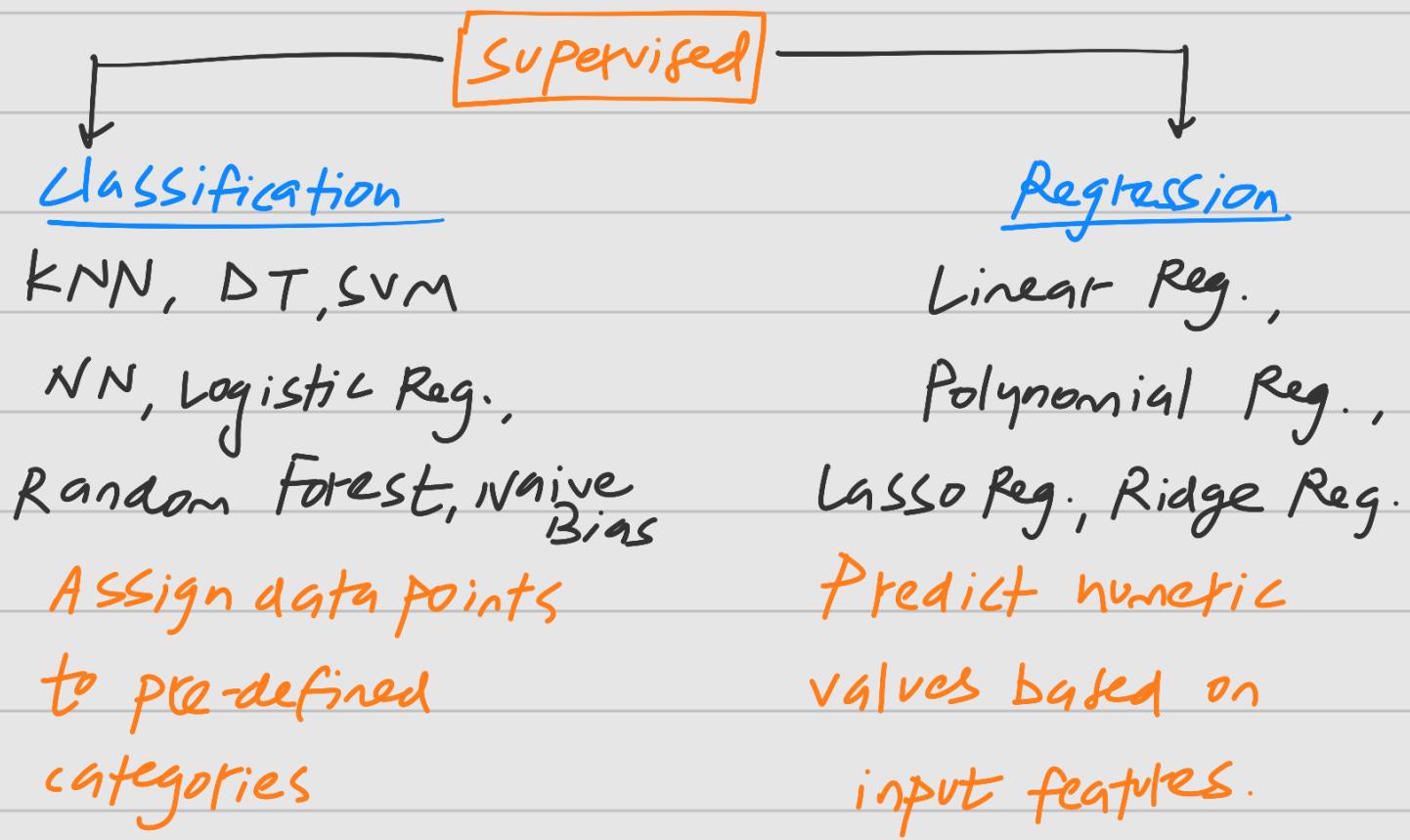


Supervised Learning

- Teaching a computer how to do something by showing it Examples.
- Model learns from labeled training data.
- Eg: E-mail spam detection

Emails would be labeled as spam/not-spam



Unsupervised Learning

- unlabeled data is used to train the model. ↓
input without labels
- Algorithms:
 - PCA
 - K-means clustering
 - Hierarchical clustering (BIRCH)
 - DBSCAN clustering

Regression

- Help in understanding & predicting relationship

b/w multiple Variables

(input) independent

The basis for
prediction.

