

# Gradient Descent

\* optimizer

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \dots$$

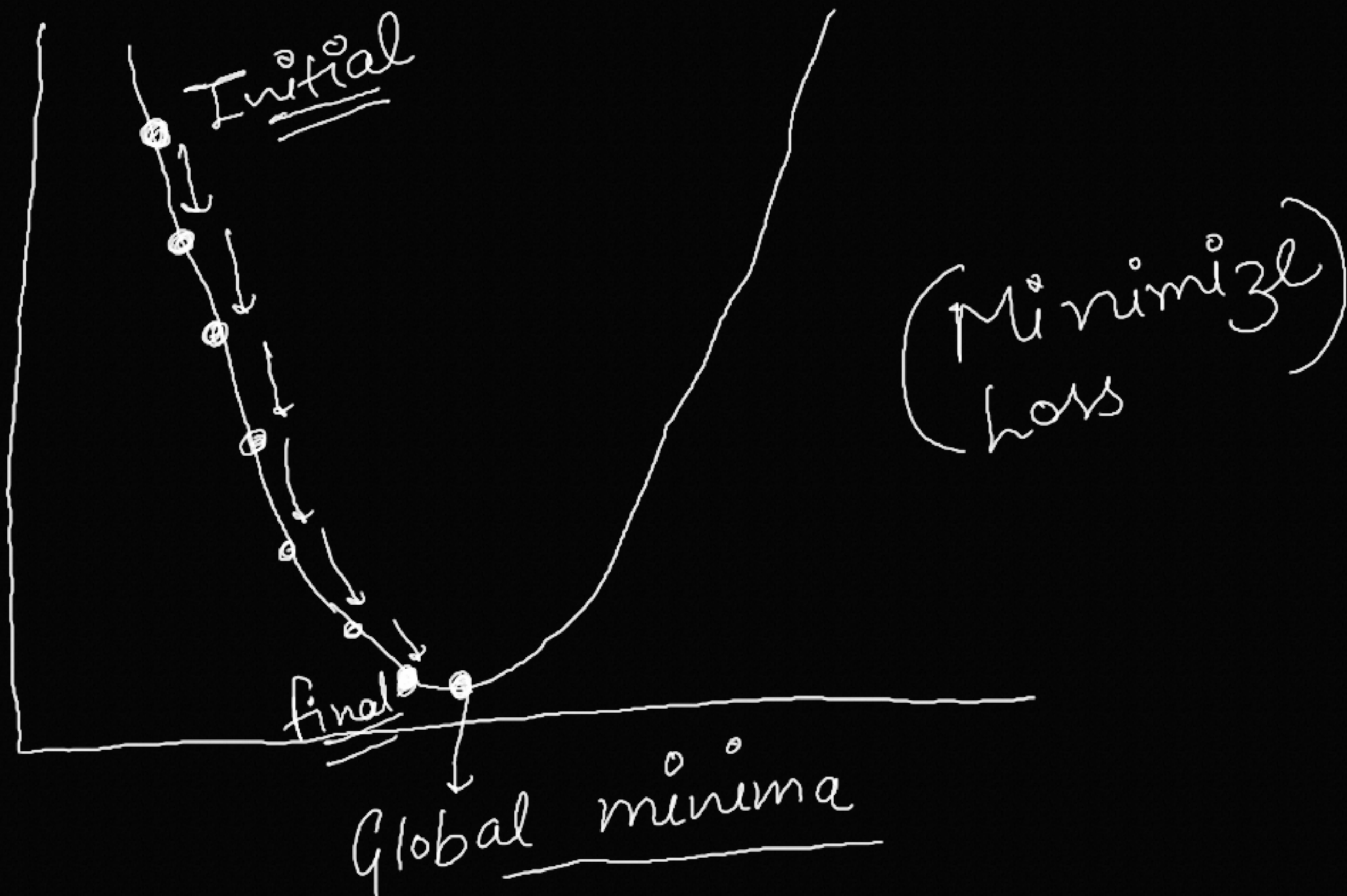
coefficient | weights

(To find best values)  
To minimize loss

$$\begin{aligned} \text{Actual} &= 5 \\ \text{Predicted} &= 4.8 \\ \text{Loss} &= \underline{\underline{0.2}} \end{aligned}$$

$$y_{\text{actual}} - y_{\text{predicted}}$$

# General



# Mathematical - 1

$$f(x) = (x-3)^2$$

(y)

