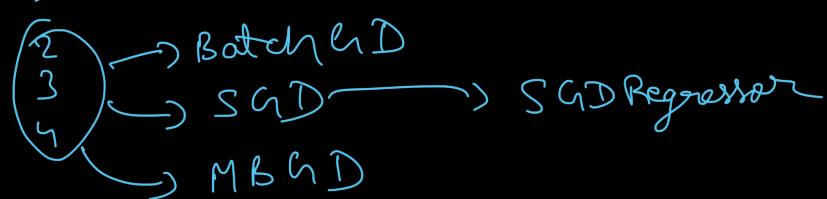


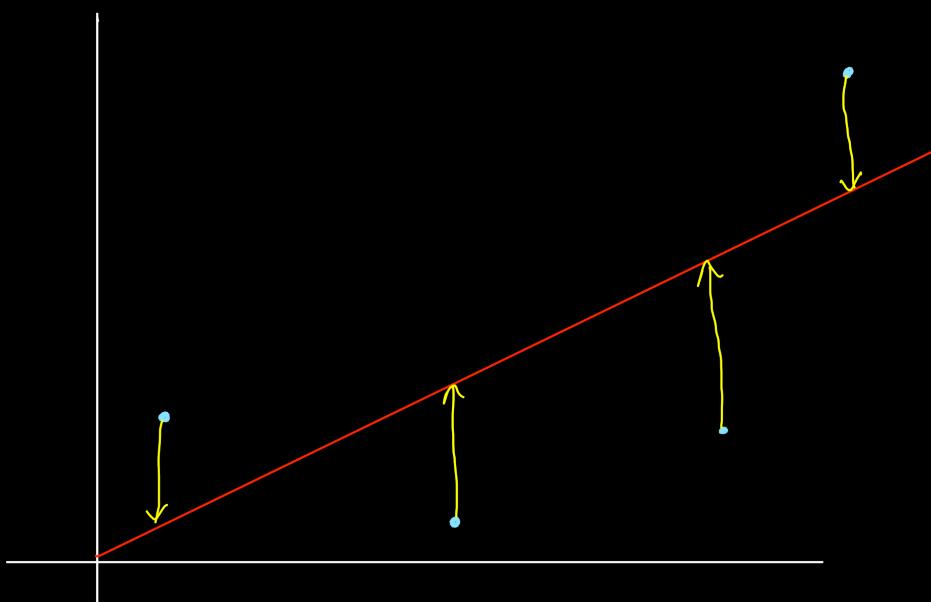
# Gradient Descent from Scratch

wikipedia : Gradient descent is a first-order iterative optimization algorithm for finding a local minimum of a differentiable function. The idea is to take repeated steps in the opposite direction of the gradient (or approximate gradient) of the function at the current point, because this is the direction of steepest descent. Conversely, stepping in the direction of the gradient will lead to a local maximum of the function; the procedure is then known as gradient ascent.

The Plan :



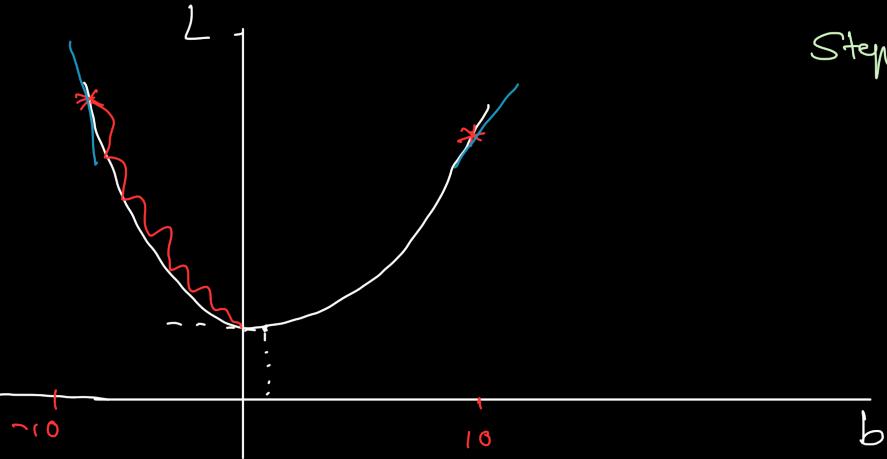
Intuition :



$$\begin{aligned}
 & \frac{\partial \mathcal{L}}{\partial m} = \frac{\partial}{\partial m} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\
 & \hat{y}_i = mx_i + b \\
 & \mathcal{L}(m, b) = \sum_{i=1}^n (y_i - mx_i - b)^2
 \end{aligned}$$

