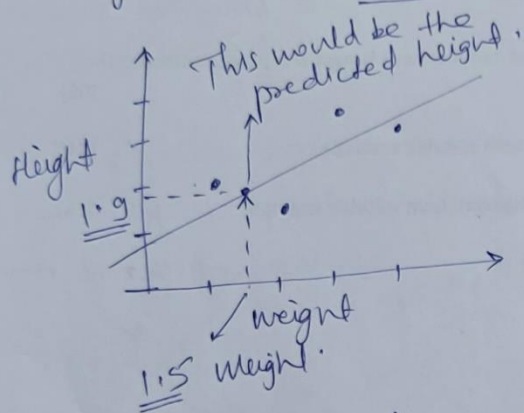


Linear Regression Using Gradient Descent (An Iterative Approach.)

⇒ When we fit a line with Linear Regression, we optimize the Intercept and slope.



$$H = \text{Intercept} + \text{slop} * \underset{\substack{\downarrow \\ \text{weight}}}{w}$$

\downarrow

$$\hat{y} = a_0 + a_1 x$$

Equation of the line.

⇒ Gradient Descent is used to optimize the slope & Intercept so that line is the best fit for given data points.

→ Statistics, Machine Learning & Data science

$$\text{Predicted Height} = \underbrace{\text{Intercept} + \text{slop} * \text{weight}}_{\text{Equation of the line}}$$

So, gradient descent can fit a line to data by finding the optimal values.

⇒ Start GD to find Intercept; then we used it to solve for Intercept & slope.

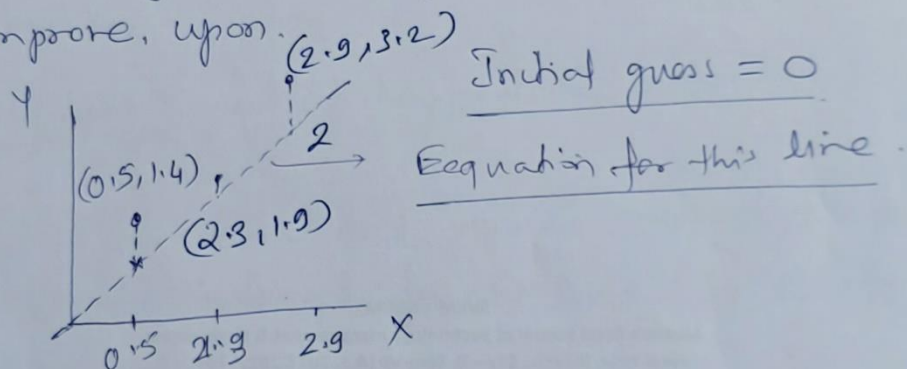
→ Currently assume that least square estimate for the slope $= 0.64$

$$\text{Predicted Height} = \text{Intercept} + 0.64 * \text{Weight}$$

& we will decide GD to find the optimal value for intercept.

→ Pick a random value for Intercept.

→ This is just initial guess that gives GD. Something to improve upon.



→ we will find out sum of residuals for this line.

In ML - sum of the squared Residuals is type of loss function - put this value in the eqn.

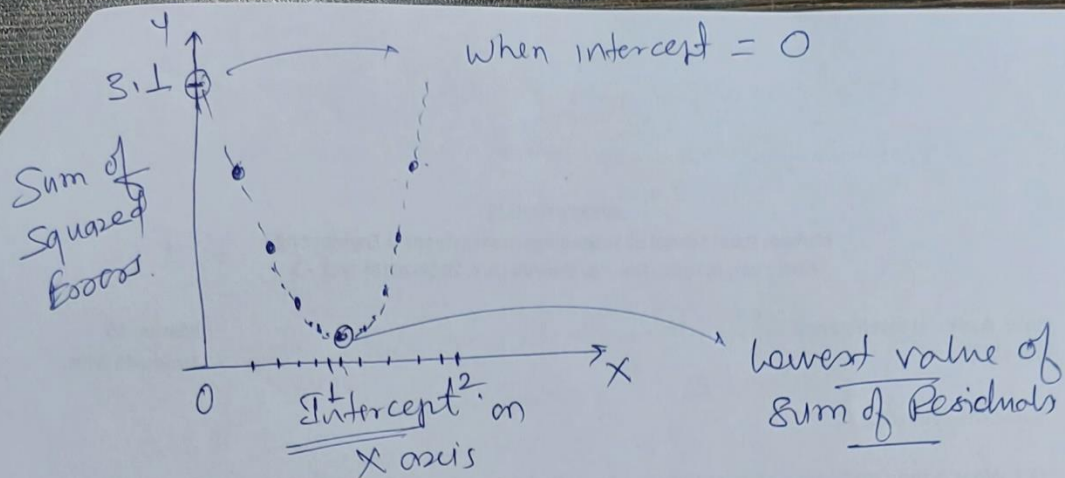
$$PH = 0 + 0.64 * 0.5 = \underline{\underline{0.32}}$$

