

## Module-1 VECTOR CALCULUS

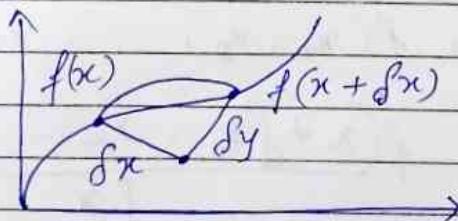
### ⇒ Function of several variables

A real valued function of real variables ' $n$ ' is a function that takes as input  $n$  real numbers, commonly represented by the variables  $\{x_1, x_2, x_3, \dots, x_n\}$  for producing another real number. The value of the function commonly denoted as  $f(x_1, x_2, x_3, \dots, x_n)$

### ⇒ Differentiation & partial differentiation

- \* Differentiation of a uni-variate function  
 The difference quotient of a uni-variate function  $y = f(x)$ ,  $x, y \in \mathbb{R}$  is

$$\boxed{\frac{dy}{dx} = \frac{f(x + \delta x) - f(x)}{\delta x}}$$



It computes the slope of secant line through two points  $f(x_0)$  and  $f(x_0 + \delta x)$

It can also be considered the average slope of the function between  $x$  and  $x + \Delta x$  if we assume the function  $f$  to be a linear function for in the limit  $\Delta x \rightarrow 0$  we obtain the tangent of the function  $f(x)$  at the point  $x$  if  $x$  is differentiable.

→ For  $h > 0$  the derivative of the function  $f(x)$  at the point  $x$  for  $h > 0$  is defined as

$$\frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{f(x+h)}$$

\* Secant line becomes tangent line when  $h \rightarrow 0$

⇒ Partial differentials

For  $x \in \mathbb{R}^n$  of  $n$  variables  $(x_1, x_2, \dots, x_n)$  we define the partial derivative of the function  $f(x_1, x_2, \dots, x_n)$

$$\frac{\partial f}{\partial x_1} = \lim_{\Delta x_1 \rightarrow 0} \frac{f(x_1 + \Delta x_1, x_2, \dots, x_n) - f(x_1, x_2, \dots, x_n)}{\Delta x_1}$$

$$\frac{\partial f}{\partial x_2} = \lim_{\Delta x_2 \rightarrow 0} \frac{f(x_1, x_2 + \Delta x_2, x_3, \dots, x_n) - f(x_1, x_2, \dots, x_n)}{\Delta x_2}$$

$$\frac{df}{dx_n} = \lim_{\delta x \rightarrow 0} \frac{f(x_1, x_2, \dots, x_n + \delta x) - f(x_1, x_2, \dots, x_n)}{\delta x}$$

Similarly,

$$1. \frac{\partial^2 f}{\partial x_1^2} = \frac{\partial}{\partial x_1} \left( \frac{\partial f}{\partial x_1} \right)$$

$$2. \frac{\partial^2 f}{\partial x_2^2} = \frac{\partial}{\partial x_2} \left( \frac{\partial f}{\partial x_2} \right)$$

$$3. \frac{\partial^2 f}{\partial x_1 \partial x_2} = \frac{\partial}{\partial x_1} \left( \frac{\partial f}{\partial x_2} \right)$$

$$4. \frac{\partial^2 f}{\partial x_2 \partial x_1} = \frac{\partial}{\partial x_2} \left( \frac{\partial f}{\partial x_1} \right)$$

$$f(x_1, x_2, \dots, x_n)$$

$$, x_2, \dots, x_m)$$

## Problems 6



1. Obtain the partial derivative for  
 $f(x, y) = (x + 2y^3)^2$

$$\begin{aligned}f(x, y) &= (x + 2y^3)^2 \\&= x^2 + (2y^3)^2 + 2(x)(2y^3) \\&= x^2 + 4y^6 + 4xy^3\end{aligned}$$

$$\begin{aligned}\frac{\partial f}{\partial x} &= \frac{d}{dx} (x^2 + 4y^6 + 4xy^3) \\&= 2x + 4y^3\end{aligned}$$

$$\begin{aligned}\frac{\partial f}{\partial y} &= \frac{d}{dy} (x^2 + 4y^6 + 4xy^3) \\&= 24y^5 + 12xy^2\end{aligned}$$

$$\begin{aligned}\frac{\partial^2 f}{\partial x^2} &= \frac{d}{dx} (2x + 4y^3) \\&= 2\end{aligned}$$

$$\frac{\partial^2 f}{\partial y^2} = \frac{d}{dy} (24y^5 + 12xy^2)$$

$$= \left( 160y^4 + 128xy^3 \right)$$

$$\frac{\partial^2 f}{\partial x \partial y} = \frac{\partial}{\partial x} \left( 24y^5 + 128x^2y^2 \right)$$

$$= 12y^2$$

$$\frac{\partial^2 f}{\partial y \partial x} = \frac{\partial}{\partial y} \left( 2x + 4y^3 \right)$$

$$= 12y^2$$

Q. If  $Z = x^2y + xy^2$ ,  $x = at$ ;  $y = 2at$ . Find  $\frac{dz}{dt}$ .

$$Z = x^2y + xy^2$$

$$x = at, y = 2at$$

$$\begin{aligned} Z &= (at)^2(2at) + (at)(2at)^2 \\ &= a^2t^2(2at) + (at)(4a^2t^2) \\ &= 2a^3t^3 + 4a^3t^3 \end{aligned}$$

$$\frac{dz}{dt} = \frac{d}{dt} (2a^3t^3 + 4a^3t^3)$$

$$\frac{dz}{dt} = [2a^3(3t^2) + 4a^3(3t^2)]$$

$$= \frac{dz}{dt} = 6a^3t^2 + 12a^3t^2$$

$$= \frac{dz}{dt} = 18a^3t^2$$

$$z \rightarrow (x, y) \rightarrow t$$

$$\frac{dz}{dt} = \frac{\partial z}{\partial x} \left( \frac{dx}{dt} \right) + \frac{\partial z}{\partial y} \left( \frac{dy}{dt} \right)$$

$$\frac{\partial z}{\partial x} = f(x^2y + xy^2)$$

$$= 2xy + y^2$$

$$\frac{\partial z}{\partial y} = \frac{d}{dy} (x^2y + xy^2)$$

$$= x^2 + 2xy$$

$$\frac{dx}{dt} = \frac{d(at)}{dt} = a$$

$$\frac{dy}{dt} = \frac{d(2at)}{dt} = 2a$$

$$\frac{dz}{dt} = [2xy + y^2](a) + [x^2 + 2xy](2a)$$

$$x = at ; y = 2at$$

$$\frac{dz}{dt} = \left[ (at)(2at) + (2at)^2 \right] (a) + \left[ a^2 t^2 + 2(at) \cdot (2at) \right] (2a)$$

$$= 4a^3 t^2 + 4a^3 t^2 + 2a^3 t^2 + 8a^3 t^2$$

$$\boxed{\frac{dz}{dt} = 18a^3 t^2}$$

4.  $Z = x^2 + y^2 \rightarrow$  where  $x = e^u \sin v$ ,  
 $y = e^u \cos v$ , find  $\frac{\partial z}{\partial u}, \frac{\partial z}{\partial v}$

$$z \rightarrow (x, y) \rightarrow (u, v)$$

$$\frac{\partial z}{\partial u} = \frac{\partial z}{\partial x} \left( \frac{dx}{du} \right) + \frac{\partial z}{\partial y} \left( \frac{dy}{du} \right)$$

$$\frac{\partial z}{\partial x} = \frac{\partial}{\partial x} (x^2 + y^2) = 2x$$

$$\frac{\partial z}{\partial y} = \frac{\partial}{\partial y} (x^2 + y^2) = 2y$$

$$\frac{dx}{du} = \frac{\partial}{\partial u} (e^u \sin v) = e^u \sin v$$

$$\frac{dy}{dv} = \frac{\partial}{\partial v} (e^u \sin v \cos v) = e^u \sin v$$

$$\frac{\partial z}{\partial u} = \partial x (e^u \sin v) + \partial y (e^u \sin v)$$

$$\text{at } x = e^u \sin v, y = e^u \cos v$$

$$\frac{\partial z}{\partial u} = \partial(e^u \sin v)(e^u \sin v) + \partial(e^u \cos v)(e^u \sin v)$$

$$= 2e^{2u} \sin^2 v + 2e^{2u} \cos^2 v$$

$$= 2e^{2u} (\sin^2 v + \cos^2 v)$$

$$= 2e^{2u}$$

$$\frac{\partial z}{\partial v} = \frac{\partial z}{\partial x} \left( \frac{\partial x}{\partial v} \right) + \frac{\partial z}{\partial y} \left( \frac{\partial y}{\partial v} \right)$$

$$\frac{\partial x}{\partial v} = \frac{\partial}{\partial v} (e^u \sin v) = e^u \cos v$$

$$\frac{\partial y}{\partial v} = \frac{\partial}{\partial v} (e^u \cos v) = -e^u \sin v$$

$$\frac{\partial z}{\partial v} = \partial x (e^u \cos v) + \partial y (-e^u \sin v)$$

$$= \partial(e^u \sin v)(e^u \cos v) + \partial(e^u \cos v)(-e^u \sin v)$$

$$= 2e^{2u} \sin v \cos v - 2e^{2u} \sin v \cos v$$

$$= 0$$

5 For  $f(x_1, x_2) = x_1^2 x_2 + x_1 x_2^3 \in R$   
 compute the partial derivatives. Also obtain  
 the gradient and mention its size.

$$\frac{\partial f}{\partial x_1} = \frac{d}{dx_1} (x_1^2 x_2 + x_1 x_2^3)$$

$$= 2x_1 x_2 + x_2^3$$

$$\frac{\partial f}{\partial x_2} = \frac{d}{dx_2} (x_1^2 x_2 + x_1 x_2^3)$$

$$= x_1^2 + 3x_1 x_2^2$$

$$\frac{\partial^2 f}{\partial x_1^2} = \frac{d}{dx_1^2} (2x_1 x_2 + x_2^3)$$

$$= 2x_2$$

$$\frac{\partial^2 f}{\partial x_2^2} = \frac{d}{dx_2^2} (x_1^2 + 3x_1 x_2^2)$$

$$= 6x_1 x_2$$

$$\frac{\partial^2 f}{\partial x_1 \partial x_2} = \frac{d}{\partial x_1} (x_1^2 + 3x_1 x_2^2)$$

$$= 2x_1 + 3x_2^2$$

$$\frac{\partial^2 f}{\partial x_2 \partial x_1} = \frac{d}{\partial x_2} (2x_1 x_2 + x_2^3)$$

$$= 2x_1 + 3x_2^2$$

## Gradient

$$\nabla = \frac{df(i)}{dx} + \frac{df(j)}{dy} + \frac{df(k)}{dz}$$

$$\nabla f = \frac{df(i)}{dx} + \frac{df(j)}{dy} + \frac{df(k)}{dz}$$

$$\nabla f = \frac{df(i)}{dx_1} + \frac{df(j)}{dx_2}$$

$$= \frac{d}{dx_1} (x_1^2 x_2 + x_1 x_2^3)(i) + \frac{d}{dx_2} (x_1^2 x_2 + x_1 x_2^3)j$$

$$\nabla f = (2x_1 x_2 + x_2^3)i + (x_1^2 + 3x_1 x_2^2)j$$

writing in matrix form,

$$\begin{bmatrix} 2x_1 x_2 + x_2^3 & , & x_1^2 + 3x_1 x_2^2 \end{bmatrix}$$

Therefore, the size of gradient is  
(1 x 2).

6. Let  $f(x_1, x_2) = x_1^2 + 2x_2$  where

$x_1 = 5\sin t$ ,  $x_2 = \cos t$ . Find

$\frac{df}{dt}$ . Also find the gradient.

$$f \rightarrow (x_1, x_2) \rightarrow (t)$$

$$\frac{df}{dt} = \frac{\partial f}{\partial x_1} \left( \frac{dx_1}{dt} \right) + \frac{\partial f}{\partial x_2} \left( \frac{dx_2}{dt} \right)$$

$$\frac{\partial f}{\partial x_1} = \frac{\partial}{\partial x_1} (x_1^2 + 2x_2) = 2x_1$$

$$\frac{\partial f}{\partial x_2} = \frac{\partial}{\partial x_2} (x_1^2 + 2x_2) = 2$$

$$\frac{dx_1}{dt} = \frac{\partial}{\partial t} (\sin t) = \cos t$$

$$\frac{dx_2}{dt} = \frac{\partial}{\partial t} (\cos t) = -\sin t$$

$$\frac{df}{dt} = 2x_1(\cos t) + 2(-\sin t)$$

$$\frac{df}{dt} = 2(\sin t)(\cos t) + 2(-\sin t)$$

$$\boxed{\frac{df}{dt} = 2\sin t \cos t - 2\sin t}$$

## Gradient

$$\nabla f = \left( \frac{\partial f}{\partial x_1} \right) i + \left( \frac{\partial f}{\partial x_2} \right) j$$

$$\nabla f = \frac{d}{dx_1} (x_1^2 + 2x_2) i + \frac{d}{dx_2} (x_1^2 + 2x_2) j$$

$$\nabla f = (2x_1) i + 2j$$

Matrix form  $\Rightarrow [2x_1, 2]$

The size of gradient is  $(1 \times 2)$ .

7.  $u = xy + yz + zx$  at  $x = t \cos t$ ,  
 $y = t \sin t$  and  $z = t$ . Find  $\frac{du}{dt}$

$$u \rightarrow (x, y, z) \rightarrow (t)$$

$$\frac{du}{dt} = \frac{\partial u}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial u}{\partial y} \cdot \frac{dy}{dt} + \frac{\partial u}{\partial z} \cdot \frac{dz}{dt}$$

$$\begin{aligned} \frac{\partial u}{\partial x} &= \frac{\partial}{\partial x} (xy + yz + zx) \\ &= y + z \end{aligned}$$

$$\begin{aligned} \frac{\partial u}{\partial y} &= \frac{\partial}{\partial y} (xy + yz + zx) \\ &= x + z \end{aligned}$$

$$\frac{\partial u}{\partial z} = \frac{\partial}{\partial z} (xy + yz + zx)$$

$$= y + zx$$

$$\frac{\partial x}{\partial t} = \frac{d}{dt} (t \cos t) = -\sin t$$

$$\frac{dy}{dt} = \frac{d}{dt} (-\sin t) = \cos t$$

$$\frac{dz}{dt} = \frac{d}{dt} (t) = 1$$

$$\frac{du}{dt} = (y+z)(-\sin t) + (x+z)(\cos t) + (y+x)(1)$$

$$= (t \sin t + t)(-\sin t) + (t \cos t + t)(\cos t)$$

$$+ (t \cos t + t \sin t)$$

$$= -t \sin^2 t - \cancel{t \sin t} + t \cos^2 t + t \cos t$$
$$+ t \cos t + \cancel{t \sin t}$$

$$= t \cos^2 t - t \sin^2 t + 2t \cos t$$
$$= t (\cos^2 t - \sin^2 t) + 2t \cos t$$

## Gradient of Vector valued functions

For a function  $f : R^n \rightarrow R^m$  and a vector  $x = (x_1, x_2, \dots, x_n)^T \in R^n$  the corresponding vector of function values is given by  $f(x) = \begin{bmatrix} f_1(x) \\ \vdots \\ f_m(x) \end{bmatrix} \in R^m$

The partial derivatives of a vector valued function with respect to  $x_i \in R$  where  $i = 1 \text{ to } n$  is given as the vector  $\frac{\partial f}{\partial x_i} =$

$$\left[ \begin{array}{c} \frac{\partial f_1}{\partial x_i} \\ \frac{\partial f_2}{\partial x_i} \\ \vdots \\ \frac{\partial f_n}{\partial x_i} \end{array} \right]$$

$$\lim_{h \rightarrow 0} \frac{f(x_1, x_2, \dots, \underline{x_i + h}, \dots) - f(x_1, x_2, \dots, \underline{x_i}, \dots)}{h}$$

$$\lim_{h \rightarrow 0} \frac{f(\underline{x_1}, x_2, \dots, \underline{x_i + h}, \dots) - f(x_1, x_2, \dots, \underline{x_i}, \dots)}{h}$$

$$\lim_{h \rightarrow 0} \frac{f(x_1, x_2, \dots, \underline{x_i + h}, \dots) - f(x_1, x_2, \dots, \underline{x_i}, \dots)}{h}$$

we know that the gradient of  $f$  with respect to a vector is the row vector of the partial derivatives. Every partial derivative  $\frac{\partial f_i}{\partial x_j}$  is itself a column

vector. We obtain the gradient of  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  w.r.t  $x \in \mathbb{R}^n$  by collecting these partial derivatives

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} + \frac{\partial f_1}{\partial x_2} + \dots + \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} + \frac{\partial f_2}{\partial x_2} + \dots + \frac{\partial f_2}{\partial x_n} \\ \vdots \\ \frac{\partial f_m}{\partial x_1} + \frac{\partial f_m}{\partial x_2} + \dots + \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

↓

$$= \frac{\partial f_1(x_1)}{\partial x_1} \quad \frac{\partial f_2(x_1)}{\partial x_2} \quad \dots \quad \frac{\partial f_m(x_1)}{\partial x_n}$$

$$\frac{\partial f_1}{\partial x_1} \quad \frac{\partial f_2}{\partial x_2} \quad \frac{\partial f_m}{\partial x_n}$$

$$\begin{array}{ccc} | & | & | \\ | & | & | \\ | & | & | \end{array} \quad \text{ER} \quad m \times n$$

$$\frac{\partial f_1}{\partial x_1} \quad \frac{\partial f_2}{\partial x_2} \quad \frac{\partial f_m}{\partial x_n}$$

The gradient of a function  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  is a matrix of size  $m \times n$  is called the 'Jacobians'.

$$\text{i.e } J = \Delta x f = \frac{\partial f(x)}{\partial x}$$

→ The gradient of matrix with respect to a vector

Let  $A \in \mathbb{R}^{4 \times 2}$  ← matrix of order  $4 \times 2$   
and  $x \in \mathbb{R}^3$  ← vector in 3 dimensional space i.e.

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

Partial derivatives  $\frac{\partial A}{\partial x_i} \in \mathbb{R}^{4 \times 2}$ ,  $\frac{\partial A}{\partial x_j} \in \mathbb{R}^{4 \times 2}$

and  $\frac{\partial A}{\partial x_3} \in \mathbb{R}^{4 \times 2}$

→ ∴ Gradient of a matrix w.r.t a vector  $\frac{\partial A}{\partial x} \in \mathbb{R}^{4 \times 2 \times 3}$

Let  $f = Ax$ ;  $f \in \mathbb{R}^n$ ;  $A \in \mathbb{R}^{m \times n}$   
and  $x \in \mathbb{R}^n$

Gradient is the collection of all partial derivatives

i.e.  $\frac{\partial f}{\partial A} = \left[ \begin{array}{c} \frac{\partial f_1}{\partial A} \\ \frac{\partial f_2}{\partial A} \\ \vdots \\ \frac{\partial f_n}{\partial A} \end{array} \right] \in \mathbb{R}^{m \times (m \times n)}$

where  $\frac{\partial f_i}{\partial A} \in \mathbb{R}^{1 \times (m \times n)}$

In general,  $f_i = \sum_{j=1}^N A_{ij} x_j$ ,  $i = 1, 2, \dots, m$ .

then, the partial derivative of  $f$  w.r.t  
a row of  $A$

$$\frac{\partial f_i}{\partial A_q} = x_q \Rightarrow \frac{\partial f_i}{\partial A_i} = x^T \in \mathbb{R}^{1 \times 1 \times N} \text{ and}$$

$$\frac{\partial f_i}{\partial A_{-k} f_i} = 0^T \in \mathbb{R}^{1 \times 1 \times N}$$

Since  $f^o$  maps onto  $\mathbb{R}$  and each row  
of  $A$  is of size  $1 \times N$  we obtain a  $(1 \times 1 \times 6)$   
sized terms as the partial derivatives  
of  $f^o$  w.r.t a row of  $A$

$$\frac{\partial f^o}{\partial A} = \begin{bmatrix} 0^T \\ \vdots \\ 0^T \\ x^T \\ 0^T \\ \vdots \\ 0^T \end{bmatrix} \in \mathbb{R}^{1 \times (m \times N)}$$

$\Rightarrow$  gradient of matrices w.r.t. matrices

Consider a matrix  $A \in R^{M \times N}$  and  
 $f = R^{M \times N} \rightarrow R^{N \times N}$  with  $f(A) = A^T A$   
=  $K \in R^{M \times N}$

The gradient  $\frac{\partial K}{\partial A}$  is given by,

$$\frac{\partial K}{\partial A} \in A^{(N \times N) \times (M \times N)}$$

$$\Rightarrow \frac{\partial K_{pq}}{\partial A} \in R^{i \times M \times N}$$

for  $p, q = 1, 2, \dots, N$  where  $K_{pq}$   
is the  $(p, q)^{th}$  entry of  $K = f(A)$

$$K_{pq} = \gamma_{pq} T_{pq} = \sum_{m=1}^M A_{mp} \cdot A_{mq}$$

then

$$\frac{\partial K_{pq}}{\partial A_{ij}} = \sum_{m=1}^M \frac{\partial}{\partial A_{ij}} (A_{mp} \cdot A_{mq}) = \gamma_{pqij}$$

where  $\gamma_{pqij} = A_{iq}$  if  $j=p$  and  $p \neq i$

$A_{ip}$  if  $j=q$  and  $p \neq i$

$\frac{\partial A}{\partial A}$  if  $j=p$  and  $p=q$

0, otherwise.

## Problems on Jacobians 6-

1. If  $x_1, y \in \mathbb{R}^2$  and  $y_1 = -2x_1 + x_2$   
 $y_2 = x_1 + x_2$ . Show that the Jacobian determinant |def (J)| = 3

Soln 6.  $y_1 = -2x_1 + x_2$        $y_2 = x_1 + x_2$

$$J = \begin{vmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} \end{vmatrix}$$

$$J = \begin{vmatrix} \frac{\partial}{\partial x_1} (-2x_1 + x_2) & \frac{\partial}{\partial x_2} (-2x_1 + x_2) \\ \frac{\partial}{\partial x_1} (x_1 + x_2) & \frac{\partial}{\partial x_2} (x_1 + x_2) \end{vmatrix}$$

$$J = \begin{vmatrix} -2 + 0 & 0 + 1 \\ 1 + 0 & 0 + 1 \end{vmatrix}$$

$$J = \begin{vmatrix} + & - \\ -2 & 1 \\ 1 & 1 \end{vmatrix}$$

$$= |-2 - 1|$$

$$= |-3|$$

$$J = 3 //$$

2. Find the Jacobian determinant of  
 $y_1 = x_1^2 + x_2^2 + x_3^2 \quad y_2 = x_1 x_2 + x_2 x_3 + x_3 x_1$   
 $y_3 = x_1 + x_2 + x_3$

Soln -  $J = \begin{vmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \frac{\partial y_1}{\partial x_3} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \frac{\partial y_2}{\partial x_3} \\ \frac{\partial y_3}{\partial x_1} & \frac{\partial y_3}{\partial x_2} & \frac{\partial y_3}{\partial x_3} \end{vmatrix}$

$J = \begin{vmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \frac{\partial y_1}{\partial x_3} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \frac{\partial y_2}{\partial x_3} \\ \frac{\partial y_3}{\partial x_1} & \frac{\partial y_3}{\partial x_2} & \frac{\partial y_3}{\partial x_3} \end{vmatrix}$

$\therefore$  Since the R<sub>3</sub> is having same elements

$\Rightarrow J = 0$  ✓.

