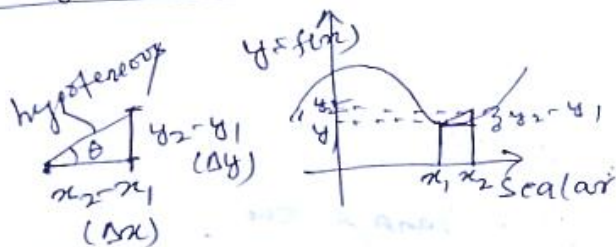


## 27. Solving optimization problem

### 27-1: Differentiation

- optimization problems :- p.c.n, log. Reg., Linear Regression
- Differentiation is used a lot in ml for optimization.

### Single Variable differentiation



$$y = f(x)$$

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

↑ Differentiation of  $y$  w.r.t.  $x$

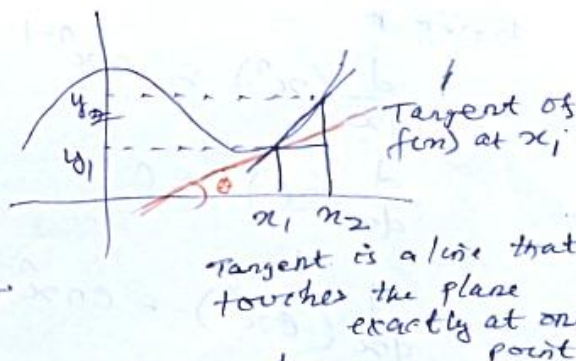
Intuitively

$\frac{dy}{dx}$  means, how much does  $y$  change as  $x$  changes.  
(rate of change of  $y$  w.r.t.  $x$ )

$$\frac{y_2 - y_1}{x_2 - x_1} = \frac{\Delta y}{\Delta x}$$

$$\boxed{\frac{dy}{dx} = \lim_{\Delta x \rightarrow 0} \frac{\Delta y}{\Delta x}}$$

Tangent is the hypotenuse that we obtain as  $\Delta x \rightarrow 0$ .



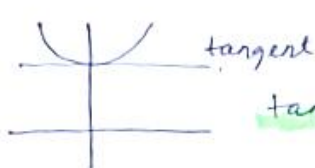
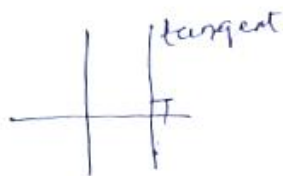
$$\tan \theta = \frac{dy}{dx} \text{ as } \Delta x \rightarrow 0; \quad \frac{\Delta y}{\Delta x} = \frac{dy}{dx}$$

$\frac{dy}{dx}$  = slope of the tangent to  $f(x)$

$\left[ \frac{dy}{dx} \right]_{x_1}$  = slope of the tangent to  $f(x)$  at  $x = x_1$

Tan  $\theta$

angle bet<sup>n</sup> Tangent & x-axis



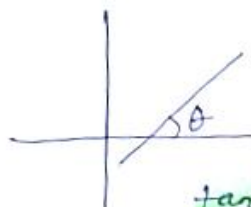
$$\tan \theta = \tan 0$$

$$\tan(90) = \text{undefined} \quad \tan \theta = \tan 0 = 0$$



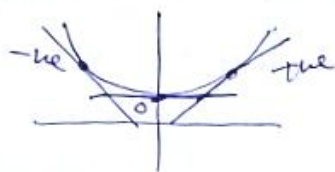
$$90 < \theta \leq 180$$

$$\tan \theta = -ve$$



$$\tan \theta = +ve$$

$$0 < \theta \leq 90$$



Slopes  $\tan \theta$

Intuitively:

differentiation is rate of change of  $y$  as  $x$  changes

Geometrically:

it is slope of the tangent to  $f(x)$

Basic

$$\frac{d}{dx} (x^n) = nx^{n-1} \quad \Rightarrow \quad \frac{d}{dx} x^2 = 2x$$

$$\frac{d}{dx} (c) = 0$$

$$\frac{d}{dx} (3) = 0$$

late 1600's  
Newton  
Leibnitz

$$\frac{d}{dx} (cx^n) = cnx^{n-1}$$

$$\frac{d}{dx} (\log(x)) = \frac{1}{x}$$

$$\frac{d}{dx} (e^x) = e^x$$

$$\frac{d}{dx} (f(x) + g(x)) = \frac{d}{dx} f(x) + \frac{d}{dx} g(x)$$

Chain Rule

$$\frac{d}{dx} f(g(x)) = \frac{df}{dg} \cdot \frac{dg}{dx}$$

### Chain Rule

$$f(g(x)) = (a - bx)^2$$

$$\begin{cases} g(x) = a - bx \\ f(x) = x^2 \end{cases}$$

$$\frac{d}{dx} f(g(x)) = \frac{df}{dg} \cdot \frac{dg}{dx}$$

$$\frac{dg}{dx} = \frac{d}{dx} (a - bx) = \frac{d}{dx} (a) - \frac{d}{dx} (bx) \\ = 0 - b = -b$$

