

## Gradient Descent Algorithm (GDA)

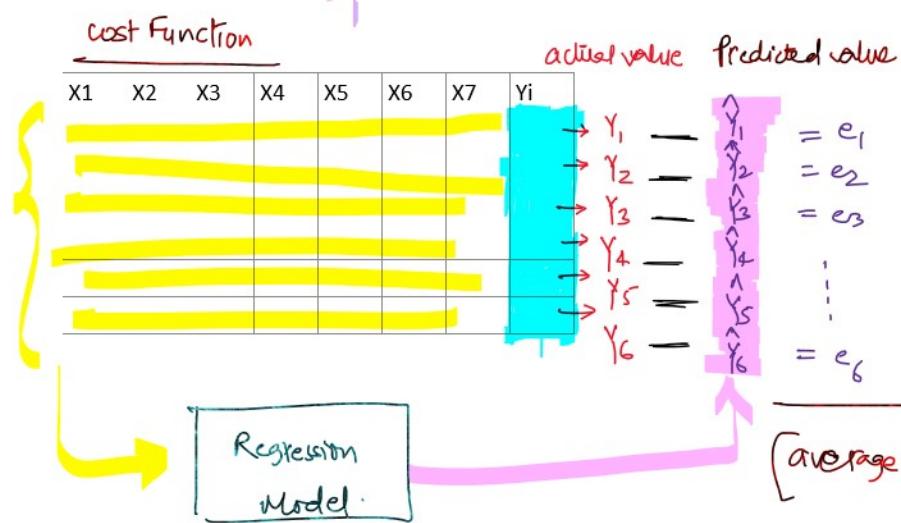
21 September 2025 11:24

### \* Pro tip Loss Function vs Cost Function

Loss Function: It is the function to capture the difference between actual and predicted value. for a single row or a single record or a single training example.

$$\text{error} = [y_i - \hat{y}_i] \text{ for some } i\text{th row}$$

actual      predicted  
value      value



for regression : MSE → Mean Squared Error

### Cost Function:

It is used to refer an average of the loss functions in some way over the entire training dataset

### Some famous cost Functions

MSE: Mean Squared Error

③ ② ①

MSE : Mean squared error



① error =  $(y_i - \hat{y}_i)$

② squaring the error =  $(y_i - \hat{y}_i)^2$

③ Mean of the squared error = 
$$\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

$$\sqrt{MSE}$$

Task: Read about different cost functions for classification problems. What is a logit???

- Logistic Regression

- Decision Trees

Support vector Machines

↳ SVM ...

Binary  
as well as  
Multi-class.

Gradient Descent Algorithm (GDA)

Gradient: is the rate of change → differentiation (partial derivatives)  
 or  
 derivative  
+ → (slope / rate of change)

Descent coming down (direction is downwards)

→ Gradient descent tells us how steep the slope is and which direction decreases the error.

$\nabla J$ : gradient of  $J \rightarrow$  it points in the direction of maximum increase  
 ↓  
 cost function

GD: Gradient descent moves in the opposite direction of the gradient because we want to minimize the error.

## Gradient descent algorithm

repeat until convergence {

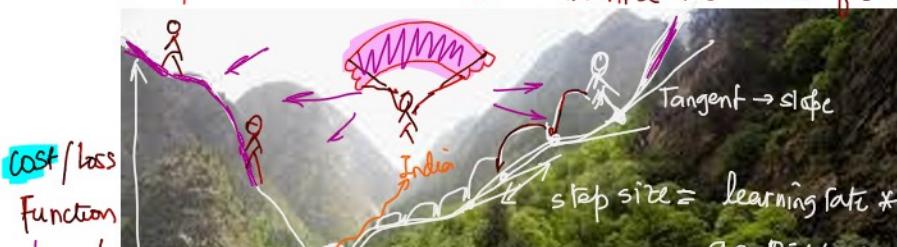
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

(for  $j = 1$  and  $j = 0$ )

}

Intuition: Abhinandan - Mig-21 Pilot got captured on the enemy's land.

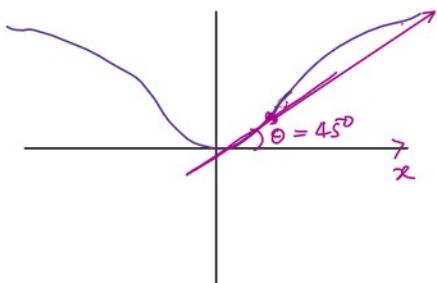
② abc: Abhinandan needs to minimize the chance of getting caught.