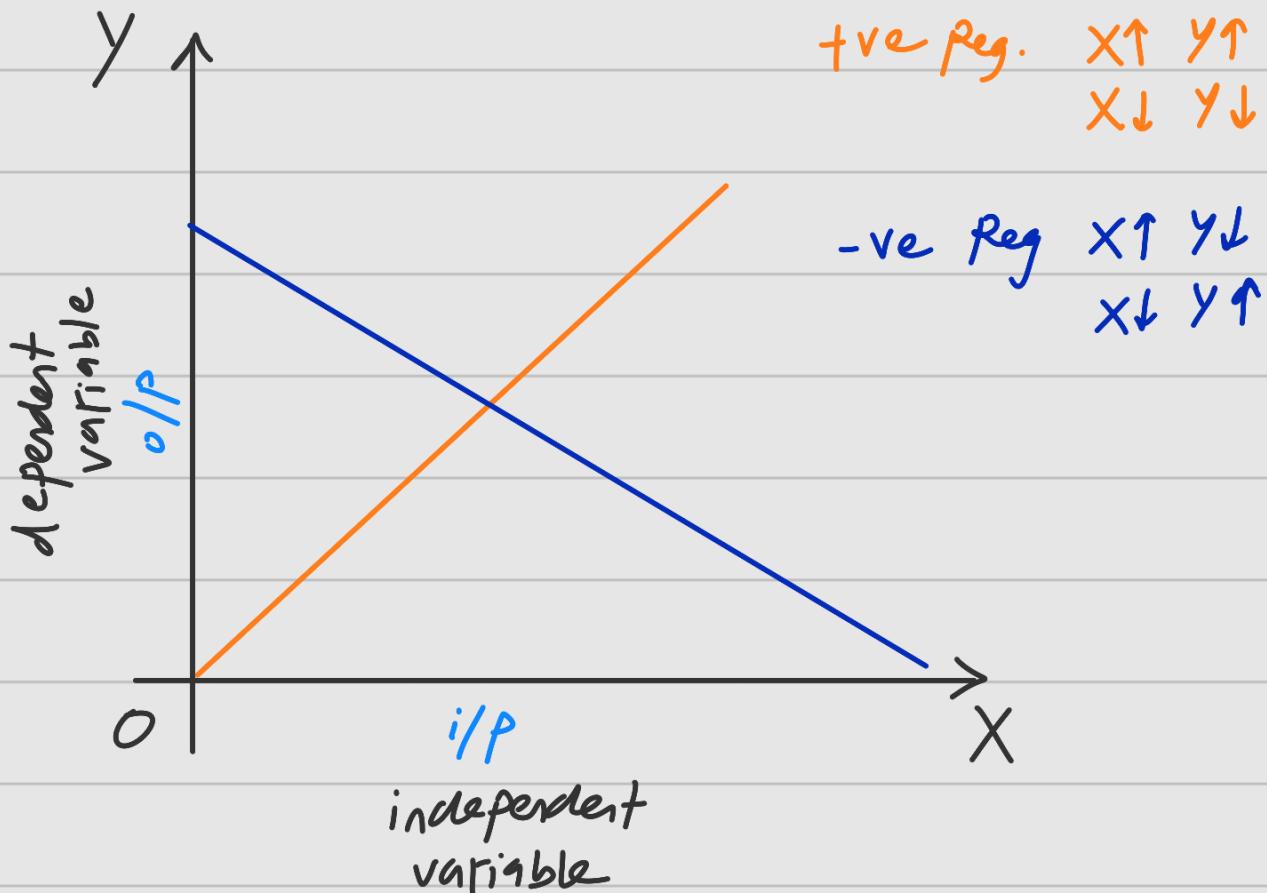
dependent (output)

to be predicted

Eg: Predicting salary based on years of exp.

