# SOLVING OPTIMIZATION PROBLEMS

# DIFFERENTIATION

Solving (Optimization) problems: differentiation

✓ optmzn. problems :- (PCA, logistic regᵖⁿ; lr. regression)

✓ (ML) — { differentiation → single-value, vector { 99% }
maxima & minima ⊐ $11^{th}$ & $12^{th}$ →
⤷ undergrad

Integration { less used }
Diff. eqns

The idea of Differentiation and Maxima,Minima is used in 99% of Machine Learning

Single-variable diffn:   $y = f(x)$

$y = f(x)$

$\frac{y_2 - y_1}{x_2 - x_1} = \frac{\Delta Y}{\Delta X} = Tan\theta$

$\left\{ \frac{dy}{dx} = \frac{df}{dx} = y' = f' \right\}$

⤷ differentiation of $y$ w.r.t $x$

scalar

"$\frac{dy}{dx}$" (intuitively) → rate of change of $y$ as $x$ changes
→ how much does $y$ change as $x$ changes

dy/dx - Rate of change of y as x changes. dy/dx = $\frac{y_2 - y_1}{x_2 - x_1}$ = $\frac{\Delta y}{\Delta x}$ = $Tan\theta$
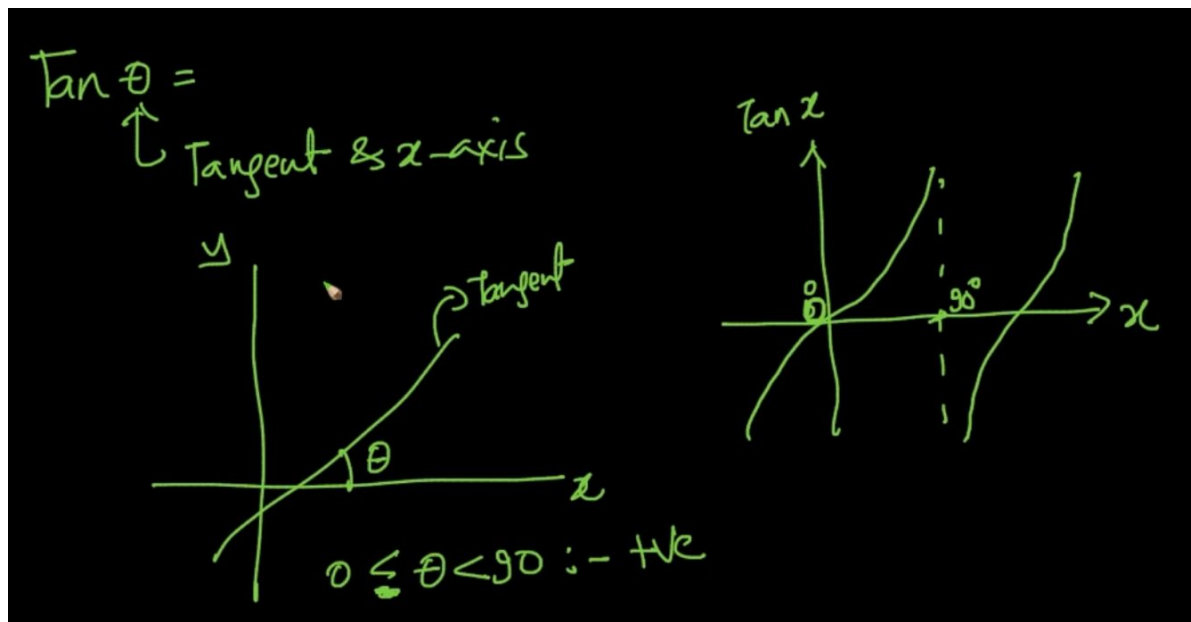
$$\left[\left(\frac{dy}{dx}\right) = \lim_{\Delta x \to 0} \frac{\Delta y}{\Delta x}\right]$$



