CMR INSTITUTE OF TECHNOLOGY			USN								CELEBRACION DE LA CONTRACTION		
Sub:	Internal Assessr	Aug 2024 Optimization Techniques Code:							BCS405C				
Date:	02/08/2024 Duration: 90 mins Max Marks: 50 Sem: IV Branch:							nch:	CSDS/CSML				
Answer any five of the following.								Marks	OB CO	E RB T			
1	¹ Use Steepest Descent method for $f(x, y) = x - y + 2x^2 + 2xy + y^2$ starting from (0,0)						0)	10	CO3	L1,L2			
	Use Newton Raphson method to find the smallest root and the second smallest positive roots of the equation $\tan x = x$ correct to 4 decimal places.							ots of	10	CO4	L2,L3		
	Define Hessian matrix. Using Hessian matrix classify the relative extreme for the function $f(x, y) = \frac{1}{3}x^3 + xy^2 - 8xy + 3$							2+8	CO3	L3			
4	Find the linear regression coefficients using Gradient Descent method.							10	CO4	L3			
5	Explain in brief 1. Adagrad optimization strategy 2. RMS prop								5+5	CO5	L2		
	What is the difference between convex optimization and non-convex optimization. Write short notes on ADAM.							5+5	CO5	L2,L3			

CMR INSTITUTE OF TECHNOLOGY			USN						CORDING WITH ALL		
Sub:	Internal Assess	Aug 2024 Optimization Techniques Code:						BCS405C			
Date:	02/08/2024 Duration: 90 mins Max Marks: 50 Sem: IV Branch:						Branch:	CSDS/CSML			
Answ	Answer any five of the following.								Marks	OE CO	BE RBT
1	Use Steepest Descent method for $f(x, y) = x - y + 2x^2 + 2xy + y^2$ starting from (0,0)							10	CO3	L1,L2	
	Use Newton Raphson method to find the smallest root and the second smallest positive roots of the equation $\tan x = x$ correct to 4 decimal places.						10	CO4	L2,L3		
	Define Hessian matrix. Using Hessian matrix classify the relative extreme for the function $f(x, y) = \frac{1}{3}x^3 + xy^2 - 8xy + 3$					2+8	CO3	L3			
4	Find the linear regression coefficients using Gradient Descent method.						10	CO4	L3		
5	Explain in brief 1. Adagrad optimization strategy 2. RMS prop						5+5	CO5	L2		
	What is the difference between convex optimization and non-convex optimization. Write short notes on ADAM.						5+5	CO5	L2,L3		

Working Procedure to Fit a Linear Regression Line to Data Using the Gradient Descent Method

(1) The following steps are used to fit a linear regression line $\hat{y} = a + bx$

for the given data:

Step 1: Define the Model and Cost Function

The linear regression model you want to fit is given by:

$$\hat{y} = a + bx$$

where \hat{y} is the predicted value, *a* and *b* are the parameters (weights) of the model and *x* is the input feature.

The cost function (error function) for linear regression is the Mean Squared Error (MSE):

$$MSE = \frac{1}{n} \times \sum_{i=1}^{n} (y_i - \hat{y})^2$$
$$MSE = \frac{1}{n} \times \sum_{i=1}^{n} (y_i - (a + bx_i))^2$$

Where *n* is the number of data points.

Step 2: Initialize Weights and Hyperparameters

Initialize the weights a and b to some arbitrary values. Start with a = 0 and b = 0. Also need to set hyperparameters:

- Learning rate (α): A small positive value that controls the step size in each iteration.
- Number of iterations: The number of times we update the weights.

Professor, Department of Artificial Intelligence & Data Science, Don Bosco Institute of Technology, Bangalore. .ب

Step 3: Gradient Descent

For each iteration, calculate the gradients of the cost function with respect to the weights (a and b) and update the weights accordingly.

The gradients are given by:

$$\Delta a = -\frac{2}{n} \sum_{i=1}^{n} (y_i - (a + bx_i))$$

$$\Delta b = -\frac{2}{n} \sum_{i=1}^{n} x_i \times (y_i - (a + bx_i))$$

Update the weights using the gradients:

$$a_{new} = a - \alpha \times \Delta a$$

 $b_{new} = b - \alpha \times \Delta b$

Repeat this process for the specified number of iterations.

Step 4: Predict

After training, use the final values of *a* and *b* to make predictions:

$$\hat{\mathbf{y}} = a + bx$$

Step 5: Evaluate and visualize

Evaluate the quality of the linear regression model by calculating the final MSE on your training data:

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

Also, visualize the linear regression line by plotting it alongside the data points.