↓
dependent

↓
independent



Linear Regression

Eqn. of linear Regression :

$$y = mx + b$$

$y \rightarrow$ dependent variable

$x \rightarrow$ independent variable

$m \rightarrow$ slope of line

(change in y for unit change in x)

$b \rightarrow$ intercept

(value of y , when x is 0)

Least Square

Eg: Predicting Pizza prices based on diameter (least square method)

<u>Diameter (X)</u>	<u>Price (Y)</u>
8	10
10	13
12	16

Sol'n $y = mx + b$

calculating m & b first

x_i	y_i	$\bar{x} - x_i$	$\bar{y} - y_i$	$(\bar{x} - x_i)^2$	$(\bar{x} - x_i)(\bar{y} - y_i)$
8	10	2	3	4	6
10	13	0	0	0	0
12	16	-2	-3	4	6
$\sum x = 30$	$\sum y = 39$	$\sum = 0$	$\sum = 0$	$\sum = 8$	$\sum = 12$
$\bar{x} = 10$	$\bar{y} = 13$				

$$m = \frac{\sum (\bar{x} - x_i)(\bar{y} - y_i)}{\sum (\bar{x} - x_i)^2} = \frac{+2^3}{+2^2} = \boxed{\frac{3}{2}} //$$

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

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

$$b = 13 - (1.5)(10)$$

$$b = 13 - 15$$

$$\boxed{b = -2} //$$

Best fit line

$$\boxed{y = 1.5x - 2} //$$

Error / loss \rightarrow MSE

$$\frac{1}{n} \sum_{i=1}^n (y_i - y_p)^2$$

↓ *predicted*
Actual

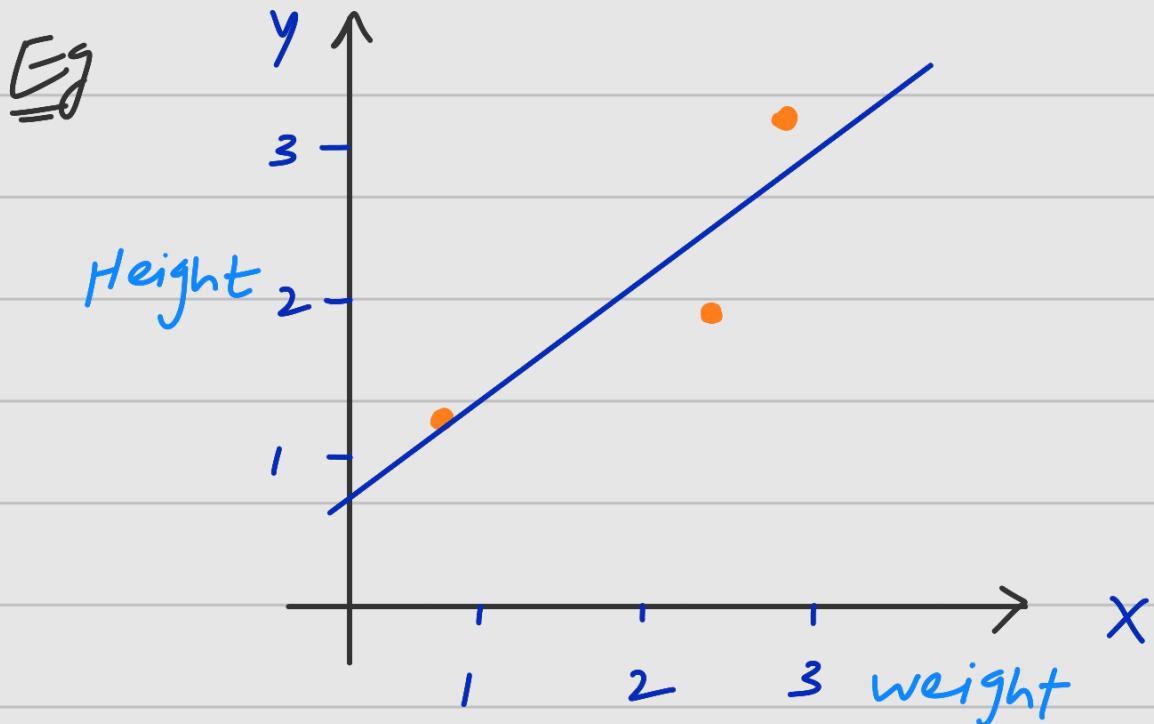
Gradient Descent

works for any differential model

Suitable for large datasets

"used for optimization"

loss $f \rightarrow$ sum of squared residuals
 $\sum (\text{observed} - \text{predicted})^2$