$$\frac{df}{dg} = \frac{d g^2}{d g} = 2g = 2(a - bx)$$

$$\frac{d}{dx} f(g(x)) = 2(a - bx) \cdot -b \\ = -2b(a - bx)$$

### Online Differentiation tools

[www.derivative-calculator.net](http://www.derivative-calculator.net)

$$\log(1 + \exp(ax))$$

$\boxed{g'}$  ↵

ln: natural logarithm

$f(x)$

$f'(x)$ : gradient function

[www.wolframalpha.com/input/](http://www.wolframalpha.com/input/)

derivative of  $x^4 \sin x$  ↵



## 27-2 Online differentiation tools

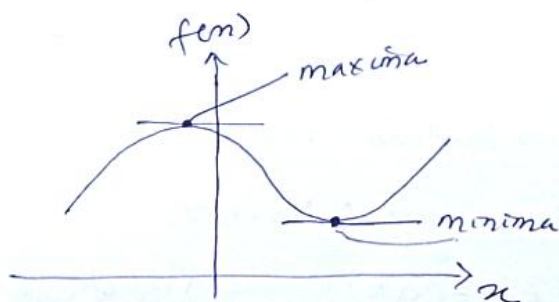
→ [www.derivative-calculator.net](http://www.derivative-calculator.net)

→ Derivative calculator with steps ([www.numberempire.com/derivative-calculator.php](http://www.numberempire.com/derivative-calculator.php))

[www.symbolab.com/...](http://www.symbolab.com/)

[www.wolframalpha.com/...](http://www.wolframalpha.com/)

## 27-3 maxima and minima



At minima and maxima, slope = 0

$f(x) = x^2 - 3x + 2$ , find maxima and minima.

$$\text{slope} = \frac{df}{dx} = 0.$$

$$\frac{df}{dx} = 2x - 3 + 0 = 2x - 3 = 0 \Rightarrow x = \underline{\frac{3}{2}} = 1.5$$

$$\text{At } x = \frac{3}{2}; \text{ slope} = 0.$$

Q// Is this a maxima or minima.

$$f(1.5) = -0.25$$

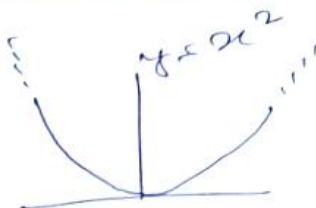
$$f(1) = 0$$

$$\Rightarrow f(1.5) < f(1) \Rightarrow 1.5 \text{ can't be maxima.}$$

This means we have minima at  $x = \underline{1.5}$ .

Example 2

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



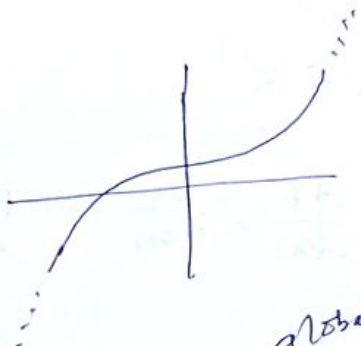
This has no maxima.

but minima at  $x=0$ .

Has no minima

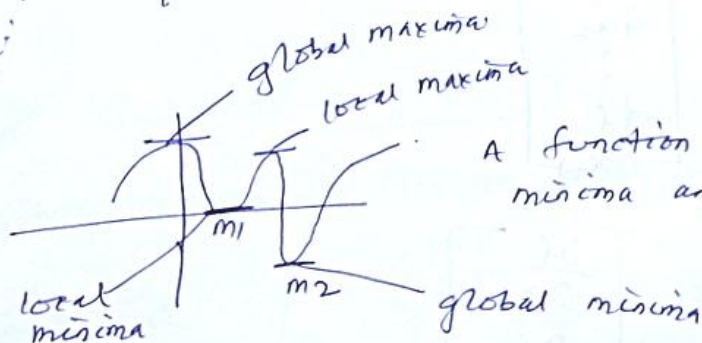
Note # The function may or may not have maxima and minima.

example -3



no minimum  
no maximum

example -4



A function has multiple minima and maxima.

example 5

$$f(x) = \log(1 + \exp(ax))$$

$$\frac{df}{dx} = \frac{a \cdot \exp(ax)}{1 + \exp(ax)} = 0 \quad \left. \vphantom{\frac{df}{dx}} \right\} \begin{array}{l} \text{solving this is} \\ \text{not trivial/easy.} \end{array}$$

$\rightarrow \frac{df}{dx} = 0$ , to find the minima or maxima ~~they~~ <sup>they may</sup> not be always possible. (for complex functions).