### \* Product Rule :-

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

### \* Sum Rule :-

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

### \* Chain Rule :-

$$\frac{d}{dx}(f \circ g)(x) = \frac{d}{dx}(g(f(x))) = \frac{\partial g}{\partial f} \cdot \frac{\partial f}{\partial x}$$

If a function  $h: \mathbb{R} \rightarrow \mathbb{R}$ ,  $RCL = f \circ g$   
 with  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$  &  $g: \mathbb{R} \rightarrow \mathbb{R}^2$  and  
 $f(x) = -e^{x_1 x_2}, x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} =$

$g(t) = t \cos t + t \sin t$ . Compute gradient  
 of  $h$  with respect to  $t$ .

Soln of  $\frac{\partial h}{\partial x} \in \mathbb{R}^{1 \times 2}$ ,  $\frac{\partial g}{\partial x} \in \mathbb{R}^{2 \times 1}$

$$f(x) = e^{x_1 x_2}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} t \cos t \\ t \sin t \end{bmatrix}$$

$$x_1 = t \cos t$$

$$x_2 = t \sin t$$

$$\frac{\partial h}{\partial t} = \frac{\partial f}{\partial x} \cdot \frac{\partial x}{\partial t}$$

$$= \left[ \begin{array}{cc} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \\ \end{array} \right] \left[ \begin{array}{c} \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial t} \end{array} \right]$$

$$\frac{\partial f}{\partial x_1} = \frac{\partial}{\partial x_1} \left( e^{x_1 x_2} \right)$$

$$= e^{x_1 x_2} \cdot x_2^2$$

$$\frac{\partial f}{\partial x_2} = \frac{\partial}{\partial x_2} (e^{x_1 x_2^2}) = 2e^{x_1 x_2^2} \cdot e^{x_1 x_2^2}$$

$$\frac{\partial x_1}{\partial t} = \frac{\partial}{\partial t} (t \cos t)$$

- product rule -

$$\frac{\partial x_1}{\partial t} = \cos(t) - (t) \sin t$$

$$\frac{\partial x_2}{\partial t} = \frac{\partial}{\partial t} ((t) \sin t)$$

$$= \sin t + (t) \cos t$$

$$\frac{\partial h}{\partial t} = \begin{bmatrix} e^{x_1 x_2^2} \cdot x_2^2 & 2e^{x_1 x_2^2} \cdot e^{x_1 x_2^2} \end{bmatrix}.$$

$$\begin{bmatrix} \cos t - t \sin t \\ \sin t + t \cos t \end{bmatrix}$$

If  $f(x) = Ax$  where  $f(x) \in \mathbb{R}^m$   
 $A \in \mathbb{R}^{m \times n}$ ,  $x \in \mathbb{R}^n$ , obtain  
gradient and mention its size

Sol:  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$   
 $y = f(x)$

Range  $\mathbb{R}^m$  - Domain  $\in \mathbb{R}^m$ .

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} f_1(x) \\ f_2(x) \\ \vdots \\ f_n(x) \end{bmatrix} = \begin{bmatrix} f_1(x_1, \dots, x_n) \\ f_2(x_1, \dots, x_n) \\ \vdots \\ f_n(x_1, \dots, x_n) \end{bmatrix}$$

gradient collection of partial derivatives

$$\begin{bmatrix} \frac{dy_1}{dx} \\ \frac{dy_2}{dx} \\ \vdots \\ \frac{dy_n}{dx} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & & \vdots \\ \frac{\partial f_n}{\partial x_1} & \frac{\partial f_n}{\partial x_2} & \cdots & \frac{\partial f_n}{\partial x_n} \end{bmatrix}$$

$$f(x) = Ax \quad ; \quad A \in \mathbb{R}^{M \times N} \quad x \in \mathbb{R}^N$$

$$f(x) = y \in \mathbb{R}^M$$

Obtain the gradient of scalar  
 $\phi = 4x_0 + 2x_1 - 3x_2 + x_3$   
 with respect to the matrix

$$x = \begin{bmatrix} x_0 & x_1 \\ x_2 & x_3 \end{bmatrix}$$

$$x^T = \begin{pmatrix} x_0 & x_1 \\ x_1 & x_2 \end{pmatrix}$$

$$\nabla \phi = \frac{\partial \phi}{\partial x}$$

$$\frac{\partial}{\partial x} = \begin{bmatrix} \frac{\partial \phi}{\partial x_0} & \frac{\partial \phi}{\partial x_1} \\ \frac{\partial \phi}{\partial x_1} & \frac{\partial \phi}{\partial x_2} \end{bmatrix}$$

$$\frac{\partial \phi}{\partial x} = \begin{bmatrix} 4 & -3 \\ 2 & 1 \end{bmatrix} = 10$$

Size of the matrix is  $2 \times 2$

Obtain the gradient of scalar

$$y = x_0 + x_1 + 4x_2 - 6x_3 \text{ w.r.t}$$

matrix  $\begin{bmatrix} x_0 & x_1 \\ x_1 & x_2 \end{bmatrix}$

$$x^T = \begin{bmatrix} x_0 & x_2 \\ x_1 & x_3 \end{bmatrix}$$

$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y}{\partial x_0} & \frac{\partial y}{\partial x_2} \\ \frac{\partial y}{\partial x_1} & \frac{\partial y}{\partial x_3} \end{bmatrix}$$

$$\frac{\partial y}{\partial x} = \begin{bmatrix} 1 & 4 \\ 1 & -6 \end{bmatrix} = 10.$$

Size of the matrix is  $2 \times 2$

obtain the gradient of vector

$$f = [e^{x_0 x_1} \cdot e^{x_0 x_3}] \text{ w.r.t}$$

the matrix  $\begin{bmatrix} x_0 & x_1 \\ x_2 & x_3 \end{bmatrix}$

$$x^T = \begin{bmatrix} x_0 & x_2 \\ x_1 & x_3 \end{bmatrix}$$

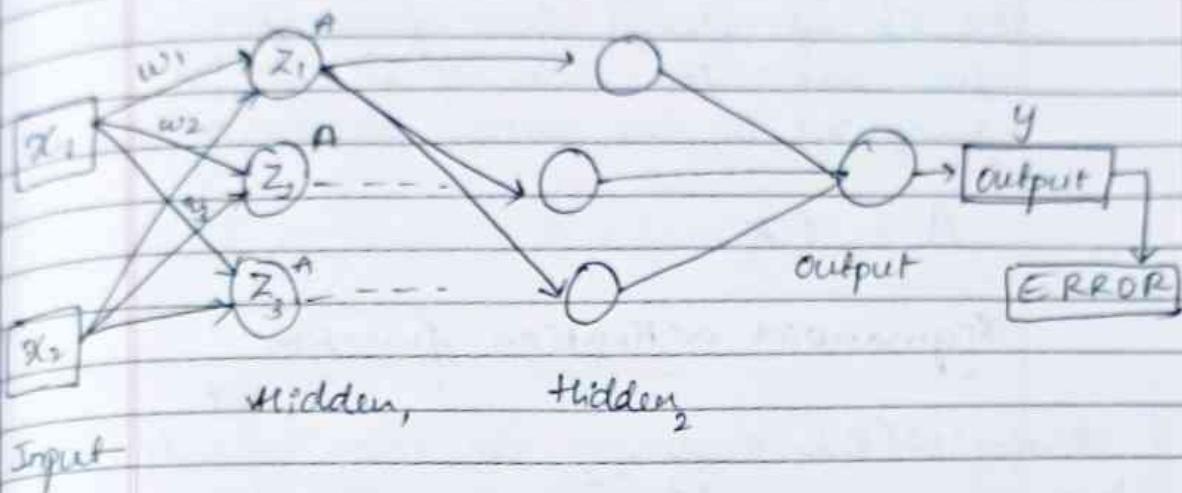
$$\frac{\vec{\partial f}}{\vec{\partial x}} = \begin{bmatrix} \frac{\vec{\partial f}}{\vec{\partial x_0}} & \frac{\vec{\partial f}}{\vec{\partial x_2}} \\ \frac{\vec{\partial f}}{\vec{\partial x_1}} & \frac{\vec{\partial f}}{\vec{\partial x_3}} \end{bmatrix}$$

$$\frac{\vec{\partial f}}{\vec{\partial x}} = \begin{bmatrix} [e^{x_0 x_1} \cdot x_1 \cdot 0] & [0 \cdot e^{x_0 x_3} \cdot x_3] \\ [e^{x_0 x_1} \cdot x_0 \cdot 0] & [0 \cdot e^{x_0 x_3} \cdot x_2] \end{bmatrix}$$

The above matrix is the gradient vector.

## Module - 2

### Back propagation



The whole learning process is divided into a recursive, forward & back -ward propagation.

### Forward propagation

Inputs are passed to neurons of hidden layer with some randomly initialized weights along biases ( $w_1, w_2, w_3$ )

as a linear transformation shown below.

$$Z = (\text{input} * \text{weight}) + \text{bias}$$

To solve complex problems, non-linear

transformation is introduced to be achieved by an activation function.

The o/p of linear transformation ( $z$ ) i.e weighted sum of inputs is supplied to the below activation function.

$$A = f(z)$$

Sigmoid activation function

$$S(z) = \frac{1}{1+e^{-z}}$$

Every layer funct applies this linear followed by non-linear equations hence, final o/p of every layer would be the o/p of activation function.

We get predicted o/p in an o/p layer ( $y$ ) doing same processes in every layer

Backward propagation

Predicted o/p ( $y$ ) may differ from the actual o/p hence loss is calculated using loss function e.g.  $J$

This tells that how diverted our prediction is from the actual.

Let's say loss using mean sum of squared loss function.

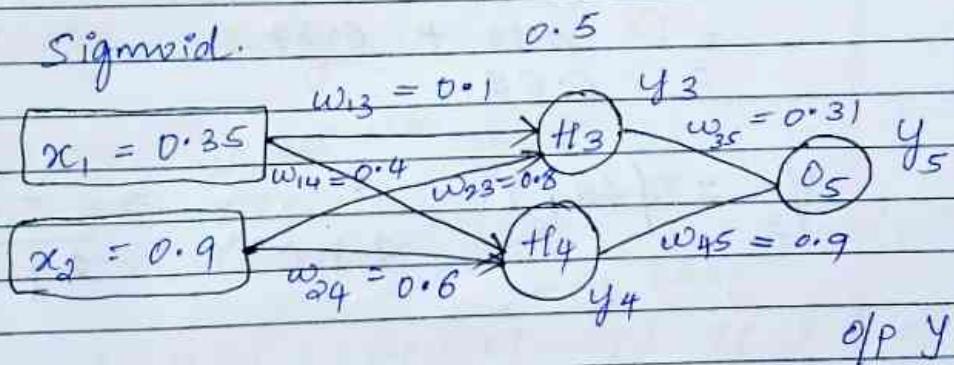
i.e defined as

$$J = \text{loss} = \sum_{i=1}^N (y - \hat{y})^2$$

↑ actual o/p  
↑ predicted o/p

### Problems:-

1. Assume that the neurons have sigmoid activation function. Perform a forward pass and a backward pass on the network.  
- e. Assume that the actual o/p of  $y$  is 0.5 and learning rate is  $\alpha = 1$ . Perform another forward pass. and write the given network also.



Input.

hidden layer

SD1^n^s Forward pass (to calculate o/p  $y_3, y_4, y_5$ )

To calculate  $y_3, y_4$

$$= a_j^o = \sum_j (w_{ij}^o \times x_i)$$

Sigmoid activation function =  $y_j^o = F(a_j^o) = \frac{1}{1+e^{-a_j^o}}$

$$\begin{aligned}y_3 &= a_1 = (w_{13} \times x_1) + (w_{23} \times x_2) \\&= (0.1 \times 0.35) + (0.8 \times 0.9) \\&= 0.035 + 0.72 \\&= 0.755\end{aligned}$$

$$\begin{aligned}y_4 &= a_2 = (w_{14} \times x_1) + (w_{24} \times x_2) \\&= (0.4 \times 0.35) + (0.6 \times 0.9) \\&= 0.14 + 0.54 \\&= 0.68\end{aligned}$$

$$\begin{aligned}y_3 &= f(a_1) = \frac{1}{1+e^{-a_1}} = \frac{1}{1+e^{-0.755}} \\&= 0.68\end{aligned}$$

$$y_4 = f(a_2) = \frac{1}{1+e^{-0.68}} = 0.6637$$

$$y_5 = a_3 = (0.38 \times 0.68) + (0.9 \times 0.6613) \\ = 0.2040 + 0.5942 \\ = 0.8013$$

$$y_5 = f(a_3) = \frac{1}{1 + e^{-0.8013}} = 0.6903$$

$$\Rightarrow \text{output} = y_5 = 0.6903 \approx 0.69$$

$$\Rightarrow \text{Error} = \frac{\text{target}}{\text{Target o/p}} - \frac{y_5}{\text{o/p obtained}} \\ = 0.5 - 0.69 \\ = -0.19$$

- Updating weights to reduce the error -

$$\Delta w_{ij} : \Delta w_{ji} = n \delta_j o_i$$

$o_i \rightarrow$  o/p at  $i^{\text{th}}$  unit

$\delta_j \rightarrow$  error at  $j^{\text{th}}$  unit

$n \rightarrow$  learning rate  $\rightarrow$  target output

$$\delta_j = o_j(1 - o_j)(t_j - o_j) \quad \text{If } J \text{ is an output unit}$$

$$\delta_j = o_j(1 - o_j) \sum_k \delta_k w_{kj} \quad \text{If } J \text{ is a hidden unit}$$

$\rightarrow k$  is the no. of outputs.

Backward pass

Computing  $d_3$   $d_4$   $d_5$  at  $y_3$   $y_4$   $y_5$

For output unit 5

$$d_5 = o_5(1-o_5)(t_5 - o_5)$$

$$= 0.6903 (1-0.6903) (0.5 - 0.6903)$$

$$= 0.6903 (0.3097) (-0.19)$$

$$= -0.0406$$

For hidden unit 3

$$d_3 = h_3(1-h_3) * w_{35} * d_5$$

$$= 0.68(1-0.68) * 0.3 * -0.0406$$

$$= -0.0027$$

$$d_4 = h_4(1-h_4) * w_{45} + d_5$$

$$= 0.6637(1-0.6637) * 0.9 * -0.0406$$

$$= -0.0082$$

$$\Delta w_{45} = n \delta_5 y_4$$

$$= 1 \times -0.0406 \times 0.6031 \\ = -0.0269$$

$$\begin{aligned}(\text{updated}) w_{45} &= \Delta w_{45} + w_{45}(\text{old}) \\ &= -0.0269 + 0.9 \\ &= 0.8731\end{aligned}$$

$$\begin{aligned}\Delta w_{14} &= n \delta_4 y_1 = 1 \times -0.0027 \times \\ &\quad 0.35 \\ &= -0.0029\end{aligned}$$

$$\begin{aligned}(\text{updated}) w_{14} &= -0.0029 + 0.4 \\ &= 0.3971\end{aligned}$$

$$\begin{aligned}\Delta w_{13} &= n \delta_3 x_1 = 1 \times -0.0027 \times 0.35 \\ &= -0.0009\end{aligned}$$

$$(\text{updated}) w_{13} = 0.0991$$

$$\begin{aligned}(\text{updated}) w_{23} &= 0.7976 \\ " w_{24} &= 0.5926 \\ " w_{35} &= 0.2724\end{aligned}$$

Forward pass.

$$\begin{aligned}a_1 &= (w_{13} \times x_1) + (w_{23} \times x_2) \\ &= (0.0991 \times 0.35) + (0.7976 \times 0.9) \\ &= 0.7525\end{aligned}$$

$$y_3 = \frac{1}{1 + e^{-0.7525}} = 0.6797$$

$$a_2 = (w_{14} \times x_1) + (w_{24} \times x_2)$$
$$= 0.6723$$

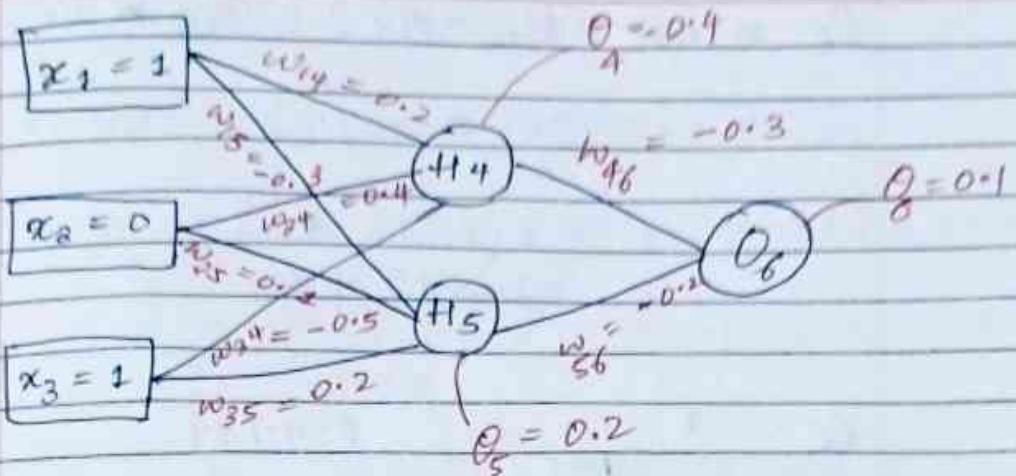
$$y_4 = f(a_2) = \frac{1}{1 + e^{-0.6723}} = 0.6620$$

$$a_3 = (w_{35} \times y_3) + (w_{45} \times y_4)$$
$$= 0.7631$$

$$y_5 = f(a_3) = \frac{1}{1 + e^{-0.7631}} = 0.6820$$

$$\text{Error} = Y_{\text{target}} - y_5$$
$$= 0.5 - 0.6820$$
$$= -0.1820$$

2. Assume that the neurons have a sigmoid activation function. Perform a forward pass & backward pass on the network. Assume that the actual opf of  $y=1$  and learning rate is 0.9. perform another forward pass.



Sol<sup>n</sup> Forward Pass &

$$\begin{aligned}
 H_4 &= a_1 = (w_{14} \times x_1) + (w_{24} \times x_2) + (w_{34} \times x_3) \\
 &= (0.2 \times 1) + (0.4 \times 0) + (-0.5 \times 1) \\
 &= 0.2 + (-0.5) \\
 &= -0.3 - 0.4 \\
 &= 0.3318
 \end{aligned}$$

$$H_4 = f(a_1) = \frac{1}{1 + e^{-a_1}} = \frac{1}{1 + e^{+0.7}} = 0.3318$$

$$\begin{aligned}
 H_5 &= a_2 = (w_{15} \times x_1) + (w_{25} \times x_2) + (w_{35} \times x_3) \\
 &= (-0.3 \times 1) + (0.1 \times 0) + (0.2 \times 1) \\
 &= -0.3 + 0.2 \\
 &= -0.1 + 0.2 \\
 &= 0.5250
 \end{aligned}$$

$$H_5 = f(a_2) = \frac{1}{1 + e^{-a_2}} = \frac{1}{1 + e^{-0.1}} = 0.5250$$

Evaluating biases  $\Rightarrow H_4 = 0.3318$   
 $\Rightarrow H_5 = 0.5250$

~~$H_4 = 0.3318 - 0.4 =$~~

$$\begin{aligned}
 O_6 &= o_3 = (w_{46} \times h_4) + (w_{56} \times h_5) \\
 &= (-0.3 \times 0.3318) + (-0.2 \times 0.5261) \\
 &= -0.0995 + -0.1050 \\
 &= -0.2045 + 0.1
 \end{aligned}$$

$$O_6 = \frac{1}{1 + e^{-0.2045}} = 0.4739$$

$\Rightarrow O_6 = 0.4739$

Error = Target o/p - Obtained o/p

$$\Rightarrow (\text{Error} = 1 - 0.4739)$$

### Backward pass

- Computing  $d_4$   $d_5$   $d_6$

For output unit 6

$$\begin{aligned}
 d_6 &= O_6(1-O_6)(t_6 - O_6) \\
 &= 0.4739(1-0.4739)(1-0.4739) \\
 &= 0.4739(0.5261)(0.5261) \\
 &= 0.1312
 \end{aligned}$$

for hidden unit 6

$$\begin{aligned}d_4 &= h_4(1-h_4) * w_{46} * d_6 \\&= 0.3318(1-0.3318) * (-0.3) * 0.1312 \\&= 0.3318(0.6682) * (-0.3) * 0.1312 \\&= -0.0087\end{aligned}$$

$$\begin{aligned}d_5 &= h_5(1-h_5) * w_{56} * d_6 \\&= 0.5250(1-0.5250) * (-0.2) * 0.1312 \\&= 0.5250(0.4750) * (-0.2) * 0.1312 \\&= -0.0065\end{aligned}$$

$$\begin{aligned}\Delta w_{46} &= n d_6 \neq d_4 \\&= 0.9 \times 0.1312 \times 0.3318 \\&= 0.0392\end{aligned}$$

$$\begin{aligned}\text{updated } w_{46} &= 0.0392 + (-0.3) \\&= -0.2608\end{aligned}$$

$$\begin{aligned}\Delta w_{14} &= n d_4 x_1 = 0.9 \times -0.0087 \times 1 \\&= -0.0078\end{aligned}$$

$$\begin{aligned}\text{updated } w_{14} &= -0.0078 + (0.2) \\&= 0.1922\end{aligned}$$

$$\begin{aligned}\Delta w_{25} &= n d_5 x_2 = 0.9 \times -0.0065 \times 0 \\&= 0\end{aligned}$$

$$\text{updated } w_{25} = 0.1$$

Similarly:

$$w_{14} = 0.192$$

$$\Delta w_{14} = -0.0078$$

$$w_{15} = -0.306$$

$$\Delta w_{15} = -0.0059$$

$$w_{24} = 0.9$$

$$\Delta w_{24} = 0$$

$$w_{25} = 0.1$$

$$\Delta w_{25} = 0$$

$$w_{34} = -0.508$$

$$\Delta w_{34} = -0.0078$$

$$w_{35} = 0.194$$

$$\Delta w_{35} = -0.0059$$

$$w_{46} = -0.261$$

$$\Delta w_{46} = 0.0392$$

$$w_{56} = -0.138$$

$$\Delta w_{56} = -0.2008 \cdot 0.0619$$

— Changing bias value.—

$$\Delta \theta_6 = n d_6 = 0.9 \times 0.1312 \\ = 0.1181$$

$$\theta_6 = \Delta \theta_6 + (\text{old}) \theta_6 \\ = 0.1181 + 0.1 \\ = 0.2181$$

$$\Delta \theta_4 = n d_4 = 0.9 \times -0.0087 \\ = -0.0078$$

$$\theta_4 = \Delta \theta_4 + (\text{old}) \theta_4 \\ = -0.0078 - 0.4 \\ = -0.4078$$

$$\Delta \theta_5 = n d_5 = 0.9 \times -0.0065 \\ = -0.0059$$

$$\begin{aligned}\theta_5 &= \Delta\theta_5 + (\text{old})\theta_5 \\ &= -0.0059 + 0.2 \\ &= 0.1942.\end{aligned}$$