Suppose  $m = 78.35$

$$\mathcal{L} = \sum_{i=1}^n (y_i - 78.35x_i - b)^2$$



Step 1: Select a random  $b$

$$[b = -10]$$

Slope = -ve then increment in  $b$  value

Slope = +ve then decrement in  $b$  value

$$\boxed{b_{\text{new}} = b_{\text{old}} - \text{slope}}$$

$$\boxed{b_{\text{new}} = b_{\text{old}} - \eta \text{slope}} \rightarrow \text{learning Rate}$$

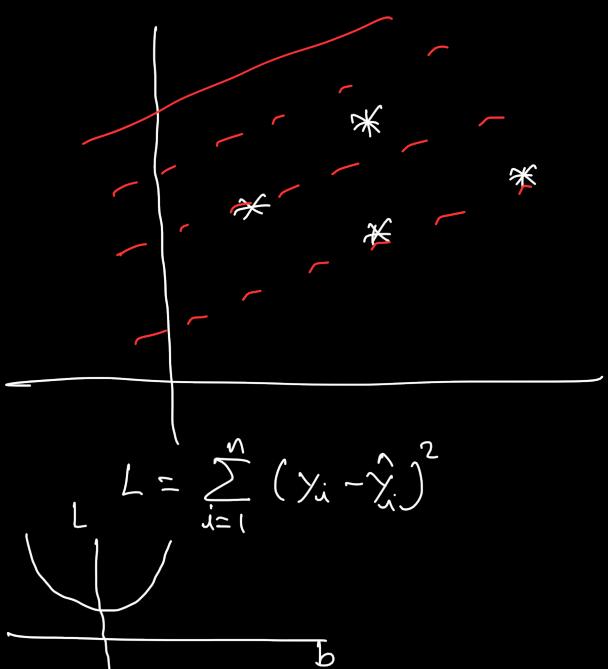
### When to Stop

1) difference b/w  $b_{\text{old}}$  and  $b_{\text{new}} \geq 0.0001$

2) fixed number of  $\frac{\text{iterations}}{\hookrightarrow \text{epochs}}$  such as 1000, 100