Predicted = (slope \times weight) + intercept
Height

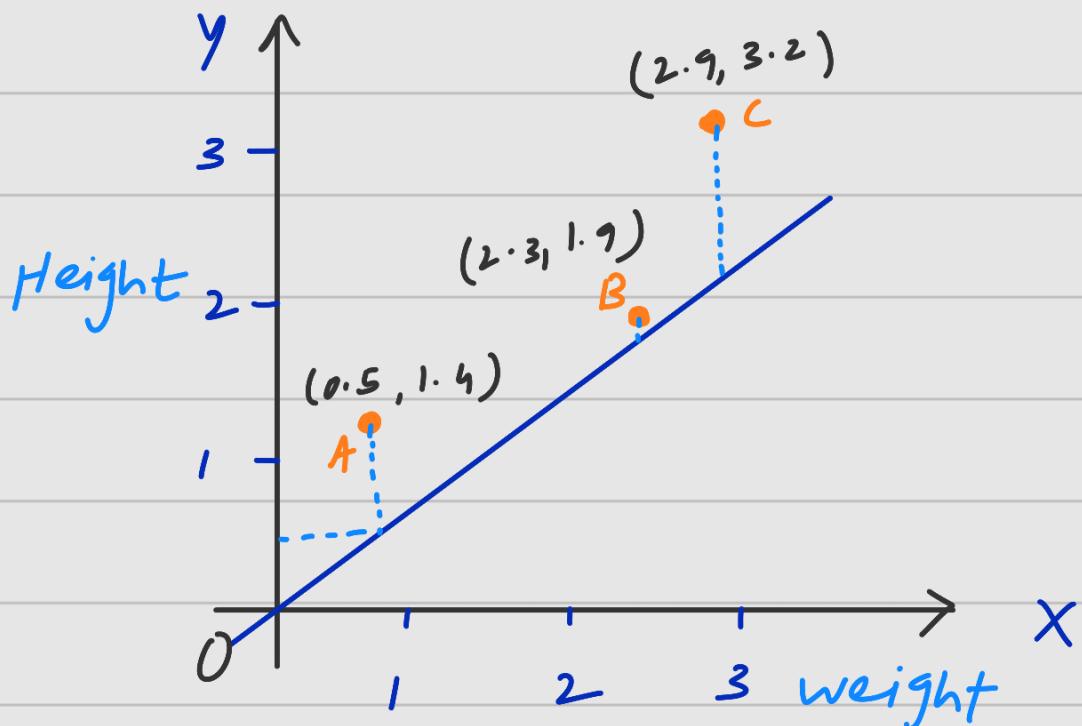
$$[y = mx + b]$$

Let's add least square estimate of
slope = 0.64

Predicted Height = $0.64 \times \text{weight} + \text{intercept}$

- we'll use gradient descent to find optimal value for intercept
- Firstly we pick random value for intercept.
- This is just initial guess that gives Gradient Descent to improve upon.

Let, intercept = 0



Point A → Actual height = 1.4
Actual weight = 0.5

$$\begin{aligned}\text{Pred. Height} &= 0.64 \times \overbrace{\text{weight}}^0.5 + 0 \\ &= 0.64 \times 0.5 = 0.32\end{aligned}$$

$$\text{Residual}(A) = \text{Actual Height} - \text{Predicted Height}$$
$$= 1.4 - 0.32 = 1.1$$

Point B → Actual height = 1.9

Actual weight \leftarrow = 2.3

$$\text{Pred. Height} = 0.64 \times \text{weight} + 0$$
$$= 0.64 \times 2.3 = 1.472$$

$$\text{Residual}(B) = \text{Actual} - \text{Pred. Height}$$
$$= 1.9 - 1.472 \approx 0.4$$