$$\text{Residual} = OH - PH.$$

$$= 1.4 - 0.32$$

$$= 1.1 \rightarrow \text{Residuals.}$$

$$\begin{aligned} \text{Sum of Square Residual} &= (1.1)^2 + (0.4)^2 + (1.3)^2 \\ &= 3.1 \end{aligned}$$



plotting the point on graph for increasing value of residuals.

Is the best? what happens if best value is in between.

So, it is very slow and tedious method to find out the intercept for the line (by least square error method)

Don't ~~disparis~~ you think GD will be more efficient method to find intercept (best).

⇒ GD only does a "few calculations" far from the optimal solⁿ. and increases the number of calculations closer to the optimal value.

→ It takes big steps when it is far away and baby steps when it is close.

Let's understand SD to find optimal value of intercept starting from Random value (RV)

$$\underline{RV = 0} \Rightarrow$$

① $RV = 0 \Rightarrow$ When we calculated SSE

$$SSE = \overset{\text{First Point}}{(OV - PV)}^2 + \overset{\text{Second Point}}{(OV - PV)}^2 + \overset{\text{Third Point}}{(OV - PV)}^2$$

$$= (1.4 - (I + 0.64 * W))^2$$

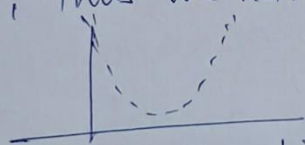
$PH = I + \overset{0.64}{S} * W$

$$\underline{SSE} = \frac{\sum (y - \bar{y})^2}{n}$$

$$= (1.4 - (I + 0.64 * 0.5))^2 + (1.9 - (I + 0.64 * 2.3))^2 + (3.2 - (I + 0.64 * 2.9))^2$$

Note: we can put any value for intercept and predict the new height.

Now we can put any values of Intercept and get SSE, Thus we have equation for this curve



\Rightarrow Now we can take derivative of this function and determine the slope at any value for the intercept.

\Rightarrow Let's take derivative of SSE with respect to the Intercept

$$\frac{d}{d(\text{Intercept})} SSE = \frac{d}{d(\text{Intercept})} (1.4 - (I + 0.64 * 0.5))^2 + \frac{d}{d(\text{Intercept})} \text{Second Point} + \frac{d}{d(\text{Intercept})} \text{Third point}$$

$$\frac{d}{d(\text{intercept})} (1.4 - (I + 0.64 * 0.5))^2 =$$

To take derivative of this, we need to apply
(Chain rule), \rightarrow moving square to the front

