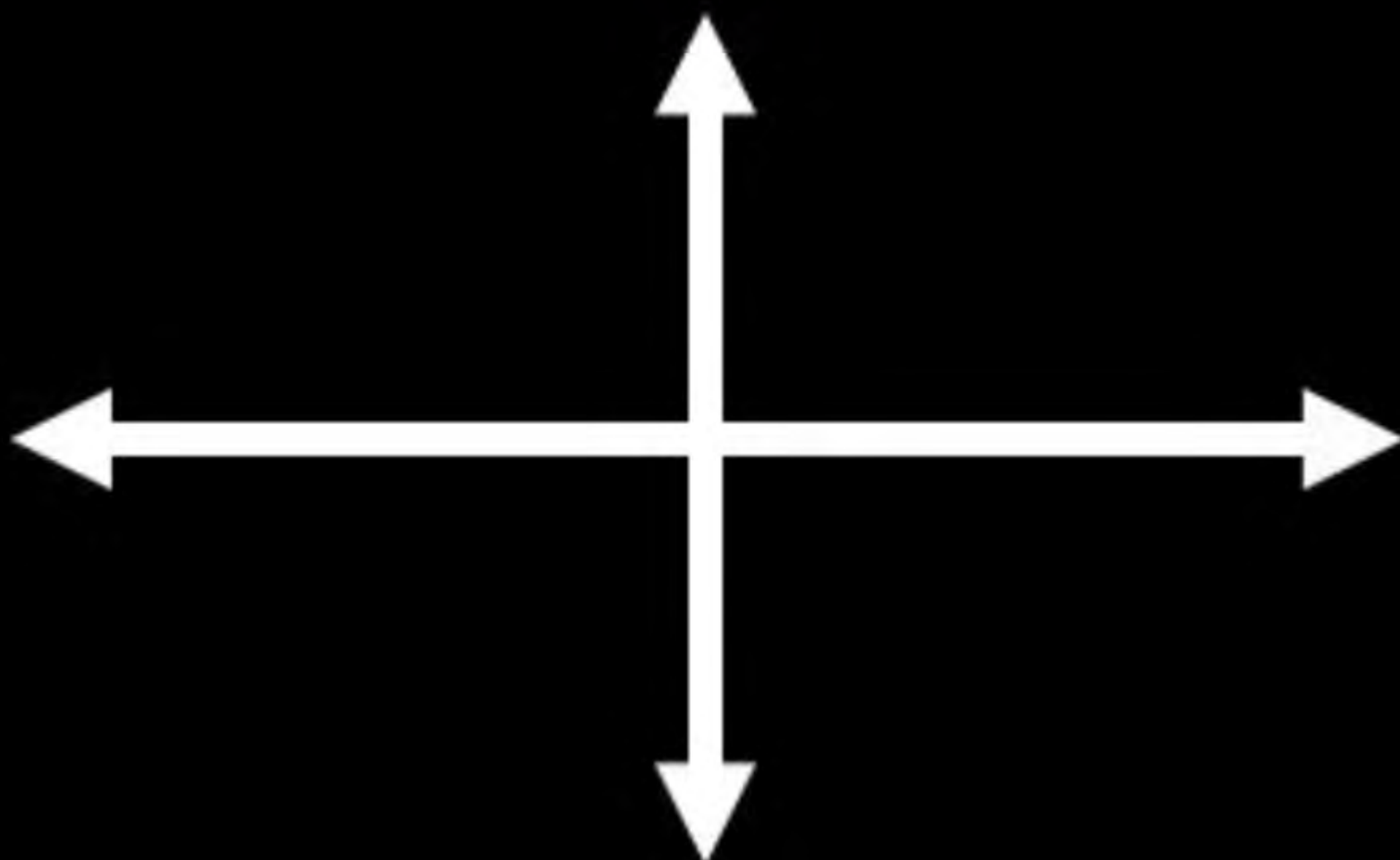$\dfrac{\partial L}{\partial c} = 0$

$\Rightarrow$ 2 equation and we get $m, c$.

**Problem 3 – Predict Sale of I-phone based on Age of customer**

We don't have any expert now, and data has only two Points. So
_____

What is the best model now ?