QD  $\rightarrow$  minimize our fn

We know,  $x=3$

$$y = (3-3)^2$$
$$= \underline{\underline{0}}$$

let see how QD works.

Step ① :- Initial guess,  
lets say our guess,  $x=0$ .

Step ② :- Calculate gradient (partial derivative)  
w.r.t to coefficient  
( $x$ ).

$$f(x) = (x-3)^2$$

$$\frac{df}{dx} = 2(x-3) \Rightarrow \text{gradient}$$

Step ③ :- Update Rule

$$x_{\text{new}} = x_{\text{old}} - \alpha \cdot \text{gradient}$$

What is alpha ( $\alpha$ )?

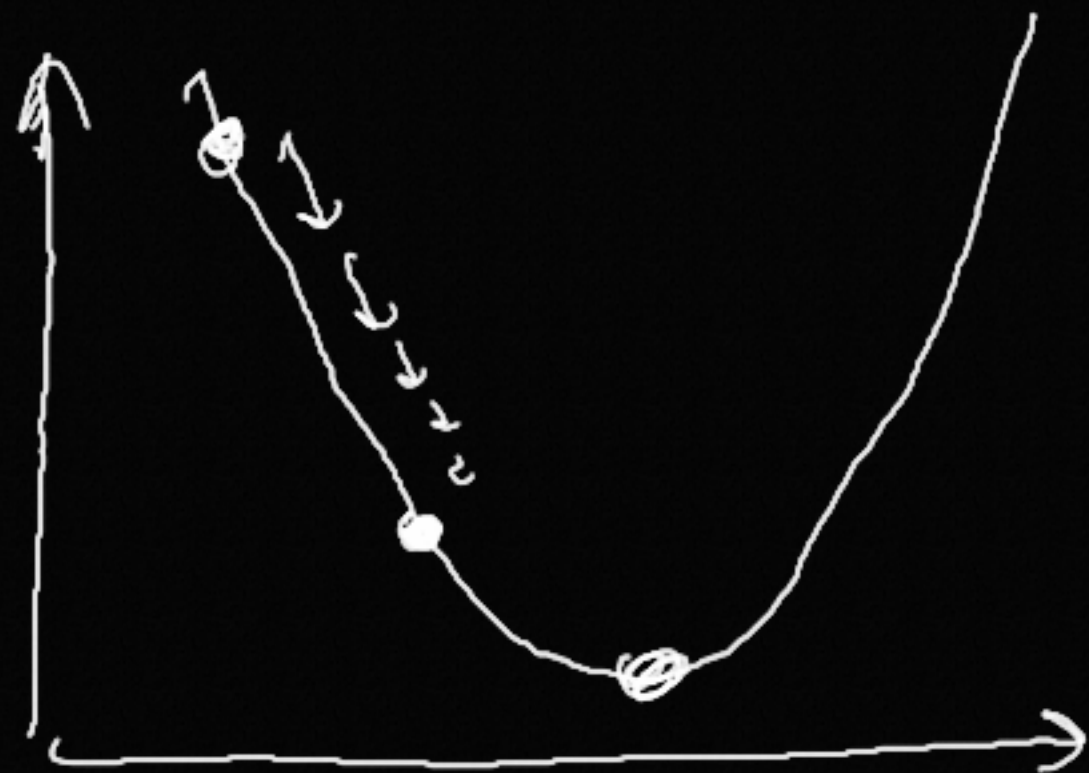
$\alpha \rightarrow$  small positive number generally  $\begin{pmatrix} 0.01, 0.001 \\ 0.1 \end{pmatrix}$