$\rightarrow$  Solution: Gradient Descent to find minima & maxima.

# 27-4 Vector Calculus: Grad

So far we have used  
 $x$ : Scalar

What if  $x$  is vector.

$$y = \sum_{i=1}^d a_i x_i = a^T x$$

Example  $f(x) = y = a^T x$        $x = \langle x_1, x_2, \dots, x_d \rangle$

$a = \langle a_1, a_2, \dots, a_d \rangle$  constant

$\frac{dy}{dx}$  vector =  $\nabla_x f$  (grad or Del of  $f$  w.r.t.  $x$ )

$\nabla_x f =$  vector  $\in \mathbb{R}^d$

$$\begin{bmatrix} \frac{df}{dx_1} \\ \frac{df}{dx_2} \\ \vdots \\ \frac{df}{dx_d} \end{bmatrix}$$

$\frac{df}{dx_i} = \frac{\partial f}{\partial x_i} \rightarrow$  partial differentiation.

$$= \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix}$$

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

$\nabla_x f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_d \end{bmatrix} = a$

$\frac{\partial f}{\partial x_1} = a_1$   
 $\frac{\partial f}{\partial x_2} = a_2$   
 $\vdots$



$$\frac{d}{dn}(an) = \underbrace{a}_{\text{Scalar}} \cdot \left[ \nabla_n (a^T n) = \underbrace{a}_{\text{vector}} \right]$$

### Example

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

$\langle n_{ei} y_i \rangle$  : constant (come from  $D_{\text{train}}$ )  
variable is  $w$ .

using chain rule

$$\nabla_{\vec{p}} \mathcal{L} = \left( \frac{-y_i \omega x_i}{1 + \exp(-y_i \omega x_i)} \right) + 2\gamma \omega = 0$$

derivation of conjugate gradient

27.5 Gradient descent algorithm

while solving maxima & minima, we were using

$$\frac{df}{dn} \stackrel{\text{scalar}}{=} 0 \quad \bigg| \quad \nabla_n f \stackrel{\text{vector}}{=} 0$$

But it is not easy for complex function.

But it is not easy to find the minimum of a function.

Alternative soln: Gradient Descent Algo. (for minima)

→ This is an iterative algorithm

$x_D \leftarrow$  latest guess of  $x^*$  (optimal)

our problem is  $x^* = \arg \min_x f(x)$

## Enthalpizatiòn

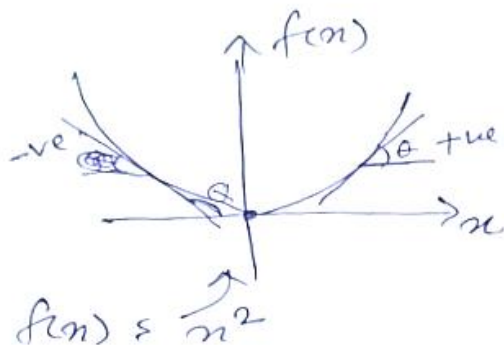
$n_0 \leftarrow$  first guess of  $n^*$

$n_1 \leftarrow \text{Iteration 1}$

$m_2 \leftarrow \text{Iteration 2}$

Gradient Ascent  
Algo. (for maxima):

$$x_k \leftarrow \text{iter. } k.$$



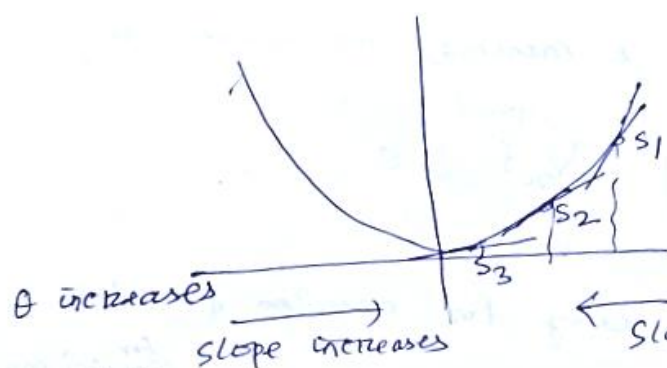
$$f(x) = x^2$$

$$x^* = \underset{x}{\operatorname{argmin}} f(x)$$

nd minima : slope become zero.  
 one side : slope is +ve  
 other side : slope is -ve

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

Geometric observation ① slope changes its sign from +ve to -ve at minima.



② As you move closer to  $x^*$ ; slope reduces.

$\theta$  decreases

How Gradient Descent algo. work?