Updated bias  $\delta$

$$\theta_4 = -0.4078$$

$$\theta_5 = 0.1942$$

$$\theta_6 = 0.2181$$

Forward pass &

$$\begin{aligned}\Rightarrow h_4' &= (w_{14} \times x_1) + (w_{24} \times x_2) + (w_{34} \times x_3) \\ &= (0.192 \times 1) + (-0.508 \times 1) \\ &= 0.192 - 0.508 \\ &= -0.3160 + (-0.4078) = -0.7238\end{aligned}$$

$$\begin{aligned}\Rightarrow h_5' &= (w_{15} \times x_1) + (w_{35} \times x_3) \\ &= -0.306 + 0.194 \\ &= -0.1120 + 0.1942 = 0.0822\end{aligned}$$

$$\theta_6' \neq w_{46} \times h_4'$$

after bias

$$\Rightarrow y_4' = \frac{1}{1+e^{-0.3160}} = 0.4217 \Rightarrow 0.3266.$$

$$\Rightarrow y_5' = \frac{1}{1+e^{-0.1120}} = 0.4720 \Rightarrow 0.5205.$$

$$\begin{aligned}\Rightarrow \theta_6' &= (w_{46} \times y_4') + (w_{56} \times y_5') \\ &= (-0.261 \times 0.4217) + (-0.138 \times 0.4720) \\ &= -0.1571 + 0.2181 \\ &= 0.0610\end{aligned}$$

$$\Rightarrow \boxed{y_6 = 0.5152.}$$

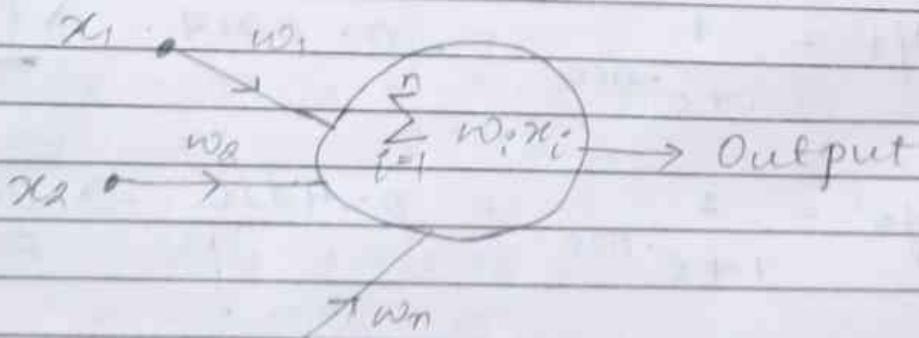
$$\underline{\text{Error}_6} = 1 - 0.5152$$

$$\therefore \boxed{\text{Error} = 0.4848.}$$

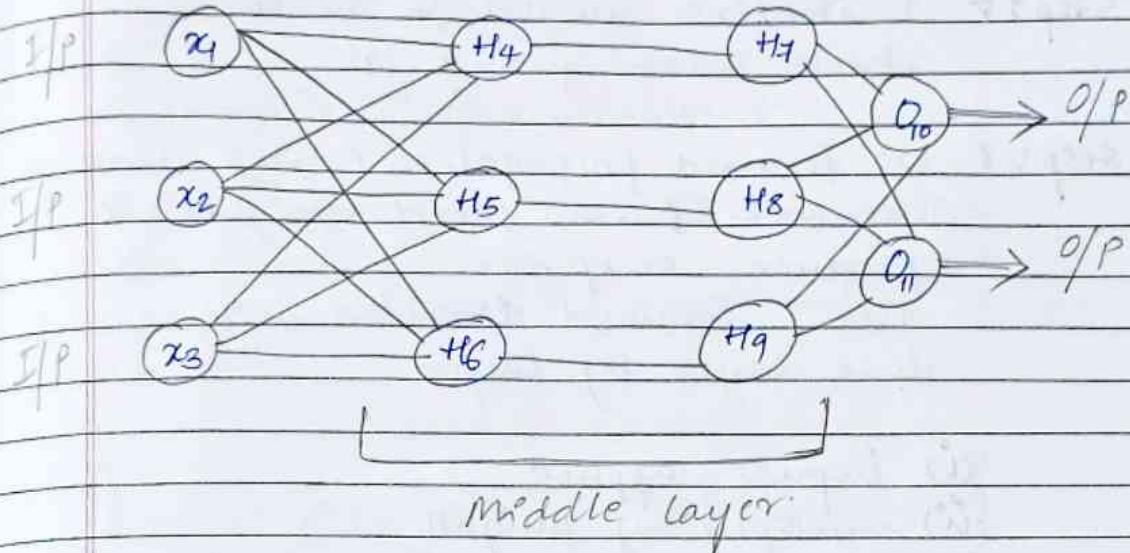
Explain Neural Networks &

- \* Neural networks are made up of many artificial neurons.
- \* Neuron can have any no. of inputs from 1 to n.  
 $(x_1, x_2, \dots, x_n)$
- \* Each input into the neuron has its own weight ( $w_1, w_2, \dots, w_n$ ) .
- \* O/p is equal to  $(x_1w_1 + x_2w_2 + \dots + x_nw_n)$

$$\text{Output} = \sum_{i=1}^n x_i w_i$$



## Multi-layer Network f



## Sigmoid Activation function f

$$1. \quad f(a_j) = \frac{1}{1 + e^{-a_j}}$$

$$2. \quad \tanh x : \quad f(x_a) = \frac{1 - e^{-a_i}}{1 + e^{-a_i}}$$

## Back propagation Algorithm

Step 1 Initialize weights @ random  
choose learning rate 'n'

Step 2 Do forward propagation (pass) through network (with fixed weights) to produce outputs.

Here, forward direction will be done layer by layer

- (i) inputs applied
- (ii) multiply by weights
- (iii) sum and squared by activation function
- (iv) Output passed to each neuron in next layer.

Repeat all the above until network op is produced

Step 3 Do Back propagation and error calculation

(i) Compute the error ( $\Delta$  or local gradient) for each output unit  $\Delta_k$ .

(ii) Compute error for each hidden unit  $\Delta_j$  by back propagation errors layer by layer.

step 4: Update all the weights  $w_{ij}$  by gradient descent and go back to step 2 and repeat until the network training completion

### Automatic differentiation

This is special case of Back propagation.

Consider,

$$f = x^2(x+y) \text{ where } c = x^2$$

$$S = x+y$$

$$f = x^2(x+y), \quad F = c(S)$$

$$\frac{\partial F}{\partial x} = \frac{\partial F}{\partial c} \cdot \frac{\partial c}{\partial x} + \frac{\partial F}{\partial S} \cdot \frac{\partial S}{\partial x}$$

$$= S \cdot (2x) + c(1)$$

$$\Rightarrow \frac{\partial F}{\partial x} = 2xs + c$$

$$\frac{\partial F}{\partial y} = \frac{\partial F}{\partial c} \cdot \frac{\partial c}{\partial y} + \frac{\partial F}{\partial S} \cdot \frac{\partial S}{\partial y}$$

$$= S(0) + c(1)$$

$$\Rightarrow \frac{\partial F}{\partial y} = c$$

# Gradients in a deep network

- \* The function value 'y' is computed as a many level function composition.

$$y = (f_k \circ f_{k-1} \circ \dots \circ f_1)(x)$$

w.r.t  $f \circ g = f(g(x))$

$$\Rightarrow \text{for e.g. } f(x) = x^2 \\ g(x) = x+1$$

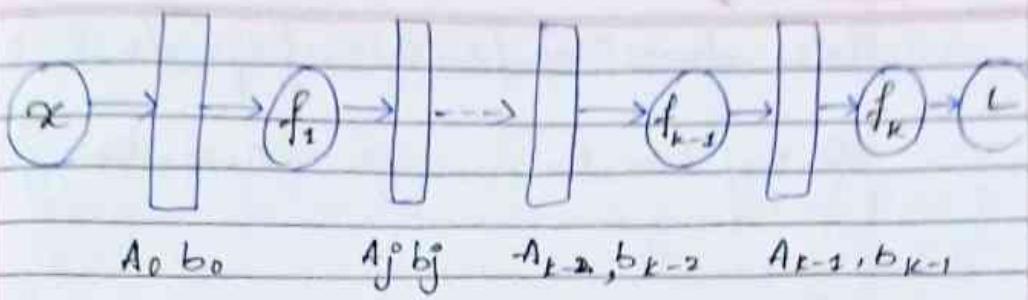
$$\begin{aligned} f \circ g &= f[g(x)] \\ &= f(x+1) \\ &= (x+1)^2 \end{aligned}$$

Hence, similarly  $y = (f_k(f_{k-1}(\dots(f_1(x))))$

- \* Here,  $x$  are inputs or images and  $y$  are observations.

$f_i ; i=1, \dots, k$ . possesses its own parameters

- \* Forward pass is a multi-layer neural network to compute the loss ' $L$ ' as a function of input  $x$  and the parameters  $A_i, B_i, b_i$



here,

$$f_i(x_{i-1}) = \sigma(A_{i-1}x_{i-1} + b_{i-1})$$

$f_i \rightarrow$  function in the  $i^{th}$  layer

$x_{i-1} \rightarrow$  output of the layer  $i-1$

$\sigma \rightarrow$  Activation function logistic  
sigmoid  $\frac{1}{1+e^{-x}}$

$$\tanh x = \frac{1-e^{-x}}{1+e^{-x}}$$

ReLU  $\rightarrow$  Rectified Linear Unit

In order to train the models, we require the gradient of loss function ' $L$ ' w.r.t to all model parameters  $A_j^T, b_j$ . where  $j = 1, 2, \dots, k$ .

$f_0 = x \rightarrow$  input

$$f_i = \sigma_i(A_{i-1}f_{i-1} + b_{i-1}) \text{ for } i = 1, 2, \dots, k$$

To find  $A_j^T, b_j$  for  $j = 0, 1, 2, \dots, k-1$ .

we have to minimize the loss function

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \|f_i(\theta, x_i) - y_i\|^2$$

Loss function  $L(\theta) = \|y - f_k(\theta, x)\|$ .

$$\theta = \{A_0, b_0, \dots, A_{k-1}, b_{k-1}\}$$

By definition of maxima or minima,  
we require the partial derivatives of  
 $L$  w.r.t parameters.

$$\theta^j = (A_j, b_j) ; j = 0, 1, \dots, k-1$$

By chain rule we have,

$$\frac{\partial L}{\partial \theta_{k-1}} = \frac{\partial L}{\partial f_k} \cdot \frac{\partial f_k}{\partial \theta_{k-1}}$$

$$\frac{\partial L}{\partial \theta_{k-2}} = \frac{\partial L}{\partial f_k} \cdot \underbrace{\frac{\partial f_k}{\partial \theta_k}}_{\downarrow} \cdot \underbrace{\frac{\partial \theta_k}{\partial \theta_{k-2}}}_{\text{partial derivative w.r.t inputs of op layer w.r.t parameters}}$$

$$\frac{\partial L}{\partial \theta_{k-3}} = \frac{\partial L}{\partial f_k} \cdot \frac{\partial f_k}{\partial \theta_{k-1}} \cdot \frac{\partial \theta_{k-1}}{\partial \theta_{k-2}} \cdot \frac{\partial \theta_{k-2}}{\partial \theta_{k-3}}$$

⋮

$$\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial f_k} \cdot \frac{\partial f_k}{\partial \theta_{k-1}} \cdots \frac{\partial \theta_{i+2}}{\partial \theta_{i+1}} \cdot \frac{\partial \theta_{i+1}}{\partial \theta_i}$$

## Optimization using Second derivative Test -

A necessary condition for  $x_0$  to be an extreme point (max or min point) of  $f(x)$  is the  $\nabla f(x_0) = 0$  where  $x \in \mathbb{R}^m$  and  $x_0 \in \mathbb{R}^n$  is called critical point/stationary point. A sufficient condition for a stationary point  $x_0$  of  $f(x)$  to be an extreme (maximum or minimum).

- $f$  has relative minimum if  $f_{xx}(x_0)f_{yy}(x_0) - f_{xy}^2(x_0)$  is greater than 0. If  $f_{xx}(x_0) > 0$
- $f$  has relative maximum if  $f_{xx}(x_0)f_{yy}(x_0) - f_{xy}^2(x_0) > 0$  and  $f_{xx}(x_0) < 0$ .
- $f$  has a saddle point (neither max nor min.) If  $f_{xx}(x_0)f_{yy}(x_0) - f_{xy}^2(x_0) < 0$ .
- The test is inconclusive if  $f_{xx}(x_0)f_{yy}(x_0) - f_{xy}^2(x_0) = 0$ .

## Hessian Matrix form &

A sufficient condition for an stationary point  $x_0$  of  $f(x_0)$  to be an extremum,  $f$  has:

- $f$  has relative minimum if  $H(x_0)$  is positive definite
- $f$  has relative maximum if  $H(x_0)$  is negative definite

- is a saddle point if  $H(x_0)$  is indefinite  
 → the test is inconclusive if otherwise

### Problems 6-

2. Using the Hessian matrix classify the relative extremum for  $f(x_1, x_2)$

$$f/3 x_1^3 + x_2 x_2^2 - 8x_1 x_2 + 3.$$

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} x_1^2 + x_2^2 - 8x_2 \\ 2x_1 x_2 - 8x_1 \end{bmatrix}$$

$$\nabla f \begin{bmatrix} \delta f / \delta x_1 \\ \delta f / \delta x_2 \end{bmatrix} = 0$$

$$= x_1^2 + x_2^2 - 8x_2 = 0 \quad 2x_1 x_2 - 8x_1 = 0.$$

$$2x_1 x_2 = 8x_1$$

$$x_2 = \frac{8x_1}{2x_1} = 4$$

$$x_1 = 4$$

$$\Rightarrow x_1 = 0$$

Taking  $x_1 = 0$  in  $x_1^2 + x_2^2 - 8x_2 = 0$ .

$$= 0^2 + x_2^2 - 8x_2 = 0$$

$$x_2^2 = 8x_2$$

$$x_2 = 8 \text{ or } x_2 = 0$$

Taking  $x_2 = 4$  in  $x_1^2 + x_2^2 + -8x_2 = 0$

$$= x_1^2 + 16 - 8(4) = 0$$

$$= x_1^2 + 16 - 32 = 0$$

$$x_1^2 = 16$$

$$\Rightarrow x_1 = 4 \text{ or } -4$$

Points are  $(0, 0), (0, 4), (4, 4), (-4, 4)$

$$\frac{\delta f}{\delta x_1, x_2} = \frac{\delta^2 f}{\delta x_1^2} = 2x_1$$

$$\frac{\delta f}{\delta x_2, x_1} = \frac{\delta^2 f}{\delta x_2^2} = 2x_2$$

$$\frac{\delta^2 f}{\delta x_1 \delta x_2} = 2x_2 - 8$$

$$\frac{\delta^2 f}{\delta x_2 \delta x_1} = 2x_1 - 8$$

$$H = \begin{bmatrix} \frac{\delta^2 f}{\delta x_1^2} & \frac{\delta^2 f}{\delta x_1 \delta x_2} \\ \frac{\delta^2 f}{\delta x_2 \delta x_1} & \frac{\delta^2 f}{\delta x_2^2} \end{bmatrix}$$

$$H = \begin{bmatrix} 2x_1 & 2x_2 - 8 \\ 2x_2 - 8 & 2x_1 \end{bmatrix}$$

H at all points is-

$$H(x_0) \text{ at } (0, 0) = \begin{bmatrix} 0 & -8 \\ -8 & 0 \end{bmatrix} = 0 - 64 = -64 \text{ is indefinite} \\ \rightarrow \text{saddle point.}$$

$$H(x_0) \text{ at } (0, 4) = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} = 0 \text{ indefinite} \\ \rightarrow \text{saddle point.}$$

$$H(x_0) \text{ at } (4, 4) = \begin{bmatrix} 8 & 0 \\ 0 & 8 \end{bmatrix} = 64 - 0 = 64 \text{ +ve definite} \\ \rightarrow \text{minimum point}$$

$$H(x_0) \text{ at } (-4, 4) = \begin{bmatrix} -8 & 0 \\ 0 & -8 \end{bmatrix} = 64 - 0 = 64 \text{ -ve definite} \\ \rightarrow \text{maximum point}$$

minimum point  $(4, 4)$

maximum point  $(-4, 4)$

a.  $f(x_1, x_2, x_3) = x_1 + 2x_3 + x_2 x_3 - x_1^2$

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \frac{\partial f}{\partial x_3} \end{bmatrix} = \begin{bmatrix} 1 - 2x_1 \\ x_3 - 2x_2 \\ 2 + x_2 - 2x_3 \end{bmatrix}$$

$$\nabla f = 0$$

$$1 - 2x_1 = 0$$

$$x_3 - 2x_2 = 0$$

$$2 + x_2 - 2x_3 = 0$$

$$2x_1 = 1$$

$$x_3 = 2x_2$$

$$2 + x_2 - 4x_2 = 0$$

$$x_1 = 1/2$$

$$2 - 3x_2 = 0$$

$$x_3 = 4/3$$

$$2 = 3x_2$$

$$\therefore x_1 = 1/2, x_2 = 2/3, x_3 = 4/3$$

Hessian Matrix  $f$ -

$$H(x_0) = \begin{bmatrix} f_{x_1 x_1} & f_{x_1 x_2} & f_{x_1 x_3} \\ f_{x_2 x_1} & f_{x_2 x_2} & f_{x_2 x_3} \\ f_{x_3 x_1} & f_{x_3 x_2} & f_{x_3 x_3} \end{bmatrix}$$

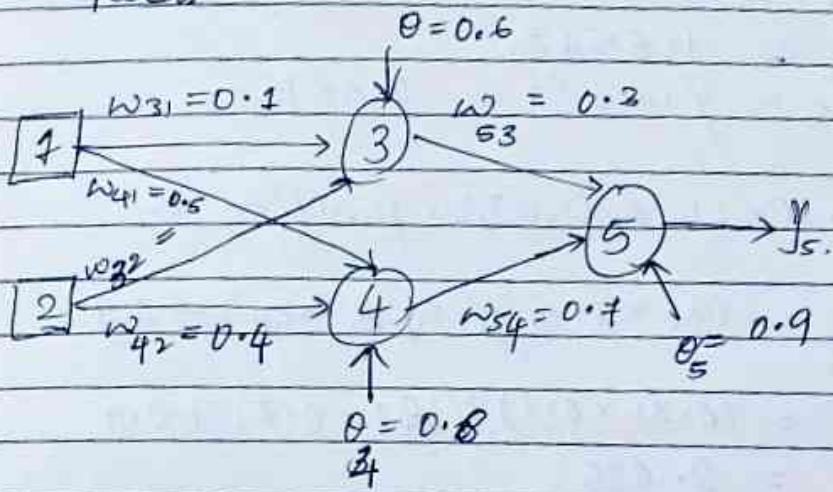
$$= \begin{bmatrix} -2 & 0 & 0 \\ 0 & -2 & 1 \\ 0 & 1 & -2 \end{bmatrix}$$

$$H(x) = \begin{bmatrix} -2 \end{bmatrix} = -2 < 0$$

$$H(x) = \begin{bmatrix} -2 & 0 \\ 0 & -2 \end{bmatrix} = 4 - 0 = 4 \quad \text{-ve definite, maximum value.}$$

$$H(x) = \begin{bmatrix} -2 & 0 & 0 \\ 0 & -2 & 1 \\ 0 & 1 & -2 \end{bmatrix} = -2(4-1) - 0(2) + 0(1) \\ = -6 > 0$$

3. Find output neuron 5, if input vector  $(0.7, 0.3)$  using the activation function ReLU.



$$x_1 = 0.7 \quad x_2 = 0.3$$

$$h_3^f = a_3 = (0.7 \times 0.1) + (0.3 \times 0.5) + 0.6 = 0.82$$

$$h_3^f = \frac{1}{1 + e^{-a_3}} = 0.6942$$

$$h_4^f = a_4 = 1.2$$

$$h_4^f = \frac{1}{1 + e^{-a_4}} = 0.7807$$

## ReLU activation function

$$\max(0, a)$$

To find ( $a_3, a_4$ )

$$a_3 = \max(0, 0.82) = 0.82$$

$$a_4 = \max(0, 1.27) = 1.27$$

— By Sigmoid Activation function  $f$  —

$$a_5 = (0.6942)(0.3) + (0.7807)(0.7) + 0.9 \\ = 1.6548$$

$$y_5 = \frac{1}{1+e^{-a_5}} = 0.8395.$$

By ReLU activation function  $f$ -

$$a_5 = (a_3 \times w_{53}) + (a_4 \times w_{54}) + 0.9 \\ = (0.82 \times 0.3) + (1.27 \times 0.7) + 0.9 \\ = 2.035$$

$$y_5 = \max(0, a_5) = (0, 2.035) = 2.035.$$

## Steepest Ascent / Descent - or - Gradient Ascent / Descent

Gradient ascent/descent is an OT algorithm used in ML to minimize or maximize the cost function by iteratively adjusting parameters in the direction of the negative gradient, aiming to find the optimal set of parameters.

The cost function represents the discrepancy b/w the predicted o/p of the model and the actual o/p.

Gradient ascent/descent aims to find the parameters that minimize/maximize the discrepancy and improve the model's performance.

### Working rule

In this method, we start from an initial trial point  $x_0$  and iteratively move along the steepest descent direction until the optimum point is found.

Step 1:- Start with an arbitrary point  $x_0$

Step 2:- Determine the optimal step length  $\lambda_k$  [minimize  $f(x + \lambda_k s)$ ,  
 $s_0 = -\nabla f$ ]

In the direction of  $-\nabla f(x_k)$  and  
 set  $x_{k+1} = x_k - \lambda_k \nabla f(x_k)$  ;  
 $k = 0, 1, 2, 3, \dots$ .

Step 5f Test the new point  $x_{k+1}$  for optimality. If  $x_{k+1}$  is optimum, we stop the procedure. Else go to the next step.

Step 4f Set  $k = k+1$  and go to Step 2.