It should be optimal acc to data



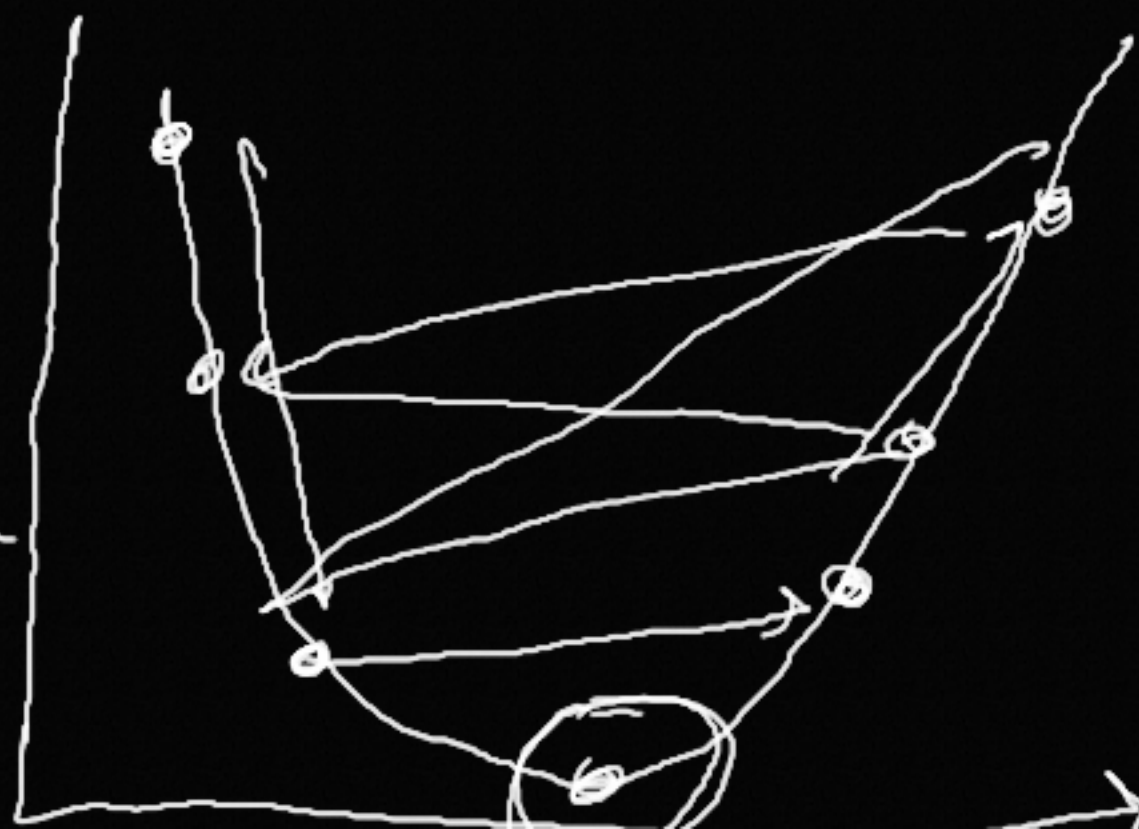
If  $\alpha \rightarrow$  too small  $\rightarrow$  Vanishing Gradient  
(0.0000001)

It takes forever to  
reach global minima



$\alpha \rightarrow$  too large  $\rightarrow$  Exploding Gradient  
(10, 100)