geom. intr

$y = f(x)$

Tangent off(x) @ $x_1$

$x_1$ $x_2$

Tangent is the hyp that

We obtain as $\delta x \to 0$

$\boxed{\text{Tan } \theta = \frac{dy}{dx}}$

$x_2 \to x_1 \Rightarrow \delta x \to 0$

as $\delta x \to 0$; $\frac{\delta y}{\delta y} = \frac{dy}{dx}$

dy/dx = $\lim_{\Delta x \to 0} \frac{\Delta y}{\Delta x}$ .Let's break down it. The derivative of y w.r.t  x (dy/dx) is equal to rate of change of y w.r.t rate of change of x i.e ( $\frac{\Delta y}{\Delta x}$ ) and $\Delta x$ i.e $x_2 - x_1$ is close to 0.
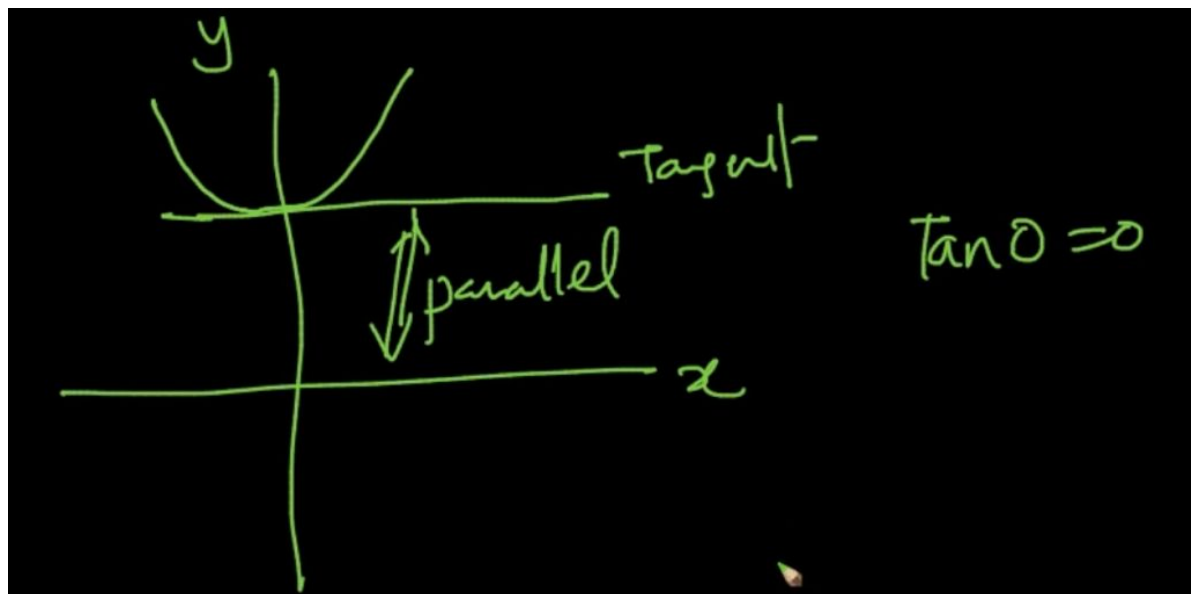
If $\Delta x$ is 0 then the tangent of f(x) at that point x is our (dy/dx) and it's $\theta$ is angle between slope and x-axis
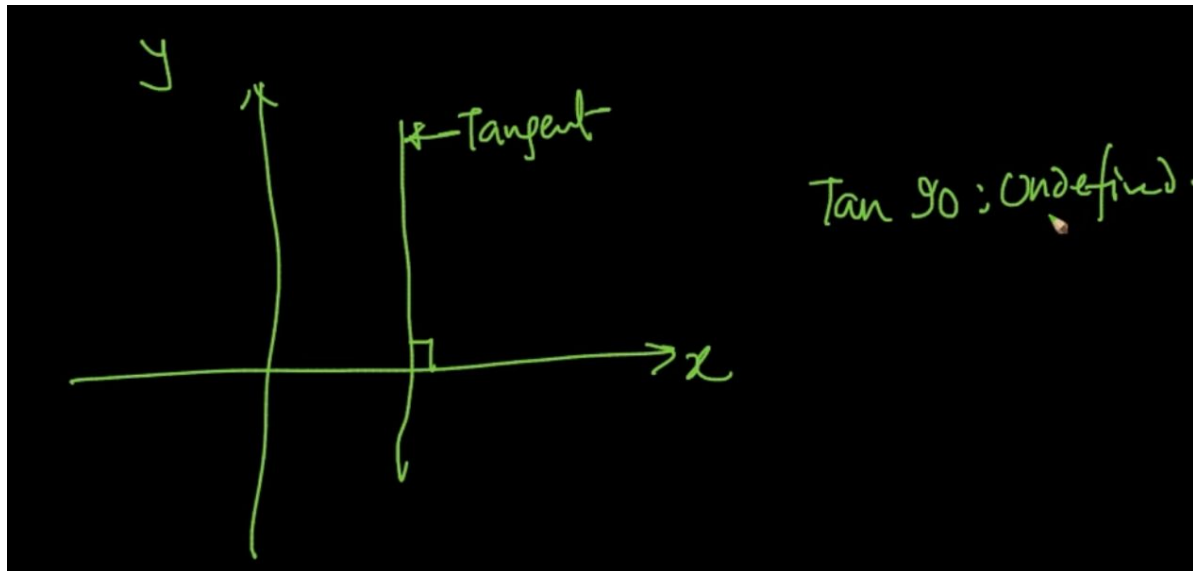
$$\frac{dy}{dx} = \text{Slope of the tangent to } f(x)$$

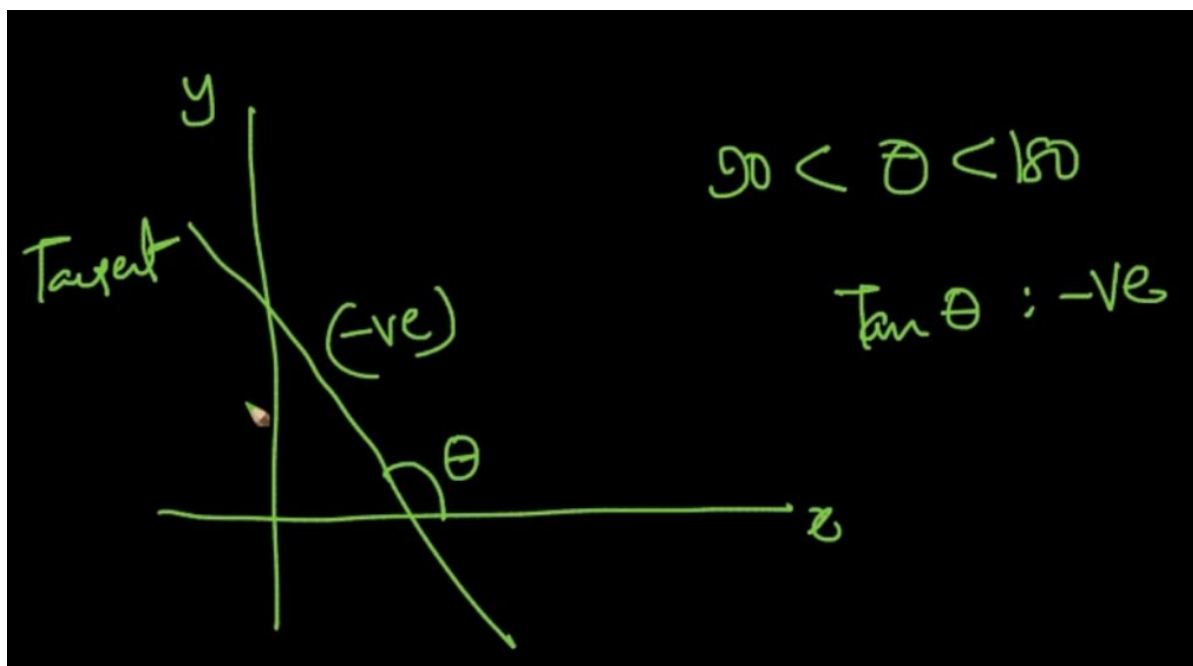$$\left[\frac{dy}{dx}\right]_{x_1} = \text{Slope of the tangent of } f(x) @ x = x_1$$

Tan θ =
↳ Tangent & x-axis

$0 \leq \theta < 90 :- +ve$

Tan x

Tangent when 0 < θ < 90 is positive

Tangent

parallel

Tan θ = 0

When θ = 0 tangent is parallel to x-axis

Tan 90 is undefined



Tan θ is negative when 90 < θ < 180

$$\text{Chain-rule:} \qquad \frac{d}{dx} f(g(x)) = \frac{df}{dg} \cdot \frac{dg}{dx}$$

$$\rightarrow f(g(x)) = (a - bx)^2 \qquad \bigg| \qquad \frac{d}{dx} f(g(x)) = \frac{df}{dg} \cdot \frac{dg}{dx}$$

$$\Bigg( \Bigg[ \begin{array}{l} g(x) = (a - bx) \\ f(x) = x^2 \end{array}$$

$$\frac{dg}{dx} = \frac{d}{dx}(a - bx) = \frac{d}{dx}(a)$$
$$- \frac{d}{dx}(bx)$$
$$\bullet = -b$$

It's the most important concept of Calculus in Deep Learning and Machine Learning.