6/6/24

Ex.  $f(x_1, x_2) = x_1 - x_2 + 2x_1^2 + 2x_1x_2 + x_2^2$   
 Starting from the point  $x_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

$$\Rightarrow f(x_1, x_2) = x_1 - x_2 + 2x_1^2 + 2x_1x_2 + x_2^2$$

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 1 + 4x_1 + 2x_2 \\ -1 + 2x_1 + 2x_2 \end{bmatrix}, (1)$$

$$\nabla f(x_0) = \nabla f \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \neq \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

To find optimal length (Iteration 1)

$$f \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 1_0 \begin{bmatrix} 1 \\ -1 \end{bmatrix} = f \begin{bmatrix} -1_0 \\ 1_0 \end{bmatrix}$$

$$f[-\lambda_0, \lambda_0] = -\lambda_0 - \lambda_0 + 2\lambda_0^2 - 2\lambda_0\lambda_0 + \lambda_0^2$$

$$= -2\lambda_0 + \lambda_0^2$$

$$\frac{\delta f}{\delta \lambda_0} = 2\lambda_0 - 2$$

$$\frac{\delta f}{\delta \lambda_0} = 0 \Rightarrow 2\lambda_0 - 2 = 0$$

$$2\lambda_0 = 2$$

$$\lambda_0 = 1$$

$$x_4 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 2 \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \end{bmatrix} \neq 0$$

- After substituting in (2) —

$$\nabla f(x_4) = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$

- (iteration 2) —

$$\text{minimize } f[x_4 - \lambda (\nabla f(x_4))]$$

$$f \left[ \begin{bmatrix} -1 \\ 1 \end{bmatrix} - \lambda \begin{bmatrix} -1 \\ -1 \end{bmatrix} \right] = f \left[ \begin{bmatrix} -1+\lambda_1 \\ 1+\lambda_1 \end{bmatrix} \right]$$

$$f(-1+\lambda_1, 1+\lambda_1) = -1+\lambda_1, -(1+\lambda_1) + 2(-1+\lambda_1)^2$$

$$+ 2(-1+\lambda_1)(1+\lambda_1) + (1+\lambda_1)^2$$

$$= 5\lambda_1^2 - 2\lambda_1 - 1$$

$$\frac{\delta f}{\delta \lambda_1} = 10\lambda_1 - 2 = 0 \Rightarrow \lambda_1 = 1/5$$

$$\chi_2 = (\chi_1 - \lambda_1 \nabla f(\chi_1))$$

$$= \left( \begin{bmatrix} -1 \\ 1 \end{bmatrix} - 1/5 \begin{bmatrix} -1 \\ -1 \end{bmatrix} \right)$$

$$= \begin{bmatrix} -0.8 \\ 1.2 \end{bmatrix} \neq \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$\chi_2$  isn't an optimal point

$$\nabla f(\chi_2) = \begin{pmatrix} 0.2 \\ -0.2 \end{pmatrix}$$

(iteration 3) —

$$\text{minimize } f \left[ \chi_2 + \lambda_2 (-\nabla f(\chi_2)) \right]$$

$$f \left[ \begin{bmatrix} -0.8 \\ 1.2 \end{bmatrix} - \lambda_2 \begin{bmatrix} 0.2 \\ -0.2 \end{bmatrix} \right] \Rightarrow f \left[ \begin{bmatrix} -0.8 - 0.2\lambda_2 \\ 1.2 + 0.2\lambda_2 \end{bmatrix} \right]$$

$$f(-0.8 - 0.2\lambda_2, 1.2 + 0.2\lambda_2)$$

$$= (-0.8 - 0.2\lambda_2)^2 + (1.2 + 0.2\lambda_2)^2 + 2(-0.8 - 0.2\lambda_2)^2 + 2(-0.8 - 0.2\lambda_2)(1.2 + 0.2\lambda_2) + (1.2 + 0.2\lambda_2)^2$$

$$= -1.2 - 0.16\lambda_2 + 0.12\lambda_2^2$$

$$\nabla f(\bar{x}_2) = 0.08\lambda_2 - 0.08$$

$$= \frac{1}{2} = 1$$

$$x_3 = \left[ x_2 - \lambda_2 \nabla f(x_2) \right]$$

$$\cancel{x_3} \left[ \begin{bmatrix} -0.8 \\ 1.2 \end{bmatrix} - 1 \begin{bmatrix} 0.2 \\ -0.2 \end{bmatrix} \right] = \begin{bmatrix} -1 \\ 1.4 \end{bmatrix} \neq \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$x_3$  is not an optimal point.

$\Rightarrow$

## Automatic Differentiation

This is a special case of back propagation.

consider  $f = x^2(x+y)$  where  $c = x^2$ ,  
 $s = x+y$

$$F \rightarrow (c, s) \rightarrow (x, y)$$

$$F \rightarrow (x, y)$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial c} * \frac{\partial c}{\partial x} + \frac{\partial f}{\partial s} * \frac{\partial s}{\partial x}$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial c} * \frac{\partial c}{\partial y} + \frac{\partial f}{\partial s} * \frac{\partial s}{\partial y}$$

$$\frac{\partial f}{\partial x} = 2[cs] * \frac{\partial (x^2)}{\partial x} + 2(cs) * \frac{\partial (x+y)}{\partial x}$$

$$\frac{\partial f}{\partial x} = (s * 2x) + (c + 1)$$

$$\boxed{\frac{\partial f}{\partial x} = 2cx + c}$$

$$\frac{\partial f}{\partial y} = \frac{\partial}{\partial c} [cs] + \frac{\partial}{\partial y} [x^2] + \frac{\partial}{\partial s} [cs] + \frac{\partial}{\partial y} [x+y]$$

$$\frac{\partial f}{\partial y} = s \cdot x(0) + c + 1$$

$$\boxed{\frac{\partial f}{\partial y} = c}$$

$\Rightarrow$  Gradients in a deep network

The Function value  $y$  is computed  
as a many level function composition

$$\text{i.e } Y = (f_k \circ f_{k-1} \circ \dots \circ f_1)(x)$$

$$f \circ g(x) = f(g(x))$$

$$\text{Ex: } f(x) = x^2$$

$$g(x) = x+1$$

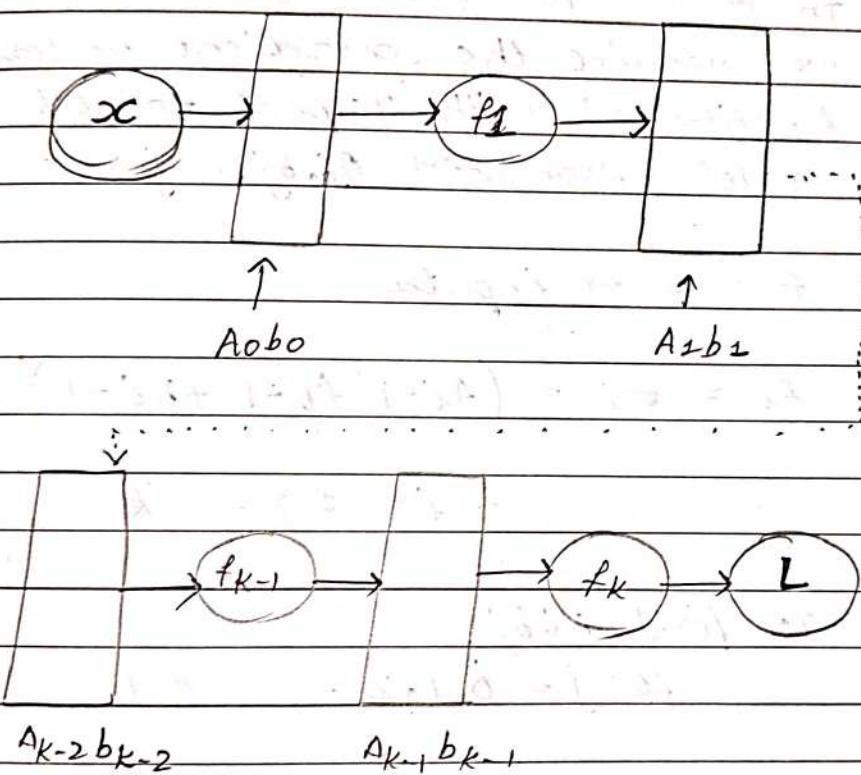
$$\begin{aligned} f \circ g &= f(g(x)) \\ &= f(x+1) \\ &= f(x+1)^2 \\ &= (x+1)^2 \end{aligned}$$

$$f_k(f_{k-1}(\dots \circ f_1(x)) \dots))$$

Here  $x$  are inputs or images.  $Y$  are observations

$f_i$ ;  $i = 1 \dots k$ , possesses its own parameters.

Forward pass is a multi-layer neural network to compute the pass ' $L$ ' as a function of the inputs  $x$  and the parameters  $A_i, b_i$



$$\text{Here } f_i(x_{i-1}) = \sigma(A_{i-1}x_{i-1} + b_{i-1})$$

$f_i \rightarrow$  function in the  $i$ th layer  
 $x_{i-1} \rightarrow$  output of layer  $i-1$

$\sigma \rightarrow$  Activation function

$$\text{logistic sigmoid} = \frac{1}{1 + e^{-\sigma}}$$

$$\tanh x = \frac{1 + e^{-x}}{1 + e^{+x}}$$

In order to train the models we require the gradient of loss function ' $L$ ' with respect to all model parameters  $A_j, b_j, j = 1 \dots k$

$$f_0 = x \rightarrow \text{inputs}$$

$$f_i = \sigma_i = (A_{i-1} f_{i-1} + b_{i-1})$$

$$= i = 1, 2, \dots, k$$

To find  $A_i, b_j$

$$\text{for } j = 0, 1, 2, \dots, k-1$$

we have to minimize the loss function

$$L(\theta) = \|y - f_k(\theta, x)\|^2$$

$$\theta = \{A_0, b_0, \dots, A_{k-1}, b_{k-1}\}$$

By definition of maxima or minima we require the partial derivative of 'L' with respect to the parameter  $\theta_j = (A_j, b_j) \Rightarrow j = 0, 1, \dots, k-1$

By chain rule we have

$$\frac{\partial L}{\partial \theta_{k-1}} = \frac{\partial L}{\partial f_k} \cdot \frac{\partial f_k}{\partial \theta_{k-1}}$$

$$\frac{\partial L}{\partial \theta_{k-2}} = \frac{\partial L}{\partial f_k} \cdot \frac{\partial f_k}{\partial f_{k-1}} \cdot \frac{\partial f_{k-1}}{\partial \theta_{k-2}} \rightarrow \begin{matrix} \text{Partial} \\ \text{derivative} \\ \downarrow \\ \text{w.r.t} \\ \text{inputs} \end{matrix}$$

w.r.t  
parameters

Similarly

$$\frac{\partial L}{\partial \theta_{k-3}} = \frac{\partial L}{\partial f_k} \cdot \frac{\partial f_k}{\partial f_{k-1}} \cdot \frac{\partial f_{k-1}}{\partial \theta_{k-2}} \cdot \frac{\partial f_{k-2}}{\partial \theta_{k-3}}$$

$$\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial f_k} \cdot \frac{\partial f_k}{\partial f_{k-1}} \cdots \frac{\partial f_{k-i}}{\partial f_{k-1}} \cdot \frac{\partial f_{k-1}}{\partial \theta_i} \cdots \frac{\partial f_{k-i+1}}{\partial \theta_i}$$

→ Optimization using second derivative test

Procedure

1. A necessary condition for  $x_0$  to be an extreme point (maximum or minimum) of  $f(x)$  is  $\nabla f(x_0) = 0$  where  $x \in \mathbb{R}^n$  and  $x_0 \in \mathbb{R}^n$  is called critical point or stationary point.
2. A sufficient condition for a stationary point  $x_0$  of  $f(x)$  to be an extreme (maximum or minimum) is that

case - 1:  $f$  has relative minimum if

$$f_{xx}(x_0) f_{yy}(x_0) - f_{xy}^2(x_0) > 0$$

and

$$f_{xx}(x_0) > 0$$

case - 2:  $f$  has relative maximum if

$$f_{xx}(x_0) f_{yy}(x_0) - f_{xy}^2(x_0) < 0$$

and

$$f_{xx}(x_0) < 0$$

case - 3:  $f$  has a saddle point

(neither maximum or minimum) if

$$f_{xx}(x_0) f_{yy}(x_0) - f_{xy}^2(x_0) < 0$$

case - 4: The test is inconclusive

$$\text{if } f_{xx}(x_0) f_{yy}(x_0) - f_{xy}^2(x_0) = 0$$

→ Hessian matrix form

A sufficient condition for a stationary point  $x_0$  of  $f(x)$  to be extremum (maximum or minimum)

case: 1 :  $f$  has relative minimum if  $f(h_0)$  is positive definite

case: 2 -  $f$  has relative maximum if  $f(h_0)$  is negative definite

case: 3 -  $f$  has a saddle point if  $f(h_0)$  is indefinite

case: 4 - The test is inconclusive otherwise

## Problems

1. Using the Hessian matrix to classify the relative extremum

$$\text{for } f(x_1, x_2) = \frac{1}{3}x_1^3 + x_1x_2^2 - 8x_1x_2 + 73$$

solution

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} x_1^2 + x_2^2 - 8x_2 \\ 2x_2x_1 - 8x_1 \end{bmatrix}$$

Equate to zero.

$$x_1^2 + x_2^2 - 8x_2 = 0 \quad \text{---(1)}$$

$$2x_2x_1 - 8x_1 = 0$$

DATE :

$$2x_2 x_1 = 8x_1$$

$$2x_2 = \frac{8x_1}{x_1}$$

$$2x_2 = 8$$

$$x_2 = \frac{8}{2}$$

$$\boxed{x_2 = 4}$$

$$x_1^2 + x_2^2 - 8x_2 = 0$$

$$x_1^2 + (4)^2 - 8(4) = 0$$

$$x_1^2 + 16 - 32 = 0$$

$$x_1^2 - 16 = 0$$

$$x_1^2 = 16$$

$$x_1 = \pm 4$$

with  $x_2 = 4$

stationary point  $= (+4, 4), (-4, 4)$   
 $(0, 0), (0, 8)$

$$2x_1x_2 - 8x_1 = 0 \quad ;$$

$$2x_1(x_2 - 4) = 0$$

$$- 2x_1 = 0$$

$$x_2 - 4 = 0$$

$$x_1 = 0, x_2 = 4.$$

$$f_{x_1 x_1} = \frac{\partial^2 f}{\partial x_1^2} = \frac{\partial}{\partial x_1} (x_1^2 + x_2^2 - 8x_2)$$

$$= 2x_1$$

$$f_{x_2 x_2} = \frac{\partial^2 f}{\partial x_2^2} = \frac{\partial}{\partial x_2} (x_1^2 + x_2^2 - 8x_2)$$

$$= 2x_2$$

$$f_{x_1 x_2} = \frac{\partial^2 f}{\partial x_1 \partial x_2} = \frac{\partial}{\partial x_1} (x_1^2 + x_2^2 - 8x_2)$$

$$= 2x_2 - 8$$

$$f_{x_2 x_1} = \frac{\partial^2 f}{\partial x_2 \partial x_1} = \frac{\partial}{\partial x_2} (x_1^2 + x_2^2 - 8x_2)$$

$$= 2x_1 - 8$$

$$H = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix}$$

$$H = \begin{bmatrix} 2x_1 & 2x_2 - 8 \\ 2x_2 - 8 & 2x_1 \end{bmatrix}$$

$H$  at  $(0, 0)$

$$H = \begin{bmatrix} 0 & -8 \\ -8 & 0 \end{bmatrix} = 0 - (-64) = -64$$

solution  
saddle point

$H$  at  $(0, 8)$

$$H = \begin{bmatrix} 0 & 8 \\ 8 & 0 \end{bmatrix} \quad 0 - (64) = -64$$

indefinite  
saddle point

$H$  at  $(4, 4)$

$$H = \begin{bmatrix} 8 & 0 \\ 0 & 8 \end{bmatrix} \quad 64 - 0 = 64$$

minimum point

$H$  at  $(4, 4)$

$$H = \begin{bmatrix} -8 & 0 \\ 0 & -8 \end{bmatrix} \quad 64 - 0 = 64$$

maximum point

2.

Using the Hessian matrix to classify the relative extremum

$$f(x_1, x_2, x_3) = x_1 + 2x_3 + x_2 x_3 - x_1^2 \\ \vdots \quad \quad \quad - x_2^2 - x_3^2$$

definition

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \frac{\partial f}{\partial x_3} \end{bmatrix} = \begin{bmatrix} 1 - 2x_1 \\ x_3 - 2x_2 \\ 2 + x_2 - 2x_3 \end{bmatrix}$$

$$\nabla f = 0$$

$$1 - 2x_1 = 0 \quad \text{--- (1)}$$

$$x_3 - 2x_2 = 0 \quad \text{--- (2)}$$

$$2 + x_2 - 2x_3 = 0 \quad \text{--- (3)}$$

$$1 - 2x_1 = 0$$

$$1 = 2x_1$$

$$\boxed{\frac{1}{2} = x_1}$$

$$\boxed{x_3 = 2x_2}$$

$$\frac{x_3}{x_2} = 2$$

$$2 + x_2 - 2(2x_2) = 0$$

$$2 + x_2 - 4x_2 = 0$$

$$2 - 3x_2 = 0$$

$$2 = 3x_2$$

$$\boxed{\frac{2}{3} = x_2}$$

$$x_3 - 2x_2 = 0$$

$$x_3 - 2\left(\frac{2}{3}\right) = 0$$

$$x_3 - \frac{4}{3} = 0$$

$$\boxed{x_3 = \frac{3}{4}}$$

$$(x_1, x_2, x_3) = (1/2, 2/3, 3/4)$$

$$f_{x_1 x_1} = \frac{\partial^2 f}{\partial x_1^2} = 2(1 - 2x_1) = -2$$

$$f_{x_2 x_2} = \frac{\partial^2 f}{\partial x_2^2} = \frac{2}{\partial x_2} (x_3 - 2x_2) = -2$$

$$f_{x_3 x_3} = \frac{\partial^2 f}{\partial x_3^2} = \frac{2}{\partial x_3} (2 + x_2 - 2x_3) = -2$$

$$f_{x_1 x_2} = f_{x_2 x_1} = \frac{\partial^2 f}{\partial x_1 \partial x_2} = \frac{2}{\partial x_1 \partial x_2} (x_3 - 2x_2) = 0$$

$$f_{x_1 x_3} = f_{x_3 x_1} = \frac{\partial^2 f}{\partial x_1 \partial x_3} = \frac{2}{\partial x_1 \partial x_3} (2 + x_2 - 3x_3) = 1$$

$$f_{x_2 x_3} = f_{x_3 x_2} = \frac{\partial^2 f}{\partial x_2 \partial x_3} = \frac{\partial}{\partial x_2} \left( 2 + x_2 \cdot 2x_3 \right) = 1.$$

$$H(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_1 \partial x_3} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \frac{\partial^2 f}{\partial x_2 \partial x_3} \\ \frac{\partial^2 f}{\partial x_3 \partial x_1} & \frac{\partial^2 f}{\partial x_3 \partial x_2} & \frac{\partial^2 f}{\partial x_3^2} \end{bmatrix}$$

$$H(x) = \begin{bmatrix} -2 & 0 & 0 \\ 0 & -2 & 1 \\ 0 & 1 & -2 \end{bmatrix}$$

~~$$H(x_2) = -2(x_2 - 1) + 0 + 0$$~~

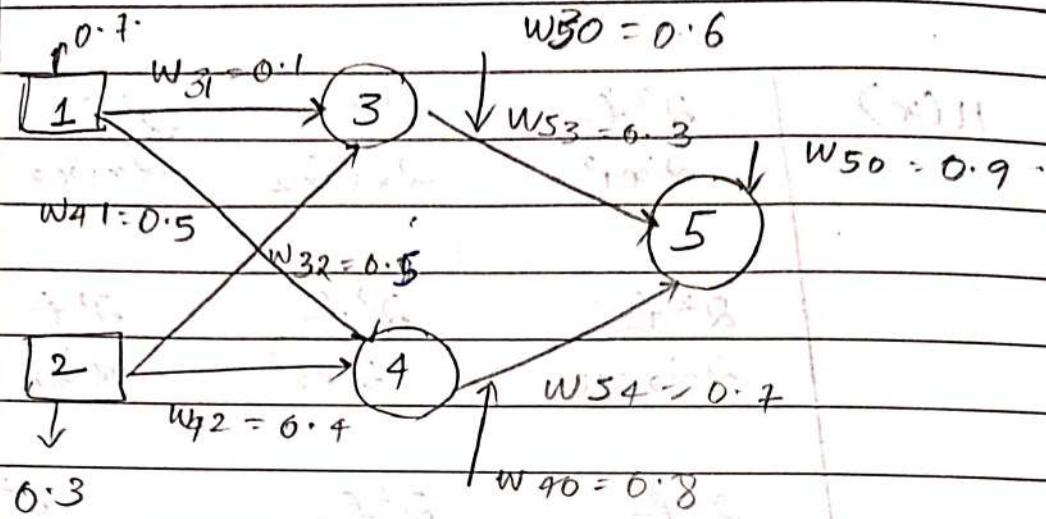
~~$$\begin{aligned} H(x) &= -2(3) \\ &= -6 \end{aligned}$$~~

$$H(x) = -2$$

$$H(x) = \begin{bmatrix} -2 & 0 \\ 0 & -2 \end{bmatrix} = -4 > 0$$

$$H(x) = \begin{bmatrix} -2 & 0 & 0 \\ 0 & -2 & 1 \\ 0 & 1 & -2 \end{bmatrix} = -6 > 0$$

Find output at neuron 5 if input vector  $[0.7, 0.3]$  using the activation function ReLU



$$\text{logistic sigmoid function} = \frac{e^{-x}}{1 + e^{-x}}$$

$$\begin{aligned} a_3 &= (w_{31} * 0.7) + (w_{32} * 0.3) + w_{30} \\ &= (0.1 * 0.7) + (0.5 * 0.3) + 0.6 \\ &= 0.07 + 0.15 + 0.6 \\ &= 0.82 \end{aligned}$$

$$\begin{aligned} y_3 &= F(a_3) \\ &= \frac{1}{1 + e^{-0.82}} \end{aligned}$$

$$y_3 = 0.6942$$

$$a_4 = (w_{41} * 0.7) + (w_{42} * 0.3) + w_{40}$$

$$a_4 = (0.5 * 0.7) + (0.4 * 0.3) + 0.8$$

$$a_4 = 0.35 + 0.12 + 0.8$$

$$a_4 = 1.27$$

$$y_4 = f(a_4)$$

$$= \frac{1}{1 + e^{-a_4}}$$

$$= \frac{1}{1 + e^{-1.27}}$$

$$y_4 = 0.781$$

### ReLU Activation Function

$$\max(0, a)$$

To find  $a_3$

$$\max(0, a_3)$$

$$= \max(0, 0.82)$$

$$a_3 = 0.82 = y_3$$

To find  $a_4$

$$\max(0, a_4)$$

$$= \max(0, 1.27)$$

$$a_4 = 1.27 = y_4$$

$$a_5 = (y_3 * w_{53}) + (y_4 * w_{54}) + w_{50}$$

$$= (0.6942 * 0.3) + ($$

$$= (0.82 * 0.3) + (1.27 * 0.7) + 0.9$$

$$= 0.2460 + 0.8890 + 0.9 = 2.0350$$

$$y_s = \max(0, a_5)$$

$$y_s = \max(0, 2.035)$$

$$y_s = 2.035$$

→ Steepest Ascent / Descent

(OR)

Gradient Ascent / Descent method

Gradient descent is an optimization algorithm used in machine learning to minimize or maximize the cost function by iteratively adjusting parameters in the direction of negative gradient aiming to find the optimal set of parameters.

