



**BITS** Pilani  
Pilani Campus

# Mathematical Foundations

MFDS Team





# **Mathematical Foundations**

# **Webinar 3**

# Agenda

---

1. Maxima and Minima
2. Gradient of a Vector-Valued Function
3. Gradient Descent
4. Gradient Descent Examples
5. Previous Year Problems

# 1. Maxima and Minima

---

- For a given function  $f(x, y)$ ,
- If  $f_{xx} < 0$  and  $f_{xx}f_{yy} - f^2_{xy} > 0$  at  $(a,b)$ , then  $f$  has a local maximum value at  $(a,b)$ . ✓
- If  $f_{xx} > 0$  and  $f_{xx}f_{yy} - f^2_{xy} > 0$  at  $(a,b)$ , then  $f$  has a local minimum value at  $(a,b)$ . ✓
- If  $f_{xx}f_{yy} - f^2_{xy} < 0$  at  $(a,b)$ , then  $f$  has a saddle point at  $(a,b)$ . ✓
- If  $f_{xx}f_{yy} - f^2_{xy} = 0$  at  $(a,b)$ , then the test is inconclusive and further investigation is required. ✓

# Hessian matrix

- ▶ The considerations of the previous slide show that determining whether there is a minimum or maximum at the point  $(a, b)$  boils down to looking at the following matrix and asking if it is positive-definite or not.

- ▶ This is the Hessian matrix is  $\begin{bmatrix} f_{xx} & f_{xy} \\ f_{yx} & f_{yy} \end{bmatrix}$

$$H = \begin{vmatrix} r & s \\ s & t \end{vmatrix} = rt - s^2$$

**Note : The discriminant  $f_{xx}f_{yy} - f_{xy}^2$  is noting but the determinant of the Hessian matrix.**

## Example

- Find the local extreme values of the function  $f(x, y) = xy - x^2 - y^2 - 2x - 2y + 4$ .

### Solution

The function is defined and differentiable for all  $x$  and  $y$  and its domain has no boundary points. The function therefore has extreme values only at the points where  $f_x$  and  $f_y$  are simultaneously zero. This leads to

$$f_x = y - 2x - 2 = 0, \quad f_y = x - 2y - 2 = 0, \quad \text{or} \quad x = y = -2.$$

Therefore, the point  $(-2, -2)$  is the only point where  $f$  may take on an extreme value. To see if it does so, we calculate

$$f_{xx} = -2, \quad f_{yy} = -2, \quad f_{xy} = 1.$$

The discriminant of  $f$  at  $(a, b) = (-2, -2)$  is  $f_{xx}f_{yy} - f_{xy}^2 = (-2)(-2) - (1)^2 = 4 - 1 = 3$ .

The combination  $f_{xx} < 0$  and  $f_{xx}f_{yy} - f_{xy}^2 > 0$  tells us that  $f$  has a local maximum at  $(-2, -2)$ .

The value of  $f$  at this point is  $f(-2, -2) = 8$

Thus, the function has a local maximum value of 8 at the point  $(-2, -2)$ , and there are no other local extrema.

## Example

2. Find the local extreme values of  $f(x, y) = xy$ .

**Solution** Since  $f$  is differentiable everywhere (Figure 14.43), it can assume extreme values only where

$$f_x = y = 0 \quad \text{and} \quad f_y = x = 0.$$

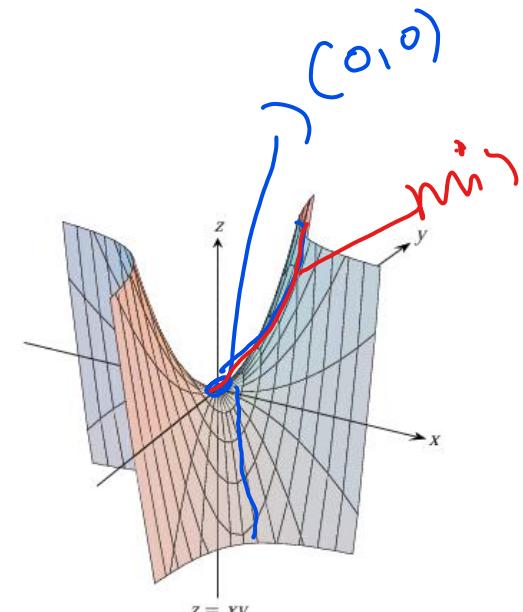
(0, 0)

Thus, the origin is the only point where  $f$  might have an extreme value. To see what happens there, we calculate

