

### Question 3

Consider the following function:

$$f(x_1, x_2, x_3) = -2x_1^2 + 4x_1 + x_2x_1 - x_2 + x_3x_2 - x_3^2$$

a) Compute  $\nabla f$

#### Solution

□ Recall, for a function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $x \mapsto f(x)$ ,  $x \in \mathbb{R}^n$  of  $n$  variables

$x_1, \dots, x_n$  we define the Partial derivatives as:

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

$\vdots$

$$\frac{\partial f}{\partial x_n} = \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_{n-1}, x_n+h) - f(x)}{h}$$

and collect them in the row vector

$$\nabla f = \text{grad } f = \frac{df}{dx} = \left[ \frac{\partial f(x)}{\partial x_1} \quad \frac{\partial f(x)}{\partial x_2} \quad \dots \quad \frac{\partial f(x)}{\partial x_n} \right] \in \mathbb{R}^{1 \times n} \dots \textcircled{A}$$

where  $n$  is the number of variables and 1 is the dimension of the image/range/codomain of  $f$ . Here we defined the Column vector  $x = [x_1, \dots, x_n]^T \in \mathbb{R}^n$ .

The row vector in  $\textcircled{A}$  is called the gradient of  $f$  or the jacobian

\* Therefore, given  $f(x_1, x_2, x_3) = -2x_1^2 + 4x_1 + x_2x_1 - x_2 + x_3x_2 - x_3^2$

$$\nabla f = \left[ \frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \frac{\partial f}{\partial x_3} \right] \in \mathbb{R}^{1 \times 3}$$

$$\square \frac{\partial f}{\partial x_1} = -4x_1 + 4 + x_2$$

$$\square \frac{\partial f}{\partial x_2} = x_1 - 1 + x_3$$

$$\square \frac{\partial f}{\partial x_3} = x_2 - 2x_3$$

$$\therefore \nabla f = [-4x_1 + 4 + x_2, \quad x_1 - 1 + x_3, \quad x_2 - 2x_3]$$

b) Find the Solutions for  $\nabla f(x) = 0$

Solution

$$\text{Recall } \nabla f = [-4x_1 + 4 + x_2, x_1 - 1 + x_3, x_2 - 2x_3]$$

$$\text{At } \nabla f = 0;$$

$$[-4x_1 + 4 + x_2, x_1 - 1 + x_3, x_2 - 2x_3] = [0, 0, 0]$$

$$\text{i.e. } -4x_1 + 4 + x_2 = 0 \Rightarrow 4x_1 - x_2 = 4 \dots (1)$$

$$x_1 - 1 + x_3 = 0 \Rightarrow x_1 + x_3 = 1 \dots (2)$$

$$x_2 - 2x_3 = 0 \Rightarrow x_2 = 2x_3 \dots (3)$$

$$\text{Recall from (2): } x_1 = 1 - x_3 \dots (4)$$

Sub for (1) and (3) in eqn... (1)

$$\Rightarrow 4(1 - x_3) - 2x_3 = 4$$

$$\text{i.e. } 6x_3 = 0 \Rightarrow x_3 = 0,$$

$$\text{Sub for } x_3 \text{ in (4): } x_1 = 1 - 0 = 1$$

$$\text{Sub for } x_3 \text{ in (3): } x_2 = 2x_3 = 2(0) = 0$$

$$\text{Therefore; } (x_1, x_2, x_3) = (1, 0, 0).$$

1c

Starts at a value  $\begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix}$ , and consider learning rates  $\{1, \frac{1}{2}\}$

Execute two steps of the gradient descent algorithm for each learning rate separately

Solution

\* Gradient Descent Algorithm:

Repeat until Convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

(for  $j=1$  and  $j=0$ )

}

□ Alpha is called the Learning rate - A tuning parameter in the Optimization Process. It decides the length of the steps.

Hence, we have:

$$x_i := x_i - \alpha \frac{\partial}{\partial x_i} f(x_1, x_2, x_3)$$

for  $i = 1, 2, 3$ .

where  $\alpha = \text{learning rate}$

$\Rightarrow$  Starting from Point  $[1, 2, 0]$ , update the values of  $x_1, x_2$  and  $x_3$  according to the gradient rule.