Point C → Actual Height = 3.2

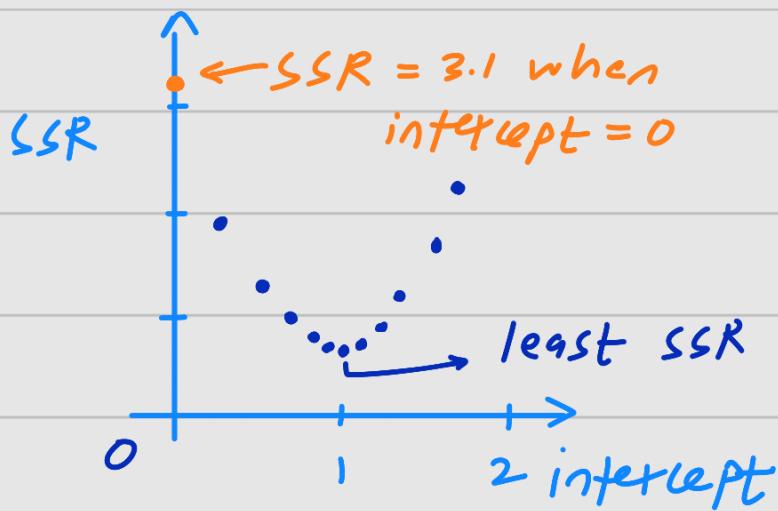
Actual weight = 2.9

$$\text{Pred. Height} = 0.64 \times \text{weight} + 0$$
$$= 0.64 \times 2.9 = 1.856$$

$$\text{Residual}(C) = \text{Actual} - \text{Pred. Height}$$
$$= 3.2 - 1.856 \approx 1.3$$

$$\text{Sum of Squared Residuals} = 1.1^2 + 0.4^2 + 1.3^2$$

$$\cong 3.1$$



For increasing values of intercept we get diff points

Notice Big steps when its far away

Sum of Squared = \sum (Actual - Predicted Height)²
Residuals

$$\begin{aligned}
 &= (1.4 - (\text{intercept} + 0.64 \times 0.5))^2 \\
 &+ (1.9 - (\text{intercept} + 0.64 \times 2.3))^2 \\
 &+ (3.2 - (\text{intercept} + 0.64 \times 2.9))^2
 \end{aligned}$$

Now we can put any value for intercept & get sum of square residuals

We can take derivative of this fn & take find slope at any value of intercept.

$$\frac{d(\text{SSR})}{d \text{ intercept}} = \frac{d[(1.4 - (\text{intercept} + 0.64 \times 0.5))^2]}{d \text{ intercept}} +$$

$$\frac{d[(1.9 - (\text{intercept} + 0.64 \times 2.3))^2]}{d \text{ intercept}} +$$

$$\frac{d[(3.2 - (\text{intercept} + 0.64 \times 2.9))^2]}{d \text{ intercept}}$$

$$\frac{d[(1.4 - (\text{intercept} + 0.64 \times 0.5))^2]}{d \text{ intercept}}$$

$$2(1.4 - (\text{intercept} + 0.64 \times 0.5))(-1)$$

$$\frac{d(1.4 - \text{intercept} - k)}{d \text{ intercept}}$$

0 - 1 - 0

(-1) —

$$-2(1.4 - (\text{intercept} + 0.64 \times 0.5))$$

$$-2(1.9 - (\text{intercept} + 0.64 \times 2.3))$$

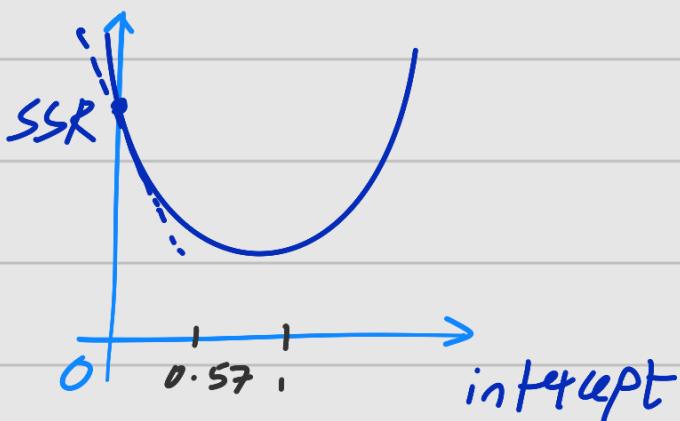
$$-2(3.2 - (\text{intercept} + 0.64 \times 2.9))$$

Grad. descent will use it to find where the SSR is lowest.