$$= 2 (1.4 - (I + 0.64 * 0.5)) * (-1)$$

$$= \cancel{2} (1.4 - I - \cancel{0.64 * 0.5}) * (-1)$$

do not content intercept
to derivative of Constant is zero

$$= -2 (1.4 - (I + 0.64 * 0.5)) +$$

Derivative of first point

$$- 2 (1.9 - (I + 0.64 * 2.3)) +$$

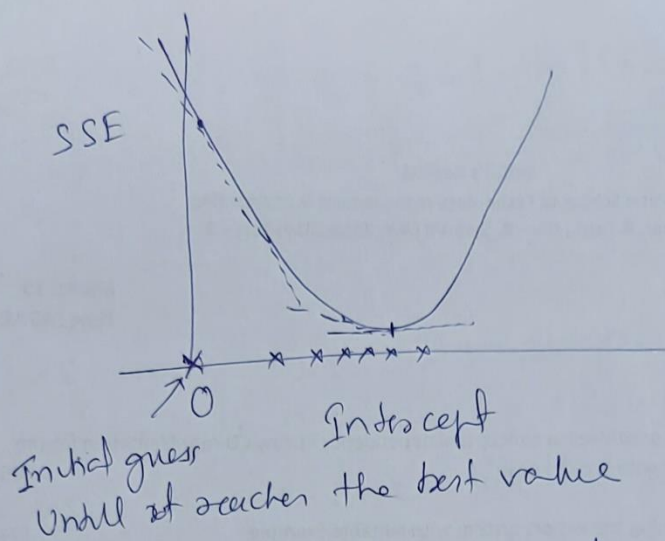
$$- 2 (3.2 - (I + 0.64 * 2.9)) \rightarrow \text{eqn. 1}$$

Now we have derivatives, so GD will use to find
where SSE is lowest.

Note
If we use LS method to solve for the optimal value
for intercept, we would simply find whether the
slope of curve = 0

— x —

In contrast GD find minimum value by taking
steps from initial guess until it reaches the
best value.



This makes GD very useful when it is not possible to solve for where the derivative = 0,
It is useful in many situations.

⇒ We started by taking intercept to the random number, in this case, it was 0.

put 0 in the equation of derivative

$$\begin{aligned} \frac{d}{d(\text{intercept})} \text{SSE} &= -2(1.4 - (0 + 0.64 \times 0.5)) + \\ &\quad -2(1.9 - (0 + 0.64 \times 2.3)) + \\ &\quad -2(3.2 - (0 + 0.64 \times 2.9)) \\ &= -5.7 \end{aligned}$$

when Intercept = 0 → slope of the curve is -5.7

Intercept =

Note - closer we get to the optimal value for the Intercept, the closer the slope of the curve gets to 0.

This means when the slope of the curve is close to 0, then we should take baby steps, because we are closer to the optimal value.

When the slope is far from zero then take bigger steps, because we are far from optimal value of 'intercept'.

\Rightarrow If we take huge step then increase the SSE so the size of the step should be related to the slope.

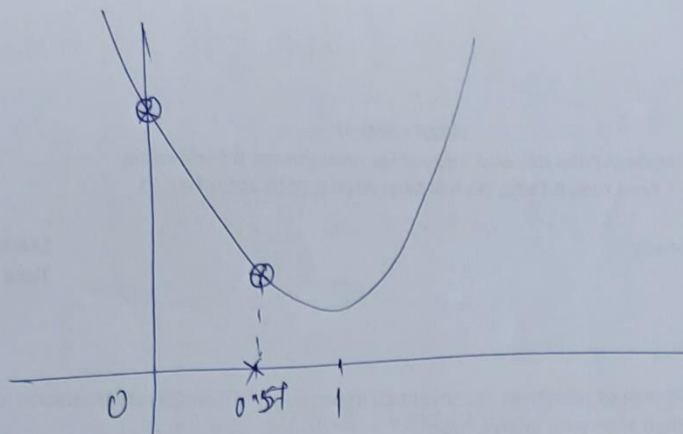
GD determine the step size by multiplying slope by a small number called "The learning Rate"

$$\text{Step Size} = -5.7 \times 0.1 = -0.57$$

When intercept was zero step size was ~~-5.7~~ -0.57

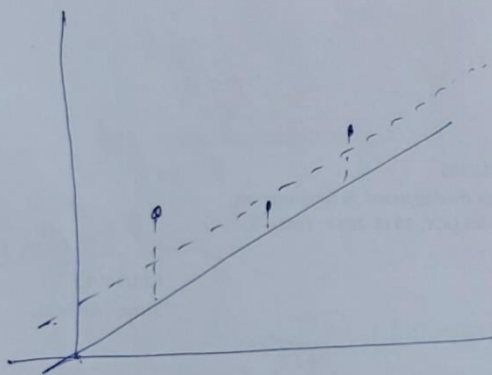
$$\begin{aligned} \text{New Intercept} &= \text{old intercept} - \text{step size} \\ &= 0 - (-0.57) = \underline{0.57} \end{aligned}$$

with step size we can calculate new "intercept"



In one big step we moved to much closer to the optimal value for the intercept.

⇒ Going back to the original data the original line with intercept 0 is moved much closer to the points.