$\frac{d}{dx}(f(g(x)) = \frac{df}{dg} \cdot \frac{dg}{dx}$ . First, we'll calculate $\frac{dg}{dx}$ .Now, calculate df/dg.
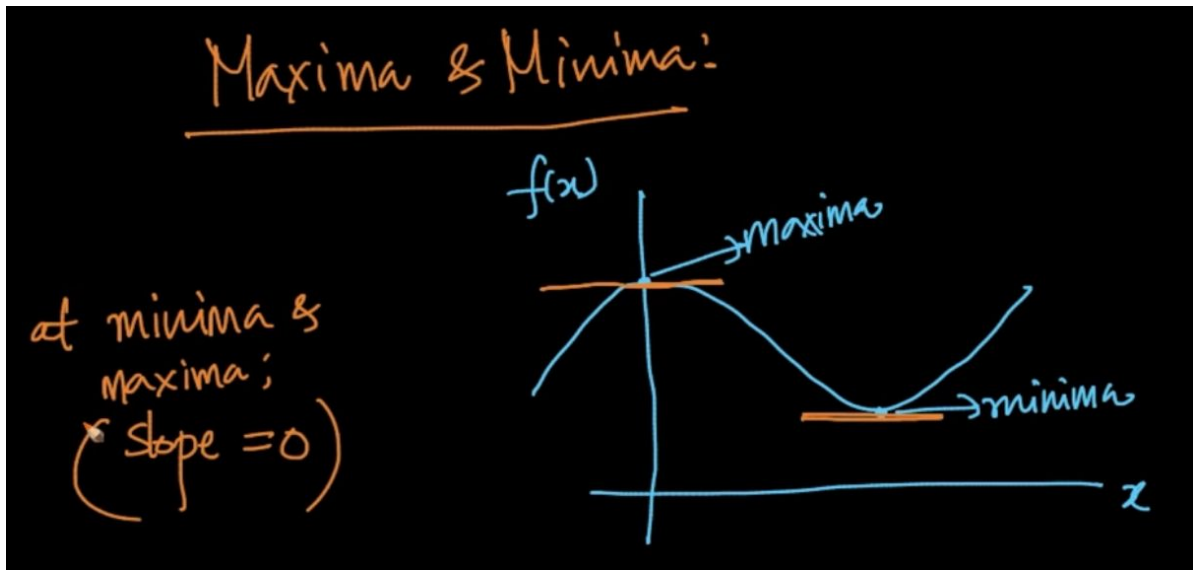
$$\frac{dg}{dx} = -b$$

$$\& \quad g(x) = z$$

$$f(g(x)) = z^2$$

$$\frac{df}{dg} = \frac{d z^2}{dz} = 2z$$

If we take g(x) = z then f(g(x)) = $z^2$ for the above f(g(x)) so df/dg = 2z

So $\frac{df}{dg} \cdot \frac{dg}{dx}$ = 2.(a-bx). -b for the above problem

# MAXIMA AND MINIMA



At minima and maxima the slope is = 0



We know that df/dx is our slope so if we make it = 0 then we can get the value of maxima/minima. We are calculating df/dx = 0 and getting the value of x. So, f(x) would tell
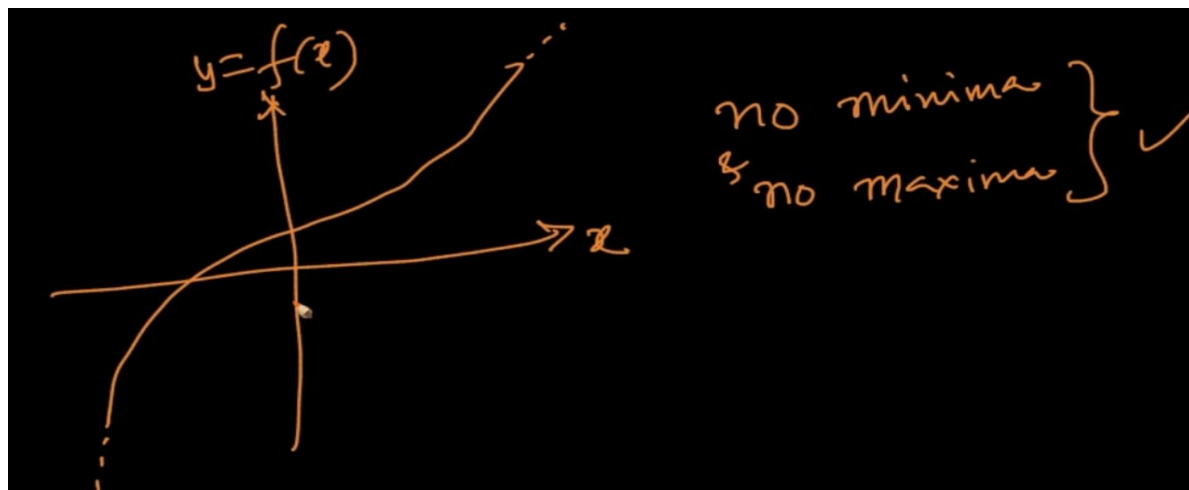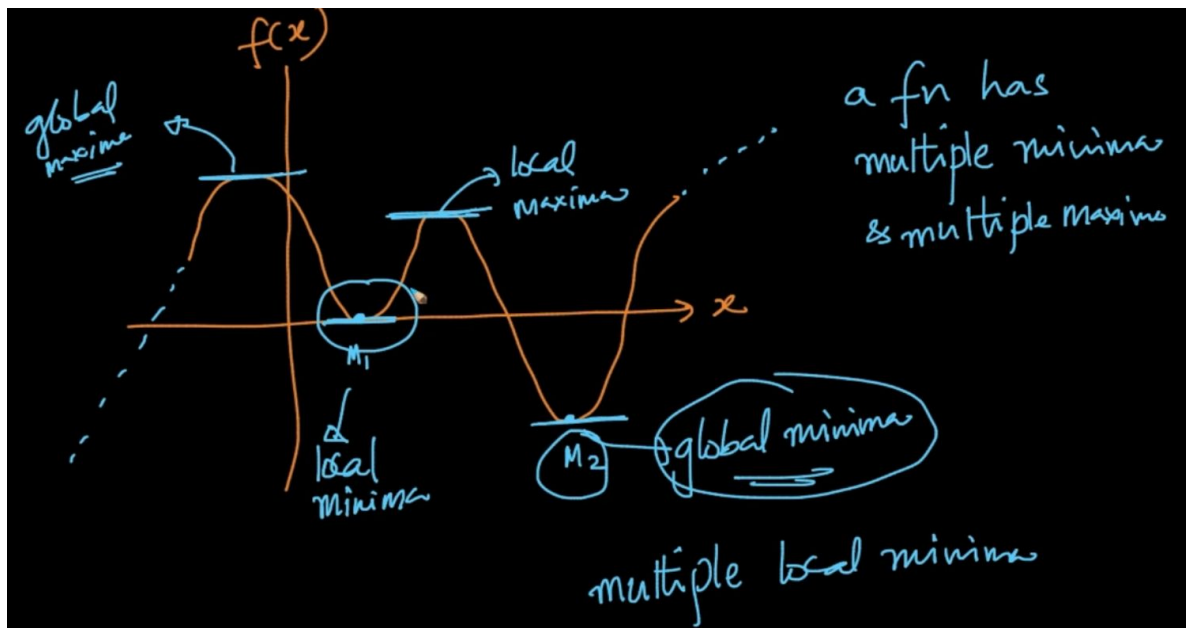
us the minima/maxima but we don't know whether it's telling maxima or minima .For this we calculate f(x) for a near value towards it. Here, df/dx = 0 at x = 1.5 and f(1.5) = -0.25 We take f(1) which is 0 , Since f(1.5) < f(1) it cannot be maxima so it'll be minima