### Mathematical Formulation

$$m = 78.35$$



Step 1 → Start with random  $b$

for  $i$  in epochs:  $\boxed{\eta = 0.01}$

$$b_{\text{new}} = b_{\text{old}} - \eta \times \text{slope}$$

$$\frac{dL}{db} = \frac{d}{db} \left( \sum_{i=1}^n (y_i - \hat{y}_i)^2 \right)$$

$$= \frac{d}{db} \sum_{i=1}^n (y_i - mx_i - b)^2$$

$$= 2 \sum_{i=1}^n (y_i - mx_i - b)(-1)$$

Slope =  $-2 \sum_{i=1}^n (y_i - mx_i - b)$

at  $b = 0$

$$= -2 \sum_{i=1}^n (y_i - 78.35x_i - 0) = \text{slope}_{\text{bold}}$$

$b_{\text{new}} = b_{\text{old}} - \eta \text{slope}_{\text{bold}}$

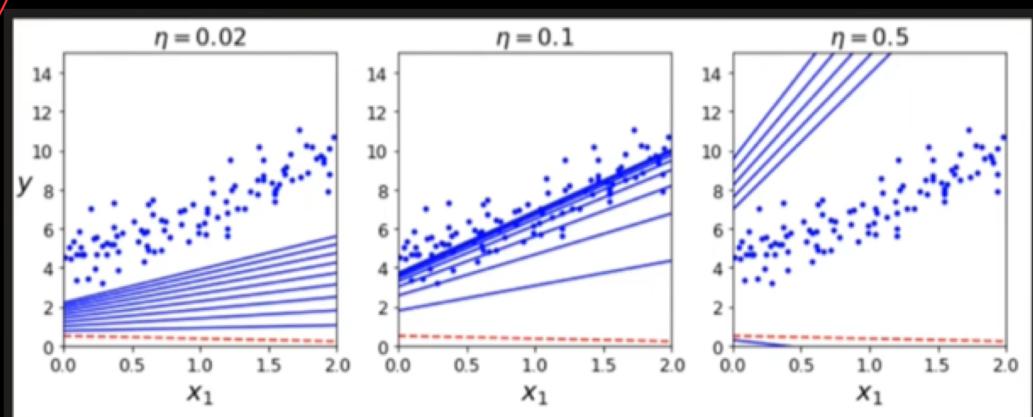
Step Size

Code from scratch

Visualization 1

## Few Discussions

1. Effect of Learning rate
2. The universality of Gradient Descent



epoch = 10

$b_{\text{new}} = b_{\text{old}} - \eta \text{slope}$

→ we only need slope to use GD. And slope comes from derivative of loss function.

Now any loss function which is differentiable at every point can be used in 1D for optimization.

Adding m into the equation

Steps

1) initialize random value for m and b

$$m=1 \text{ and } b=0$$

2) epochs = 100 ,  $\eta = 0.01$

for i in epochs:

$$b = b - \eta \text{ Slope}_b$$

$$m = m - \eta \text{ Slope}_m$$

$$\left. \begin{array}{l} L = \sum (y_i - \hat{y}_i)^2 \\ L = \sum (y_i - mx_i - b)^2 \\ L(m, b) \end{array} \right\} \begin{array}{l} \text{Loss f'm is a} \\ \text{function of} \\ m \text{ & } b \end{array}$$