Now put new intercept in derivative eqn.

$$\frac{d}{d(\text{intercept})} \text{SSE} = -2(1.4 - (0.57 + 0.64 \times 0.5)) +$$

$$-2(1.9 - (0.57 + 0.64 \times 2.3)) +$$

$$-2(3.2 - (0.57 + 0.64 \times 2.9)).$$

$$= \boxed{-2.3} \text{ Now this is the slope of the line}$$

Calculate the step size

$$\text{Step size} = \text{slope} \times \text{Learning Rate}$$

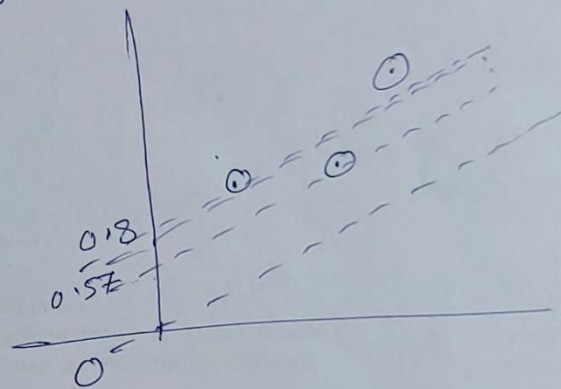
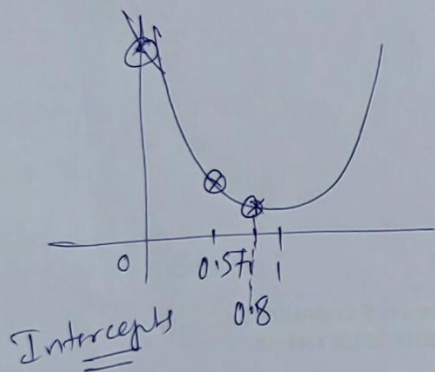
$$= -2.3 \times 0.1 = -0.23$$

$$\text{Step size} = -0.23$$

$$\text{New Intercept} = \text{old Intercept} - \text{step size}$$

$$= 0.57 - (-0.23)$$

$$= 0.8$$



Overall SSE getting smaller with the optimum value of intercept

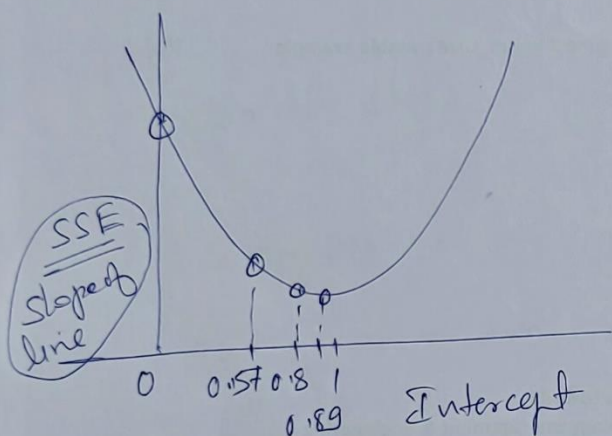
→ Now calculate derivative again -

$$\frac{d}{d(\text{intercept})} \text{SSE} = (-2(1.4 - (0.8 + 0.64 \times 0.59))) + (-2(1.9 - (0.8 + 0.64 \times 2.3))) + (-2(3.2 - (0.8 + 0.64 \times 2.9)))$$

$$= \boxed{-0.9} \text{ we get}$$

$$\begin{aligned}\text{Step Size} &= \text{slope} \times \text{LR} \\ &= -0.9 \times 0.1 = -0.09 \\ \text{Step Size} &= -0.09\end{aligned}$$

$$\begin{aligned}\text{New Intercept} &= 0.8 - (-0.09) \\ &= \boxed{0.89}\end{aligned}$$



$$\text{New Intercept} = 0.92$$

$$\text{New I} = 0.94$$

$$\text{New I} = 0.95$$

Step gets smaller and smaller when we get closer at the bottom.

After 6 steps GD estimate for the Intercept is 0.95

Note: The least Squares estimate for the intercept is also 0.95

So, GD has done a job, but without comparing it to gold standard, how does GD know to stop taking step?

GD stops when the step size is very close to 0.

⇒ Step size will be very close to 0, when the slope is very close to 0.