$$f_{xx} = 0, \quad f_{yy} = 0, \quad f_{xy} = 1.$$

$$\sqrt{t - s}$$

$$f_{xx}f_{yy} - f_{xy}^2 = -1, \quad < 0$$

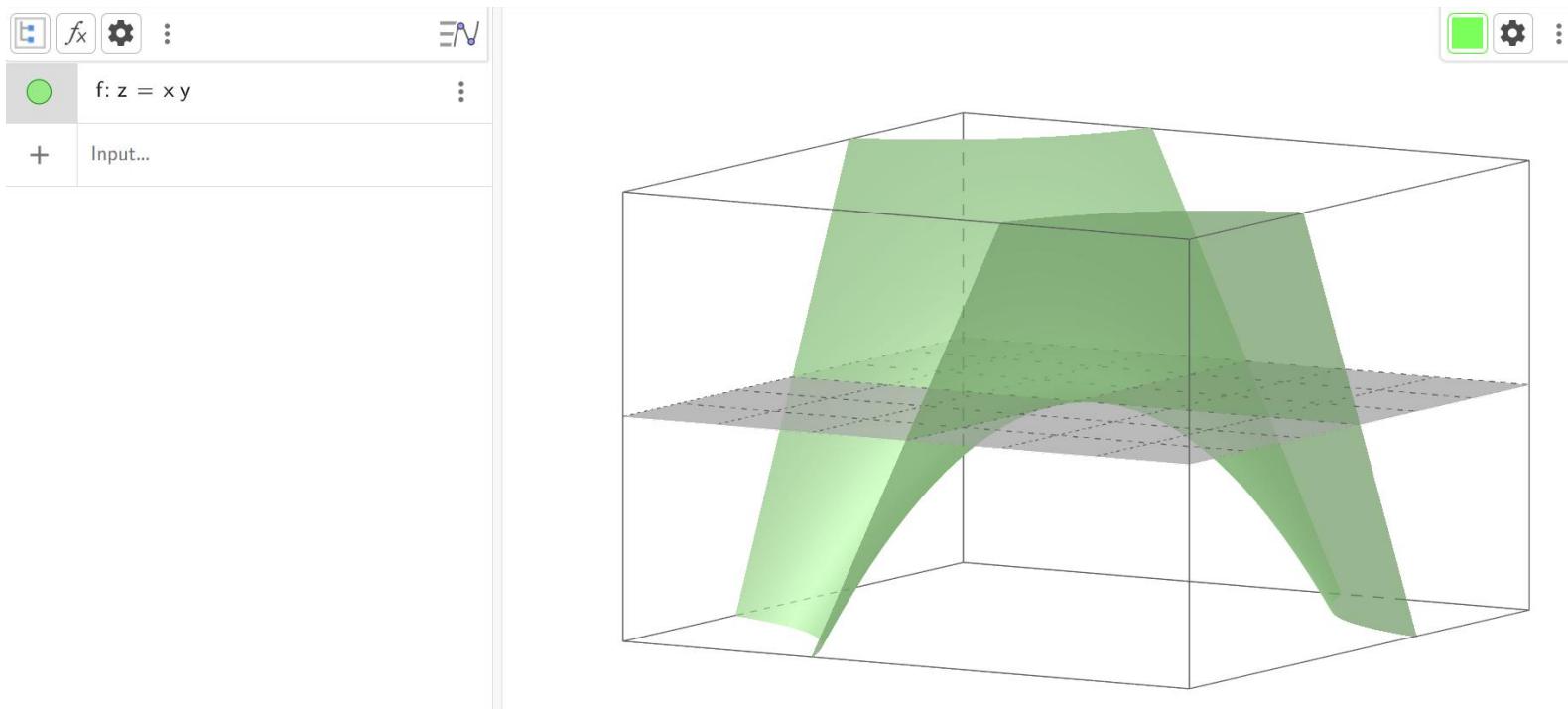


The discriminant,

1) max

2) saddle

3)  $\sqrt{t - s} > 0$  &  $\sqrt{s} > 0$  max  
 4)  $f_{xx} \quad \sqrt{t - s^2} = 0$



## Maxima and Minima of a function of three variables using the Hessian.

Consider a function  $f(x, y, z)$  of three variables.

### Step 1: Critical Points

To find critical points, we need to solve the system of equations given by the

$$\text{gradient of } f \text{ and is given by } \nabla f = \left[ \frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z} \right]^T = [0 \quad 0 \quad 0]^T$$

The solutions to this system represent critical points

### Step 2: Hessian Matrix

The Hessian matrix, denoted as  $H$ , is a square matrix of second-order partial derivatives of  $f$  and is given by

$$H = \begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} & \frac{\partial^2 f}{\partial x \partial z} \\ \frac{\partial^2 f}{\partial y \partial x} & \frac{\partial^2 f}{\partial y^2} & \frac{\partial^2 f}{\partial y \partial z} \\ \frac{\partial^2 f}{\partial z \partial x} & \frac{\partial^2 f}{\partial z \partial y} & \frac{\partial^2 f}{\partial z^2} \end{bmatrix}$$

$$\begin{aligned} f_x &= 0 \\ f_y &\neq 0 \\ f_z &= 0 \end{aligned}$$

### Step 3: Second Derivative Test

The determinant of the Hessian matrix,  $\det(H)$ , is used to determine the nature of critical points:

Let  $x_0$  be a critical point of  $f$ .

$$\det(H) = (\lambda_1, \lambda_2, \lambda_3) = 0$$

- i) If the Hessian matrix  $H$  is positive definite (i.e. all eigenvalues are strictly positive) then  $\delta^T H(x_0)\delta > 0$ . So  $f(x) > f(x_0)$  and hence  $x_0$  is a point of local minima.
- ii) If the Hessian matrix  $H$  is negative definite (i.e. all eigenvalues are strictly negative) then  $\delta^T H(x_0)\delta < 0$ . So,  $f(x) < f(x_0)$  and hence  $x_0$  is a point of local maxima.
- iii) If the Hessian matrix  $H$  has both positive and negative eigenvalues then  $x_0$  is a saddle point.
- iv) If the determinant of Hessian matrix  $H$  is zero then the test is inconclusive.

This method helps classify critical points as maxima, minima, or saddle points in functions of three variables

## Example

3. Discuss the maximum and minimum value of the function  $f(x,y,z)=x^2+y^2+z^2+x-2z-xy$ .

**Solution:**  $f(x,y,z) = x^2 + y^2 + z^2 + x - 2z - xy$

$$\nabla_{x,y,z} f(x,y,z) = [f_x \quad f_y \quad f_z] = [2x+1-y \quad 2y-x \quad 2z-2]$$

$$f_x = 2x+1-y, \quad f_y = 2y-x, \quad f_z = 2z-2$$

$$f_{xx} = 2$$

To find critical points, we need to solve the system of equations given by the gradient of  $f$  is equal to zero. i.e

$$\nabla f = \left[ \frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z} \right]^T = [0 \quad 0 \quad 0]^T$$

$$2x+1-y=0 \quad \Rightarrow y = 2x+1 \quad (1)$$

$$2y-x=0 \quad \Rightarrow x=2y \quad (2)$$

$$\Rightarrow z=1$$

Substituting (2) in (1) gives  $y = -1/3$ ,  $x = -2/3$  and  $z = 1$

$$(x, y, z) = \left( -\frac{2}{3}, -\frac{1}{3}, 1 \right)$$

$$f_x = 2n+1-y, \quad f_y = 2y-x \quad f_z = 2z-2$$

$$f_{xx} = 2, \quad f_{yy} = 2, \quad f_{zz} = 2$$

$$f_{xy} = -1, \quad f_{yz} = 0, \quad f_{xz} = 0$$

The only critical point is  $(-2/3, -1/3, 1)$  where the function will be either minima or maxima.