The cost function represents the discrepancy between predicted output of the model and the actual output.

Gradient Descent or ascent aims to find the parameters that minimize or maximize this discrepancy and improve the model performance.

## Working Principle

In this method we start from a initial trial point  $x_0$  and iteratively move along the steepest descent direction until the optimum point is found.

1. start with an arbitrary point  $x_0$
2. determine the optimal step length  $\lambda_k$   
 [minimize or maximize]  $f(x + \lambda_k s_i, s_i = -\nabla f)$   
 in the direction of  $-\nabla f(z_k)$  and  
 set  $x_{k+1} = x_k - \lambda_k \nabla f(x_k)$ ;  $k = 0, 1, 2, \dots$
3. Test the new point  $x_{k+1}$  for optimouility. If  $x_{k+1}$  is optimum we stop the process otherwise go to the next step.
4. set  $k = k+1$  and go to step 2

Problems

1.

$$f(x_1, x_2) = x_1 - x_2 + 2x_1^2 + 2x_1x_2 + x_2^2$$

starting from the point  $x_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

solution:

$$f(x_1, x_2) = x_1 - x_2 + 2x_1^2 + 2x_1x_2 + x_2^2$$

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 1 + 4x_1 + 2x_2 \\ -1 + 2x_2 + 2x_1 \end{bmatrix}$$

$$\nabla f(x_0) = \nabla f \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} 1 + 4(0) + 2(0) \\ -1 + 2(0) + 2(0) \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ -1 \end{bmatrix} \neq \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

To find optimal length

$$\text{minimize } f[x_0 + \lambda_0(-\nabla f(x_0))]$$

$$f \left[ \begin{bmatrix} 0 \\ 0 \end{bmatrix} + -\lambda_0 \begin{bmatrix} 1 \\ -1 \end{bmatrix} \right]$$

$$f = \left[ \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \lambda_0 \begin{bmatrix} 1 \\ -1 \end{bmatrix} \right]$$

$$f = \begin{bmatrix} [0] & [x_0] \\ 0 & x_0 \end{bmatrix}$$

$$f = \begin{bmatrix} -x_0 \\ x_0 \end{bmatrix}$$

Substitution in  $f(x_1, x_2)$ .

$$\begin{aligned} f(-x_0, x_0) &= -x_0 - x_0 + 2x_0 \\ &\quad + 2x_0 x_0 + x_0^2 \\ f(-x_0, x_0) &= -2x_0 + x_0^2 \end{aligned}$$

$$\frac{\partial f}{\partial x_0} = 2x_0 - 2$$

$$\frac{\partial f}{\partial x_0} = 0 \quad \rightarrow \text{minimum}$$

$$2x_0 - 2 = 0$$

$$2x_0 = 2$$

$$x_0 = 1$$

$$x_1 = \begin{bmatrix} [0] & [-1] \\ 0 & 1 \end{bmatrix}$$

$$x_1 = \begin{bmatrix} [0] & [1] \\ 0 & 1 \end{bmatrix}$$

$$x_1 = \begin{bmatrix} -1 \\ 1 \end{bmatrix} \neq \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Substitute in ① and ②

$$\nabla f(x_1) = \begin{bmatrix} 1 + 4(-1) + 2(1) \\ -1 + 2(-1) + 2(1) \end{bmatrix}$$

$$\nabla f(x_1) = \begin{bmatrix} 1 - 4 + 2 \\ -1 - 2 + 2 \end{bmatrix}$$

$$\nabla f(x_1) = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$

Iteration 2

$$\text{minimize } f\left[x_1 - \lambda_1(\nabla f(x_1))\right]$$

$$f\left[\begin{bmatrix} -1 \\ 1 \end{bmatrix} - \lambda \begin{bmatrix} -1 \\ -1 \end{bmatrix}\right]$$

$$f\left[\begin{bmatrix} -1 \\ 1 \end{bmatrix} + \begin{bmatrix} \lambda_1 \\ \lambda_1 \end{bmatrix}\right]$$

$$f\left[\begin{bmatrix} -1 + \lambda_1 \\ 1 + \lambda_1 \end{bmatrix}\right]$$

$$f(-1 + \lambda_1, 1 + \lambda_1) = -1 + x_1 - 1 + x_1$$

$$\dots + 2(-1 + \lambda_1)^2 + 2(-1 + \lambda_1)$$

$$\dots (1 + \lambda_1) + (1 + \lambda_1)^2$$

$$= -2 + 2(\lambda_1^2 + 1 - 2\lambda_1)$$

$$+ 2(-1 + x_1(1 + \lambda_1))$$

$$+ 1 + \lambda_1^2 + 2\lambda_1$$

~~$$= -2 + 2(\lambda_1^2 + 1 - 2\lambda_1)$$~~

$$= -2 + 2\lambda_1^2 + 2 - 4\lambda_1 + 1 + \lambda_1^2 + 2\lambda_1$$

$$+ 2(-1 - \lambda_1 + \lambda_1 + \lambda_1^2)$$

$$= 3\lambda_1^2 - 2\lambda_1 - 2 - 2\lambda_1 + 2x_1 + 1$$

$$+ 2\lambda_1^2$$

$$= 5\lambda_1^2 - 2\lambda_1 - 1$$

$$\frac{\partial f}{\partial x_1} = 10\lambda_1 - 2$$

$$\frac{\partial f}{\partial \lambda_1} = 0$$

$$10\lambda_1 - 2 = 0$$

$$10\lambda_1 = 2$$

$$\lambda_1 = \frac{2}{10}$$

$$x_1 = \frac{1}{5}$$

$$\lambda_1 = 1/5$$

$$OC_2 = \left[ \begin{bmatrix} -1 \\ 1 \end{bmatrix} - \frac{1}{5} \begin{bmatrix} -1 \\ 1 \end{bmatrix} \right]$$

$$\nabla f(x_1 - \lambda x_2) = [1 + 4(-0.8) + 2(1.2)]$$

$$\nabla f(x_2) = \begin{bmatrix} [-1] \\ [1.2] \end{bmatrix} + \begin{bmatrix} [1] \\ [1.2] \end{bmatrix}$$

$$x_2 = \begin{bmatrix} [-0.8] \\ [1.2] \end{bmatrix} \neq \begin{bmatrix} [0] \\ [0] \end{bmatrix}$$

$x_2$  is not optimal

Substitute  $x_2$  value in ① and ②

$$\nabla f(x_2) = \begin{bmatrix} 1 + 4(-0.8) + 2(1.2) \\ -1 + 2(-0.8) + 2(1.2) \end{bmatrix}$$

$$\nabla f(x_2) = \begin{bmatrix} 1 - 3.2 + 2.4 \\ -1 - 1.6 + 2.4 \end{bmatrix}$$

$$\nabla f(x_2) = \begin{bmatrix} 0.2 \\ -0.2 \end{bmatrix}$$

Iteration 3.

$$\text{minimize } f[x_2 + \lambda_2 (-\nabla f(x_2))]$$

$$f \left[ \begin{bmatrix} [-0.8] \\ [1.2] \end{bmatrix} - \lambda_2 \begin{bmatrix} [0.2] \\ [-0.2] \end{bmatrix} \right]$$

$$f(-0.8 + \lambda_2 0.2, 1.2 + \lambda_2 0.2) =$$

$$f(-0.8 - \lambda_2 0.2, 1.2 + \lambda_2 0.2)$$

$$f(-0.8 - \lambda_2 0.2, 1.2 + \lambda_2 0.2)$$

Substitute in main equation.

$$= (-0.8 - \lambda_2 0.2) - (1.2 + \lambda_2 0.2) \\ + 2(-0.8 - \lambda_2 0.2)^2 \\ + 2(-0.8 - \lambda_2 0.2)(1.2 + \lambda_2 0.2) \\ + (1.2 + \lambda_2 0.2)^2$$

$$= (-0.8 - \lambda_2 0.2 - 1.2 - \lambda_2 0.2) \\ + 2(0.64 + 0.04\lambda_2^2 - 0.32\lambda_2) \\ + 2(-0.96 - 0.16\lambda_2 - 0.24\lambda_2 \\ - 0.2\lambda_2^2) \\ + (1.44 + 0.4\lambda_2 0.04\lambda_2^2 + 0.48\lambda_2)$$

~~$$= (-0.8 + \lambda_2 0.2 - 1.2 - \lambda_2 0.2 + 1.28) \\ + 0.08\lambda_2^2 - 0.64\lambda_2 - 1.92 - \\ 0.32\lambda_2 - 0.48\lambda_2 - 0.4\lambda_2^2 \\ + 1.44 + 0.04\lambda_2^2 + 0.48\lambda_2$$~~

~~$$= -1.2 - 1.36\lambda_2 - 0.28\lambda_2^2$$~~

~~$$2x_1 = -1.36 - 0.56\lambda_2 \\ -1.36 = 0.56\lambda_2$$~~

$$\begin{aligned} & -0.8 - 0.2\lambda_2 \div (1.2 + 0.2\lambda_2) \\ & - 2(0.8 + 0.2\lambda_2)^2 \\ & + 2(-0.8 - 0.2\lambda_2)(1.2 + \lambda_2 + 0.2) \\ & + (1.2 + 0.2\lambda_2)^2 \end{aligned}$$

$$\frac{\partial f}{\partial \lambda_2} = 0.08\lambda_2 - 0.08$$

$$\lambda_2 = 1$$

$$x_3 = (x_2 - \lambda_2 \nabla f(x_2))$$

$$x_3 = \left[ \begin{matrix} -0.8 \\ 1.2 \end{matrix} \right] - 1 \left[ \begin{matrix} 0.2 \\ -0.2 \end{matrix} \right]$$

=

$$\left[ \begin{matrix} -1 \\ 1.4 \end{matrix} \right] \neq \left[ \begin{matrix} 0 \\ 0 \end{matrix} \right]$$

$x_3$  is not optimal point.

2.  $f(x, y) = 4x^2 - 4xy + 2y^2$

$$x^{(0)}, y^{(0)} = (2, 3)$$

tion  $f(x, y) = 4x^2 - 4xy + 2y^2$

$$\nabla f = \left[ \begin{matrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{matrix} \right] = \left[ \begin{matrix} 8x - 4y \\ -4x + 4y \end{matrix} \right]$$

Substituting  $(x, y) = (2, 3)$

$$\nabla f(x_0) = \begin{bmatrix} (g \circ 2) - f(3) \\ -f(2) + f(3) \end{bmatrix} = \begin{bmatrix} 4 \\ 4 \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

(Iteration 1)

$$\text{minimize } f \cdot [x_0 + \lambda_0 (-\nabla f(x_0))]$$

$$= f \left[ \begin{bmatrix} 2 \\ 3 \end{bmatrix} - \lambda_0 \begin{bmatrix} 4 \\ 4 \end{bmatrix} \right]$$

$$= f \left[ \begin{bmatrix} 2 - 4\lambda_0 \\ 3 - 4\lambda_0 \end{bmatrix} \right]$$

$$= f(2 - 4\lambda_0, 3 - 4\lambda_0)$$

$$= 4(2 - 4\lambda_0)^2 - 4(2 - 4\lambda_0)(3 - 4\lambda_0) + 4(3 - 4\lambda_0)^2$$

$$= 4(4 + 16\lambda_0^2 - 16\lambda_0) - 4(6 - 8\lambda_0 - 12\lambda_0 + 16\lambda_0^2) + 2(9 + 16\lambda_0^2 - 24\lambda_0)$$

$$= 10 + 32\lambda_0^2 - 32\lambda_0$$

$$\frac{\partial f}{\partial \lambda_0} = 64\lambda_0 - 32$$

equate to 0

$$\frac{\partial f}{\partial \lambda_0} = 0$$

$$2\lambda_0$$

$$64\lambda_0 - 32 = 0$$

$$64\lambda_0 = 32$$

$$\lambda_0 = \frac{32}{64}$$

$$\lambda_0 = \frac{1}{2}$$

$$x_0 = 0.5 \text{ (initial)}$$

$$x_1 = x_0 + \lambda_1 (-\nabla f(x_0))$$

$$= \begin{bmatrix} 2 \\ 3 \end{bmatrix} - \lambda_1 \begin{bmatrix} 4 \\ 9 \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix} - \frac{1}{2} \begin{bmatrix} 4 \\ 9 \end{bmatrix}$$

$$= \begin{bmatrix} 2 \\ 3 \end{bmatrix} - \frac{1}{2} \begin{bmatrix} 4 \\ 9 \end{bmatrix}$$

$$= \begin{bmatrix} 2 \\ 3 \end{bmatrix} - \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

$$= \begin{bmatrix} 0 \\ 1 \end{bmatrix} \neq \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

$$\nabla f(x_1) = \begin{bmatrix} 8(0) - 4(1) \\ 0 - 4(0) + 4(1) \end{bmatrix} = \begin{bmatrix} 4 \\ 4 \end{bmatrix}$$

$$= \begin{bmatrix} -4 \\ 4 \end{bmatrix}$$

Iteration (2)

$$\text{minimize } f(x_1 - \lambda)$$

$$\text{minimize } f(x_1 + \lambda_1 (-\nabla f(x_1)))$$

$$= f \left[ \begin{bmatrix} 0 \\ 1 \end{bmatrix} - \lambda_1 \begin{bmatrix} -4 \\ 4 \end{bmatrix} \right]$$

$$= f \left[ \begin{bmatrix} 0 \\ 1 \end{bmatrix} + \begin{bmatrix} 4\lambda_1 \\ -4\lambda_1 \end{bmatrix} \right]$$

$$= \begin{bmatrix} 4\lambda_1 \\ 1 - 4\lambda_1 \end{bmatrix}$$

$$f(4\lambda_1, 1 - 4\lambda_1)$$

$$= 4(4\lambda_1)^2 - 4(4\lambda_1)(1 - 4\lambda_1) + 2(1 - 4\lambda_1)^2$$

$$= 4(16\lambda_1^2) - 4(4\lambda_1 - 16\lambda_1^2) + 2(1 + 16\lambda_1^2 - 8\lambda_1)$$

$$= 64\lambda_1^2 - 16\lambda_1 + 64\lambda_1^2 + 2 + 32\lambda_1^2 - 16\lambda_1$$

$$= 160\lambda_1^2 - 32\lambda_1 + 2$$

$$\frac{\partial f}{\partial \lambda_1} = 320\lambda_1 - 32$$

$\lambda_1$

equate it to zero.

$$\frac{\partial f}{\partial \lambda_1} = 0$$

$$320\lambda_1 - 32 = 0$$

$$320\lambda_1 = 32$$

$$\lambda_1 = \frac{32}{320}$$

$$\boxed{\lambda_1 = 1/10}$$

$$x_2 = x_1 - \gamma_1 \nabla f(x_1)$$

$$x_2 = \left[ \begin{bmatrix} 0 \\ 1 \end{bmatrix} - \frac{1}{10} \begin{bmatrix} -4 \\ 9 \end{bmatrix} \right]$$

$$x_2 = \begin{bmatrix} 2/5 \\ 3/5 \end{bmatrix} + \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

$x_2$  is not optimal

$$f(x_2) = 2/5 + 3/5 + 2/5 \times (-4) + 3/5 \times 9 = 16.8$$

# Problems On Automatic Differentiation

1. draw a computational graph of the function  $f(x) = \sqrt{x^2 + e^{x^2} + \cos(x^2 + e^{x^2})}$   
 also find  $\frac{df}{dx}$  by using automatic differentiation.

solution Let  $x^2 = a$        $e^{x^2} = b$   
 $e^{x^2} = e^a = b$

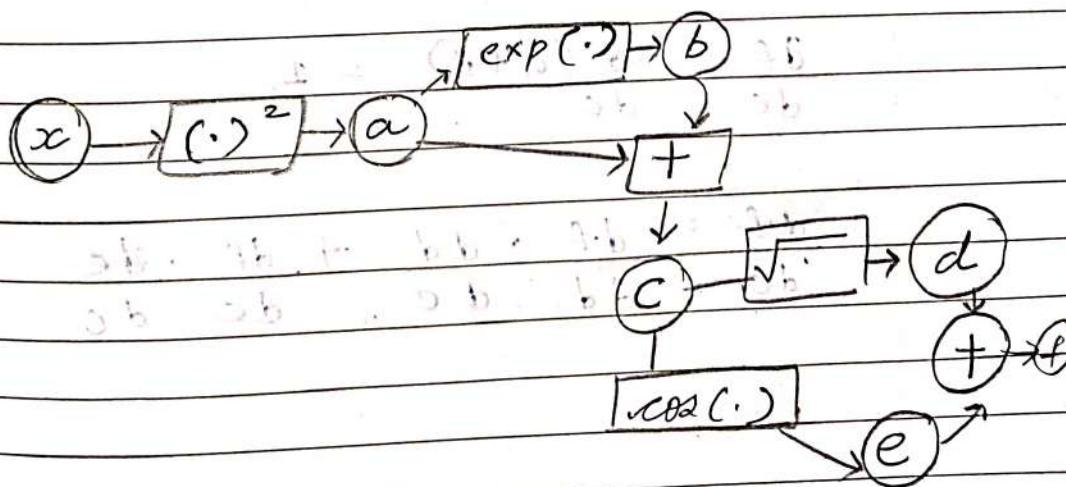
$$c = a + b$$

$$d = \sqrt{c}$$

$$e = \cos(c)$$

$$f(x) = d + e$$

computational Graph



To find  $\frac{df}{dx}$

$$\frac{da}{dx} = \frac{d(x^2)}{dx} = 2x$$

$$\frac{db}{da} = \frac{d(e^a)}{da} = e^a$$

$$\frac{dc}{da} = \frac{d(a+b)}{da} = 1$$

$$\frac{dc}{db} = \frac{d(a+b)}{db} = 1$$

$$\frac{dd}{dc} = \frac{d\sqrt{c}}{dc} = \frac{1}{2\sqrt{c}}$$

$$\frac{de}{dc} = \frac{d(\cos c)}{dc} = -\sin(c)$$

$$\frac{df}{dd} = \frac{d(d+e)}{dd} = 1$$

$$\frac{df}{de} = \frac{d(d+e)}{de} = 1$$

$$\frac{df}{dc} = \frac{df}{dd} \cdot \frac{dd}{dc} + \frac{df}{de} \cdot \frac{de}{dc}$$

$$\frac{df}{dc} = \frac{1}{2\sqrt{c}} + \frac{1}{2} \sin(c)$$

$$\frac{df}{db} = \frac{df}{dc} \cdot \frac{dc}{db} =$$

$$\frac{df}{db} = \frac{df}{dc} \cdot \times$$

$$\frac{df}{da} = \frac{df}{db} \cdot \frac{db}{da} + \frac{df}{dc} \cdot \frac{dc}{da}$$

$$= \frac{df}{db} \cdot e^a + \frac{df}{dc} \cdot f'(1)$$

$$= \frac{df}{dc} [e^a + 1] \cdot \times$$

$$\frac{df}{dx} = \frac{df}{da} \cdot \frac{da}{dx}$$

$$= \frac{df}{da} (2x)$$

$$= \frac{df}{dc} (e^a + 1) (2x)$$

$$= \left[ \frac{1}{2\sqrt{c}} + \int \sin(c) \right] [e^a + 1] [2x]$$

(2)

$$\text{Let } f(x_1, x_2) = -\log x_1 + x_1 x_2 - \sin x_2$$

1. Draw a computational graph of  $f(x_1, x_2)$

2. Evaluate  $f$  at  $x_1, x_2 = (2, 5)$  by Forward trace

Solution

$$v_0 = 2$$

$$v_1 = \log v_0 = \log 2 = 0.693$$

$$v_2 = v_1 * x_1 = 0.693 * 2 = 1.386$$

$$\begin{aligned} v_3 &= v_0 * v_2 \\ &= 2 * 1.386 = 2.772 \end{aligned}$$

$$v_4 = 10$$

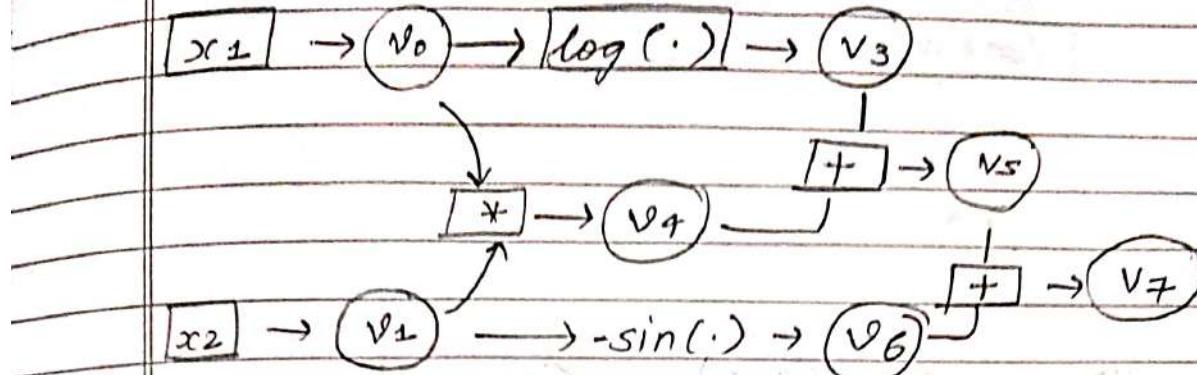
$$\rightarrow v_5 = v_3 + v_4 = 10.693$$

$$\begin{aligned} v_6 &= -\sin(v_1) \\ &= -\sin(0.693) \\ &= 0.9589 \end{aligned}$$

$$v_7 = v_5 + v_6$$

$$v_7 = 10.693 + 0.9589$$

$$v_7 = 11.643$$



3. If  $f = x^2 (x+y)$  where  $C = x^2$ ;  $S = x+y$

(i) draw a computational graph.