Professor, Department of Artificial Intelligence & Data Science, Don Bosco Institute of Technology, Bangalore. .ب

Step 6: Iterate as Needed

We need to adjust the learning rate and the number of iterations to find the best-fitting line. If the cost is not converging or fluctuating, you may need to modify the hyperparameters.

This process allows you to iteratively update the weights to minimize the cost function, resulting in a linear regression line that best fits the given data.

Examples

1. We have recorded the weekly average price of a stock over 6 consecutive days. Y shows the weekly average price of the stock and x shows the number of the days. Try to fit the best possible function ' f ' to establish the relationship between the number of the day and conversion rate. (Applying Gradient descent) where f(x) = y = a + b * x.

X	1	2	3	4	5	6
Y	10	14	18	22	25	33

The initial values of a & b are a = 4.9 & b = 4.401. The learning rate is mentioned as .05. The error rate of a & b should be less than .01. Plot the predicted and actual data in a graph.

Solution: Given data: $X = x_i = 1, 2, 3, 4, 5, 6$ $Y = y_i = 10, 14, 18, 22, 25, 33$ n = 6Initialization: a = 4.9 and b = 4.401Learning rate $\alpha = 0.05$ Maximum allowable error for a and b = 0.01The goal is to minimize the MSE (mean squared error) defined as

Professor, Department of Artificial Intelligence & Data Science, Don Bosco Institute of Technology, Bangalore.

Batch Gradient Descent	Stochastic Gradient Descent
No random shuffling of points are required.	The data sample should be in a random order, and this is why we want to shuffle the training set for every epoch.
Can't escape shallow local minima easily.	SGD can escape shallow local minima more easily.
Convergence is slow.	Reaches the convergence much faster.
It updates the model parameters only after processing the entire training set.	It updates the parameters after each individual data point.
The learning rate is fixed and cannot be changed during training.	The learning rate can be adjusted dynamically.
It typically converges to the global minimum for convex loss functions.	It may converge to a local minimum or saddle point.
It may suffer from overfitting if the model is too complex for the dataset.	It can help reduce overfitting by updating the model parameters more frequently.

Q9b What is the difference between convex optimization and non convex optimization?

Convex optimization and non-convex optimization are both optimization problems, but they differ in the number of optimal solutions they can have:

• Convex optimization

In convex optimization, there can only be one globally optimal solution, or it may be possible to prove that there is no feasible solution. Convex optimization is easier and more reliable because convex functions have a unique global minimum. Convex problems can also be solved efficiently, even when they are very large. Examples of convex optimization problems include multi-period processor speed scheduling, minimum time optimal control, and grasp force optimization.

Non-convex optimization

In non-convex optimization, the objective or some of the constraints are non-convex, which can lead to multiple feasible regions and multiple locally optimal points within each region. This can make optimization more challenging. Non-convex optimization can still be a good choice if the optimization scheme doesn't get stuck in a local minimum. It can also be used to implement more accurate state dynamics. However, even simple-looking non-convex optimization problems with only ten variables can be very challenging, and problems with hundreds of variables can be intractable.

called bias correction, we obtain the corrected first and second impulses respectively.

These correction 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.

9a)AdaGrad optimization strategy

Another optimization strategy I would like to introduce is called 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 , which is **the sum of squared gradients** up to the time .

SGD mit Impuls

AdaGrad

 $v_{t+1} \leftarrow \rho v_t + \nabla_{\theta} \mathcal{L}(\theta) \qquad g_0 = 0$ $\theta_j \leftarrow \theta_j - \epsilon v_{t+1} \qquad g_{t+1} \leftarrow g_t + \nabla_{\theta} \mathcal{L}(\theta)^2$ $\theta_j \leftarrow \theta_j - \epsilon \frac{\nabla_{\theta} \mathcal{L}}{\sqrt{g_{t+1}} + 1e^{-5}}$

When we optimize a weights θ_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**.

Disadvantages: there is a problem with this optimization algorithm.

If the training takes too long. Over time, this term the sum of squared gradients 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.

9a) &10c)RMSProp

There is a slight modification of AdaGrad called "RMSProp". This modification is intended to solve the previously described problem that can occur with AdaGrad. In RMSProp, the running sum of squared gradients g_{t+1} is maintained. However, instead of allowing this sum to increase continuously over the training period, we allow the sum to decrease.