In practice minimum
step size = 0.001 or smaller.

So the slope = 0.009.

Step size = 0.009 × 0.1
= 0.0009 which is smaller than
0.1 (tolR).

So GD stops.

—x—

⇒ GD also includes a limit on the number of steps it will take before giving up.

In practice, ~~step size~~ =
maximum no. of steps = 1,000 or
greater

① → we have decided to use sum of squared Error SSE as loss function, to evaluate how well a line fits the data.

② Then we took the derivative of the SSE, In other words we took derivative of loss function.

③ Picked Random value for the intercept = 0

④ Calculate derivative when 'intercept was 0'
 (Slope)

⑤ put that slope into the step size calculations

$$\text{Step size} = \text{slope} \times LR$$

⑥ Calculated New Intercept

$$NI = OI - \text{Step size}$$

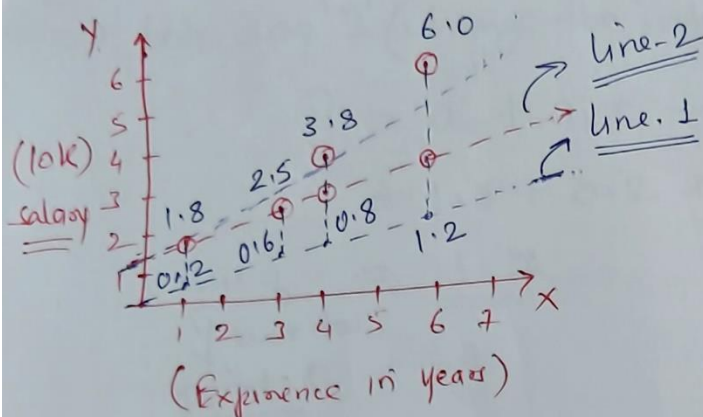
⑦ we put new intercept into the derivative and repeated everything until step size was close to 0.

LR using Gradient Descent Method (Example)

Data

(1, 1.8), (3, 2.5), (4, 3.8), (6, 6.0)

(1, 1.8), (3, 2.5), (4, 3.8), (6, 6.0)



Now we have to set LR model for this
→ straight line

eqn. $y = mx + c$

we can rewrite it as

$$y = a_0 + a_1 x$$

↓ ↓
Intercept slope

lets assume that $a_0 = 0$ & $a_1 = 0.2$

what will be the equation of line?

$$\bar{y} = 0 + 0.2x \quad \text{--- line-1}$$

Need to calculate the "errors" with respect to this line.

Formula for error = Observed value - Estimated value

$$\text{Error with point 1} = (y - \bar{y}) = (1.8 - 0.2)^2 = 2.56$$

$$\text{Sum of squared Error} = (1.8 - 0.2)^2 + (2.5 - 0.6)^2$$

$$(\text{Loss function}) + (3.8 - 0.8)^2 + (6.0 - 1.2)^2$$

$$= 2.56 + 3.61 + 9 + 28.84$$

$$= 38.21$$

Value of loss function with respect to line-1

$$SSE = \boxed{38.21} \Rightarrow \text{when } a_0 = 0 \text{ \& } a_1 = 0.2$$

\Rightarrow let's say 'I' (change the intercept) as 1.5

$$\bar{y} = a_0 + a_1 x \rightarrow \text{line-2}$$

$$\bar{y} = 1.5 + 0.2x \quad \text{--- eqn of the line.}$$

Value of \bar{y} w.r to point-1
1.7

$$\begin{aligned} \text{Sum of Squared Error} &= (y_1 - \bar{y}_1)^2 + (y_2 - \bar{y}_2)^2 + \\ (\text{loss function}) &\quad (y_3 - \bar{y}_3)^2 + (y_4 - \bar{y}_4)^2 \end{aligned}$$

$$= \cancel{(1.7)}^2$$

$$= (1.8 - 1.7)^2 + (2.5 - 2.1)^2 + (3.8 - 2.3)^2 + (6 - 2.9)^2$$

$$= 0.01 + 0.16 + 2.25 + 9.61$$

$$SSE = \boxed{12.03} \Rightarrow \text{when } a_0 = 1.5 \text{ \& } a_1 = 0.2$$

Now we will change the slope from 0.2 to 0.6

$$\bar{y} = 1.5 + 0.6x \quad \text{--- line-3}$$

$$\begin{aligned}\text{loss function} &= (y_1 - \bar{y}_1)^2 + \dots \\ &= (1.8 - 2.1)^2 + (2.5 - 3.3)^2 \\ &\quad + (3.8 - 3.9)^2 + (6 - 5.1)^2 \\ &= (0.09) + (0.64) + (0.01) + (0.81)\end{aligned}$$