① pick an initial point (at random)  $x_0$

② find  $x_1$  s.t.  $x_1$  is closer  $x^*$  than  $x_0$

Update function make

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

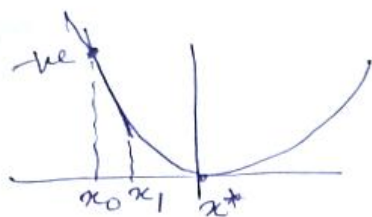
step size  
let step size  $\leq 1$



$$\Rightarrow x_1 = x_0 - 1 \cdot [\text{Some +ve value}]$$

$$\Rightarrow x_1 < x_0$$





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

$$= x_0 - 1 * (-ve \text{ value})$$

$$= x_0 + 1$$

∴  $x_1 > x_0$

③  $x_2 = x_1 - \gamma \cdot \left[ \frac{df}{dx} \right]_{x_1}$



At any iteration,

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

update function of Gradient descent

$x_0, x_1, x_2 \dots x_k$

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

If  $(x_{k+1} - x_k)$  is very small then

terminate the loop. and

return  $x^* = x_k$

If  $x$  is a vector instead of scalar, replace

$\frac{df}{dx}$  with  $\nabla_x f$

$$x_{k+1} = x_k - \gamma \left( \nabla_x f \right)_{x_k}$$

## Gradient Descent Algo.

Initialize  $x^{(0)}$ ,  $\eta > 0$

until convergence do

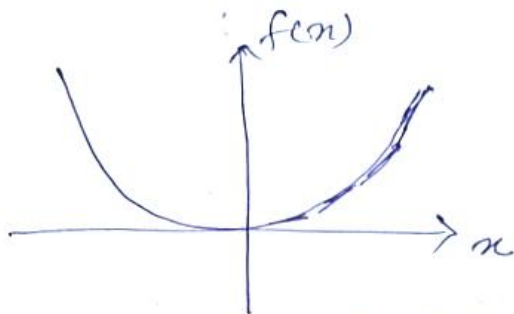
$$x^{(k+1)} = x^{(k)} - \eta \nabla f(x_k)$$

$$x^* = x^{(k+1)}$$

return  $x^*$

problem:

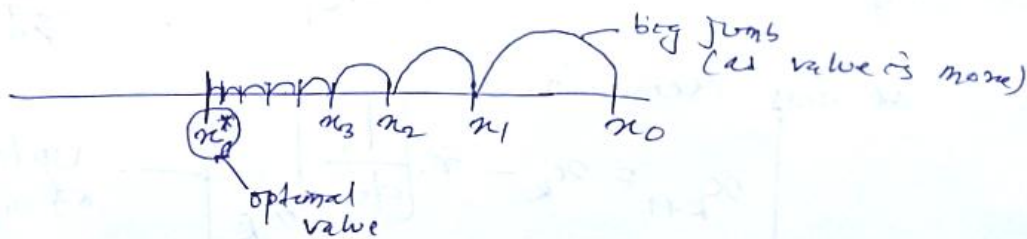
$$x^* = \underset{x}{\operatorname{argmin}} f(x)$$



$$x_{k+1} = x_k - r \cdot \left[ \frac{df}{dx} \right]_{x_k}$$

At every iteration the slope is decreasing.

$$\left[ \frac{df}{dx} \right]_{x_0} \geq \left[ \frac{df}{dx} \right]_{x_1} \geq \left[ \frac{df}{dx} \right]_{x_2} > \dots$$



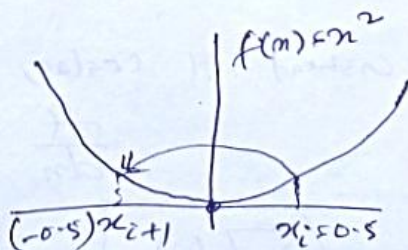
27-6

Learning rate or step size

update equation

$$x_i = x_{i-1} - r \cdot \left[ \frac{df}{dx} \right]_{x_{i-1}}$$

step size / learning rate.



let  $r = 1$

$$\frac{df}{dx} = 2x$$

$$x_{i+1} = x_i - r \cdot \left[ \frac{df}{dx} \right]_{x_i}$$

$$x_{i+1} = 0.5 - 1 \cdot (2 \cdot 0.5) = -0.5$$

We simply jumped over  $x^*$ .

$$x_{i+2} = -0.5 - 1 \cdot (2 \times -0.5) = -0.5 - 1(-1) = 0.5$$

### Oscillation problem

Oscillating between  $x = 0.5$  and  $x = -0.5$ .

This happens because we made  $\alpha$  as constant.

⑩ we will never converge to the optimal point.

### Remedy for oscillation

① one technique is to reduce  $\alpha$  with each iteration.

( $\alpha$  is function of iteration number)

$$\alpha = h(i) \quad \text{s.t. as } i \uparrow; \alpha \downarrow$$

|  
iteration