The Hessian matrix  $H = \begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} & \frac{\partial^2 f}{\partial x \partial z} \\ \frac{\partial^2 f}{\partial y \partial x} & \frac{\partial^2 f}{\partial y^2} & \frac{\partial^2 f}{\partial y \partial z} \\ \frac{\partial^2 f}{\partial z \partial x} & \frac{\partial^2 f}{\partial z \partial y} & \frac{\partial^2 f}{\partial z^2} \end{bmatrix} = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix}$

Further, The determinant of Hessian  $H$  is ,  $\det(H) = \begin{vmatrix} 2 & -1 & 0 \\ -1 & 2 & 0 \\ 0 & 0 & 2 \end{vmatrix} = 6 > 0$

The determinants of all the leading principal minors of  $H$  are given by .

$$M_1 = 2 > 0, \quad M_2 = \begin{vmatrix} 2 & -1 \\ -1 & 2 \end{vmatrix} = 3 > 0; \quad M_3 = \begin{vmatrix} 2 & -1 & 0 \\ -1 & 2 & 0 \\ 0 & 0 & 2 \end{vmatrix} = 6 > 0$$

If  $\det(H) > 0$  and all the leading principal minors of  $H$  are positive, then the critical point is a local minimum. Thus we have a minimum at  $(-2/3, -1/3, 1)$  and the minimum value is  $-4/3$ .

## 2. Gradient of a Vector-Valued Function



### a) Gradient of a vector with respect to the input vector

For a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and a vector valued function  $x = [x_1, x_2, \dots, x_n]^T \in \mathbb{R}^n$ , the

corresponding vector function is given as  $f(x) = \begin{bmatrix} f_1(x) \\ \vdots \\ f_m(x) \end{bmatrix} \in \mathbb{R}^m$ . The gradient of

$f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  with respect to  $x \in \mathbb{R}^n$  is given by

$$\frac{df(x)}{dx} = \left[ \begin{array}{c} \frac{\partial f_1(x)}{\partial x_1} \dots \frac{\partial f_1(x)}{\partial x_n} \\ \vdots \\ \frac{\partial f_m(x)}{\partial x_1} \dots \frac{\partial f_m(x)}{\partial x_n} \end{array} \right]$$

$$= \left[ \begin{array}{cc} \frac{\partial f_1(x)}{\partial x_1} & \dots & \frac{\partial f_1(x)}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m(x)}{\partial x_1} & \dots & \frac{\partial f_m(x)}{\partial x_n} \end{array} \right] \in \mathbb{R}^{m \times n}$$

#### Example

4. Let  $f = [3x_1^2x_2 \quad 2x_1 + x_2^8]^T$ , then find the Gradient of  $f$ .

Solution:

$$\begin{aligned} \frac{df}{dx} &= \begin{bmatrix} \frac{\partial f_1(x)}{\partial x_1} & \frac{\partial f_1(x)}{\partial x_2} \\ \frac{\partial f_2(x)}{\partial x_1} & \frac{\partial f_2(x)}{\partial x_2} \end{bmatrix} \\ &= \begin{bmatrix} 6x_1x_2 & 3x_1^2 \\ 2 & 8x_2^7 \end{bmatrix} \end{aligned}$$

$$f = \begin{pmatrix} f_1 \\ f_2 \end{pmatrix} \quad x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

## b) Gradient of a scalar with respect to a matrix

If  $f : \mathbb{R}^{M \times N} \rightarrow \mathbb{R}$  is a scalar function of a matrix  $X \in \mathbb{R}^{M \times N}$ , then the gradient of  $f$  with respect

to  $X$  is defined as the  $M \times N$  matrix given by  $\nabla_X f = \begin{bmatrix} \frac{\partial f}{\partial X_{11}} & \frac{\partial f}{\partial X_{12}} & \dots & \frac{\partial f}{\partial X_{1N}} \\ \frac{\partial f}{\partial X_{21}} & \frac{\partial f}{\partial X_{22}} & \dots & \frac{\partial f}{\partial X_{2N}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial X_{M1}} & \frac{\partial f}{\partial X_{M2}} & \dots & \frac{\partial f}{\partial X_{MN}} \end{bmatrix} \in \mathbb{R}^{M \times N}$ . A = \begin{bmatrix} xy^T & yw^T \end{bmatrix}

### Example

5. Let  $f = 2x + 4y + 3z + 4w$  and the matrix  $A = \begin{bmatrix} x & y \\ z & w \end{bmatrix}$  the find the gradient of  $f$  w.r.t matrix  $A$ .

~~+ 3t~~

**Solution:** The gradient of  $f$  is  $\frac{df}{dA} = \begin{bmatrix} \frac{\partial f}{\partial x} & \frac{\partial f}{\partial y} \\ \frac{\partial f}{\partial z} & \frac{\partial f}{\partial w} \end{bmatrix} = \begin{bmatrix} 2 & 4 \\ 3 & 4 \end{bmatrix}$

~~$\frac{\partial f}{\partial A}$~~

### c) Proof one of the gradient identities involving trace.

#### Example

6. Consider the function  $f(x) = \text{Tr}(Ax)$ , where  $A \in \mathbb{R}^{2 \times 2}$  and  $x \in \mathbb{R}^{2 \times 2}$ . Determine the gradient of 'f' with respect to x, denoted as  $\nabla_x f(x)$ .

**Solution:**

$$A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \text{ and } x = \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{pmatrix}$$

$$Ax = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{pmatrix}$$

$$Ax = \begin{pmatrix} a_{11}x_{11} + a_{12}x_{21} & a_{11}x_{12} + a_{12}x_{22} \\ a_{21}x_{11} + a_{22}x_{21} & a_{21}x_{12} + a_{22}x_{22} \end{pmatrix}_{2 \times 2}$$

$$f(x) = \text{Tr}(Ax)$$

$$f(x) = (a_{11}x_{11} + a_{12}x_{21} + a_{21}x_{12} + a_{22}x_{22}) \rightarrow R$$

$$\nabla_x f(x) = \begin{pmatrix} \frac{\partial f(x)}{\partial x_{11}} & \frac{\partial f(x)}{\partial x_{12}} \\ \frac{\partial f(x)}{\partial x_{21}} & \frac{\partial f(x)}{\partial x_{22}} \end{pmatrix} = \begin{pmatrix} a_{11} & a_{21} \\ a_{12} & a_{22} \end{pmatrix} = A^T$$

$$\begin{aligned}
 f: & \quad \rightarrow R \\
 f: & \quad \rightarrow R \\
 f: & \quad \rightarrow R \\
 \frac{\partial f}{\partial x} = & \begin{bmatrix} f_{x_{11}} & f_{x_{12}} \\ f_{x_{21}} & f_{x_{22}} \end{bmatrix}
 \end{aligned}$$

$$f = \text{Tr}(A\hat{x})$$

$$\frac{df}{dx} = A^T$$



### 3. Gradient Descent

---

**Definition:** Gradient Descent is an optimization algorithm used to minimize a function by iteratively moving in the direction of steepest descent.