There can be some functions which might not have maxima/minima . Ex: In the above orange figure there's only minima not maxima and vice-versa for the blue figure



There can be some functions which might not have maxima/ minima

A function can have multiple maxima and minima as seen above.

**Local Minima/ Maxima** - It's the Minimum /Maximum in the neighborhood as seen above

**Global Maxima/Minima** - It's the Minimum/ Maximum among all the Minimums/ Maximums



Some functions cannot be solved by simple df/dx = 0. So, we'll use Gradient Descent for that

**VECTOR CALCULUS GRAD**



We've dealt with scalars in Differentiation but in ML things are in vector. So we'll deal with vector calculus. Our , y = f(x) = $a^T x$ where x is vector and a is a vector of constants



df/dx = $\nabla_x f$ which is Del ( $\nabla$ ) of f w.r.t x as seen. This is differentiation of vector

As seen above it's $\nabla_x f$ = [df/dx1, df/dx2.....df/dxd] but df/dx1 is partial derivation and we write it as $\delta f / \delta x1$ . It means that x has <x1,x2,x3......xd> but we are derivating the function with only one component x1 at a time not the whole x

$$f(x) = y = a^T x = \sum_{i=1}^{d} a_i x_i = a_1 x_1 + a_2 x_2 + \cdots + a_n x_n$$

$$\nabla_x f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_d \end{bmatrix} = a$$

$$\frac{\partial f}{\partial x_1} = a_1$$

$$\frac{\partial f}{\partial x_2} = a_2$$

$$\nabla_x (a^T x) = a$$

$$\frac{d}{dx}(ax) = a$$

$y = a^T x = a_1 x_1 + a_2 x_2 + ......  a_n x_n$ . So when we do the derivation of above . We'll get 'a'

$$\mathcal{L}(w) = \sum_{i=1}^{n} \log\left(1 + \exp\left(-y_i w^T x_i\right)\right) + \lambda w^T w$$

$\langle x_i, y_i \rangle \rightarrow$ constants $\rightarrow$ $\lambda$ train

$$\frac{d \lambda w^2}{w} = 2\lambda w$$

$$\nabla_w \mathcal{L} = \frac{(-y_i x_i)(\exp(-y_i w^T x_i))}{1 + \exp(-y_i w^T x_i)} + 2\lambda w = 0$$

(Gradient descent)

Finding minima.maxima with $\square_x L$ is very hard because the equation is complex af.

So, we'll apply computational techniques like Gradient Descent to solve it

## GRADIENT DESCENT: GEOMETRIC INTUITION



Gradient descent is an iterative algorithm. In an Iterative algorithm first we'll guess the value

of $x^*$ (Point where we want to reach)which is $x_0$. Now by applying our iterative algorithm Gradient Descent we'll get $x_1$ which is closer to our targeted $x^*$ than $x_o$

geom.

slope changes f(x)
-its sign from
+ve to -ve @ minima

(-ve)  minima  (+ve)

$x$

minima
one side:- slope : +ve
other side:- slope : -ve

✓ $x^* = \underset{x}{\arg\min} \; f(x)$

$\min f(x)$
$\Rightarrow \max -f(x)$

$\max f(x)$
$\Rightarrow \min -f(x)$

We want to find $x^*$ = argmin f(x) where we want to minimize f(x). min f(x) = max(-f(x)) and vice versa. One important observation here is that the slope change it's sign from +ve to -ve After reaching minima
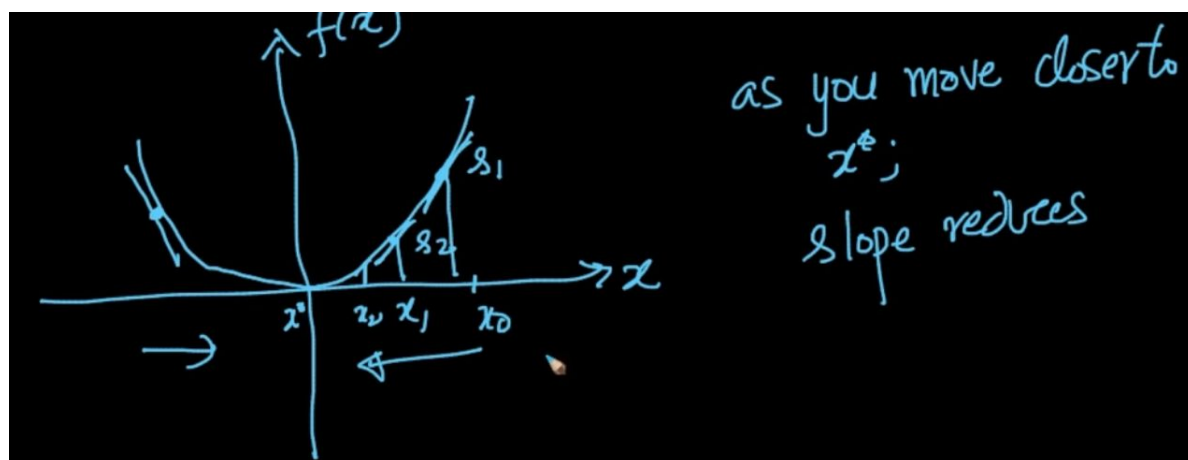