(ii)  $\frac{df}{dx}$ ,  $\frac{df}{dy}$  using chain rule

(iii) draw a computational graph in forward mode to show the results of

solution

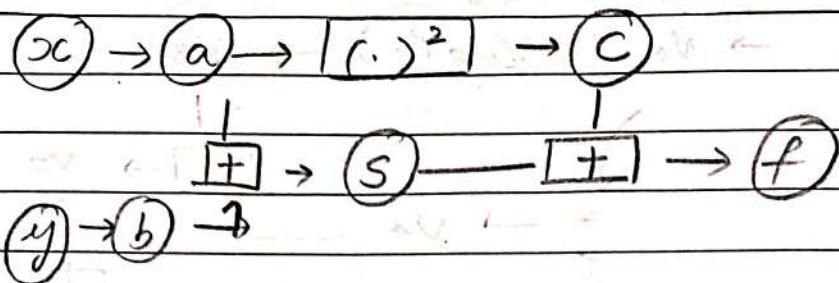
$$x = a$$

$$y = b$$

$$C = x^2 = a^2$$

$$S = a+b = x+y$$

$$f = C(S)$$



(ii)  $f \rightarrow (cs) \rightarrow (x \ y)$

$$f \rightarrow (x \ y)$$

$$\frac{df}{dx} = \frac{df}{dc} \cdot \frac{dc}{dx} + \frac{df}{ds} \cdot \frac{ds}{dx}$$

$$\frac{df}{dx} = \frac{d(cs)}{dc} \cdot \frac{d(x^2)}{dx} + \frac{d(cs)}{ds} \cdot \frac{d(x+y)}{dx}$$

$$\frac{df}{dx} = s(2x) + c(\cdot 1)$$

$$\frac{df}{dx} = s(2x) + c$$

$$\frac{df}{dx} = (x+y)2x + x^2$$

$$\frac{df}{dx} = 2x^2 + 2xy + x^2$$

$$\frac{df}{dx} = 3x^2 + 2xy$$

$$\frac{df}{dy} = \frac{df}{dc} \cdot \frac{dc}{dy} + \frac{df}{ds} \cdot \frac{ds}{dy}$$

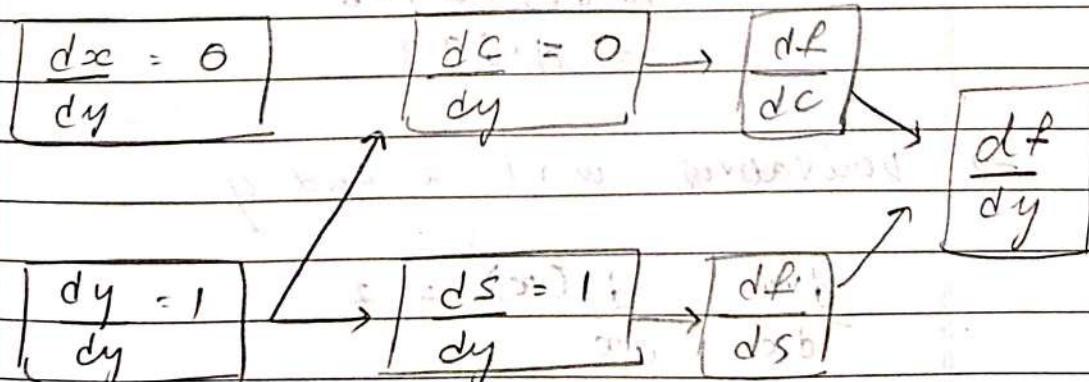
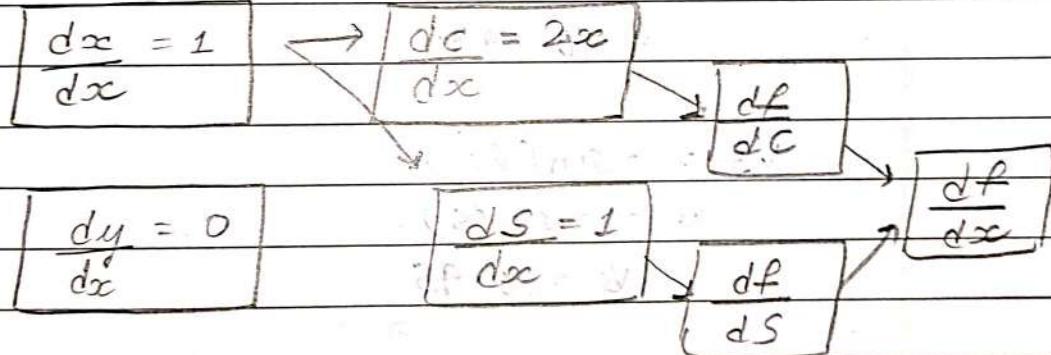
$$\frac{df}{dy} = \frac{d(cs)}{dc} \cdot \frac{d(x^2)}{dy} + \frac{d(cs)}{ds} \cdot \frac{d(x+y)}{dy}$$

$$\frac{df}{dy} = s(0) + c \cdot 1$$

$$\frac{df}{dy} = 1 \cdot c = x^2.$$

iii Computational graph to show results

a) Graph of  $x$  and  $y$



4. If  $f(x, y) = \log x + xy - \sin y$  compute the forward trace derivative of  $f$  w.r.t  $x$ , at  $(2, 5)$

Solution

$$x = v_0 = 2$$

$$y = v_1 = 5$$

$$v_2 = \log(v_0)$$

$$= \log(2)$$

$$v_2 = 0.6931$$

$$v_3 = v_0 * v_1$$

$$v_3 = 2 * 5$$

$$v_3 = 10$$

$$v_4 = v_2 + v_3$$

$$= 0.6931 + 10$$

$$= 10.693$$

$$v_5 = -\sin(v_2)$$

$$v_5 = -\sin(5)$$

$$v_5 = 0.95$$

$$v_6 = v_4 + v_5$$

$$= 10.693 + 0.958$$

$$= 11.643$$

→ Derivatives w.r.t  $x$  and  $y$

$$\frac{d v_0}{d x} = \frac{d}{d x}(x) = 1$$

$$\frac{dv_2}{dx} = \frac{d(\log v_0)}{dx} = 0$$

$$\frac{dv_2}{dx} = \frac{d(\log v_0)}{dx} = \frac{d(\log x)}{dx} = \frac{1}{x} = \frac{1}{2} = 0.5$$

$$\frac{dv_3}{dx} = \frac{d(v_0 v_1)}{dx}$$

$$= v_0 \cdot \frac{d(v_1)}{dx} + v_1 \cdot \frac{d(v_0)}{dx}$$

$$= v_0(0) + v_1(1)$$

$$= 2(0) + 5(1)$$

$$= 5$$

$$\frac{dv_4}{dx} = \frac{d(v_2 + v_3)}{dx}$$

$$= \frac{dv_2}{dx} + \frac{dv_3}{dx}$$

$$= 0.5 + 5$$

$$= 5.5$$

$$\frac{dv_5}{dx} = \frac{d(C \cdot \sin(v_1))}{dx}$$

$$= -\cos(v_1) \cdot \frac{dv_1}{dx}$$

$$= -\cos(v_1) \cdot 0$$

$$= 0$$

$$\frac{dV_G}{dx} = \frac{d}{dx}(V_A + V_S) \quad (\text{Given})$$

$$= \frac{dV_A}{dx} + \frac{dV_S}{dx}$$

$$= 5.5 + 0.1 \times 5 \quad (\text{Given})$$

$$= 5.5$$

Derivative with respect to  $y$ .

$$\frac{dV_0}{dy} = \frac{d(x)}{dy} = 0$$

$$\frac{dV_1}{dy} = \frac{d(y)}{dy} = 1$$

$$\frac{dV_2}{dy} = \frac{d(\log V_0)}{dy}$$

$$= \frac{1}{\log(V_0)} \cdot \frac{dV_0}{dy}$$

$$= \frac{1}{\log 2} \cdot 0 \quad (\text{Given})$$

$$\frac{dV_2}{dy} = 0$$

$$\frac{dv_3}{dy} = \frac{d}{dy}(v_0 \cdot v_1)$$

$$= v_0 \cdot \frac{dv_1}{dy} + v_1 \cdot \frac{dv_0}{dy}$$

$$= v_0 \cdot (1) + v_1 (0)$$

$$= 2(1) + 5(0)$$

$$= 2$$

$$\frac{dv_4}{dy} = \frac{d}{dy}(v_2 + v_3)$$

$$= \frac{dv_2}{dy} + \frac{dv_3}{dy}$$

$$= 0 + 2$$

$$= 2$$

$$\frac{dv_5}{dy} = \frac{d}{dy}(-\sin(v_1))$$

$$= -\cos(v_1) \cdot \frac{d(v_1)}{dy}$$

$$= -\cos(s) \cdot 1$$

$$= -0.2837$$

$$\frac{dv_6}{dy} = \frac{d}{dy}(v_4 + v_5)$$

$$= \frac{dv_4}{dy} + \frac{dv_5}{dy} = 2 - 0.2837 \\ = 1.7163$$

## Gradients of quadratic cost

obtain the quadratic cost of function

$$C = (\hat{y} - y)^2$$

$$C = (\hat{y} - y)^2$$

$$u = \hat{y} - y$$

$$C = u^2$$

$\hat{y} \rightarrow$  Estimated value

$y \rightarrow$  True value

$$\frac{dc}{du} = 2u = 2(\hat{y} - y)$$

$$\frac{dc}{du} = 2(\hat{y} - y) = \frac{\partial h}{\partial b} - \frac{\partial h}{\partial b} = \frac{\partial h}{\partial b}$$

$$u = \hat{y} - y$$

$$\frac{dy}{d\hat{y}} = 1$$

$$\frac{dy}{d\hat{y}}$$

$$\frac{dc}{d\hat{y}} = \frac{dc}{du} \cdot \frac{du}{d\hat{y}}$$

$$= 2(\hat{y} - y) \cdot 1$$

Estimated  $y$  w.r.t model parameters

$$\hat{y} = mx + b$$

$$\text{d}y = q(mx + b)$$

$$= x$$

$$\frac{dy}{db} = \frac{d(mx+b)}{db}$$

$$= 1$$

Quadratic cost w.r.t model parameters

$$\frac{dc}{db} = \frac{dc}{dy} \cdot \frac{dy}{db}$$

$$\frac{dc}{dm} = \frac{dc}{dy} \cdot \frac{dy}{dm}$$

$$\nabla c = \begin{bmatrix} \frac{dc}{db} \\ \frac{dc}{dm} \end{bmatrix} = \begin{bmatrix} 2(y - \hat{y}) \\ 2(\hat{y} - y)x \end{bmatrix}$$

Q1.

Obtain the gradient of mean squared error [average quadratic cost]

solution

mean Estimated Error

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

$$u = \hat{y}_i - y_i$$

$$C = \frac{1}{n} \sum u^2$$

$\hat{y}_i \rightarrow$  Estimated value

$y_i \rightarrow$  True value

$$\frac{dc}{du} = \frac{1}{n} \sum 2u$$

$$= \frac{1}{n} \sum 2(\hat{y}_i - y_i)$$

$$u = \hat{y}_i - y_i$$



## Fibonacci Sequence

1 1 2 3 5 8 13 21 34 55  
 $F_0 F_1 F_2 F_3 F_4 F_5 F_6 F_7 F_8 F_9$

$$F_0 = 1 = F_1, \quad F_n = F_{n-1} + F_{n-2}, \quad n \geq 2$$



## Fibonacci Search method

### Algorithm / method

Step 1. Determine the number of iterations

Step 2. Define  $L_2^* = \frac{F_{n-2}}{F_n} * l_0$

where  $l_0 = b - a$

$$\Rightarrow L_2^* = \frac{F_{n-2}}{F_n} (b - a) \rightarrow \begin{array}{l} \text{length of} \\ \text{interval} \\ \text{uncertainty} \end{array}$$

choose the two initial points  $x_1$  and  $x_2$  be such that  $x_1 = a + L_2^*$

$$\Rightarrow a + \frac{F_{n-2}}{F_n} * l_0$$

$$x_2 = a - L_2^* = b - \frac{F_{n-2}}{F_n} * l_0$$

Step 3 compute  $F(x_1)$  and  $F(x_2)$  and apply unimodal principle

process is repeated again and again

Unimodal Principle

set initial interval  $[a, b]$  be given  
and the number of iterations be  $n$

1. set  $L = a$ ,  $R = b$ ,  $k = 1$  (iteration number)

$\downarrow$   
left value

2. compute  $x_1, x_2$  and  $f(x_1)$  and

$f(x_2)$

If  $f(x_2) > f(x_1)$  then  $L_{k+1} = L_k$

$R_{k+1} = x_2$

If  $x_1$  is minimum preserve  $x_1$   
and take left and right value

$(L = a, R = x_2)$

If  $f(x_2) < f(x_1)$  then  $R_{k+1} = x_1$

$R_{k+1} = x_2$

If  $x_2$  is minimum preserve  $x_2$   
and take left and right value as  
 $L = x_2, R = b$ .

## Problems

PAGE NO.:

DATE:

1.

minimize  $f(x) = x^2$ , over  $[-5, 15]$  using Fibonacci search by taking  $n = 7$ .

Solution

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

$$a = -5, b = 15$$

$$L_0 = b - a$$

$$L_0 = 15 - (-5)$$

$$L = a = -5, R = b = 15 \Rightarrow L_0 = 20$$

$$L_2^* = \frac{F_{n-2}}{F_n} L_0 = \frac{F(5)_{(10)}}{F_7} = \frac{8}{21} (20) \\ \Rightarrow L_2^* = 7.619$$

$$x_1 = a + L_2^*$$

$$x_2 = b - L_2^*$$

L	R	$x_1$	$x_2$	$f(x_1)$	$f(x_2)$	$L \text{ or } R$
-5	15	2.619	7.3810	6.8592	59.47	L
		(L+R-x_2)	(L+R-x_1)			
-5	7.3810	-0.2380	2.619	0.0566	6.8592	L
-5	2.619	-2.1429	-0.2381	4.5920	0.0567	R
-2.1429	2.619	-0.2381	0.7142	0.0567	0.5101	L
-2.1429	0.7142	-1.1906	-0.2381	1.4175	0.0567	R
-1.1906	0.7142	-0.2381	-0.2383	0.0567	0.0568	R
-0.2381	0.7142	-0.2383	0.7144	0.0568	0.5104	
$L = -0.2381, R = 0.7142$						

Preserve the minimum out of  $x_1$  and  $x_2$ . L if  $f(x)$  is minimum  
R if  $f(x)$  is minimum.

$$x^* = -0.2381 + 0.7192$$

$$x^* = 0.2381$$

$$f(x^*) = 0.0567$$

solution

Let  $L$  = fixed  $x_2$   
 $x_2$  is  $R$

PAGE NO.:

DATE:

2. minimize  $f(x) = x^2 + 54$ ,  $[0, 5]$  using

Fibonacci Search by taking  $n = 3$

*solution*  $f(x) = x^2 + 54$

$$a = 0, b = 5, n = 3$$

$$L = a = 0$$

$$R = b = 5$$

$$L_0 = b - a$$

$$L_0 = 5 - 0$$

$$L_0 = 5$$

$$L_2^* = \frac{F_{n-2}(L_0)}{F_n} = \frac{F_1(5)}{F_3} = \frac{1}{3}(5)$$

$$L_2^* = 2.5 \quad 1.6667$$

$$x_1 = a + L_2^*$$

$$x_2 = b - L_2^*$$

$L$	$R$	$x_1$	$x_2$	$f(x_1)$	$f(x_2)$	$L^{(0)}$ $R$
0	5	2.567	3.333	35.1724	27.3105	$R$
1.667	5	3.333	3.334	27.3105	27.3123	$L$
1.667	3.334	1.668	3.333	35.1563	27.3105	$L = 1.667, b = 3.334$ $x^* = 2.5605 f(x^*) = 27.8481$

4. Minimize  $f(x) = x(x-1.5)$  in  $[0, 1]$  within the interval of uncertainty 0.25 of the initial interval of uncertainty

Solution

$$\Delta_n \leq 0.25 \Delta_0$$

$$\frac{\Delta_n}{\Delta_0} \leq 0.25$$

$$\frac{1}{F_n} \leq 0.25$$

$$F_n \geq 4$$

$$1 \ 1 \ 2 \ 3 \ 5$$

$$n=4$$

$$L = a = 0$$

$$R = b = 1$$

$$\Delta_0 = b - a$$

$$\Delta_0 = 1 - 0$$

$$\Delta_0 = 1$$

$$L_2 = \frac{F_{n-2}}{F_n} (\Delta_0)$$

$$= \frac{F_{n-2}}{F_n} (1)$$

$$= \frac{F_2}{F_n} (1)$$

$$= \frac{2}{5} (1)$$

$$L_2^* = 0.4$$

$$x_1 = a + L_2^*$$

$$x_2 = b - L_2^*$$

L	R	$x_1$	$x_2$	$f(x_1)$	$f(x_2)$	$L/R$
0	1	0.4	0.6	-0.94	-0.54	R
0.4	1	0.6	0.8	-0.54	-0.56	R
0.6	1	0.8	0.8	-0.56	-0.56	L
0.6	0.8					

$$x^* = \frac{0.6 + 0.8}{2}$$

$$x^* = 0.70$$

$$f(x^*) = 0.70(0.7 - 1.5)$$

$$= -0.56$$

5. minimize  $f(x) = \begin{cases} 1-x & , x < 1/2 \\ x^2 & , x \geq 1/2 \end{cases}$

over  $[-1, 1]$  perform six iterations.

solution  $L = a = -1, n = 6$

$$R = b = 1$$

$$\ell_0 = b - a$$

$$= 1 - (-1)$$

$$= 1 + 1$$

$$\ell_0 = 2$$

$$L_2^* = \frac{F_{n-2}}{F_n} (\ell_0)$$

$$= \frac{F_4}{F_6} (\ell_0)$$

$$= \frac{5}{8} (2)$$

$$= 1.3$$

$$(2 - 1.3) \cdot 1.3 = 0.91$$

$$L_2^* = 0.7692$$

$$x_1 = a + L_2^*$$

$$x_2 = b - L_2^*$$

L	R	$x_1$	$x_2$	$f(x_1)$ (1st func)	$f(x_2)$ (1st func)	$4/R$
-1	1	-0.2308	0.2308	0.6154	0.3846	R
-0.2308	1	0.2308	0.5384	0.3846	0.2898	R
0.2308	1	0.5384	0.6924	0.2898	0.1794	L
0.2308	0.6924	0.3848	0.5384	0.3076	0.2898	R
0.3848	0.6924	0.5384	0.5388	0.2898	0.2903	L
0.3848	0.5388					

$$x^* = \frac{0.3848 + 0.5388}{2}$$

$$\therefore x^* = 0.4618$$

$\left[ \begin{array}{l} \text{1st} \\ \text{func} \end{array} \right] f(x^*) = 0.2691$

$\Rightarrow$  local and global optima:

Optimization refers to finding the set of inputs to an objective function that results in the maximum or minimum output.

\* General form of optimization:

$$\text{minimize } x \text{ (or maximum) } f_0(x)$$

objective function.

subject to the constraints  $f_i(x) \leq 0$  (or =)  
(OR  $\geq 0$ ) for  $i = 1, 2, \dots, m$ .

\* Local optima:

A local optimum refers to a point within the domain of a function where the function attains lowest or highest value in its local neighbourhood.

mathematically a point  $x^*$  is said to be a local minimum or maximum if there exists a neighbourhood around  $x^*$

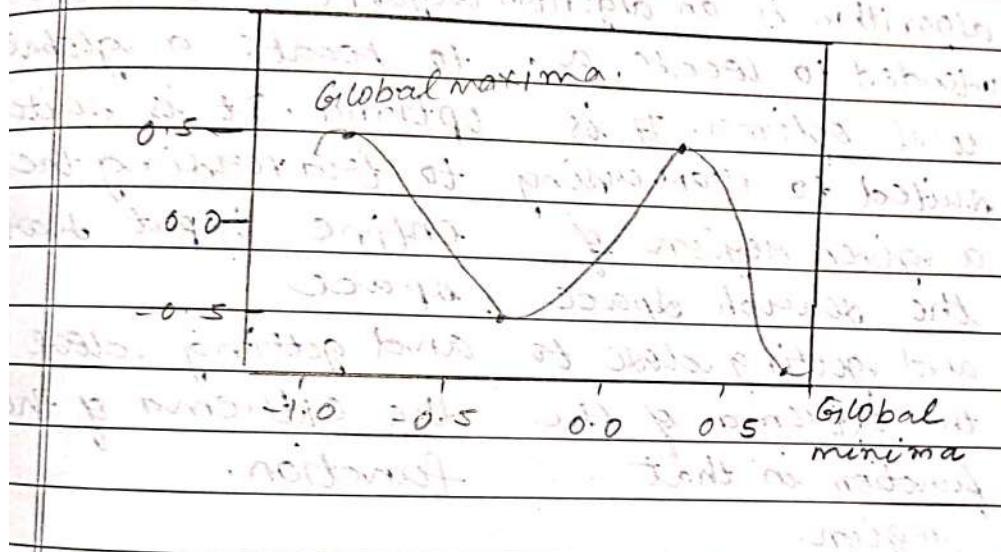
such that  $f(x^*) \leq f(x) \rightarrow$  local minimum  
 $f(x^*) \geq f(x) \rightarrow$  maximum

For all ( $\forall$ )  $x$  within that neighbourhood

In other words a local optimum is a point where the function is lower or higher than its neighbouring points.

### i. Global optima

The Global optimum of a function refers to the point within its domain where the function attains the lowest or highest value over the entire function. In other words it is the overall best solution for the optimization problem.



### Difference between local optima and Global optima

A local optimization algorithm should be used when you know that you are in the region of global optima or your objective function contains single optima.

A Global optimization algorithm should be used when you know very little about structure of objective function, response surface or when you know that the function contains local optima.

Pmax or min	<p>2. It is the extrema of the objective function for a given region of the input space.</p>	<p>A Global optima is the extrema of objective function for the entire input search space</p>
	<p>3. A local optimization algorithm also called a local search algorithm is an algorithm intended to locate a local optima. It is suited to traversing a given region of the search space and getting close to the extrema of the function in that region.</p>	<p>A global optimization algorithm (also called a global search algorithm) is intended to locate a global optima. It is suited to traversing the entire input search space and getting close to the extrema of the function.</p>
	<p>4. Local optimization methods are widely used in applications where there is a value in finding a good point if not the very best</p>	<p>Global optimization is used for problems with a small number of variables where computing time is not critical and the value of finding the true global solution is very high.</p>

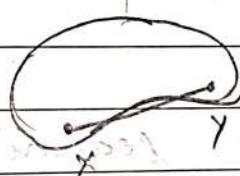
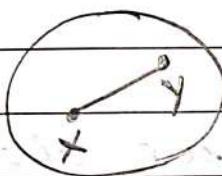
5. 3 Examples of local search algorithm using the definitions include:

- a) Nelder Mead Algorithm
- b) BFGS algorithm
- c) Hill Climbing algorithm.

3 Examples of global search algorithm using the definitions include:

- a) Genetic algorithm
- b) simulated-annealing algorithm
- c) Particle-swarm optimization.