**Purpose:** The primary goal of Gradient Descent is to find the local minimum of a function. It's widely used in machine learning for optimizing models and finding optimal parameters.

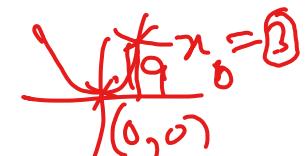
#### How it Works:

- At each iteration, the algorithm computes the gradient of the objective function with respect to the parameters.
- It then updates the parameters in the opposite direction of the gradient to reduce the value of the function.
- This process continues until convergence is achieved or a stopping criterion is met.

## Algorithm

$$f(w) = 3^2 \\ = 9$$

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



Data :  $x_0 \in \mathbb{R}^n$

$$f'(x) = 2x$$

$$x_{n+1} = x_n - \alpha \cdot f'(x_n)$$

Step 0: Set  $i = 0$

$$= 3 - (0.9) \cdot 2(3)$$

$$x_0 = 3$$

$$\alpha = (0, 1)$$

Step 1: If  $\nabla f(x_i) = 0$  then stop  $x_1 = 3 - \frac{1}{9} \cdot 6^3$

$$\alpha = 0.25$$

else, compute search direction  $h_i = -\nabla f(x_i)$ .

Step 2: Compute the size  $\lambda_i \in \arg \min_{\lambda \geq 0} f(x_i + \lambda \cdot h_i)$

Step 3: Set  $x_{i+1} = x_i + \lambda_i h_i$  go to step 1.

$$x_0 = 3 \\ x_1 = 1.5 \checkmark$$

## Derivation of Step size ( $\tau$ )

- Suppose we want to minimize the **quadratic function**:  $f(x) = \frac{1}{2}x^T Qx - b^T x$
- Let the iterative formula for minimum point be  $x_{k+1} = x_k - \tau_k \nabla f(x_k)$

Where  $\tau_k$  be the step size

- Let  $g(\tau) = f(x_{k+1}) = f(x_k - \tau_k \nabla f(x_k))$

$$g(\tau) = \frac{1}{2} [x_k - \tau_k \nabla f(x_k)]^T Q [x_k - \tau_k \nabla f(x_k)] - b^T [x_k - \tau_k \nabla f(x_k)]$$

$$g(\tau) = a\tau^2 + d\tau + c$$

$$a = \frac{1}{2} [\nabla f(x_k)]^T Q \nabla f(x_k)$$

$$d = [b^T - x_k^T Q] \nabla f(x_k) = -[\nabla f(x_k)]^T \nabla f(x_k)$$

Here  $g(\tau)$  is quadratic and concave function

- Condition for extrema of  $\frac{dg}{d\tau} = 0$  is

$$\frac{dg}{d\tau} = 2a\tau + d = 0 \Rightarrow \tau = -\frac{d}{2a}$$

$$\tau_k = \frac{\left[ \nabla f(x_k) \right]^T \nabla f(x_k)}{\left[ \nabla f(x_k) \right]^T Q \nabla f(x_k)}$$

*Therefore*

$$\tau = \frac{S^T S}{S^T Q S}$$

## 4. Gradient Descent Examples

### Example:

7. Find the minimum of  $f(x,y) = 3x^2+y^2$  by

- Computing the gradient of  $f$  and  $\tau$
- Solve with initial values  $x_0 = 1$  and  $y_0 = 3$

**Solution:**  $f(x,y) = 3x^2 + y^2$  with  $(x_0, y_0) = (1, 3)$

$$f(1,3) = 3 \cdot 1^2 + 3^2 = 12$$

**Step1:** We find  $\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix}$

$$s = \nabla f = \begin{bmatrix} 6x \\ 2y \end{bmatrix} \text{ and } H = \begin{bmatrix} 6 & 0 \\ 0 & 2 \end{bmatrix}$$

## Step2: Compute

$\tau_i = \frac{s^T s}{s^T H s}$  is the step size.

$$\begin{aligned}x_0 &= 3 \\y_0 &= 1\end{aligned}$$

$$\tau_i = \frac{\begin{bmatrix} 6x_i & 2y_i \end{bmatrix} \begin{bmatrix} 6x_i \\ 2y_i \end{bmatrix}}{\begin{bmatrix} 6x_i & 2y_i \end{bmatrix} \begin{bmatrix} 6 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 6x_i \\ 2y_i \end{bmatrix}}$$

$$\tau_i = \frac{36x_i^2 + 4y_i^2}{216x_i^2 + 8y_i^2}$$

$$\begin{aligned}\mathcal{L} &= \frac{(36)^2 + 4 \cdot 1^2}{216 \cdot 3^2 + 8 \cdot 1^2} \\&= \frac{1296 + 4}{1944 + 8} \\&= \frac{1300}{1952} \\&= 0.666\ldots\end{aligned}$$

## Step:3 Iterate the minimum point

$$x_{i+1} = x_i - \tau_i \nabla f_i$$

$$\begin{bmatrix} x_{i+1} \\ y_{i+1} \end{bmatrix} = \begin{bmatrix} x_i \\ y_i \end{bmatrix} - \left( \frac{36x_i^2 + 4y_i^2}{216x_i^2 + 8y_i^2} \right) \begin{bmatrix} 6x_i \\ 2y_i \end{bmatrix}$$