$f(x)$

$S_1$

$S_2$

$x^*$  $x_2 \, x_1$  $x_0$

as you move closer to $x^*$;
slope reduces

Another observation is as we move closer to $x^*$ from the right side the slope reduces
And it increases if we move from left to right. So , let's see how Gradient Descent works

## Gradient descent:-

① pick an initial pt $= x_0$ @ random

$y = f(x)$

② $x_1 = x_0 - \gamma_* \left[\dfrac{df}{dx}\right]_{x_0}$ ← (+ve)

Let stepsize $= 1$
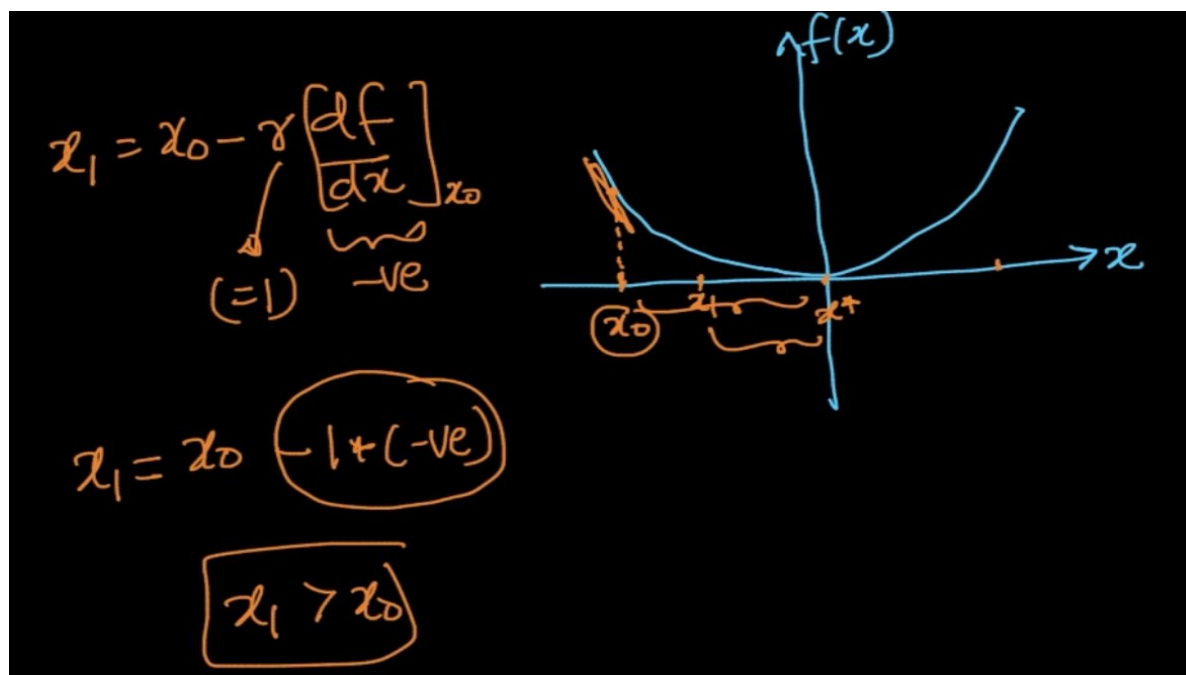
$x_1 = x_0 - 1*(+ve) \implies x_1 < x_0$

(+ve)

$x^* \leftarrow x_1 \quad x_0$

$x_1$ is closer to $x^*$ than $x_0$
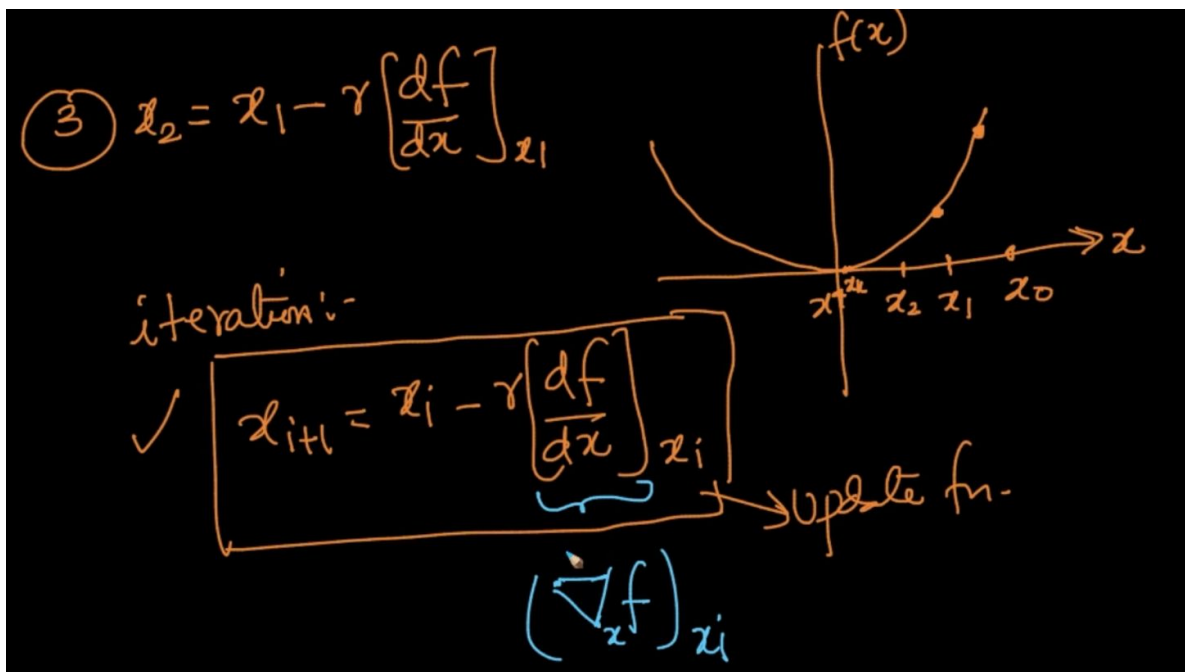
1) We'll pick a random initial point $= x_o$