$\Rightarrow$  Convex Sets



CONVEX SET  $\rightarrow$  NON-CONVEX SET

Convexity is a measure of the curvature and due to its curvature we can define two different concepts in convexity which are convex sets and convex functions.

A convex set is a collection of points in which the line XY connecting any two points X, Y in the set lies completely within the set. Ex: Solid Cube

A set is a convex set if we can draw a line segment between any two points in the set and all the points on that line segment are within that set.

Mathematically If  $x, y \in C$  is convex  
 $\rightarrow \alpha x + (1-\alpha)y \in C \forall \alpha \in [0,1]$

$$\text{Ex: } \{x | Ax = b\} = C$$

$$\{x | Ax \leq b\} = C$$

Question: why convexity is important in machine learning.

Answer In machine learning we always want to minimize our objective function which is a loss function. convex functions are important in optimization because local minima of a convex function is also a global minima therefore once we find a local minimum our minimization problem is solved and we have found minimum loss.

### CONVEX FUNCTION

A function is a convex function if its domain is a convex set and it is always below its secant segment.

PAGE NO:  
DATE:

CONVEX

NON - CONVEX

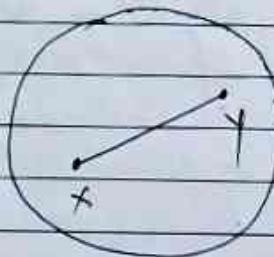
→ Three examples of local search algorithms using the definitions include;

- (i) Nelder Mead algorithm
- (ii) BFGS algorithm
- (iii) hill climbing algorithm.

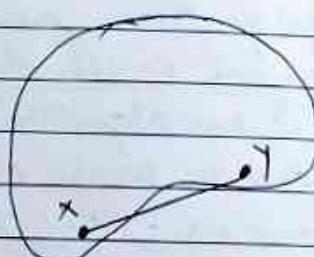
Three examples of global search algorithms using the definitions include;

- Genetic algorithm
- simulated annealing algorithm
- particle swarm optimization

### \* Convex Sets &



Convex Set



Non-convex Set

'Convexity' is a measure of the curvature and due to its curvature we can define two different concepts in convexity which are

- Convex sets
- Convex functions.

A Convex Set is a collection of points in which the line XY connecting any two points X,Y in this set lies completely within the set.

e.g. a  $\text{Kg}$  solid cube

A set is a convex set if we can draw a line segment b/w any two points in the set and all the points on that line segment are within that set.

Mathematically, if  $x, y \in C$  is convex  
 $\Rightarrow \theta x + (1-\theta)y \in C \quad \forall \theta \in [0,1]$

Ex:  $\begin{cases} x/Ax = B \\ y = C \end{cases}$

$$\begin{cases} x/Ax = B \\ y = C \end{cases}$$

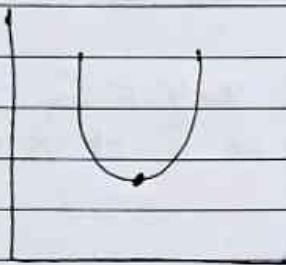
\* Why Convexity is important in machine Learning?

In ML, we always want to minimize our objective function (which is a loss function.)

Convex func. are important in optimization bcoz a local minima of a convex func. is also a global minima.

Therefore, once we find a local minima, our minimization problem is solved and we've found minimum loss.

Convex function is function where its domain is a convex set and is always below its secant segment.



Convex func.

Non-Convex func.

### Three-point Interval Search

24/06/24

This method is used for un-constraint Optimization. In this method we divide the interval in 4 equal parts. We select the central point as a functional value which contains maxima or minima.

Then, we select its interval around its central functional value.

We repeat this process until we reach such a value that shows less than epsilon ( $\epsilon$ ) value.

$$|f(x_c) - f(x_{\text{curr}})| < \epsilon$$

The above step is called stopping criteria.

1. By Using Three Point interval search [TPIS] method to find maximum of  $f(x) = x(5\pi - x)$  on  $[0, 20]$  with  $\epsilon = 0.1$ .

SOL:

Given function is  $f(x) = x(5\pi - x)$

Interval is  $[0, 20]$

$$A = 0$$

$$B = 20$$

$$n = 4$$

$$= \frac{B-A}{n} = \frac{20-0}{4} = 5$$

a	$x_0$	$x_1$	$x_2$	b
0	0.05	10	15	20

$$= f(a) = f(0) = 0(5\pi - 0) = 0,$$

$$= f(x_0) = f(5) = 5(5\pi - 5) = 25\pi - 25 \\ = 25(\pi - 1) \\ = 25(3.1416) \\ = 53.5398 //$$

$$= f(x_1) = f(10) = 10(5\pi - 10) = 57.0796,$$

$$= f(x_2) = f(15) = 15(5\pi - 15) = 10.6194,$$

$$= f(b) = f(20) = 20(5\pi - 20) = -85.8407,$$

The function has

maximum value at  $f(x_1) = 57.0796$

Central value  $\Rightarrow f(x_1) = 57.0796$ .

Sub-interval  $\Rightarrow [5, 15]$ .

$$\Rightarrow |f(x_c) - f(x_{\text{cm}})| \dots$$

$$|57.0796 - 53.5398|$$

$$= 3.5372 \neq \epsilon \quad (\because \epsilon = 0.1)$$

### Iteration 1

Now interval  $[a, b] = [5, 15]$   
 $n = 4$ .

$$\frac{b-a}{n} = \frac{15-5}{4} = \frac{10}{4} = 2.5.$$

$x_0$	$x_1$	$x_2$	$x_3$	$x_4$
5	7.5	10	12.5	15

$$f(x_3) = f(7.5) = 61.5597$$

$$f(x_4) = f(12.5) = 40.0995$$

$$f(x_0) = f(5) = 53.5398$$

$$f(x_1) = f(10) = 57.0796$$

$$f(x_2) = f(15) = 10.6194.$$

$$\text{Central value} \Rightarrow f(x_3) = 61.5597 \\ x_3 = 7.5$$

$$\text{Sub-interval} = [x_0, x_1] = [5, 10]$$

$$|f(x_c) - f(x_{\text{cmm}})|$$

$$|61.5597 - 57.0796|$$

$$= 4.4801 \neq \epsilon$$

Iteration 2

$$\text{New interval } [x_0, x_4] = [5, 10] \\ n = 4$$

$$\frac{b-a}{n} = 1.25.$$



$$f(x_5) = 59.1123$$

$$f(x_6) = 60.8822$$

~~$f(x_3) =$~~

Here,  $f(x_3) = f(7.5) = 61.5597$  is maximum value.

$$\text{Central value} = x_3 = 7.5$$

$$\text{Sub-interval} = [x_5, x_6] = [6.25, 8.75]$$

$$|f(x_c) - f(x_{cm})|$$

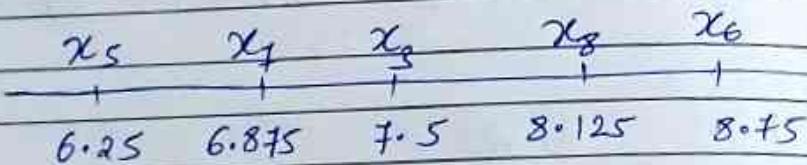
$$|61.5597 - 60.8322|$$

$$= 0.6775 \neq \epsilon$$

Iteration 3

Interval is  $[x_5, x_6] = (6.25 \rightarrow 8.75)$   
 $n = 4$

$$\frac{b-a}{n} = 0.6250$$



$$f(x_7) = 60.7266$$

$$f(x_8) = 61.6116$$

Central maximum value

$$= f(x_8) = 61.6116$$

Sub-interval =  $[x_3, x_6] = [7.5, 8.75]$

$$|f(x_c) - f(x_{cm})|$$

$$= |61.6116 - 61.5597|$$

$$= 0.0519 < \epsilon$$

∴ Maximum value of the function  $f(x)$  is  $61.6116$  at  $x_8 = 8.125$ .

2. By using TPI's method to find max. of  $f(x) = x + \frac{4}{x}$  on  $[0, 2]$  with  $\epsilon = 0.001$ .
3. By using TPI's method to find max. of  $f(x) = x \sin 4x$  on  $[0, 3]$  with  $\epsilon = 0.002$ .

1/07/2020

## Lagrangian Multipliers Method.

\* General problem :

Maximize or minimize  $f(x)$  ;

$$x = (x_1, x_2, x_3, \dots, x_m)^T = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

Subject to  $g_j(x) = 0, j = 1, 2, \dots, m$ .

Then the Lagrangian function is given

by  $L(x, \lambda_j) = f(x) + \sum_{j=1}^m \lambda_j g_j(x)$

- OR -

$$L(x_1, x_2, \dots, x_n, \lambda_1, \lambda_2, \dots, \lambda_m) = f(x_1, x_2, \dots, x_n) + \lambda_1 g_1(x) + \lambda_2 g_2(x) + \dots + \lambda_m g_m(x)$$

- OR -

$$L(x_i, \lambda_j) = f(x_i) + \sum_{j=1}^m \lambda_j g_j(x_i), i = 1, 2, \dots, n$$

$$j = 1, 2, \dots, m.$$

where  $\lambda_j \geq 0$  are called Lagrangian Multipliers.

## \* Necessary Conditions

For a function  $f(x)$  to be max. or min.

$$\frac{dL}{dx_i} = 0 \Rightarrow \frac{df}{dx_i} + \sum_{j=1}^m \lambda_j \frac{dg_j(x)}{dx_i} = 0 \quad (1)$$

and

$$\frac{dL}{d\lambda_j} = g_j(x) = 0 \quad \rightarrow (2)$$

## \* Sufficient Conditions

Let  $x_0$  be a stationary point which can be obtained from solving equations (1) and (2)

Construct the bordered Hessian matrix

$$H_B = \begin{bmatrix} 0 & \nabla g \\ (\nabla g)^T & L_H \end{bmatrix}_{(m+n) \times (m+n)}$$

order

where,

$$\nabla g = \left( \frac{dg_i}{dx_i} \right)_{m \times n} \quad * \quad L_H = \begin{bmatrix} d^2 L(x, \lambda) \\ \frac{\partial f}{\partial x_i \partial x_j} \end{bmatrix}_{m \times n}$$

$n \rightarrow$  no. of variables

$m \rightarrow$  no. of constraints.

NOTE 6 The stationary point  $x_0$  is a min. point if starting with the principle minor determinant of order  $(2m+1)$ . The last  $(n-m)$  principle minors of  $L_H$  have the sign of  $(-1)^m$ .

The stationary point  $x_0$  is the max. point if starting with the principle minor determinant of order  $(2m+1)$ . The last  $(n-m)$  principle minors of  $P$  have the alternating sign pattern starting with  $(-1)^{m+1}$ .

### Problem 6-

Optimize

Minimize  $f(x) = 3x_1^2 + 4x_2^2 + 5x_3^2$  subject to  $x_1 + x_2 + x_3 = 10$ .

Representing the given function  $f(x)$  in the std. Lagrangian form.

$$L(x, \lambda) = 3x_1^2 + 4x_2^2 + 5x_3^2 + \lambda(x_1 + x_2 + x_3 - 10)$$

$$n = 3 \quad (x_1, x_2, x_3)$$

$$m = 1 \quad (\text{subject to } x_1 + x_2 + x_3 = 10)$$

Necessary Conditions:

(Equating to zero)

$$\frac{\partial L}{\partial x_1} = 6x_1 + \lambda \Rightarrow x_1 = -\lambda/6$$

$$\frac{\partial L}{\partial x_2} = 8x_2 + \lambda \Rightarrow x_2 = -\lambda/8$$

$$\frac{\partial L}{\partial x_3} = 10x_3 + \lambda \Rightarrow x_3 = -\lambda/10$$

$$\frac{\partial L}{\partial \lambda} = x_1 + x_2 + x_3 - 10 \Rightarrow x_1 + x_2 + x_3 = 10$$

Substituting values of  $x_1, x_2, x_3$ .

$$\frac{(-\lambda)}{6} + \frac{(-\lambda)}{8} + \frac{(-\lambda)}{10} = 10$$

$$\frac{-20\lambda - 15\lambda - 12\lambda}{120}$$

$$\frac{-47\lambda}{120} = 10$$

$$\Rightarrow \lambda = \frac{-1200}{47}$$

$$x_1 = \frac{-\lambda}{6} = \frac{1200}{47 \times 6} = \frac{200}{47}$$

$$x_2 = \frac{-\lambda}{8} = \frac{1200}{47 \times 8} = \frac{150}{47}$$

$$x_3 = \frac{-\lambda}{10} = \frac{1200}{47 \times 10} = \frac{120}{47}$$

$$x_0 = x^* = (x_1, x_2, x_3)$$

$$x_0 = x^* = \left( \frac{200}{47}, \frac{150}{47}, \frac{120}{47} \right).$$

Sufficient conditions.

$$Dg = \begin{bmatrix} \frac{\partial g}{\partial x_1} & \frac{\partial g}{\partial x_2} & \frac{\partial g}{\partial x_3} \end{bmatrix}$$

$$g(x) = x_1 + x_2 + x_3 - 10.$$

for question. Optimize  $\rightarrow$  sufficient conditions to be done  
minimize  $\rightarrow$  not reqd  
maximize  $\rightarrow$  not reqd

PAGE NO.  
DATE:

$$\nabla g = [1, 1, 1]$$

$$\nabla g \text{ at } x_0 = [1, 1, 1]$$

$$(\nabla g)^T = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \quad \nabla g \text{ at } x_0 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$$L_H = \begin{bmatrix} \frac{\partial^2 L}{\partial x_1^2} & \frac{\partial^2 L}{\partial x_1 \partial x_2} & \frac{\partial^2 L}{\partial x_1 \partial x_3} \\ \frac{\partial^2 L}{\partial x_2 \partial x_1} & \frac{\partial^2 L}{\partial x_2^2} & \frac{\partial^2 L}{\partial x_2 \partial x_3} \\ \frac{\partial^2 L}{\partial x_3 \partial x_1} & \frac{\partial^2 L}{\partial x_3 \partial x_2} & \frac{\partial^2 L}{\partial x_3^2} \end{bmatrix}$$

$$L_H = \begin{bmatrix} 6 & 0 & 0 \\ 0 & 8 & 0 \\ 0 & 0 & 10 \end{bmatrix} \quad L_H \text{ at } x_0 = \begin{bmatrix} 6 & 0 & 0 \\ 0 & 8 & 0 \\ 0 & 0 & 10 \end{bmatrix}$$

$$H_B = \begin{bmatrix} 0 & \nabla g \\ (\nabla g)^T & L_H \end{bmatrix} = \begin{bmatrix} 0 & 1 & 1 & 1 \\ 1 & 6 & 0 & 0 \\ 1 & 0 & 8 & 0 \\ 1 & 0 & 0 & 10 \end{bmatrix}$$

$$\Delta_3 = \begin{vmatrix} 0 & 1 & 1 \\ 1 & 6 & 0 \\ 1 & 0 & 8 \end{vmatrix} = 0(1) - (1)(8-0) + 1(0-6) = -14$$

$\Delta_2 \Delta_3 \rightarrow$  Last two (n-m) principle minors.

Page No.:  
Date:

$$\Delta_4 = \begin{vmatrix} 6 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 8 & 0 & -1 & 1 & 8 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \end{vmatrix} +$$
$$1 \begin{vmatrix} 1 & 6 & 0 & 1 & 6 & 0 \\ 1 & 0 & 0 & -1 & 1 & 0 & 8 \\ 1 & 0 & 1 & 0 & 1 & 0 & 0 \end{vmatrix}$$
$$= 0( ) - 1 [ 1(80-10) - 0( ) + 0( ) ]$$
$$+ 1 [ 1(0-0) - 6(10-10) + 0( ) ] - 1$$
$$[ 2(0-0) - 6(0-8) + 0( ) ]$$
$$= -80 + 60 - 48$$
$$= -188$$

$$f(\mathbf{x}) = 3x_1^2 + 4x_2^2 + 5x_3^2$$
$$= 3\left(\frac{200}{47}\right)^2 + 4\left(\frac{150}{47}\right)^2 + 5\left(\frac{120}{47}\right)^2$$
$$= 54.3232 + 40.7424 + 32.5939$$
$$= 127.6595$$

2 Minimize  $f(x) = x_1^2 + x_2^2$  subject to  
 $a_1 x_1 + a_2 x_2 = b$

$$n = 2 \quad m$$

Std. Lagrangian form.

$$L(x, \lambda) = x_1^2 + x_2^2 + \lambda \left( \frac{a_1 x_1 + a_2 x_2 - b}{\alpha} \right)$$

Necessary conditions.

$$\frac{dL}{dx_1} = 2x_1 + \lambda a_1 \quad x_1 = -\frac{\lambda}{2}$$

$$\frac{dL}{dx_2} = 2x_2 + \lambda a_2 \quad x_2 = -\frac{\lambda}{2}$$

$$\frac{dL}{d\lambda} = a_1 x_1 + a_2 x_2 - b \quad a_1 x_1 + a_2 x_2 = b$$

$$= a_1 \left( -\frac{\lambda}{2} \right) + a_2 \left( -\frac{\lambda}{2} \right) = b$$

$$= -\frac{\lambda a_1^2}{2} - \frac{\lambda a_2^2}{2} = b$$

$$-\lambda \frac{(a_1^2 + a_2^2)}{2} = b$$

$$\boxed{-\lambda = \frac{2b}{a_1^2 + a_2^2}}$$

$$x_1 = \frac{-\lambda}{2} = \frac{2b}{(a_1^2 + a_2^2)} \frac{(a_1)}{2} = \frac{b a_1}{a_1^2 + a_2^2}$$

$$\lambda_2 = \frac{\partial f(b)}{\partial (a_1^2 + a_2^2)} = \frac{b a_2}{a_1^2 + a_2^2}$$

Min. point or stationary point  $x_0 = 1^*$

$$= \left( \frac{a_1 b}{a_1^2 + a_2^2}, \frac{a_2 b}{a_1^2 + a_2^2} \right)$$

$$f(x^*) = f(x_0) = \frac{b^2}{a_1^2 + a_2^2} \rightarrow \frac{a_1^2 b^2 + a_2^2 b^2}{(a_1^2 + a_2^2)^2} = \frac{(a_1^2 + a_2^2) b^2}{(a_1^2 + a_2^2)^2}$$

3. Minimize  $Z = x_1^2 + x_2^2 + x_3^2$  subject to  
conditions

$$x_1 + x_2 + 3x_3 = 2,$$

$$5x_1 + 2x_2 + 2x_3 = 5.$$

$$L(x, \lambda) = x_1^2 + x_2^2 + x_3^2 + \lambda_1(x_1 + x_2 + 3x_3 - 2) + \lambda_2(5x_1 + 2x_2 + 2x_3 - 5)$$

$$\frac{dL}{dx_1} = 2x_1 + \lambda_1 + 5\lambda_2 = 0 \Rightarrow 2x_1 = -5\lambda_2 - \lambda_1$$

$$\lambda_2 = \frac{2x_1}{-5\lambda_2}$$

$$x_1 = \frac{-\lambda_1 - 5\lambda_2}{2}$$

$$\lambda_2 = -\frac{x_1}{5\lambda_2}$$

$$x_1 = -\lambda_1 - 5\lambda_2 / 2$$

$$\frac{dL}{dx_2} = 2x_2 + \lambda_1 + 2\lambda_2$$

$$x_2 = \frac{-\lambda_1 - 2\lambda_2}{2}$$

$$\frac{dL}{dx_3} = 2x_3 + 3\lambda_1 + \lambda_2$$

$$x_3 = \frac{-3\lambda_1 - \lambda_2}{2}$$

$$\frac{dL}{d\lambda_1} = x_1 + x_2 + 3x_3 - 2 = 0$$

$$= \left( \frac{-\lambda_1 - 5\lambda_2}{2} \right) + \left( \frac{-\lambda_1 - 2\lambda_2}{2} \right) + 3 \left( \frac{-3\lambda_1 - \lambda_2}{2} \right) - 2 = 0$$

$$= -\lambda_1 - 5\lambda_2 + -\lambda_1 - 2\lambda_2 (-3\lambda_1 - \lambda_2)^3 - 4 = 0$$

$$= -11\lambda_1 - 10\lambda_2 - 4 = 0 \quad -(1)$$

$$\frac{dL}{d\lambda_2} = 5 \left( \frac{-\lambda_1 - 5\lambda_2}{2} \right) + 2 \left( \frac{-\lambda_1 - 2\lambda_2}{2} \right) +$$

$$\left( \frac{-3\lambda_1 - \lambda_2}{2} \right) = -0.5 = 0.$$

$$= -10\lambda_1 - 30\lambda_2 - 10 = 0. \quad -(2)$$

Solving;

$$\lambda_1 = -2/23 \quad \lambda_2 = -7/23$$

$$x_1 = 37/46$$

$$x_2 = 16/46$$

$$x_3 = 13/46$$

$$x^* = \left( \frac{37}{46}, \frac{16}{46}, \frac{13}{46} \right)$$

$$f(x^*) = \frac{39}{46}$$

4. Maximize  $Z = 2x_1^2 + x_2^2 + 3x_3^2 + 10x_4 +$   
 $8x_2 + 6x_3 - 100.$

subject to  $x_1 + x_2 + x_3 - 20 = 0.$

Explain Stochastic Gradient Descent method.

It is an iterative method for optimising an objective function by updating the model parameters to minimise a loss function.

It is a stochastic approximation of gradient descent optimization which replaces the actual gradient with an estimate.

Stochastic gradient is computationally efficient and can converge faster than batch gradient descent. But it can be noisy and may not converge to global minimum.

Stochastic Gradient randomly divides the set of observations into minibatches. For each minibatch the gradient is computed and the vector is moved. Once all minibatches are used we say that the iteration is finished and start the next one.

Algorithm:

1. Start with an initial set of parameters
2. Shuffle the training data set to introduce randomness
3. Randomly select training from the shuffled dataset
4. Compute the gradient of the cost function w.r.t the model parameters using the current training data set
5. Update the model parameters by taking a step in the direction of the negative gradient scaled by the learning rate.
6. Repeat the process for multiple iterations until the model converges or reaches a predefined stopping criterion

Difference b/w stochastic gradient descent and minibatch gradient descent method.

Stochastic gradient descent

1. Sgd uses a single data point or small batch
2. Updates parameters after each observation using a small batch
3. Faster than gradient descent method especially for large data set
4. Sgd forms zig-zag convergence graph
5. Suitable for large data set

Minibatch gradient descent

1. Minibatch uses a small subset of data
2. Updates parameters after each subset of data
3. More efficient than sgd and less computationally intensive than batch gradient descent method.
4. It forms a smoother convergence graph than sgd, but not as much as batch gradient descent
5. Often used in practise and considered a good compromise b/w Sgd and batch gradient descent.