### 1<sup>st</sup> Iteration

$$i=0 \quad \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} x_0 \\ y_0 \end{bmatrix} - \left( \frac{36x_0^2 + 4y_0^2}{216x_0^2 + 8y_0^2} \right) \begin{bmatrix} 6x_0 \\ 2y_0 \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix} - \left( \frac{36.1^2 + 4.3^2}{216.1^2 + 8.3^2} \right) \begin{bmatrix} 6.1 \\ 2.3 \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix} - \left( \frac{72}{288} \right) \begin{bmatrix} 6 \\ 6 \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 1.5 \end{bmatrix}$$

$$f(-0.5, 1.5) = 3(-0.5)^2 + (1.5)^2 = 3$$

$$\frac{1}{288} = \frac{1}{4}$$

### 2<sup>nd</sup> Iteration

$$i=1 \quad \begin{bmatrix} x_2 \\ y_2 \end{bmatrix} = \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} - \left( \frac{36x_1^2 + 4y_1^2}{216x_1^2 + 8y_1^2} \right) \begin{bmatrix} 6x_1 \\ 2y_1 \end{bmatrix}$$

$$\begin{bmatrix} x_2 \\ y_2 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 1.5 \end{bmatrix} - \left( \frac{36(-0.5)^2 + 4(1.5)^2}{216(-0.5)^2 + 8(1.5)^2} \right) \begin{bmatrix} 6(-0.5) \\ 2(1.5) \end{bmatrix}$$

$$\begin{bmatrix} x_2 \\ y_2 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 1.5 \end{bmatrix} - \left( \frac{18}{72} \right) \begin{bmatrix} -3 \\ 3 \end{bmatrix} \quad \frac{1}{4} = 0.25$$

$$\begin{bmatrix} x_2 \\ y_2 \end{bmatrix} = \begin{bmatrix} 0.25 \\ 0.75 \end{bmatrix}$$

$$f(0.25, 0.75) = 3(0.25)^2 + (0.75)^2 = 0.75$$

**3<sup>rd</sup> Iteration**

$$i=2 \quad \begin{bmatrix} x_3 \\ y_3 \end{bmatrix} = \begin{bmatrix} x_2 \\ y_2 \end{bmatrix} - \left( \frac{36x_2^2 + 4y_2^2}{216x_2^2 + 8y_2^2} \right) \begin{bmatrix} 6x_2 \\ 2y_2 \end{bmatrix}$$

$$\begin{bmatrix} x_3 \\ y_3 \end{bmatrix} = \begin{bmatrix} 0.25 \\ 0.75 \end{bmatrix} - \left( \frac{36.(0.25)^2 + 4.(0.75)^2}{216.(0.25)^2 + 8.(0.75)^2} \right) \begin{bmatrix} 6.(0.25) \\ 2.(0.75) \end{bmatrix}$$

$$\begin{bmatrix} x_3 \\ y_3 \end{bmatrix} = \begin{bmatrix} 0.25 \\ 0.75 \end{bmatrix} - \left( \frac{4.5}{18} \right) \begin{bmatrix} 1.5 \\ 4.5 \end{bmatrix}$$

$$\begin{bmatrix} x_3 \\ y_3 \end{bmatrix} = \begin{bmatrix} -0.125 \\ -0.375 \end{bmatrix}$$

$$f(-0.125, -0.375) = 3(-0.125)^2 + (-0.375)^2 = 0.1875$$

**4<sup>th</sup> Iteration**

$$i=3 \quad \begin{bmatrix} x_4 \\ y_4 \end{bmatrix} = \begin{bmatrix} x_3 \\ y_3 \end{bmatrix} - \left( \frac{36x_3^2 + 4y_3^2}{216x_3^2 + 8y_3^2} \right) \begin{bmatrix} 6x_3 \\ 2y_3 \end{bmatrix}$$

$$\begin{bmatrix} x_4 \\ y_4 \end{bmatrix} = \begin{bmatrix} -0.125 \\ 0.375 \end{bmatrix} - \left( \frac{36.(-0.125)^2 + 4.(0.375)^2}{216.(-0.125)^2 + 8.(0.375)^2} \right) \begin{bmatrix} 6.(-0.125) \\ 2.(0.375) \end{bmatrix}$$

$$\begin{bmatrix} x_4 \\ y_4 \end{bmatrix} = \begin{bmatrix} -0.125 \\ 0.375 \end{bmatrix} - \left( \frac{1}{4} \right) \begin{bmatrix} -0.75 \\ 0.75 \end{bmatrix}$$

$$\begin{bmatrix} x_4 \\ y_4 \end{bmatrix} = \begin{bmatrix} 0.0625 \\ 0.1875 \end{bmatrix}$$

$$f(0.0625, 0.1875) = 3(0.0625)^2 + (0.1875)^2 = 0.046875$$

At the end of 10<sup>th</sup> iteration, we get  $\begin{bmatrix} x_{10} \\ y_{10} \end{bmatrix} = \begin{bmatrix} -0.001953 \\ 0.005859 \end{bmatrix}$  and the minimum value is

$$f(-0.001953, 0.005859) = 3(-0.001953)^2 + (0.005859)^2 = 0.000046$$

## Gradient Descent with Exponential Learning Rate Decay

### Example

8. Consider the one-dimensional cost function  $J(w) = (w - 4)^2$

The gradient of the cost function is given by  $\frac{dJ}{dw} = 2(w - 4)$

You are asked to minimize  $J(w)$  using **Gradient Descent with an exponentially decaying learning rate**.

The learning rate at iteration  $t$  is defined as  $\alpha_t = \alpha_n e^{-kt}$ , where initial learning rate  $\alpha_n = 0.5$ , decay constant  $k = 0.1$  and initial weight  $w_0 = 0$ . The gradient descent update rule is  $w_{t+1} = w_t - \alpha_t \cdot \frac{dJ}{dw}$ . Find

(a) Compute the learning rate  $\alpha_t$  for  $t=0,1$  and  $2$ .

(b) Perform three iterations of gradient descent and compute  $w_1, w_2$ , and  $w_3$ .