2) $x_1 = x_0 - r * [\dfrac{df}{dx}]_{x0}$ where 'r' is step-size and $\dfrac{df}{dx}$ is the slope . Also, $x_1 < x_0$

$x_1 = x_0 - \gamma \left[\dfrac{df}{dx}\right]_{x_0}$

$(=1) \quad -ve$

$\uparrow f(x)$

$x_0 \quad x \leftarrow \quad x^*$

$x_1 = x_0 \; \boxed{-1*(-ve)}$

$\boxed{x_1 > x_0}$

Case 2: When $x_0$ is at the other side

$$\text{(3)} \quad x_2 = x_1 - \gamma \left[ \frac{df}{dx} \right]_{x_1}$$

iteration:-

$$\checkmark \quad \boxed{x_{i+1} = x_i - \gamma \left[ \frac{df}{dx} \right]_{x_i}} \quad \rightarrow \text{update fn.}$$

$$\underbrace{\left( \nabla_x f \right)_{x_i}}$$

We'll calculate $x_2$ from $x_1$. Our algorithm will do iteration. When will it stop but?

$$x_0, \; x_1, \; x_2, \; x_3, \; \cdots \cdots \; \textcircled{$x_k$} \; \textcircled{$x_{k+1}$}$$

$$\checkmark \quad x_k = x_{k-1} - \gamma \left[ \frac{df}{dx} \right]_{x_{k-1}}$$

$$\text{if} \; \left( x_{k+1} - x_k \right) \text{ is v-small}$$

$$\text{then terminate @ } \boxed{x^* = x_k}$$

We'll terminate the loop once the difference between two iterations are very small

$[\frac{df}{dx}]_{x0} > \quad [\frac{df}{dx}]_{x1} > \quad [\frac{df}{dx}]_{x2}$ .... As we can see since the derivatives are getting smaller and smaller by each iteration. We are taking a huge step and as we are moving to the convergence our step -sizes are getting smaller and smaller . This is what happens in Gradient Descent
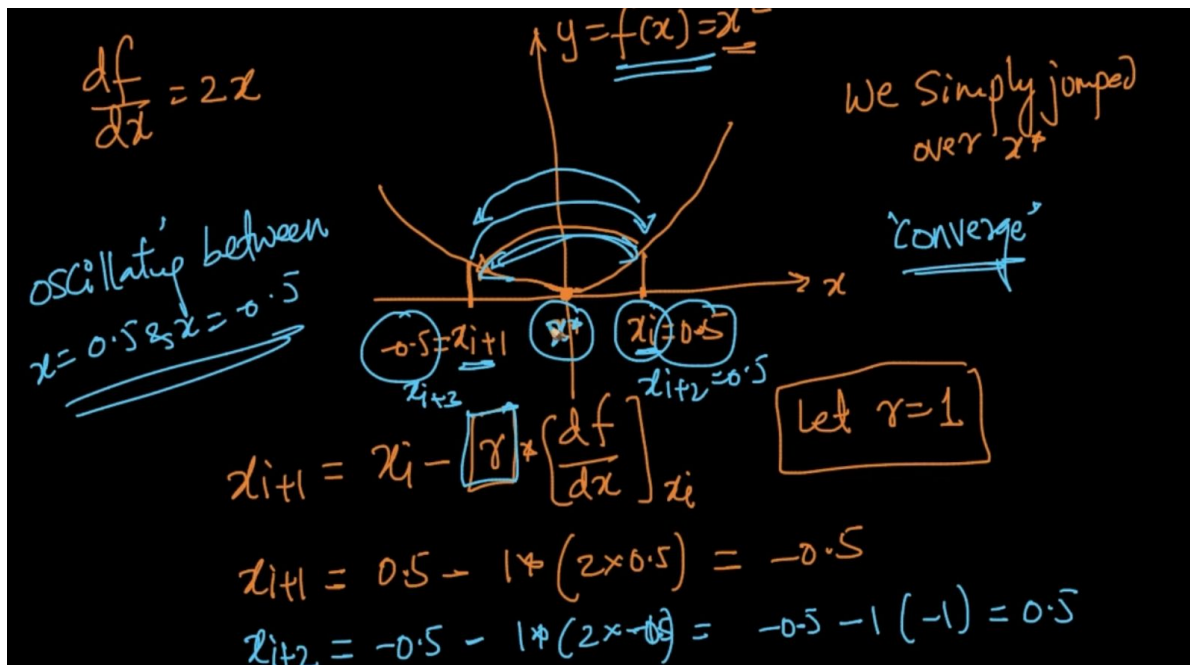
# LEARNING RATE



'r' is basically learning rate or step-size



As we can see at $x_{i+1}$ = -0.5  so we are jumping over $x^*$ and at  $x_{i+2}$ = 0.5 so we are oscillating between -0.5 and +0.5 but never converging towards $x^*$  as our r = 1

remedy for oscillation:- ✓

change $\boxed{r}$ with each iteration

→ one tech:- reduce $\boxed{r}$ with each $\boxed{\text{iteration}}$

$\left( r = \text{function of iteration number} \right)$

$r = h(i)$ ↓ iteration  s.t $i\uparrow ; r\downarrow$

One method is to reduce r with each iteration. So we want to make r = h(i) where h is a function such that as i increases r should be decreased