Gradient Descent → follows the path of the **steepest descent**  
 taking steps in the direction that  
 decreases the slope and brings Abhinandan  
 closer to the local minimum (coming to plains)

$$\tan\theta = \text{slope} = m = \frac{dy}{dx}$$



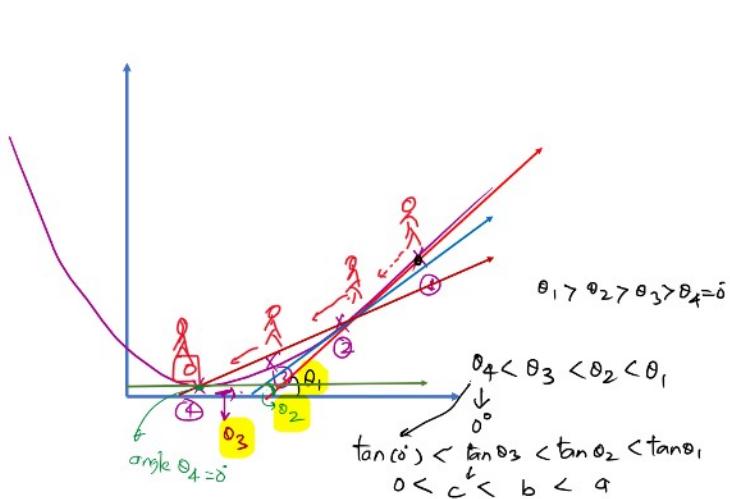
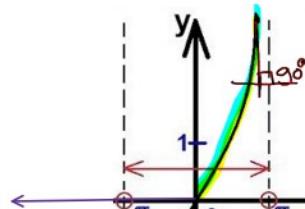
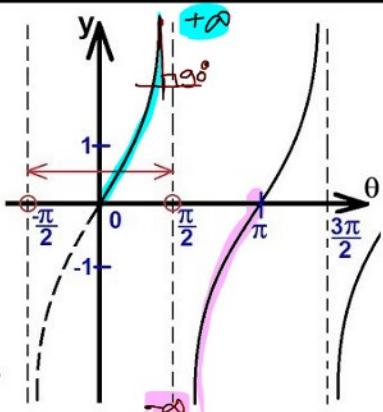
$$\tan(\theta) = \tan 45^\circ = 1$$