G8

G1

$$f(x) = (x - 4)^2$$

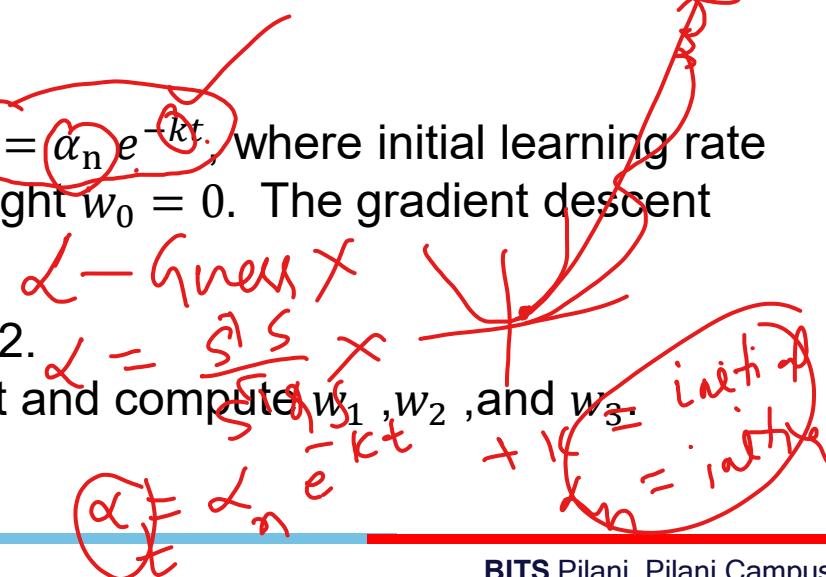
$$x_t = x_0 - \alpha_t f'(x_t)$$

$$\alpha = 0.1$$

$$\alpha = 0.09$$

$$\alpha = 0.0897$$

$$\alpha = 0.0$$



$$\alpha = \text{greedy} \quad \left| \quad \alpha = \frac{\nabla J S}{\nabla J + S} \right. \quad \left| \quad \alpha_t = f(t) e^{-kt} \right. \quad \left| \quad \alpha_t = \frac{\ln}{1+kt} \right. \quad \begin{array}{c} \text{innovate} \\ \text{achieve} \\ \text{lead} \end{array} \quad \begin{array}{c} \text{binary} \\ \text{holder} \end{array}$$

## Solution

(a)

Iteration $t$	Learning rate $\alpha_t = \alpha_n e^{-kt}$
0	$\alpha_0 = 0.5 e^{-0.1(0)} = 0.5000$
1	$\alpha_1 = 0.5 e^{-0.1(1)} \approx 0.4524$
2	$\alpha_2 = 0.5 e^{-0.1(2)} \approx 0.4094$

(b)

Iteration $t$	$\nabla J(w_t)$	$w_{t+1} = w_t - \alpha_t \cdot \frac{dJ}{dw}$
0	$\nabla J(w_0) = 2(0 - 4) = -8$	$w_1 = 0 - (0.5)(-8) = 4.0$
1	$\nabla J(w_1) = 2(4 - 4) = 0$	$w_2 = 4 - (0.4524)(0) = 4.0$
2	$\nabla J(w_2) = 0$	$w_3 = 4 - (0.4094)(0) = 4.0$

# Gradient Descent with Inverse Decay Rate.

**Example:** The same previous problem is solved using Inverse Decay

9. Given a Cost function:  $J(w) = (w - 4)^2$  and Gradient is  $\frac{dJ}{dw} = 2(w - 4)$ .

**Inverse (time-based) learning rate decay:**  $\alpha_t = \frac{\alpha_n}{1+kt}$ . Parameters:

initial learning rate  $\alpha_n = 0.5$ , Decay constant  $k = 0.1$  and initial weight  $w_0 = 0$ .

**Solution**

(a)

Iteration $t$	Learning rate $\alpha_t = \frac{\alpha_0}{1+kt}$
0	$\alpha_0 = \frac{0.5}{1 + (0.1)(0)} = 0.5000$
1	$\alpha_1 = \frac{0.5}{1 + (0.1)(1)} \approx 0.4545$
2	$\alpha_2 = \frac{0.5}{1 + (0.1)(2)} \approx 0.4167$

(b)

Iteration $t$	$\nabla J(w_t)$	$w_{t+1} = w_t - \alpha_t \cdot \frac{dJ}{dw}$
0	$\nabla J(w_0) = 2(0 - 4) = -8$	$w_1 = 0 - (0.5)(-8) = 4.0$
1	$\nabla J(w_1) = 2(4 - 4) = 0$	$w_2 = 4 - (0.4545)(0) = 4.0$
2	$\nabla J(w_2) = 0$	$w_3 = 4 - (0.4167)(0) = 4.0$

## Determining Step Size in Gradient Descent Using Binary Search

In machine learning optimisation, choosing an appropriate **step size (learning rate)** is crucial for fast and stable convergence. One way to determine the step size is by performing a **binary search (line search)** over a range of possible learning rates.

### Example:

10. Consider the cost function:  $J(w) = (w - 3)^2$ . The gradient of the cost function is  $\nabla J(w) = 2(w - 3)$ . You are currently at the parameter value  $w = 0$ . Assume that the learning rate  $\alpha$  is known to lie in the interval  $\alpha \in [0, 1]$ . The goal is to choose  $\alpha$  such that the **updated cost**  $J(w - \alpha \nabla J(w))$  is **minimized**.

- Write the expression for the updated weight  $w'$  in terms of  $\alpha$ .
- Write the cost function  $J(w')$  as a function of  $\alpha$ .
- Use **binary search** to determine the optimal step size  $\alpha$  up to **two decimal places**. Stop when the interval width is less than 0.05.
- Compute the updated weight using the optimal  $\alpha$ .

## Solution (a) Updated weight as a function of $\alpha$

At  $w = 0$ ,

$$w' = w - \alpha \nabla J(w)$$

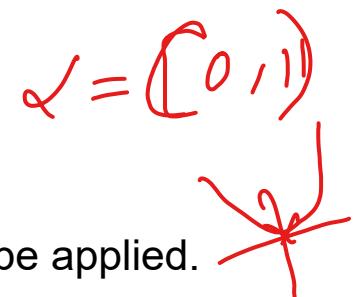
$$\nabla J(0) = 2(0 - 3) = -6$$

$$w' = 0 - \alpha(-6) = 6\alpha$$

(b) Cost as a function of  $\alpha$   $J(w') = (6\alpha - 3)^2$

(See next slide for why it is convex quadratic function )

This is a convex quadratic function in  $\alpha$ , so binary search can be applied.