$\Rightarrow$  STEP 1: Considering learning rate ( $\alpha$ ) of 1 [Recall the learning rates are  $\{1, \frac{1}{2}\}$ ]

$$* x_1 = x_1 - \alpha \frac{\partial f}{\partial x_1} = 1 - 1[-4(1) + 4 + 2] = -1$$

$$* x_2 = x_2 - \alpha \frac{\partial f}{\partial x_2} = 2 - 1[1 - 1 + 0] = 2$$

$$* x_3 = x_3 - \alpha \frac{\partial f}{\partial x_3} = 0 - 1[2 - 2(0)] = -2$$

The new Point is  $[-1, 2, -2]$

$\Rightarrow$  Step 2

$$* x_1 = x_1 - \alpha \frac{\partial f}{\partial x_1} = -1 - 1[-4(-1) + 4 + 2] = -11$$

$$* x_2 = x_2 - \alpha \frac{\partial f}{\partial x_2} = 2 - 1[-1 - 1 + (-2)] = 6$$

$$* x_3 = x_3 - \alpha \frac{\partial f}{\partial x_3} = -2 - 1[2 - 2(-1)] = -8$$

The final Position of learning rate 1 is:  $(-11, 6, -8)$ .

Using learning rate  $\frac{1}{2}$  [Starting Point  $[1, 2, 0]$ ]

STEP 2:

$$* x_1 = x_1 - \alpha \frac{\partial f}{\partial x_1} = 1 - \frac{1}{2}[2] = 0$$

$$* x_2 = x_2 - \alpha \frac{\partial f}{\partial x_2} = 2 - \frac{1}{2}[0] = 2$$

$$* x_3 = x_3 - \alpha \frac{\partial f}{\partial x_3} = 0 - \frac{1}{2}[0] = 0$$

New Point is  $[0, 2, 0]$

$\Rightarrow$  Step 2

$$* x_1 = x_1 - \alpha \frac{\partial f}{\partial x_1} = 0 - \frac{1}{2}[6] = -3$$

$$* x_2 = x_2 - \alpha \frac{\partial f}{\partial x_2} = 2 - \frac{1}{2}[-1] = \frac{5}{2}$$

$$* x_3 = x_3 - \alpha \frac{\partial f}{\partial x_3} = 0 - \frac{1}{2}[2] = -1$$

The final Position of learning rate 2 is:  $(-3, \frac{5}{2}, -1)$ .

ii) Discuss the Outcome and Compare the Outcome with Item (b)

whether they are Distances from Points in (b) or getting closer?



## Solution

Result, the learning rate is a hyperparameter that controls the step size of each iteration while moving toward a minimum of a loss function.

A smaller learning rate may lead to slower convergence but can avoid overshooting the minimum, while a larger learning rate may lead to faster convergence but may overshoot the minimum.

\* Now let's consider the given final positions of the gradient descent algorithm with two different learning rates:

Learning rate 1:  $[-11, 6, -8]$

Learning rate 2:  $[-3, \frac{5}{2}, -1]$

⇒ For the learning rate 1, the final position is quite far from the point  $[1, 0, 0]$  indicating that the algorithm overshoot the minimum. It may have taken larger steps at each iteration, causing it to miss the minimum.

$$d = \sqrt{(1 - (-11))^2 + (0 - 6)^2 + (0 - (-8))^2} = 15.62 \quad \left[ \begin{array}{l} \text{distance between two points} \\ [1, 0, 0] \text{ and } [-11, 6, -8] \end{array} \right]$$

⇒ For the learning rate  $\frac{1}{2}$ , the final position is closer to the point  $[1, 0, 0]$ . This result suggests that the algorithm may have taken smaller steps at each iteration, leading to slower convergence but a more accurate result.

$$d = \sqrt{(1 - (-3))^2 + (0 - \frac{5}{2})^2 + (0 - (-1))^2} = 4.82 \quad \left[ \begin{array}{l} \text{distance between two points} \\ [1, 0, 0] \text{ and } [-3, \frac{5}{2}, -1] \end{array} \right]$$

⇒ Overall, we can say the algorithm with a learning rate of  $\frac{1}{2}$  is better than the algorithm with a learning rate of 1, as it gets closer to the desired point.

⇒ Note: Choosing an appropriate learning rate and a good starting point are crucial for the success of the gradient descent algorithm.