$\downarrow$

$$\text{slope} = \frac{dy}{dx} = m = 1$$

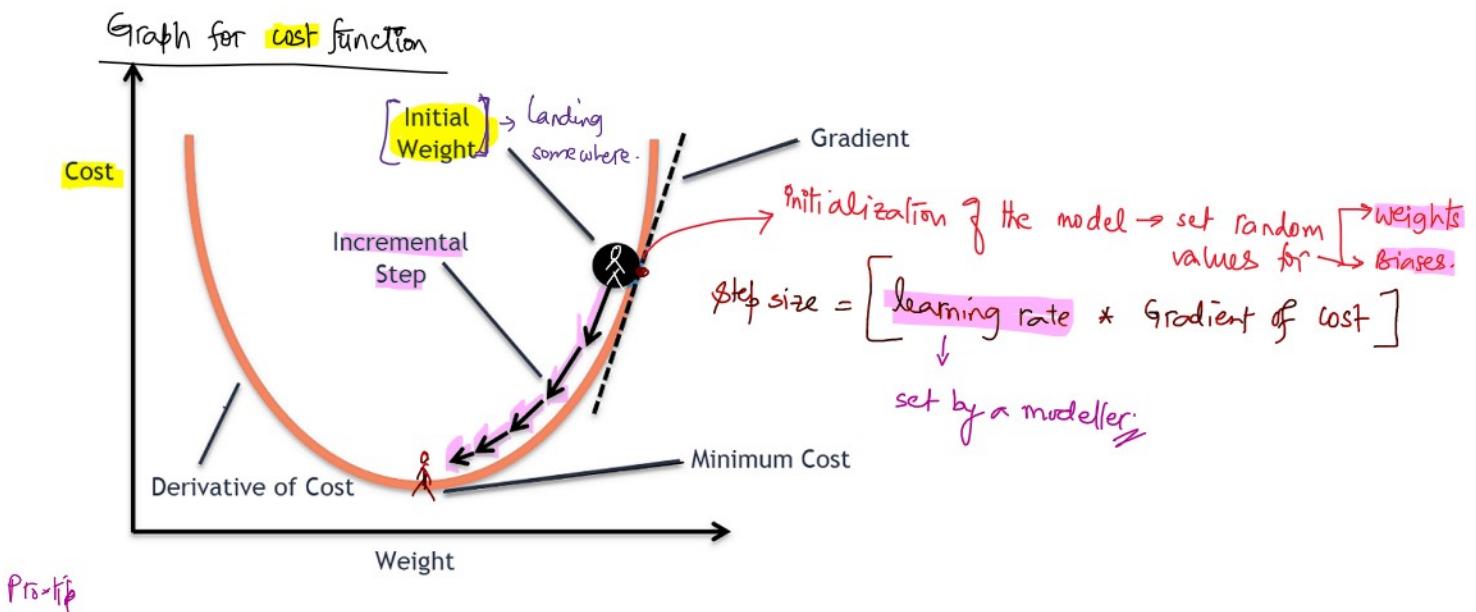
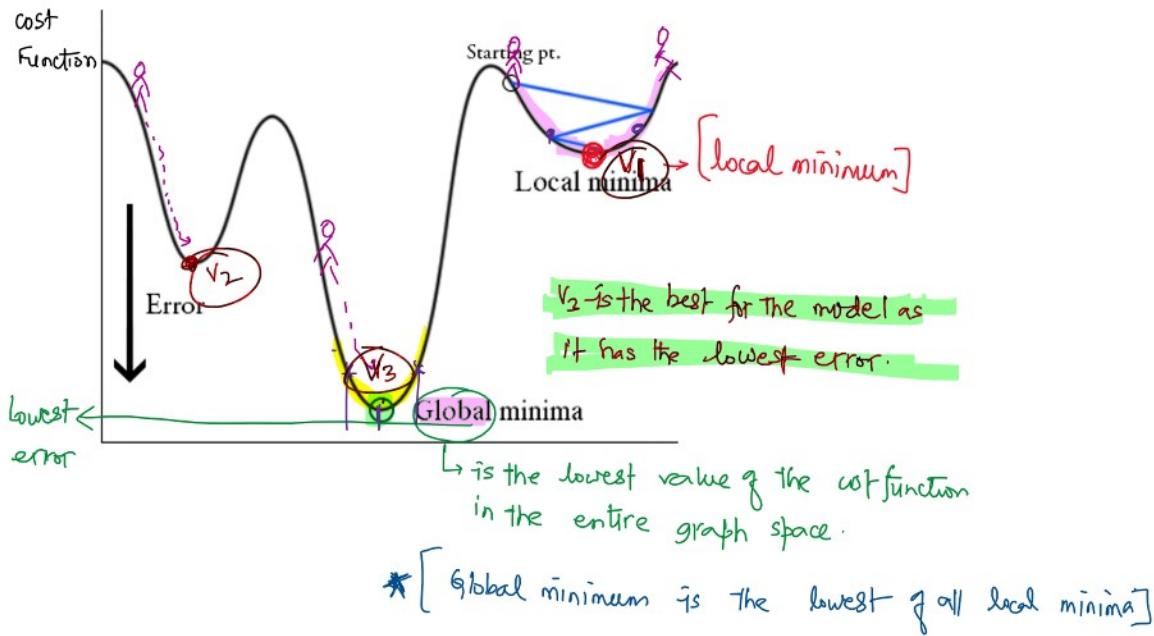
## Graph for tangent function

$\theta$ (degree)	$\theta$ (radian)	$\tan\theta$
0	0	0 0
30	$\frac{\pi}{6}$	$\frac{1}{\sqrt{3}}$ 0.577
45	$\frac{\pi}{4}$	1 1
60	$\frac{\pi}{3}$	$\sqrt{3}$ 1.732
90	$\frac{\pi}{2}$	$\infty$ $\infty$
120	$\frac{2\pi}{3}$	$-\sqrt{3}$ -1.732
135	$\frac{3\pi}{4}$	-1 -1
150	$\frac{5\pi}{6}$	$-\frac{1}{\sqrt{3}}$ -0.577



- ① As ' $\theta$ ' increased from  $0^\circ$  till  $90^\circ \rightarrow \tan\theta \rightarrow$  slope also increases
- ② and as  $\theta$  increases from  $90^\circ$  till  $180^\circ \rightarrow \tan\theta \rightarrow$  slope also decreases

# Local Minimum Vs Global Minimum



## Gradient descent algorithm

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

(for  $j = 1$  and  $j = 0$ )

cost Function :  $J(\theta_0, \theta_1)$

partial derivative of  $J(\theta_0, \theta_1)$  w.r.t.  $\theta_0, \theta_1$  partially.

gradient of cost functrn.