# GRADIENT DESCENT FOR LINEAR REGRESSION

$$L(w) = \sum_{i=1}^{n} (y_i - w^T x_i)^2 \qquad \langle x_i, y_i \rangle \rightarrow D_{train}$$

$$\nabla_w L = \sum_{i=1}^{n} \left\{ 2(y_i - w^T x_i)(-x_i) \right\}$$

① pick a random vector $w_0 = \langle \cdots \cdots \rangle$

② $w_1 = w_0 - \gamma * \sum_{i=1}^{n} (-2x_i)(y_i - w_0^T x_i)$

③ $w_2 = w_1 - \gamma * \sum_{i=1}^{n} (-2x_i)(y_i - w_1^T x_i)$

L(w) is a function for 'w' and we want to find the best w not x as x,y are constants as given in our dataset. So we need $\Box_w L$. We pick a random vector $w_o$ and apply it to get the next $w_1$. $w_1$ will be used to get $w_2$ and this step will be repeated till difference between $w_k$ and $w_{k+1}$ is very small . If that's the case then $w^*$ = $w_k$

Gradient desc

$$w_j = w_{j-1} - \gamma \sum_{i=1}^{n} (-2x_i)(y_i - w_{j-1}^T x_i)$$
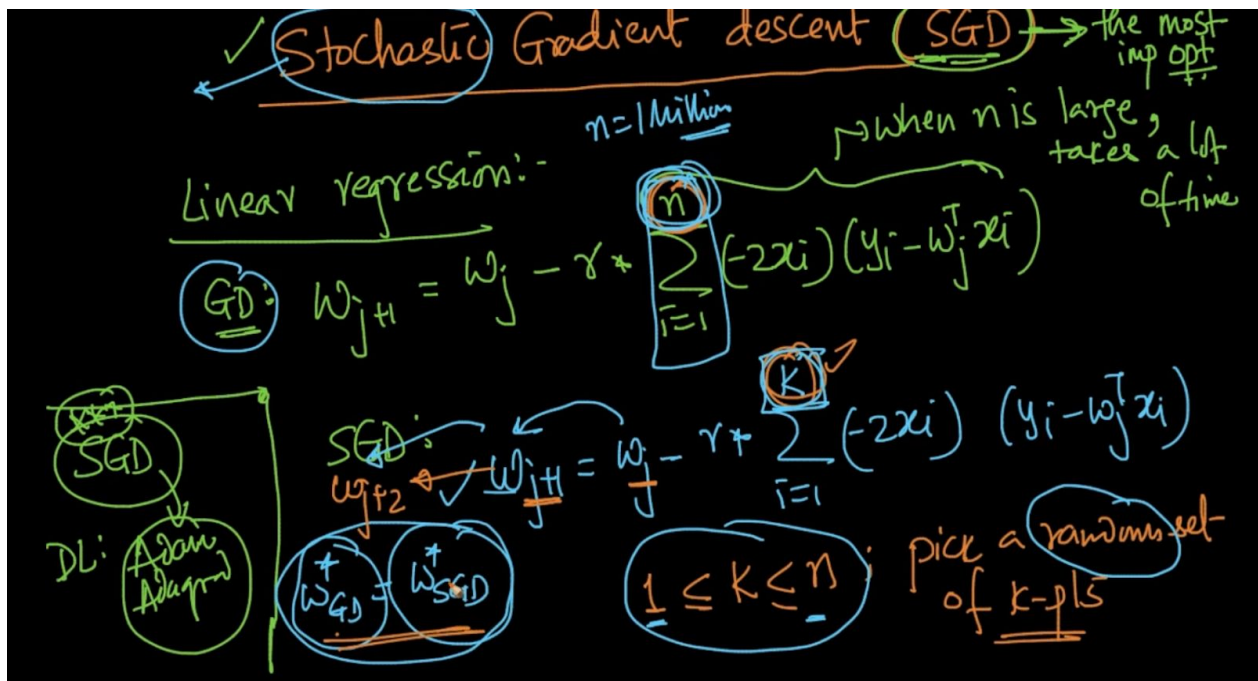
Problem:

$$w_{j-1} \text{ to } w_j$$

$$w_0, w_1, w_2$$

if n is large
n = 1 Million

if ⓝ is large
computing this Σ is very expensive

There's one problem here i.e for every iteration in w i.e w0,w1.. we need to calculate
The summation thing and if 'n' is very large then we need to do it n times for every iteration

# STOCHASTIC GRADIENT DESCENT (SGD)



Stochastic Gradient Descent is the most important optimization algorithm in ML

In Gradient Descent , we'd the problem for w as for every iteration we'd to go through the whole dataset and if n-No. Of inputs is large we are fucked.

We can use SGD (Stochastic Gradient Descent) in which we do the same operation but instead of 'n' we use 'k' where k << n and it is random set of points from the dataset

We can infer that $w^*_{GD} = w^*_{SGD}$ .

NOTE - With every iteration of w i.e $w_{j+1}$ => $w_{j+2}$ we are changing the value of k

GD:- (n=1Mill) [100] iterations to Converge $x^+$

SGD:- K = 1000; if >100 iter$^{ns}$ $x^p$ (SGD)

SGD: K = 10; (1000) ; $x^p$

K = # of random pls that you @ each iterati fu updati

If Gradient Descent, n = 1 Million points and it takes 100 iterations to converge $x^*$ but in SGD, if k = 1000, then it'll take >100 iterations to converge. Let's say 500.Same for k = 10



K : batch size in SGD)

↓

batch of random pls.

often times: (K=1); [SGD] ✓

n << $\boxed{K}$ >1 → batch SGD with 1Mill batch size =10

K: is also called batch size as we are using batch of random points

**CONSTRAINTS AND PCA**



We want to maximize a function f(x) with a constraint i.e g(x) = c. How will we tackle problems when dealing with situations like this



When we need to find a max of function with constraints we add Lagrangian Multipliers

$\lambda$, $\mu$ as seen above. Now take the derivatives w.r.t x, $\lambda$, $\mu$ as seen above. We'll get $\tilde{x}$ , which will be equal to x* that we wanted



$$\text{PCA:-} \quad \max_u \; u^T S u$$
$$\text{s.t } u^T u = 1$$

$$S = \frac{1}{n} \sum_{i=1}^{n} x_i x_i^T$$

$\rightarrow$ Co-var Matrix of X

$$\mathcal{L}(u, \lambda) = u^T S u - \lambda (u^T u - 1)$$

$$\frac{\partial \mathcal{L}}{\partial u} = 0 \implies \frac{\partial}{\partial u}(u^T S u - \lambda u^T u - \lambda) = 0$$

$$\implies \boxed{S u - \lambda u = 0} \implies$$

In PCA, we want the max of $u^T S u$ such that $u^T u = 1$ where S is covariance matrix
We take the derivative of L w.r.t u = 0 and we get Su - $\lambda$u = 0



$$\boxed{S u = \lambda u}$$ $\rightarrow$ definition of eigen-val & eigenvec

$\downarrow$vecti

Cov-Matrix

$$A u = \lambda u$$

eig vec

$\uparrow$eigen val

{ u is the eigen vec of S  
$\lambda$ is the eigen val of S }

u is the eigen-vector of S and $\lambda$ is the eigen-value of S as it's the same we'd seen in PCA