As intercept = 0, push it in derivative

$$\left. \begin{aligned} & -2(1.4 - (0 + 0.64 \times 0.5)) \\ & -2(1.4 - (0 + 0.64 \times 2.3)) \\ & -2(3.2 - (0 + 0.64 \times 2.9)) \end{aligned} \right\} = -5.7$$

for intercept = 0 ; slope = -5.7



When slope is far from zero take big steps.
When slope is near to zero take small steps

[Step size = Slope \times learning Rate]

When intercept = 0

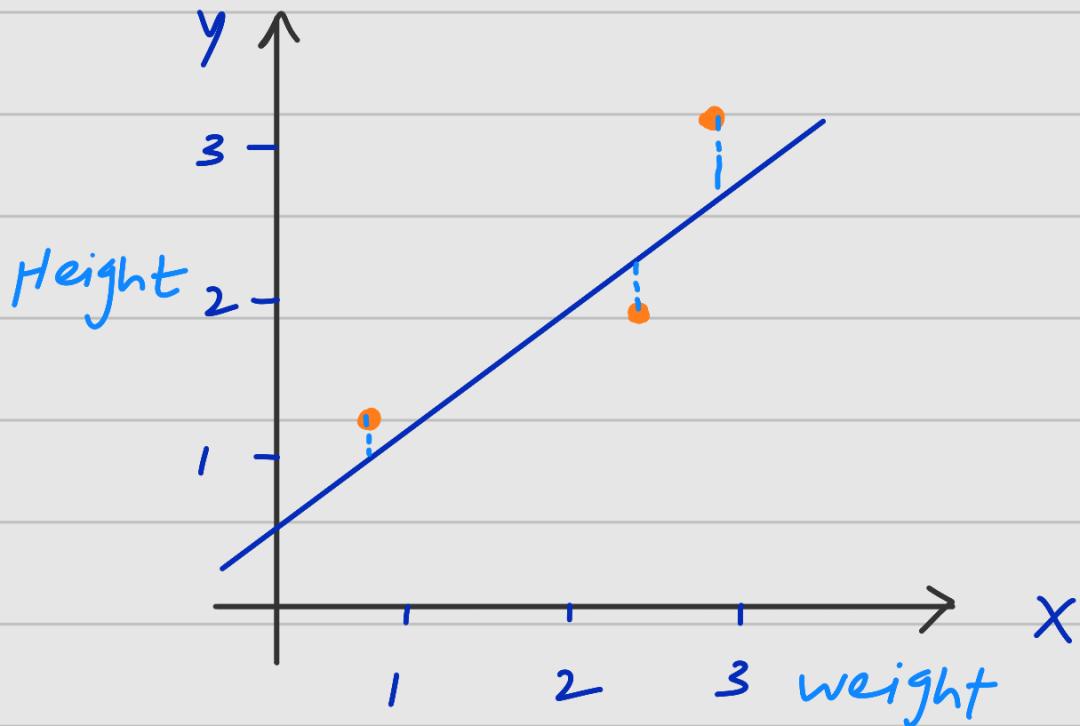
Slope = -5.7

$$\begin{aligned} \text{Step Size} &= -5.7 \times 0.1 \\ &= -0.57 \end{aligned}$$

[New intercept = old intercept - step size]

$$\text{"} \quad = 0 - (-0.57)$$

New intercept = 0.57



At intercept = 0.57, the residuals shrunked

Push new intercept = 0.57 in derivative

$$\left. \begin{aligned} & -2(1.4 - (0.57 + 0.64 \times 0.5)) \\ & -2(1.9 - (0.57 + 0.64 \times 2.3)) \\ & -2(3.2 - (0.57 + 0.64 \times 2.9)) \end{aligned} \right\} = -2.3$$

Slope = -2.3 at intercept = 0.57

$$\begin{aligned} \text{step size} &= \text{learning rate} \times \text{slope} \\ &= 0.1 \times (-2.3) \\ &= -0.23 \end{aligned}$$

$$\begin{aligned}\text{new intercept} &= \text{old intercept} - \text{step size} \\ &= 0.57 - (-0.23)\end{aligned}$$

$$\text{new intercept} = 0.80$$

Residuals shrunked even more at

$$\text{intercept} = 0.80$$

Push intercept = 0.8 in derivative

$$\left. \begin{array}{l} -2(1.4 - (0.8 + 0.64 \times 0.5)) \\ -2(1.9 - (0.8 + 0.64 \times 2.3)) \\ -2(3.2 - (0.8 + 0.64 \times 2.9)) \end{array} \right\} = -0.9$$