}

$\{$   $J = \text{J}$  and  $J = 0$  }  $\rightarrow$   $w, w_0$  minimally.  
 $\rightarrow \alpha \nabla J$ : learning rate.  
 -  $\text{learns in } 1 - e + \text{len} \quad \overline{w}$ ,  
 $\rightarrow$  direction is opposite to the gradient hence  
 it's called gradient descent.

$$\theta_1 \text{ new} = \theta_1 \text{ old} - (\text{learning rate} * \text{gradient of cost})$$

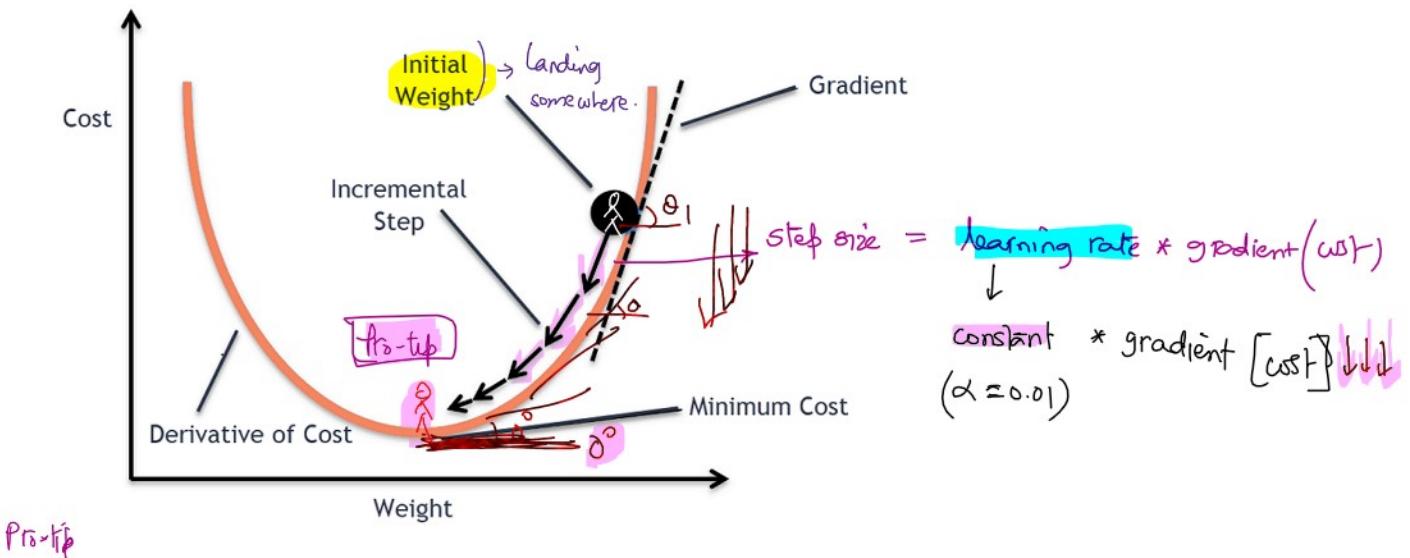
$\theta_1 \text{ new} \triangleq \theta_1 \text{ old}$

$$\theta_1 := \theta_1 - \alpha \cdot \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \leftarrow \text{updating the } \theta_1 \text{ value.}$$

$$\theta_1 \text{ new} = \theta_1 \text{ old} - \left| \begin{array}{c} \alpha \cdot \frac{\partial}{\partial \theta_1} \\ \text{gradient of cost function} \end{array} \right| \left| \begin{array}{c} \text{bias} \\ \text{weight} \\ \text{cost function} \\ J(\theta_0, \theta_1) \end{array} \right|$$

$\leftarrow \text{Learning Rate} \times \text{Gradient of Cost Function}$

$\longleftrightarrow \text{step size} \longleftrightarrow$



$$\theta_1)_{\text{new}} = \theta_1)_{\text{old}} - \alpha \cdot \frac{\partial}{\partial \theta_1} [J(\theta_0, \theta_1)]$$

gradient becomes zero  
when  $\theta = \theta^*$

$\boxed{\theta_1)_{\text{new}} \triangleq \theta_1)_{\text{old}}} \rightarrow \text{GDA converges.}$