**Problem 3 – Predict Sale of I-phone based on Age of customer**

We don't have any expert now, and data has only two Points. So

_____

But we will try to find the linear model only.
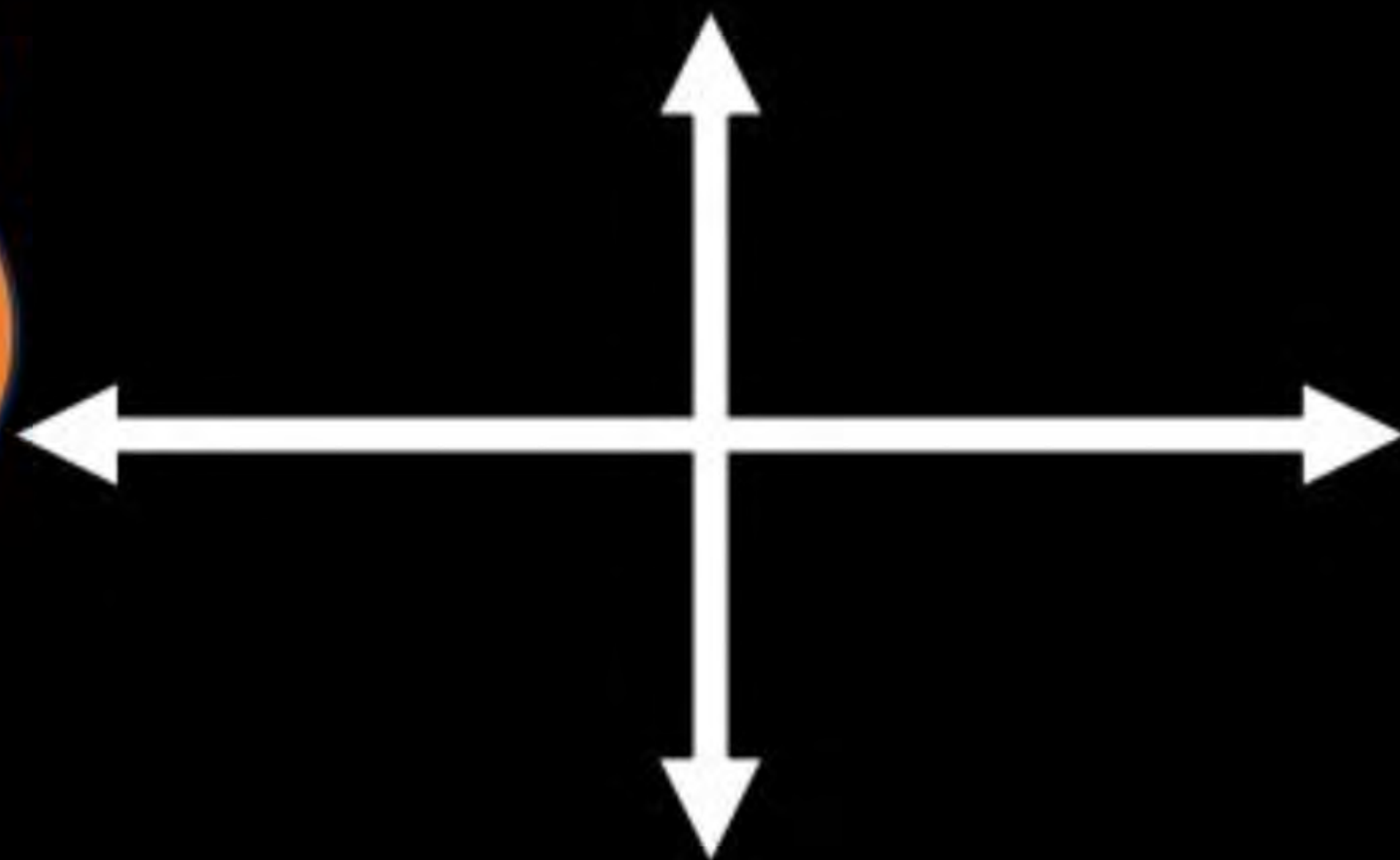
So, we must find the model that try to

## Problem 3 – Predict Sale of I-phone based on Age of customer

## Creating the best model

Now we have that is called the predicted value of input.

**Problem 3 – Predict Sale of I-phone based on Age of customer**

**Creating the best model**

Loss Functions ?? (RSS- Residual Sum of Squares)

(done)

*Reading*

The residual sum of squares (RSS), also known as the sum of squared residuals (SSR) or the sum of squared estimate of errors (SSE), is the sum of the squares of residuals

**Now how to find the best parameters ??**

Variance and mean...

mean of any variable

$x_i \circ \Rightarrow$ are some values

$\overline{x_i\circ} \Rightarrow \dfrac{\sum\limits_{i=1}^{N} x_i\circ}{\text{Number of Values}}$

$x_p = 3, 5, 7, 8$

$\overline{x_i\circ} = \dfrac{3+5+7+8}{4}$

**Now how to find the best parameters ??**

Variance of a variable

$$\left(\sigma_x^2\right) = \dfrac{\displaystyle\sum_{i=1}^{N}\left(x_i^0 - \overline{x_i^0}\right)^2}{\text{Number of values}}$$

$\sigma_{x_i} \Rightarrow$ Standard deviation.

Variance and mean...

$x_i^0 = ③, 8, 2, 7.$

$\text{mean } \overline{x_i^0} \Rightarrow \dfrac{3+8+2+7}{4} \Rightarrow ⑤$

$$\sigma_{x_i^0}^2 \Rightarrow \dfrac{(3-5)^2 + (8-5)^2 + (2-5)^2 + (7-5)^2}{4}$$

$$\Rightarrow 6.5$$

**Now how to find the best parameters ??**

If X and Y are two variables

$$\Rightarrow \quad Cov(XY) \Rightarrow \left\{ \dfrac{\displaystyle\sum_{i=1}^{N}(x_i-\bar{x}_i)(y_i-\bar{y}_i)}{N} \right\}$$

Variance and mean...