SSE with respect to all the points = 1.55 \Rightarrow When $a_0 = 1.5$ & $a_1 = 0.6$

That means we require both (Intercept & slope)

How do we get this line directly by

GD method

What will be the General eqn of loss function

$$\bar{y} = I + Sx$$

$$\begin{aligned}\text{SSE} &= (y_1 - \bar{y}_1)^2 + \dots + (y_n - \bar{y}_n)^2 \\ &= (y_1 - (I + Sx))^2 + (y_2 - (I + Sx))^2 \\ &\quad + (y_3 - (I + Sx))^2 + (y_4 - (I + Sx))^2\end{aligned}$$

Now we need to take derivative of loss function with r.t. to I & S separately.

$$\begin{aligned}\text{SSE} &= (1.8 - (I + S(1)))^2 + (2.5 - (I + S \times 3))^2 \\ &\quad + (3.8 - (I + S \times 4))^2 + \\ &\quad (6 - (I + S \times 6))^2\end{aligned}$$

step 1

Derivative w.r. to Intercept when Slope is Constant.

$$\frac{d(SSE)}{d(\text{Intercept})} = -2(1.8 - (I + S \times 1)) + \\ -2(2.5 - (I + S \times 3)) + \\ -2(3.8 - (I + S \times 4)) + \\ -2(6 - (I + S \times 6)).$$

$$\frac{d}{d(\text{Slope})} = -2 \times (1.8 - (I + S)) + \\ -2 \times 3(2.5 - (I + S \times 3)) + \\ -2 \times 4(3.8 - (I + S \times 4)) + \\ -2 \times 6(6 - (I + S \times 6))$$

step 2

Start with random values of I & S.

$$I = 0 \text{ \& } S = 0.2$$

$$\frac{d(SSE)}{d(I=0)} = ?$$

we will get
Intercept (I)
(please calculate)

$$\frac{d(SSE)}{d(S=0.2)} = ?$$

we will get
slope (S),
(please calculate)

step 3 Need to calculate step size for
Intercept & slope

$$\text{Step Size (I)} = \text{value obtained in step 2} \times \text{LR}$$

$$\text{Step Size (S)} = \text{---} - \text{---}$$

LR is
random

0.01

Need to
Calculate.

Step 4

LR \rightarrow Learning rate that
define the step size.

Calculate New Intercept

$$\text{New (I)} = \text{old (I)} - \text{step size (I)}$$

$$\text{New (I)} = 0 -$$

$$\text{New (S)} = \text{old (S)} - \text{step size (S)}$$

Step 5 put new values of I & S in
eqn of derivative of I & S again
and repeat the steps still we reach.

SSE = 0 \approx or Some random number
of iterations.

Use of Loss Function to calculate LR;

→ we used SSE as loss function. (Example),

Data

X	Y
0.5	1.4
2.3	1.9
2.9	3.2

$$SSE = (1.4 - (I + S \times 0.5))^2 +$$

$$(1.9 - (I + S \times 2.3))^2 +$$

$$(3.2 - (I + S \times 2.9))^2 \quad \text{--- (2)}$$

We need to find out value of the Intercept & slope, such that, SSE would be minimum.

⇒ Also, we need to take derivative with respect to slope.

$$\begin{aligned} \frac{d}{d(\text{Intercept})} SSE &= (1.4 - (I + S \times 0.5))^2 \\ &= 2(1.4 - (I + S \times 0.5)) \times (-1) \end{aligned}$$

Since we are taking derivative with respect to Intercept, slope would be considered as constant and derivative of Constant is 0.