Example:

$\theta_0$ : bias

$\theta_1$ : weight

$\alpha = 0.1$

$\theta_1)_{\text{initial}} = 2.1$  ←

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = 1.5 \quad \text{slope} | \text{gradient}$$

$$\theta_1)_{\text{new}} = \theta_1)_{\text{initial}} - 0.1 * 1.5$$

$$\theta_1)_{\text{new}} = 2.1 - 0.1 * 1.5 = 2.1 - (0.1 * 1.5) = 1.95$$

$$\theta_1)_{\text{new}} = 1.95$$

$$\theta_1)_{\text{new}_2} = 1.95 - 0.1 * [ ]$$

↙



## Learning Rate ( $\alpha$ or $\eta$ )

What is a hyper parameter?

→ A hyper parameter is a setting chosen before training a machine learning model.

→ it controls how the learning process happens but it is not learned from the data itself

which is not a hyperparameter



### Examples of hyperparameters

#### ① Linear / Logistic Regression

→ learning rate (gradient descent)

→ regularization  
[ L1 → Lasso ]  
[ L2 → Ridge ]

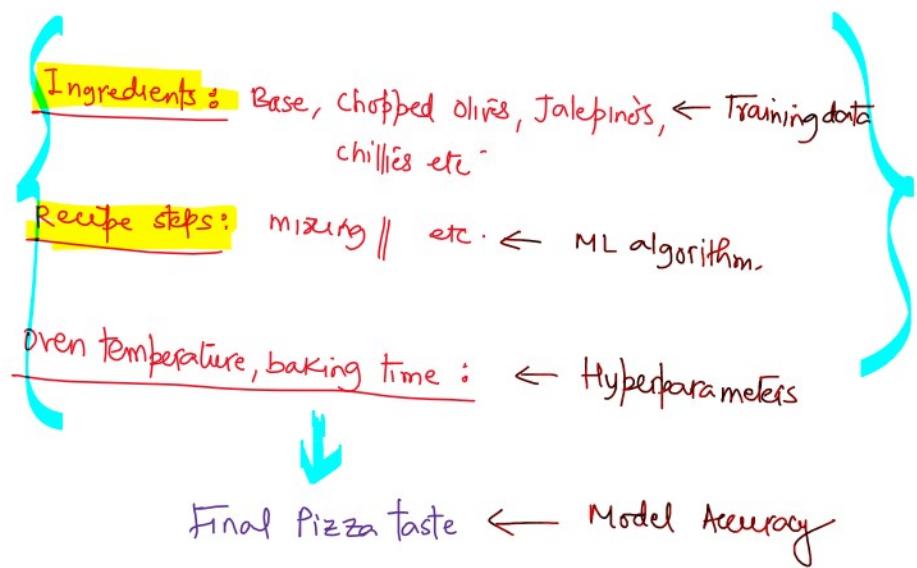
#### ② Decision Trees

→ minimum samples [ leaf ]

→ max\_depth of the tree

### ③ Neural Networks

- learning rate
- number of hidden layers
- number of neurons per layer
- batch size
- dropout rate
- optimizer choice (SGD, MBGD, Adam etc.)