Advantages and disadvantages of stochastic gradient descent method :-

Advantages :-

\* Speed : It is faster than other variants of gradient descent such as batch gradient descent, and minibatch gradient descent since it uses only one example to update the parameters  
\* Memory efficiency : Since it updates the parameters for each training example one at a time. It is memory efficient and can

handle large data set that cannot fit into memory.

\* Avoidance of Local Minima :- Due to the noisy updates in SGD it has the ability to escape from local minima and converges to a global minima.

### Disadvantages:

\* Noisy updates :- The updates in SGD are noisy and have a high variance which can be the optimization progress.

less stable and lead to oscillations among the minimum.

\* Slow convergence :- SGD makes may require more iterations to converge to the minimum since it updates the parameters for each training example one at a time.

\* Sensitivity to learning rate :- The choice of learning rate can be critical in SGD since using high learning rate can cause the algorithm to overshoot the minimum while a low learning rate can make the algorithm converge slowly.

\* Inaccurate :- Due to the noisy update SGD may not converge to the exact global minimum and can result in suboptimal solution. This can be overcome by using techniques such as learning rate scheduling and moment based updates.

### Newton-Rapson Method:

1. Choose initial value  $x_0$  and a small number  $\epsilon$  if given, if not, consider  $\epsilon = 10^{-4}$ , then compute  $f'(x_0)$
2. Compute  $f''(x_0)$
3. Calculate  $x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)}$
4. Compute  $f'(x_{i+1})$ .
5. If  $|f'(x_{i+1})| < \epsilon$ , then terminate the process  
else set  $i = i + 1$ , go to step 2.

PRO

1. Use NR method to approximate  $f(x) = 2\sin x - \frac{x^2}{10}$  with initial guess of 2.5  $\epsilon, \epsilon = 10^{-4}$ .

Given:

$$f(x) = 2\sin x - \frac{x^2}{10}, \quad x_0 = 2.5, \quad \epsilon = 10^{-4}$$

\* calc must be

$$f'(x) = 2\cos x - \frac{2x}{10} = 2\cos x - \frac{x}{5}$$

in R

$$f''(x) = -2\sin x - \frac{1}{5}$$

$$\text{i. 1st iteration: } i=0 \quad -2.10229$$

$$x_1 = x_0 - \frac{f(x_0)}{f'(x_0)} \Rightarrow 2.5 - \frac{-1.4981}{-0.28724} = 0.99502$$

$$-1.39694$$

$$x_1 = 0.995$$

$$f'(x_1) = 10.889971 \times 10^{-4}$$

ii. 2nd iteration:  $i = 1$

$$x_2 = x_1 - \frac{f'(x_1)}{f''(x_1)} \Rightarrow 0.889971 - \frac{0.99502}{-0.889971} = -1.87754$$

$$x_2 = 1.46903$$

$$f'(x_2) = -0.090621 \times 10^{-4}$$

iii. 3rd iteration:  $i = 2$

$$x_3 = x_2 - \frac{f'(x_2)}{f''(x_2)} \Rightarrow 1.46903 - \frac{-0.090621}{-2.18965}$$

$$x_3 = 1.42764$$

$$f'(x_3) = -0.000191 < 10^{-4}$$

iv. 4th iteration:  $i = 3$

$$x_4 = x_3 - \frac{f'(x_3)}{f''(x_3)} \Rightarrow 1.42764 - \frac{-0.000191}{-2.17954}$$

$$x_4 = 1.42755$$

$$f'(x_4) = 0.000003877 < 10^{-4}$$

$$\therefore f'(x_4) = 0.000003877 \text{ for } x_4 = 1.42755$$

$$\therefore f(x_4) = 1.77573$$

2. Use NR method to approximate  $f(x) = \frac{x^2 + 54}{x}$ ,  $x_0 = 2.8$   
 $\epsilon = 10^{-4}$

Given:

$$f(x) = \frac{x^2 + 54}{x}; x_0 = 2.8, \epsilon = 10^{-4}$$

$$f'(x) = \frac{2x^2 - 54}{x^2}; f''(x) = \frac{2 + 54}{x^3}$$

$$f(x) = 2x - 54 ; f'(x) = \frac{2+108}{x^3} = \frac{108}{x^3}$$

i. iteration 1 :  $i=0$

$$x_1 = x_0 - \frac{f(x_0)}{f'(x_0)} \Rightarrow 2.8 - \frac{(-40.56)}{\frac{6.91983}{6.1535}} = 11.8944$$

$$\Rightarrow 2.8 - \frac{(-1.28776)}{6.91983} = 2.9861$$

$$x_1 = 2.9861$$

$$|f'(x_1)| = 1 - 0.08379 < 10^{-4}$$

ii. iteration 2 :  $i=1$

$$x_2 = x_1 - \frac{f'(x_1)}{f''(x_1)} \Rightarrow 2.9861 - \frac{(-0.08379)}{6.05612}$$

$$x_2 = 2.999936$$

$$|f'(x_2)| = 1 - 0.000384 < 10^{-4}$$

iii. iteration 3 :  $i=2$

$$x_3 = x_2 - \frac{f'(x_2)}{f''(x_2)} \Rightarrow 2.999936 - \frac{(-0.000384)}{6.000256}$$

$$x_3 = 2.999999997$$

$$|f'(x_3)| = 1 - 0.00000016 < 10^{-4}$$

$$\therefore |f'(x_3)| = 1 - 0.00000016 < 10^{-4} \text{ for } x_3 = 2.999999997$$

$$\therefore f(x) = 27 //$$

3. Find the extrema of the function using NR method

Given :  $f(x) = 5\sin(2x) - 2x^2 - 4x$ ,

$\epsilon = 10^{-4}$ ,  $x_0 = \text{---}$  (is not given, take 0)

$$f'(x) = 25\cos(2x) - 4x - 4$$

$$= 10\cos(2x) - 4x - 4$$

$$f''(x) = -20\sin(2x) - 4$$

Final

Given:  $f(x) = 5 \sin(2x) - 2x^2 - 4x$   
 $x_0 = -2, \epsilon = 10^{-4}$

$$f'(x) = 10 \cos(2x) - 4x - 4$$

$$f''(x) = -20 \sin(2x) - 4$$

i. iteration 1 : i=0

$$x_1 = x_0 - \frac{f'(x_0)}{f''(x_0)} \Rightarrow -2 - \frac{(-2.5364)}{-19.13605}$$

$$x_1 = -2.1325$$

$$f'(x_1) = 10.20387 \times 10^{-4}$$

ii. iteration 2 : i=1

$$x_2 = x_1 - \frac{f'(x_1)}{f''(x_1)} \Rightarrow -2.1325 - \frac{(0.20387)}{-22.03159}$$

$$x_2 = -2.1232$$

$$f'(x_2) = 10.00192 \times 10^{-4}$$

iii. iteration 3 : i=2

$$x_3 = x_2 - \frac{f'(x_2)}{f''(x_2)} \Rightarrow -2.1232 - \frac{(0.00192)}{-21.86935}$$

$$x_3 = -2.1232$$

$$f'(x_3) = 10.0002655 \times 10^{-4}$$

iii: iteration 4 : i = 3

$$x_4 = x_3 - \frac{f'(x_3)}{f''(x_3)} \Rightarrow -2.1232 - \frac{(-0.0002655)}{-21.86755}$$

$$x_4 = -2.12321$$

$$|f'(x_4)| = |-0.0000001| < 10^{-4} //$$

# Optimization Techniques (BCS405C)

## Module-5: Advanced Optimization

### RMSprop:

RMSProp (Root Mean Squared Propagation) is an adaptive learning rate optimization algorithm. It is an extension of the popular Adaptive Gradient Algorithm and is designed to dramatically reduce the amount of computational effort used in training neural networks. This algorithm works by exponentially decaying the learning rate every time the squared gradient is less than a certain threshold. This helps reduce the learning rate more quickly when the gradients become small. In this way, RMSProp is able to smoothly adjust the learning rate for each of the parameters in the network, providing a better performance than regular Gradient Descent alone.

The RMSprop algorithm utilizes exponentially weighted moving averages of squared gradients to update the parameters. Here is the mathematical equation for RMSprop:

1. Initialize parameters:
  - Learning rate:  $\alpha$
  - Exponential decay rate for averaging:  $\gamma$
  - Small constant for numerical stability:  $\epsilon$
  - Initial parameter values:  $\theta$
2. Initialize accumulated gradients (Exponentially weighted average):
  - Accumulated squared gradient for each parameter:  $E_t = 0$
3. Repeat until convergence or maximum iterations:
  - Compute the gradient of the objective function with respect to the parameters:  
$$g_t = \nabla_{\theta} J(\theta_t)$$
  - Update the exponentially weighted average of the squared gradients:  
$$E_t = \gamma E_{t-1} + (1 - \gamma) g_t^2$$
  - Update the parameters:  
$$\theta_{t+1} = \theta_t - \alpha \frac{g_t}{\sqrt{E_t + \epsilon}}$$

where,

- $g_t$  is the gradient of the loss function with respect to the parameters at time  $t$
- $\gamma$  is a decay factor
- $E_t$  is the exponentially weighted average of the squared gradients
- $\alpha$  is the learning rate
- $\epsilon$  is a small constant to prevent division by zero

This process is repeated for each parameter in the optimization problem, and it helps adjust the learning rate for each parameter based on the historical gradients. The exponential moving average allows the algorithm to give more importance to recent gradients and dampen the effect of older gradients, providing stability during optimization.

### Stochastic Gradient Descent with Momentum:

SGD

SGD with Impulse

$$\theta_j \leftarrow \theta_j - \epsilon \nabla_{\theta_j} L(\theta)$$

$$v_{t+1} \leftarrow \rho v_t + \nabla_{\theta} L(\theta)$$

$$\theta_j \leftarrow \theta_j - \epsilon v_{t+1}$$

GL. 2 Stochastic GD (left), SGD with momentum (right).

On the left side in GL. 2 is the formula for the weight updates according to the regular stochastic gradient descent (SGD for short). The equation on the right represents the rule for

the updates of the weights according to the SGD with momentum. **Momentum appears here as an additional term**  $\rho \cdot v$ , which is added to the regular update rule.

Intuitively speaking, by adding this impulse term, we let our **gradient build up some sort of velocity**  $v$  during training. The velocity is the running sum of the gradients weighted by  $\rho$ .

The parameter  $\rho$  can be thought of as friction that “slows” the velocity down a bit. In general, velocity can be seen to increase with time. By using the momentum term, **saddle points and local minima become less dangerous** for the gradient. This is because the step size toward the global minimum now depends not only on the slope of the loss function at the current point, but also on the velocity that has built up over time.

If you would like a physical representation of stochastic gradient descent with momentum, please imagine a ball rolling down a hill, increasing in velocity with time. If this ball encounters an obstacle along the way, such as a hole or flat ground with no slope, its built-up velocity  $v$  would give the ball enough force to roll over this obstacle. In this case, the **flat ground represents a saddle point** and the **hole represents a local minima** of a loss function.

## Adagrad:

Another optimization strategy is called as AdaGrad. The idea behind AdaGrad is that you keep a running sum of squared gradients during optimization. In this case, we don't have a momentum term, but an expression  $g_{t+1}$ , which is **the sum of squared gradients** up to the time ( $t + 1$ ).

SGD with Impulse

$$v_{t+1} \leftarrow \rho v_t + \nabla_\theta L(\theta)$$

$$\theta_j \leftarrow \theta_j - \epsilon v_{t+1}$$

Adagrad

$$g_0 = 0$$

$$g_{t+1} \leftarrow g_t + \nabla_\theta L(\theta)^2$$

$$\theta_j \leftarrow \theta_j - \epsilon \frac{\nabla_\theta L}{1e^{-5} + \sqrt{g_{t+1}}}$$

### The updating rule of parameters at AdaGrad

When we optimize a weights  $\theta_j$ , we divide the current gradient  $\nabla_j L$  by the root of the term  $g_{t+1}$ . To understand the intuition behind AdaGrad, please imagine a loss function in a two-dimensional space. In this space, the gradient of the loss function **increases very weakly in one direction and very strongly in the other direction**. If we now sum up the gradients along the axis in which the gradients increase weakly, the squared sum of these gradients becomes even smaller.

If during the update step we divide the current gradient  $\nabla_j L$  by a very small sum of the squared gradients  $g_{t+1}$ , the quotient becomes very high. For the other axis, along which the gradients increase sharply, exactly the opposite is true.

This means that we speed up the updating process **along the axis with weak gradients** by increasing these gradients along this axis. On the other hand, we slow down the updates of the weights **along the axis with large gradients**.

However, there is a problem with this optimization algorithm.

Imagine what would happen to the sum of squared gradients if the training takes too long. Over time, this term would **grow larger**. When the current gradient is divided by this large number, the update step for the weights becomes very small. It is as if we were using **a very low learning rate**, which becomes even lower the longer the training takes. In the worst case, we would get stuck at AdaGrad and the training would go on forever.

## ADAM:

So far, we have used the momentum term to determine the velocity of the gradient and update the weight parameter in the direction of that velocity. In the case of AdaGrad and RMSProp, we used the sum of squared gradients to scale the current gradient so that we could update the weights in each space dimension with the same ratio.

These two methods seemed like pretty good ideas. Why don't we just **take the best of both worlds** and combine these ideas into a single algorithm? That's exactly the concept behind the final optimization algorithm I'd like to introduce. This algorithm is called ADAM.

The main part of this optimization algorithm consists of the following three equations. These equations may seem complicated at first glance, but if you look closely, you will see some similarities with the last three optimization algorithms.

The main part of this optimization algorithm consists of the following three equations. These equations may seem complicated at first glance, but if you look closely, you will see some similarities with the last three optimization algorithms.

$$\begin{aligned}
 m_0 &= 0 \text{ and } v_0 = 0 \\
 m_{t+1} &\leftarrow \beta_1 m_t + (1 - \beta_1) \nabla_\theta L(\theta) \quad \text{----- Impulse} \\
 v_{t+1} &\leftarrow \beta_2 v_t + (1 - \beta_2) \nabla_\theta L(\theta)^2 \quad \text{----- RMS prop} \\
 \theta_j &\leftarrow \theta_j - \frac{\epsilon}{\sqrt{v_{t+1}} + 1e^{-5}} m_{t+1} \quad \text{----- RMSprop+ Impulse}
 \end{aligned}$$

Eq. 1 Parameter update rule for ADAM Optimizer

The first expression in Eq.1 looks a bit like SGD with momentum. In this case, the term  $m_t$  would be the velocity and the term  $\beta_1$  would be the friction term. In the case of ADAM, we refer to  $m_t$  as the “first momentum”,  $\beta_1$ , on the other hand, is just a hyperparameter. However, the difference with SGD with momentum is the factor  $(1 - \beta_1)$  multiplied by the current gradient.

The second expression in Eq.1 **can be considered as RMSProp**, where we keep the running sum of squared gradients. Also in this case, there is the factor  $(1 - \beta_2)$ , which is multiplied by the squared gradient.

The term  $v_t$  in Eq.1 is called the “second momentum” and is also just a hyperparameter. The final update equation (third term in the equation) can be viewed as **a combination of RMSProp and SGD with momentum**.

ADAM has incorporated the nice features of the previous two optimization algorithms. However, there is a small problem with ADAM, and that is what happens at the beginning of training....

At the very first time step = 0 , the first and second pulse terms  $m_0$  and  $v_0$  are set to zero. After the first update of the second momentum  $v_1$  , this term is still very close to zero. When

we update the weight parameters in the last expression in equation, we divide by a very small second momentum  $v_1$ . This leads to a very large first update step.

This very large first update step is not the result of the geometry of the problem, but an artefact of the fact that we initialized the first and second momentum to zero. To address the problem of large update steps happening at the beginning of training, ADAM includes a correction clause:

$$\begin{aligned} m_{t+1} &\leftarrow \beta_1 m_t + (1 - \beta_1) \nabla_\theta L(\theta) \quad \text{----- Impulse} \\ v_{t+1} &\leftarrow \beta_2 v_t + (1 - \beta_2) \nabla_\theta L(\theta)^2 \quad \text{----- RMS prop} \\ \hat{m}_{t+1} &\leftarrow \frac{m_{t+1}}{1 - \beta_1^t} \quad \text{and} \quad \hat{v}_{t+1} \leftarrow \frac{v_{t+1}}{1 - \beta_2^t} \quad \text{----- Bias corrector} \\ \theta_j &\leftarrow \theta_j - \frac{\epsilon}{\sqrt{v_{t+1}} + 1e^{-5}} m_{t+1} \quad \text{----- RMSprop+ Impulse} \end{aligned}$$

### Eq. 2 Bias Correction for ADAM Optimizer

As you can see from Eq.2 after the initial update of the first and second pulses, we **make an unbiased estimate of these pulses** by considering the current time step. With the so-called bias correction, we obtain the corrected first and second impulses with  $\hat{m}_{t+1}$  and  $\hat{v}_{t+1}$  respectively.

These corrections cause the values of the first and second impulse to be higher at the beginning of the training than without this correction. As a result, the first update step of the neural network weight parameters does not become too large. Thus, the training is not already messed up at the very beginning.

With the additional bias corrections, we obtain the complete form of the ADAM optimizer.

## The best optimization algorithm for Machine Learning:

Now, we can discuss the question of which algorithm is the best to train a neural network.

In general, a normal gradient descent algorithm is **more than sufficient for simpler tasks**. If you are not satisfied with the accuracy of your model, you can try RMSprop or add a momentum term to your gradient descent algorithms.

However, **ADAM is the best neural network optimization algorithm** available today. This optimization algorithm is excellent for almost any deep learning problem you will ever encounter in practice. Especially if you set ADAM's hyperparameters to the following values:

- $\beta_1 = 0.9$
- $\beta_2 = 0.999$
- learning rate =  $0.001 - 0.0001$

ADAM and the above values of the hyperparameters are a very good starting point for any problem and virtually any type of neural network architecture.

For this reason, **ADAM is the default optimization algorithm** for any Deep Learning problem. Only in very few cases we can switch to other optimization algorithms like RMSprop, Adagrad etc..

It is recommended that we always start with ADAM, regardless of the neural network architecture of the problem domain.

## What is saddle point problem in Machine Learning?

The term “**saddle point**” in the context of machine learning refers to a specific point in the optimization landscape of a cost function where the gradient is zero, but the point is neither a minimum nor a maximum. Instead, it’s a point where the surface of the cost function resembles a saddle, with some dimensions curving upward and others downward.

### Key Characteristics of a Saddle Point:

#### 1. Zero gradient:

At a saddle point, the gradient of the cost function is zero. This means that the partial derivatives with respect to each parameter are all equal to zero.

#### 2. Neither Minimum nor Maximum:

Unlike a local minimum or maximum, a saddle point is a point where the cost function neither reaches a minimum nor a maximum value.

#### 3. Flat in Some Dimensions, Steep in Others:

The surface of the cost function is flat in certain dimensions (where the partial derivatives are zero) and steep in others. It creates a saddle-like shape.

#### 4. Challenge for Optimization Algorithms:

Optimization algorithms, such as gradient descent, can get stuck or converge very slowly near saddle points because the gradient is zero, and the algorithm may struggle to determine the right direction to move.

### Dealing with Saddle Points:

#### 1. Higher-Order Optimization:

Using higher-order optimization methods, such as Newton’s method or quasi-Newton methods, can help in navigating around saddle points more effectively.

#### 2. Random Initialization:

Random initialization of parameters can help algorithms escape saddle points and find better solutions.

#### 3. Noise Injection:

Injecting noise into the optimization process can help the algorithm overcome saddle points by introducing randomness.

#### 4. Gradient Clipping:

Gradient clipping can prevent the gradients from becoming too small, allowing the optimization algorithm to make more meaningful updates.

#### 5. Use of Momentum:

Techniques like momentum in optimization algorithms can help carry the algorithm through flat regions, reducing the chances of getting stuck at saddle points.

### In Summary:

Saddle points present challenges in the optimization of cost functions in machine learning. While not as problematic as local minima, they can still hinder the convergence of optimization algorithms. Researchers and practitioners employ various strategies, such as higher-order optimization, random initialization, and noise injection, to mitigate the impact of saddle points and improve the efficiency of optimization processes in high-dimensional spaces.

## Difference between convex optimization and non-convex optimization:

Convex optimization	Non-convex optimization
A <i>convex optimization problem</i> is a problem where all of the constraints are <u>convex functions</u> , and the objective is a convex function if minimizing, or a concave function if maximizing. Linear functions are convex, so linear programming problems are convex problems.	A <i>non-convex optimization problem</i> is any problem where the objective or any of the constraints are non-convex. Such a problem may have multiple feasible regions and multiple locally optimal points within each region.
Convex functions have a unique global minimum, making optimization easier and more reliable.	Non-convex functions, on the other hand, can have multiple local minima, making optimization more challenging.
Convex optimization problems are generally easier to solve computationally due to the properties of convex functions.	Non-convex optimization problems are more computationally challenging as they require more sophisticated algorithms to explore the solution space efficiently and avoid getting trapped in local minima.
Convex optimization algorithms are guaranteed to converge to the global minimum given certain conditions, such as convexity and continuity of the objective function.	Non-convex optimization algorithms may not always converge to the global minimum, especially when dealing with complex functions with multiple local minima.
Convex optimization is widely used in various fields such as machine learning, engineering, economics, and operations research due to its well-understood properties and efficient algorithms.	Non-convex optimization is essential in many real-world problems where the objective function is inherently non-convex, such as deep learning, molecular modeling, and financial modeling.

## Advantages of RMSprop over Adagrad:

**Adagrad** (short for Adaptive Gradient) is an adaptive learning rate optimization algorithm that adapts the learning rate for each parameter based on the historical gradients. It is well-suited for sparse data and large-scale problems. Adagrad computes the learning rate for each parameter individually, which allows it to converge quickly on sparse features. However, Adagrad's learning rate can become too small as training progresses, leading to slow convergence.

### Advantages:

- Adapts the learning rate for each parameter based on the sum of the squares of its past gradients, making it suitable for sparse datasets
- Performs well for problems with low learning rates

### Disadvantages:

- The learning rate can become too small over time, hindering convergence
- Requires more memory than SGD due to the need to store past gradients' sums

RMSProp (short for Root Mean Square Propagation) is an adaptive learning rate optimization algorithm that divides the learning rate by a moving average of the square root of the accumulated gradients. It is well-suited for non-convex problems and allows the model to

adapt the learning rate to different features in the data. However, RMSProp can suffer from slow convergence in some cases.

### **Advantages:**

- Adapts the learning rate for each parameter based on an exponentially decaying average of past squared gradients, making it suitable for deep neural networks with many layers
- Performs well in various domains

### **Disadvantages:**

- May converge slowly compared to other optimizers like Adam
- Not suitable for problems with very sparse data

### **Conclusion:**

- **RMSProp** adapts the learning rates based on recent gradient magnitudes, making it effective for non-stationary objectives and neural networks.
- **Adagrad** is particularly effective for sparse data, adapting the learning rate to different parameters, but it may suffer from a constantly decreasing learning rate.