$$\text{slope} = -0.9 ; \text{ at intercept} = 0.8$$

$$\text{Step size} = \text{slope} \times \text{l. rate}$$

$$= -0.9 \times 0.1 = -0.09$$

$$\begin{aligned}\text{new intercept} &= \text{old intercept} - \text{step size} \\ &= 0.8 - (-0.09) \\ &= 0.8 + 0.09\end{aligned}$$

$$\text{new intercept} = 0.89$$

↓
going on

$\text{new intercept} = 0.92$
 $\text{new intercept} = 0.94$
 $\text{new intercept} = 0.95$

After 6 steps
Gradient descent estimate
for intercept = 0.95

The least square estimate for intercept is also 0.95

[Gradient descent will stop when step size is very close to zero. i.e when slope is very close to zero]

In practice min step size = 0.001

Optimization for 2 or more Variables

- We take derivative of SSR wrt intercept as earlier (slope as K)
- Then we take derivative of SSR wrt slope (considering intercept as K)

$$\frac{d(\text{SSR})}{d(\text{slope})} = -2(0.5)(1.4 - (\text{intercept} + \text{slope} \times 0.5))$$

$$- 2(2.3)(1.9 - (\text{intercept} + \text{slope} \times 2.3))$$

$$- 2\cancel{(2.9)}(3.2 - (\text{intercept} + \text{slope} \times 2.9))$$

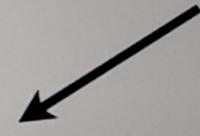
↓
 weight of
 a point

$$\frac{d}{d \text{ intercept}} \text{ Sum of squared residuals} =$$

$$-2(1.4 - (\text{intercept} + \text{slope} \times 0.5))$$

$$+ -2(1.9 - (\text{intercept} + \text{slope} \times 2.3))$$

$$+ -2(3.2 - (\text{intercept} + \text{slope} \times 2.9))$$



$$\frac{d}{d \text{ slope}} \text{ Sum of squared residuals} =$$

$$-2 \times 0.5(1.4 - (\text{intercept} + \text{slope} \times 0.5))$$

$$+ -2 \times 2.9(3.2 - (\text{intercept} + \text{slope} \times 2.9))$$

$$+ -2 \times 2.3(1.9 - (\text{intercept} + \text{slope} \times 2.3))$$

NOTE: when we have 2 or more derivative of a same $f''(\text{SSR})$, they are called gradient

Again we'll pick intercept = 0 & slope = 1
& push it in derivative

$$\frac{d}{\text{intercept}} \text{ Sum of squared residuals} =$$
$$-2(1.4 - (0 + 1 \times 0.5))$$
$$+ -2(1.9 - (0 + 1 \times 2.3))$$
$$+ -2(3.2 - (0 + 1 \times 2.9))$$
$$= -1.6$$

Slope

$$\frac{d}{\text{slope}} \text{ Sum of squared residuals} =$$
$$-2 \times 0.5(1.4 - (0 + 1 \times 0.5))$$
$$+ -2 \times 2.9(3.2 - (0 + 1 \times 2.9))$$
$$+ -2 \times 2.3(1.9 - (0 + 1 \times 2.3))$$
$$= -0.8$$

$$LR = 0.01$$

$$\text{step size (intercept)} = LR \times \text{slope}$$
$$= 0.01 \times (-1.6) = -0.016$$

$$\text{step size (slope)} = LR \times \text{slope}$$
$$= 0.01 \times (-0.8) = -0.008$$

$$\begin{aligned}\text{new intercept} &= \text{old} - \text{step size} \\ &= 0 - (-0.016) = 0.016\end{aligned}$$

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

Repeat this until max no. of steps reached OR step sizes becomes 0.001

Here the best-fit line, with intercept = 0.95
slope = 0.64, the same values we get from least square method

Tip: Standardize the data before implementing Gradient descent otherwise it can cause huge/unstable steps during training. This makes GD Explode (diverge) instead of converging

Converge: GD keeps improving model step by step until it reaches minimum loss(error)

Diverge: GD goes in wrong direction / takes too long steps, making error/loss worse