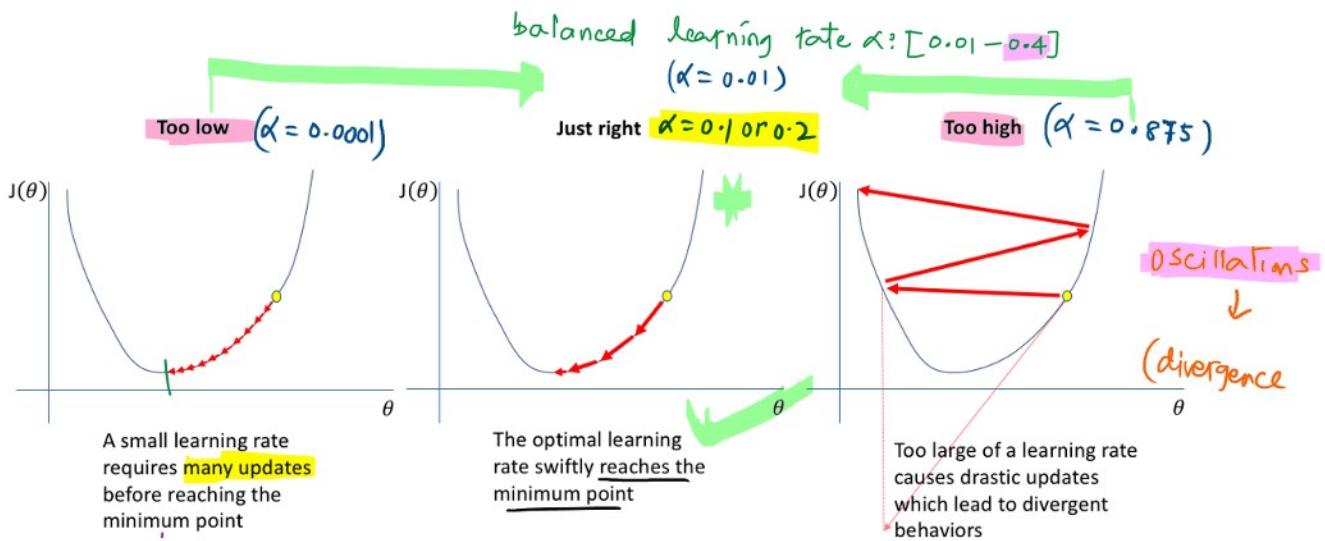
## # Learning Rate ( $\alpha$ or $\eta$ )

Learning rate is a hyperparameter in a gradient based optimization technique that determines the step size at each iteration while moving towards the minimum of the cost function.

$$\text{Step size} = \alpha * (\text{gradient of the cost function})$$

### Right choice for learning rate

$$\alpha \in [0, 1]$$



A small learning rate requires many updates before reaching the minimum point

↓  
slow training  
↓

more no. of epochs  
↓

Project budget ↑↑

Training runtime ↑↑

The optimal learning rate swiftly reaches the minimum point



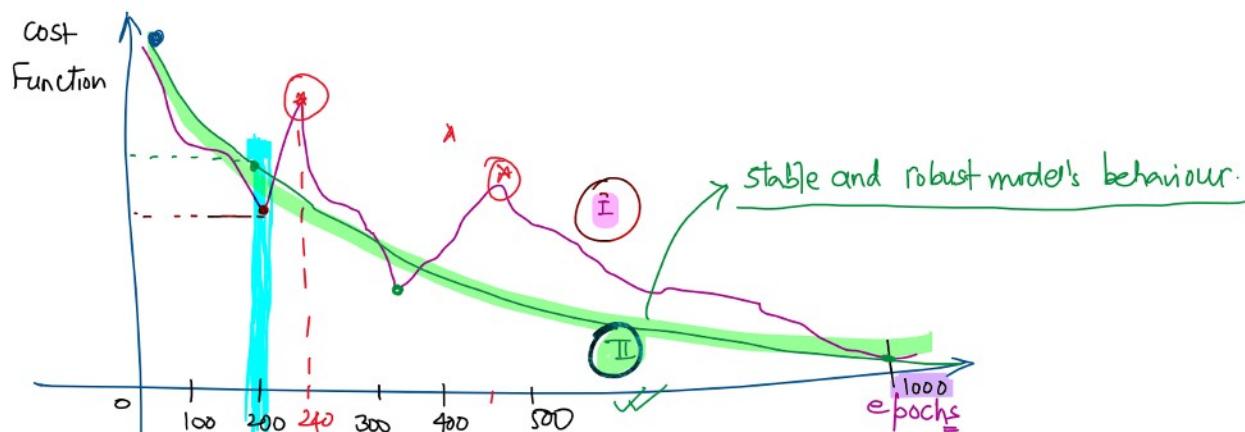
Too large of a learning rate causes drastic updates which lead to divergent behaviors

Pro-tip Why is setting up the learning rate so critical ??.

Learning rate affects the speed of convergence

• Accuracy: Whether the GDA can find the optimal solution or overshoot it, it depends on  $\alpha$ .

• stability: preventing divergence or oscillations during training is a critical step to ensure robustness of the model.



Note: As per the best practices, start with a default

learning rate  $\alpha$  or  $\eta \rightarrow [0.01 - 0.2]$

$\alpha = 0.01$

$$\alpha = [0.01 - 0.4]$$
$$[0.01, 0.03, 0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.4]$$

(10) ↴