## in Machine learning ⇒

### mean of Variable X

$$\frac{\sum_{i=1}^{N} x_i^0}{\text{No of Values}} = \frac{\sum_{i=1}^{N} x_i^0}{N}$$

### Variance of X ⇒

$$\sigma_X^2 = \frac{\sum_{i=1}^{N}\left(x_i^0 - \bar{x}_i^0\right)^2}{(N-1)}$$

### Covariance of X, Y

$$\text{Cov}(X,Y) = \frac{\sum_{i=1}^{N}\left(x_i^0 - \bar{x}_i^0\right)\left(y_i^0 - \bar{y}_i\right)}{(N-1)}$$

**Now how to find the best parameters ??**

Formulae to find direct value of m and c

So if we have any data

$(x_1, y_1) (x_2, y_2) (x_3, y_3) (x_4, y_4) ---$

So let $\boxed{y = mx + c}$ is our model

So $x = x_1, x_2, x_3, x_4 ---$

$y = y_1, y_2, y_3, y_4 ----$

$$m = \frac{Cov(x, y)}{Var\ x}$$

after Calculating m

$$\bar{y} = m\bar{x} + c$$

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

**Now how to find the best parameters ??**

Formulae to find direct value of m and c

**Example**

Obtain a linear regression for the data in below table assuming that y is the independent variable.

| $x =$ | 1 | 2 | 3 | 4 | 5 |
|-------|----|----|----|----|----|
| $y =$ | 10 | 15 | 18 | 20 | 25 |

$$\Rightarrow y = mx + c$$

$$m = \frac{Cov(x,y)}{Var(x)}$$

$$\bar{x} \Rightarrow 3 \qquad \bar{y} = 17.6$$

$$Var(x) = \frac{\sum(x_i^\circ - \bar{x})^2}{4} \Rightarrow 2.5$$

PRACTICE DOES NOT MAKE PERFECT

PRACTICE MAKES PROGRESS!

$$\text{Cov}(x,y) = \frac{\sum\limits_{i=1}^{5}(x-\bar{x})(y-\bar{y})}{4}$$

$$\Rightarrow \frac{(1-3)(10-17.6)+(2-3)(15-17.6)+(3-3)(18-17.6)+(4-3)(20-17.6)+(5-3)(25-17.6)}{4}$$

$$\Rightarrow \boxed{35/4}$$

$$m = \frac{35/4}{2.5} \Rightarrow 3.5 \checkmark$$

$$C = \bar{y} - m\bar{x}$$
$$C = 17.6 - 3.5 \times 3$$
$$\boxed{C = 7.1}$$

A set of observations of independent variable (x) and the corresponding dependent variable (y) is given below.

| x | 5 | 2 | 4 | 3 |
|---|----|----|----|----|
| y | 16 | 10 | 13 | 12 |

Based on the data, the coefficient a of the linear regression model

y = a + bx is estimated as 6.1

The coefficient b is _____ . (round off to one decimal place)

For a bivariate data set on $(x, y)$, if the means, standard deviations and correlation coefficient are

$\bar{x} = 1.0, \bar{y} = 2.0, s_x = 3.0, s_y = 9.0, r = 0.8$

Then the regression line of y on x is:

1. $y = 1 + 2.4(x - 1)$

2. $y = 2 + 0.27(x - 1)$

3. $y = 2 + 2.4(x - 1)$

4. $y = 1 + 0.27(x - 2)$

In the regression model ($y = a + bx$) where $\bar{x} = 2.50$, $\bar{y} = 5.50$ and a $= 1.50$ ($\bar{x}$ and $\bar{y}$ denote mean of variables x and y and a is a constant), which one of the following values of parameter 'b' of the model is correct?

1. 1.75

2. 1.60

3. 2.00

4. 2.50

There is no value of x that can simultaneously satisfy both the given equations. Therefore, find the 'least squares error' solution to the two equations, i.e., find the value of x that minimizes the sum of squares of the errors in the two equations. _____

$2x = 3$

$4x = 1$

We can expect one Question from here in GATE exam

**Considering data of 2 Dimensions**

Attributes,
Features,
Dimensions...

Till now we have seen a simple case of 1 D data,
now let's see 2 D Data

| Income (LPA) | Age | Sale of I-Phone (in a month) |
|---|---|---|
| 20 | 30 | 300 |
| 50 | 40 | 400 |
| 70 | 50 | 300 |

**We have N Data points**

**Now the input data is 2 D (age and income)**

The representation of D dimensional data

# 2 mins Summary

$$\text{Loss function} \Rightarrow \text{RSS} \Rightarrow \sum_{i=1}^{N} (y_i - y_{pred})^2 \Rightarrow \frac{\partial L}{\partial m}, \frac{\partial L}{\partial c}$$

$$y_{pred} = mx + c$$

$$\cdot \ m = \frac{Cov(x, y)}{Var(x)}$$

$$\cdot \ c = \bar{y} - m\bar{x}$$

THANK - YOU