# LOGISTIC REGRESSION REVISITED



We want to find w* such that $w^T w = 1$



When applying Lagrangian we get the value same as w* in which regularization was also there. So regularization can also be thought as imposing an equality constraint

# WHY L1-REGULARIZATION CREATES SPARSITY?



$L_1$ regularization creates sparsity in 'w' as compared to $L_2$ regularization. Let's see why it happens



In both the Regularizations, we are ignoring the Loss and $\lambda$. So, our goal is just to minimize 'w'. In $L_2$, we want to minimize $(w_1^2 + w_2^2 + w_3^2 + .... + w_d^2)$ and in $L_1$ we want to minimize $(|w_1| + |w_2| + ...... |w_d|)$. Let's se contribution of just $w_1$ to our regularization in L i.e $(L(w_1))$
For $L_2$, min of $w_1^2$ is the contribution of $w_1$ to $L_2$ regularization

$L_2$

$L_2(w_1) = w_1^2$

$L_2(w_1) = w_1^2$

$\dfrac{\partial L_2}{\partial w_1} = 2*w_1$

GD:
$(w_1)_0, \; (w_1)_1, \; (w_4)_2 \cdots \quad (w_1)_k$

$L_1$

$L_1(w_1) = |w_1|$

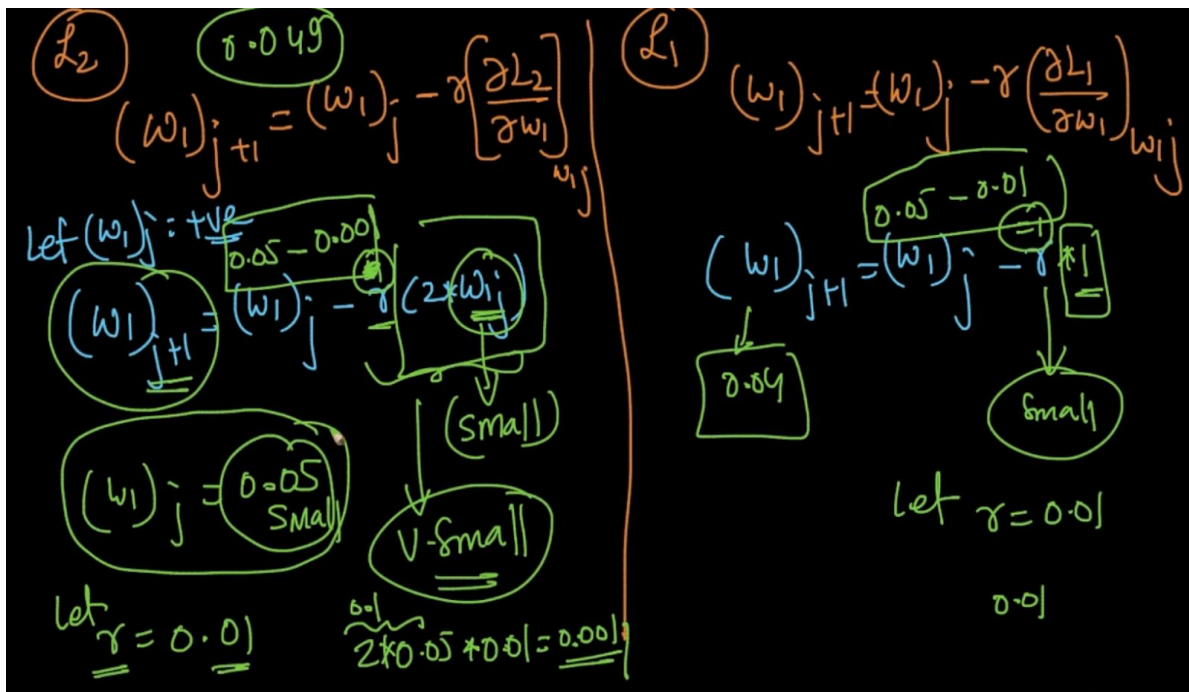$L_1(w_1) = |w_1| = \begin{cases} w_1 & \text{if } w_1 > 0 \\ -w_1 & \text{if } w_1 \le 0 \end{cases}$

$\dfrac{\partial L_1}{\partial w_1} = \begin{cases} +1 & \text{if } w_1 > 0 \\ -1 & \text{if } w_1 < 0 \end{cases}$

We are plotting the graphs of $w_1 \, v.s \, L_2(w_1)$, $w_1 \, v.s \, L_1(w_1)$ and its derivative graph.

Our gradient Descent calculates it for each iteration of $w_1$ till it converges



$L_2$    $0.049$

$(w_1)_{j+1} = (w_1)_j - \gamma \left[\dfrac{\partial L_2}{\partial w_1}\right]_{w_{ij}}$

Let $(w_1)_j : +ve$

$(w_1)_{j+1} = (w_1)_j - \gamma (2*w_1_j)$

$0.05 - 0.001$

$(w_1)_j = 0.05$   Small

Let $\gamma = 0.01$

$2*0.05 * 0.01 = 0.001$

(small)

(V-small)

$L_1$

$(w_1)_{j+1} = (w_1)_j - \gamma \left(\dfrac{\partial L_1}{\partial w_1}\right)_{wij}$

$(w_1)_{j+1} = (w_1)_j - \gamma * 1$

$0.05 - 0.01$

$0.04$

Let $\gamma = 0.01$

$0.01$

(small)

So, our $w_1$ at j'th iteration = 0.05 and r = 0.05. In $L_2$, $(w_1)_{j+1}$ = 0.05 - 0.001 = 0.049 and in

$L_1 , (w_1)_{j+1}$ = 0.05 - 0.01= 0.04 so our $L_1$ is already close to 0 as compared to $L_2$



It can be seen that $L_2$ regularization is slow at converging with each iteration whereas $L_1$ constantly changes as it has a constant slope/gradient so it converges faster towards $w_1^*$