### (c). Iteration 1:

Initial interval:  $a = 0, b = 1$



$$\alpha = \frac{0+1}{2}$$

$$\text{Midpoint: } m = \frac{0+1}{2} = 0.5$$

Derivative of  $J$  with respect to  $\alpha$ :  $J'(\alpha) = 72\alpha - 36$



$$\alpha = 0.5$$

$$\text{Evaluate at midpoint: } J'(0.5) = 72 \times 0.5 - 36 = 36 - 36 = 0$$

*constant*

Since the derivative is exactly zero, the stationary point (minimum) is found immediately.  
So optimal  $\alpha = 0.5$

(d) Updated Weight Using Optimal Step Size  $w' = 6 \times 0.5 = 3$ .

$$J(\alpha) = 36\alpha^2 - 36\alpha + 9$$

The coefficient of  $\alpha^2$  is  $36 > 0$ , so the parabola opens upwards.

For any quadratic function  $f(x) = ax^2 + bx + c$ :

- If  $a > 0$ , it is **convex** (also called strictly convex if  $a > 0$ ).
- If  $a < 0$ , it is concave.
- If  $a = 0$ , it is linear (both convex and concave).

Since the leading coefficient  $36 > 0$ ,  $J(\alpha)$  is **strictly convex** in  $\alpha$ .

A convex function has the property that any local minimum is also the global minimum, and the function has at most one minimum.

## Determining Step Size in Gradient Descent Using Golden-Section Search

Solve the same previous problem by the Golden-Section search to answer the following:

- (a) Write the expression for the updated weight  $w'$  in terms of  $\alpha$ .
- (b) Write the cost function  $J(w')$  as a function of  $\alpha$ .
- (c) Use **binary search** to determine the optimal step size  $\alpha$  up to **two decimal places**. Stop when the interval width is less than 0.05.
- (d) Compute the updated weight using the optimal  $\alpha$ .

## 5. Previous Year Problems

2

**Q1:** Consider the function  $f(x) = \cancel{ax^3} + bx^2 + cx + d$  where  $a > 0$ .

(a) Find the condition on  $a, b, c$  such that the function has two distinct critical points. Calculate the critical points in terms of  $a, b, c$ . Identify the nature of each critical point (i.e. maxima or minima). [3 Marks]

(b) Define a zone of attraction around each local minimum to be the region around it such that if gradient descent starts at any point in the region, it would end up at the given local minimum. Find the zone of attraction for each local minimum, if any, of the critical points. Justify your answer mathematically. [3 Marks]

$$f(x) = ax^3 + bx^2 + cx + d$$

### Q1a) Solution :

To find the critical points we take  $\frac{df}{dx} = 3ax^2 + 2bx + c = 0$ , to obtain two roots

$$x_1 = \frac{-b - \sqrt{b^2 - 3ac}}{3a}$$

$$\checkmark$$

$$\text{and } x_2 = \frac{-b + \sqrt{b^2 - 3ac}}{3a}.$$

In order for the two roots to be real and distinct, the quantity under the square-root sign needs to be strictly positive, i.e.

$$b^2 - 3ac > 0.$$

We see that the second derivative  $\frac{d^2f}{dx^2} = 6ax + 2b$ .

$$a x^2 + b x + c = 0$$

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

The second derivative is negative for the critical point  $x_1$  and positive for  $x_2$  which means  $x_1$  is a maxima and  $x_2$  is a minima.

$$f''(x_1) = \cancel{6ac} \left[ \frac{-b - \sqrt{b^2 - 3ac}}{3a} \right] + 2b$$

**Q1b) Solution :**  ~~$= -2\sqrt{b^2 - 3ac} + 2b = -2\sqrt{b^2 - 3ac}$~~

As identified in part (a), there is a single local minimum  $x_2 = \frac{-b + \sqrt{b^2 - 3ac}}{3a} < 0$ .

We can rewrite  $\frac{df}{dx} = 3a \left( x + \frac{b}{3a} \right)^2 + \left( c - \frac{b^2}{3a} \right)$ .

Solving for  $\frac{df}{dx} = 0$

we see that  $\frac{df}{dx} < 0$  when  $x > \frac{-b - \sqrt{b^2 - 3ac}}{3a}$  and  $x < \frac{-b + \sqrt{b^2 - 3ac}}{3a}$ .

When  $x > \frac{-b + \sqrt{b^2 - 3ac}}{3a}$ ,  $\frac{df}{dx} > 0$ .

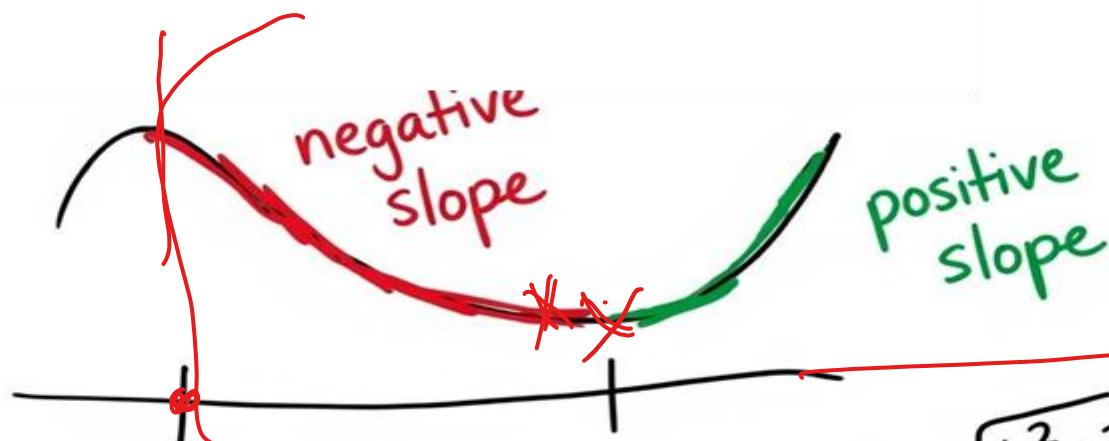
Gradient descent will take any point on the left of the local minimum  $x_2$  but greater than  $\frac{-b - \sqrt{b^2 - 3ac}}{3a}$  to  $x_2$  in a sufficiently large number of steps with a suitable step-size.

Similarly since the derivative is positive on the right of  $x_2$ , gradient descent will take any point on the right of  $x_2$  to  $x_2$  in a sufficiently large number of steps with a suitable step-size.