⇒ So we get (-1) just like before

$$\begin{aligned} \frac{d}{d(\text{Intercept})} SSE &= -2(1.4 - (I + S \times 0.5)) \\ &\quad + (-2(1.9 - (I + S \times 2.3))) + \\ &\quad (-2(3.2 - (I + S \times 2.9))). \end{aligned}$$

$$\frac{d}{d(\text{slope})} \text{SSE} = \frac{d}{d(\text{slope})} (1.4 - (I + 5 \times 0.5))^2$$

$$+ \frac{d}{d(\text{slope})} (1.9 - (I + 5 \times 2.3))^2$$

$$+ \dots$$

~~$$= 2(1.4 - (I + 5 \times 0.5)) \times 0.5$$~~

$$\frac{d(\text{SSE})}{d(\text{slope})} = -2 \times 0.5 (1.4 - (I + 5 \times 0.5))$$

$$+ -2 \times 2.3 (1.9 - (I + 5 \times 2.3))$$

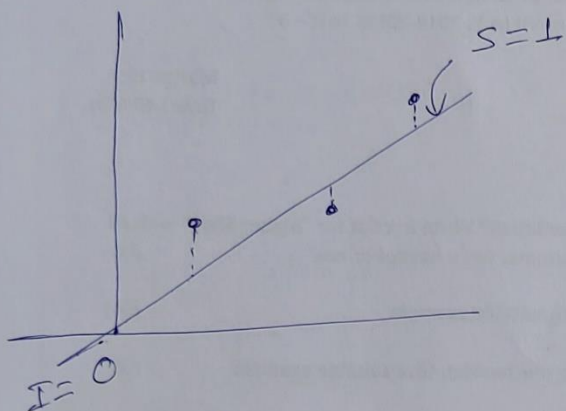
$$+ -2 \times 2.9 (1.9 - (I + 5 \times 2.9))$$

Note : When you have two or more derivatives of the same function, they are called as Gradient

We use this Gradient to Descend to lowest point in the loss function, (SSE).

That's why it is called as Gradient Descent

Now we start with Random value of I & S
 $I = 0$ & $S = 1$



Now put the values of
 $I = 0$ & $S = 1$ in
 Both the equations.

$$\frac{d}{d(I)} SSE = -1.6 \quad \text{we get} \quad (\text{Intercept})$$

$$\frac{d}{d(S)} SSE = -0.8 \quad \text{we get} \quad (\text{slope})$$

$$\Rightarrow \text{Step Size}(I) = -1.6 \times LR$$

$$\Rightarrow \text{Step Size}(S) = -0.8 \times LR \quad \text{LR} = \underline{\underline{0.01}}$$

\Rightarrow Large LR will not work in this case.

$$SS(I) = -1.6 \times 0.01 = -0.016$$

$$SS(S) = -0.8 \times 0.01 = -0.008$$

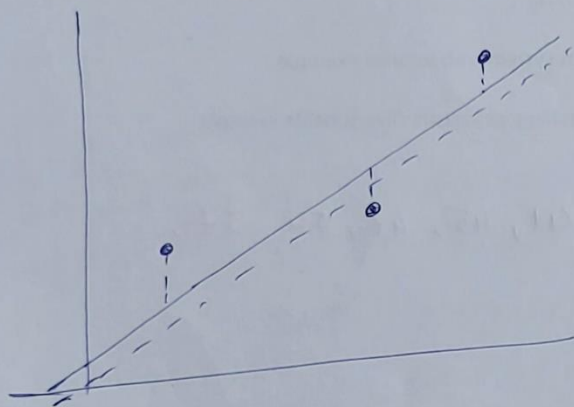
GD is very sensitive to the LR.

In practical situation LR can be large Initially and getting smaller at each step.

Now, we calculate,

$$\begin{aligned}\text{New Intercept} &= \text{old } I - \text{step size} \\ &= 0 - (-0.016) = 0.016\end{aligned}$$

$$\begin{aligned}\text{New slope} &= \text{old } S - \text{step size} \\ &= 1 - (-0.008) \\ &= 1.008\end{aligned}$$



Now we have to repeat the same process
until all the steps sizes are very small, or
we reach the maximum number of steps

we have the fit out best fitting line

$$\left. \begin{aligned}I &= \underline{\underline{0.95}} \\ S &= \underline{\underline{0.64}}\end{aligned} \right\}$$

Same value we
get from
least squares.

Now, we have understood how GD will optimize.
two parameters the slope and Intercept.