It will jump here &  
there but will not reach  
global minima



We select  $\alpha = \underline{\underline{0.1}}$

④ Repeat process until we minimize cost  
for

$$\text{loss} = \underline{\underline{3}}$$

Iteration 1 :-  $x = 0$

$$\frac{df}{dx} = 2(x-3) = 2(0-3) = \underline{\underline{-6}}$$

$$x_{\text{new}} = 0 - 0.1 \times (-6)$$
$$= \underline{\underline{0.6}}$$

$$\text{loss} = \underline{\underline{2.4}}$$

Iteration 2  $\Rightarrow x = 0.6$

$$\frac{df}{dx} = 2(x-3) = 2(0.6-3) = -\underline{\underline{4.8}}$$

$$x_{\text{new}} = 0.6 - 0.1 \times (-4.8)$$
$$= \underline{\underline{1.08}}$$

Loss = 2

Iteration 3  $\Rightarrow x = 1.08$

$$\frac{df}{dx} = -\underline{\underline{3.84}}$$

$$x_{\text{new}} = 1.08 - 0.1(-3.84) = \underline{\underline{1.464}}$$

Loss = 1.5



Repeat these process  
until  $x=3$ . or  $\checkmark$   
close to 3

OR  
Your loss will be  $\checkmark$   
closer to zero

Data

<u><math>X_1</math></u>	<u><math>X_2</math></u>	<u><math>y</math></u>
1	2	5

ML Mathematics  
 we set  $\alpha = \underline{0.01}$

Step ①  $\Rightarrow$  linear model.  

$$y = \beta_1 X_1 + \beta_2 X_2$$

Our  
GD start

Step 2 Initial guess

$$\beta_1 = 0, \beta_2 = 0$$

Do prediction

$$\hat{y}_1 = 0 \times 1 + 0 \times 2$$
$$= \underline{\underline{0}}$$

$$\text{Actual} = 5$$

$$\text{Loss} = 5 - 0$$
$$= \underline{\underline{5}}$$

Not good  
Weights are not the best  
Loss should be close to zero

Step ③ :- Calculate gradient  
•  $\beta_1$ ,  $\beta_2$   $\rightarrow$  two gradient

$\omega^{st.}$   
 $\beta_1$   $\Rightarrow \frac{dh}{d\beta_1} = \underbrace{\left(-\frac{2}{n}\right)}_{\text{normalization factor}} \sum (y - \hat{y}) x_i$

$$= \underline{-2} \times (5 - 0) \times 1$$

$$= \underline{\underline{-10}}$$

w.r.t  $\underline{\underline{\beta_2}} \rightarrow \frac{dL}{d\beta_2} = \frac{-2}{n} \sum (y - \hat{y}) x_2$

$$= -2 (5 - 0) \times 2$$

$$= \underline{\underline{-20}}$$

Step 4 Update Rule

$$\beta_{1 \text{ new}} = 0 - 0.01 \times (-10) = \underline{0.1}$$

$$\beta_{2 \text{ new}} = 0 - 0.01 \times (-20) = \underline{0.2}$$

$\beta_1 = 0.1$ ,  $\beta_2 = 0.2 \rightarrow \underline{\text{Repeat the process}}$



$$\beta_1 = 0.1, \beta_2 = 0.2$$

Do prediction

$$\begin{aligned}\hat{y} &= 0.1 \times 1 + 0.2 \times 2 \\ &= \underline{0.5}\end{aligned}$$

$$\text{Actual} = 5$$

$$\begin{aligned}\text{loss} &= 5 - 0.5 \\ &= \underline{\underline{4.5}}\end{aligned}$$

{ Now weights are 'improving' }

Repeat till you get best values for  $\beta_1$  &  $\beta_2$

## • Advantages

- \* Handle large datasets or complex models.
- \* Can help us find best values for our model  
↓  
(coefficient/weights)
- \* Works well with high-dimensional data
- \* Best for convex functions.  
(linear reg, lasso reg, Ridge reg, etc)  
(Global minima is guaranteed)

## • Disadvantages

\* Not well with non-convex function  
(Neural Networks (Deep learning))

It converges to local minima instead of global minima.



To handle this  $\rightarrow$  SGD, Adam, RMS Prop  
(stochastic)  
GD

\* find best  $\alpha$  value  $\rightarrow$  It can be hectic.

Soln  $\rightarrow$  Adam optimizer

$\alpha \rightarrow$  learning rate

$\Rightarrow$  Automatically adjust  $\alpha$   
value acc to data