$$\boxed{\text{Slope}_b = \frac{\partial L}{\partial b}}$$

$$\boxed{\text{Slope}_m = \frac{\partial L}{\partial m}}$$

$$\boxed{\frac{\partial L}{\partial b} = -2 \sum (y_i - mx_i - b)}$$

$$\therefore \text{Slope}_b \text{ at } b=0$$

$$\boxed{\frac{\partial L}{\partial m} = -2 \sum (y_i - mx_i - b)x_i}$$

$$\therefore \text{Slope}_m \text{ at } m=1$$

Visualization 2

## Effect of learning rate (hyperparameter)

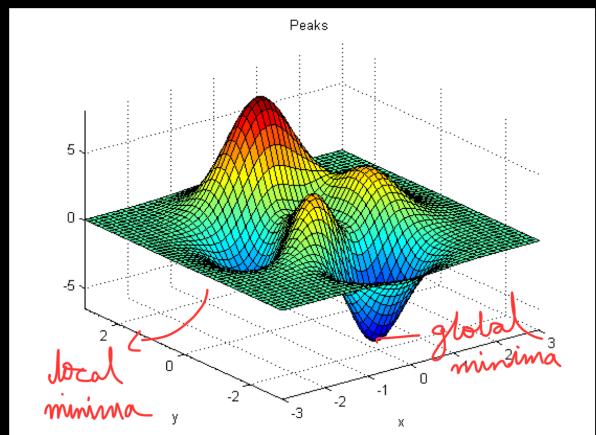
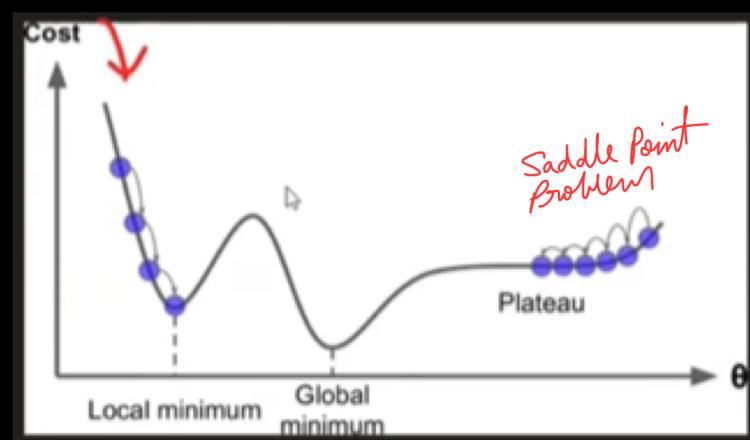
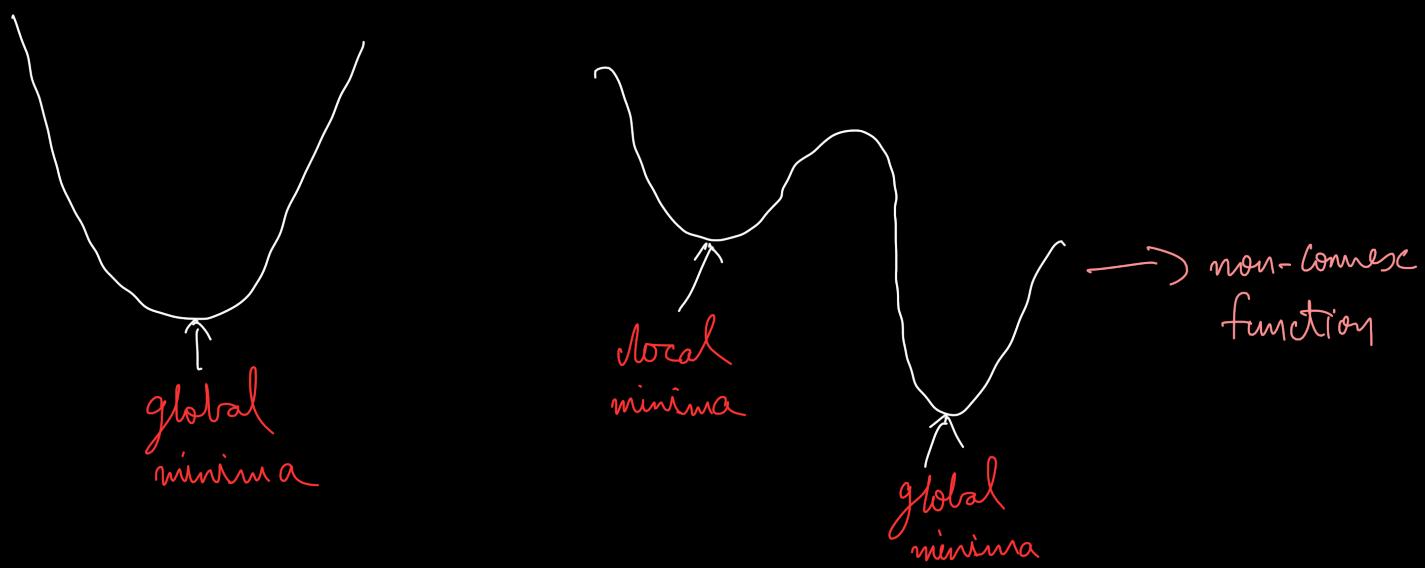
If learning rate is too small : Gradient descent will take too long to reach the minimum loss.

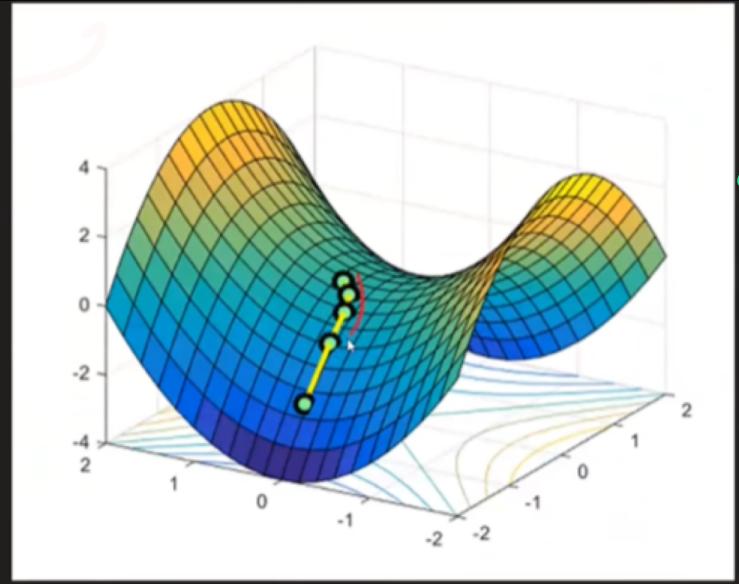
optimal learning rate ; Model reach minimal loss with fewer steps

If learning rate is too big : Gradient descent will take too long or never reach the minimum loss.

## Effect of Loss Function

$$L = \sum (Y_i - \hat{Y}_i)^2 \rightarrow \text{Convex function}$$





Saddle point Problem: Slope changes very slowly.

### Effect of data

If features in dataset are all in same scale then contour plot is circular and we descent very fast.