Thus the zone of attraction for the local minimum  $x_2 = \frac{-b + \sqrt{b^2 - 3ac}}{3a}$  is  $\frac{-b - \sqrt{b^2 - 3ac}}{3a} < x < \infty$ .

$$f''(x_2) = 2\sqrt{b^2 - 3ac} > 0.$$

i) a)



$$x_2 = \frac{-b + \sqrt{b^2 - 3ac}}{3a}$$

min

$$x_1 = \frac{-b - \sqrt{b^2 - 3ac}}{3a}$$

max

$$x_1 < x < -\infty$$

$$x =$$

$-b$

**Q2:**

$$f(x, y) = x^2 + \beta y^2$$

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

Consider a quadratic function  $f(x, y) = x^2 + \beta y^2$  where  $\beta \in \mathbb{R}$  is an unknown constant. Also assume that  $\beta > 0$ . Consider the problem of minimizing this function using gradient descent algorithm:

- (i) Derive a closed form expression (involving  $\beta$ ) for the optimal step-size  $\alpha$  for the first iteration of gradient descent if the initial point is

$$\begin{bmatrix} x_0 \\ y_0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\beta = 1$$

- (ii) If it is given to you that the optimal step-size  $\alpha = 0.5$ , derive the value of constant  $\beta$  using the formula derived in (i).

$$f(x, y) = x^2 + \beta y^2$$

## Q2)a) Solution:

The gradient of  $f(x, y)$  at  $(x_0, y_0)$  is as follows  $\nabla f(x_0, y_0) = \begin{bmatrix} 2x_0 \\ 2\beta y_0 \end{bmatrix} = \begin{bmatrix} 2 \\ 2\beta \end{bmatrix}$ .

$$\begin{aligned} f_x &= 2x \\ f_y &= 2\beta y \end{aligned}$$

The optimal step size is obtained as the solution of the following one-dimensional optimization

problem  $\underset{\alpha}{\operatorname{argmin}} f\left(\begin{bmatrix} x_0 \\ y_0 \end{bmatrix} - \alpha \nabla f(x_0, y_0)\right)$ .

$$f((x_0, y_0) - \alpha \nabla f(x_0, y_0))$$

After substituting  $x_0 = 1$  and  $y_0 = 1$ , this is equivalent to solving  $\underset{\alpha}{\operatorname{argmin}} f(1 - 2\alpha, 1 - 2\alpha\beta)$ .

In other words, we need to minimise the function  $g(\alpha) = (1 - 2\alpha)^2 + \beta(1 - 2\alpha\beta)^2$ .

$$f((1, 1) - \alpha(2, 2\beta))$$

To find the minimum of  $g(\alpha)$  with respect to  $\alpha$ , we compute the derivative  $g'(\alpha)$  and set it to zero:

$$g'(\alpha) = -4(1 - 2\alpha) + 2\beta(1 - 2\alpha\beta)(-2\beta) = 0.$$

$$f(1 - 2\alpha, 1 - 2\alpha\beta)$$

Expanding:  $g'(\alpha) = -4 + 8\alpha - 4\beta^2 + 8\alpha\beta^3 = 8\alpha(1 + \beta^3) - 4(1 + \beta^2) = 0$ .

$$2\alpha(1 + \beta^3) - (1 + \beta^2) = 0$$

$$\alpha = \frac{1 + \beta^2}{2(1 + \beta^3)}$$

Solving for  $\alpha$ :  $8\alpha(1 + \beta^3) = 4(1 + \beta^2) \Rightarrow \alpha = \frac{1 + \beta^2}{2(1 + \beta^3)}$ .

Hence the closed form expression for the optimal step-size  $\alpha$  is  $\alpha = \frac{1 + \beta^2}{2(1 + \beta^3)}$ .

## Q2)b) Solution:

$$\alpha \approx \frac{1}{2}$$

If it is given that  $\alpha = 0.5$ , substituting in the previous expression we get the equation

$$0.5 = \frac{1 + \beta^2}{2 + 2\beta^3}.$$

Rearranging we get

$$1 + \beta^2 = 0.5(2 + 2\beta^3) = 1 + \beta^3.$$

$$\beta^2 = \beta^3 \Rightarrow \beta^3 - \beta^2 = 0 \Rightarrow \beta^2(\beta - 1) = 0.$$

Hence potential values of  $\beta$  are 0 or 1.

It is given in the question that  $\beta > 0$  and (implicitly)  $\beta \neq 0$  since it is an unknown positive constant.

Hence the final answer is  $\beta = 1$ .

---

### Q3:

A data science intern arrived at a loss function given by

$$f(x) = 3x^4 - 20x^3 + 36x^2 + 10.$$

- i) Help the intern to find the stationary points and classify and hence the global minima.  
[3 marks]
- ii) Suggest the intern whether  $x = 0.5$  or  $x = 3.5$  is a better initial condition to find global minima using simple gradient descent method with reasons. [1 mark]

### **Q3)i) Solution:**

Given  $f(x) = 3x^4 - 20x^3 + 36x^2 + 10$ . So, we have

$$f'(x) = 12x^3 - 60x^2 + 72x = 0$$

$$\Rightarrow 12x(x-2)(x-3) = 0 \Rightarrow x = 0, 2, 3.$$

$$f''(x) = 36x^2 - 120x + 72$$

$$\Rightarrow f''(0) = 72 > 0 \Rightarrow \text{Point of local minima}$$

$$\Rightarrow f''(2) = -24 < 0 \Rightarrow \text{Point of local maxima}$$

$$\Rightarrow f''(3) = 36 > 0 \Rightarrow \text{Point of local minima}$$

Now evaluate the function values:

$$f(0) = 10, \quad f(3) = 3(81) - 20(27) + 36(9) + 10 = 243 - 540 + 324 + 10 = 37.$$

Since  $f(0) = 10 < f(3) = 37$ , and there is a local maximum at  $x = 2$ , the global minimum is at  $x = 0$  with  $f(0) = 10$ .

Therefore 10 is the answer (the value of the global minimum).

### **Q3)ii) Solution:**

Clearly, 0.5 is closer to 0 (global minima) and 3.5 is closer to 3 (local minima). Starting near the global minimum helps gradient descent converge to it faster and avoids getting stuck in the local minimum at  $x=3$ .

So,  $x = 0.5$  is a better initial condition.



---

# THANK YOU