 $\begin{array}{ll} \mbox{AdaGrad} & \mbox{RMS Prop} \\ g_0 = 0 & g_0 = 0, \alpha \simeq 0.9 \\ g_{t+1} \leftarrow g_t + \nabla_{\theta} \mathcal{L}(\theta)^2 & g_{t+1} \leftarrow \alpha \cdot g_t + (1-\alpha) \nabla_{\theta} \mathcal{L}(\theta)^2 \\ \theta_j \leftarrow \theta_j - \epsilon \frac{\nabla_{\theta} \mathcal{L}}{\sqrt{g_{t+1}} + 1e^{-5}} & \theta_j \leftarrow \theta_j - \epsilon \frac{\nabla_{\theta} \mathcal{L}}{\sqrt{g_{t+1}} + 1e^{-5}} \end{array}$

For RMSProp, the sum of squared gradients is multiplied by a decay rate α and the current gradient – weighted by (1- α) – is added. The update step in the case of RMSProp looks the same as in AdaGrad. Here we divide the current gradient by the sum of the squared gradients to get the nice property

of speeding up the updating of the weights along one dimension and slowing down the motion along the other.

Although SGD with momentum is able to find the global minimum faster, this algorithm takes a much longer path that could be dangerous. This is because a longer path means **more potential saddle points and local minima** of the loss function that could lie along that path. RMSProp, on the other hand, goes straight to the global minimum of the loss function without taking a detour.

1. Handling Non-stationary Objectives:

• **RMSProp** is particularly well-suited for non-stationary objectives (where the data distribution changes over time), as it can adjust more dynamically to the changes compared to Adagrad.

2. Empirical Performance:

 In practice, RMSProp often performs better than Adagrad on a variety of machine learning tasks. It tends to converge faster and reach better solutions, especially when dealing with deep learning models.

Overall, RMSProp is generally preferred for its ability to maintain a more stable and effective learning rate throughout training, leading to better performance on many complex tasks.

9c) Describle the saddle point problem in machine learning.

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:

10a)Stochastic Gradient Descent with Momentum

The first of the four algorithms I would like to introduce is called "Stochastic Gradient Descent with Momentum":

SGDSGD mit Impuls
$$\theta_j \leftarrow \theta_j - \epsilon \nabla_{\theta_j} \mathcal{L}(\theta)$$
 $v_{t+1} \leftarrow \rho v_t + \nabla_{\theta} \mathcal{L}(\theta)$ $\theta_j \leftarrow \theta_j - \epsilon v_{t+1}$

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**, 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 p.

The parameter p 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.

For a physical representation of stochastic gradient descent with momentum, 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.

Both algorithms try to reach the global minimum of the loss function, which is in a 3D space. Momentum term results in the individual gradients having less variance and thus less zig-zagging.

10a)ii)ADAM

We take the best of Adagrad and RMS prop and combine these ideas into a single algorithm called as 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.

$$\begin{split} m_0 &= 0, v_0 = 0 \\ m_{t+1} \leftarrow \beta_1 m_t + (1 - \beta_1) \nabla_{\theta} \mathcal{L}(\theta) & \text{Impuls} \\ v_{t+1} \leftarrow \beta_2 v_t + (1 - \beta_2) \nabla_{\theta} \mathcal{L}(\theta)^2 & \text{RMS Prop} \\ \theta_j \leftarrow \theta_j - \frac{\epsilon}{\sqrt{v_{t+1}} + 1e^{-5}} m_{t+1} & \text{RMS Prop + Impuls} \end{split}$$

The first expression looks a bit like SGD with momentum. In this case, the term m_t would be the velocity and the term β_1 would be the friction term. In the case of ADAM, we refer to m_t as the "first momentum." On the other hand, β_1 is just a hyperparameter. However, the difference with SGD with momentum is the factor $1 - \beta_1$ multiplied by the current gradient.

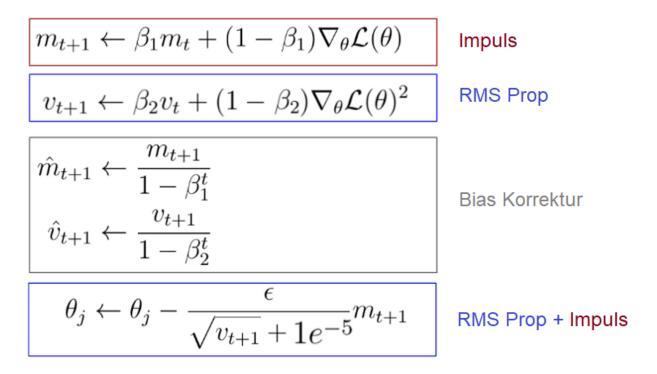
The second expression **can be considered as RMSProp**, where we keep the running sum of squared gradients. Also in this case, there is the factor $1-\beta_2$, which is multiplied by the squared gradient.

The term v_t in the equation is called the "second momentum" and is also just a hyperparameter. The final update equation can be viewed as **a combination** of RMSProp and SGD with momentum.

Disadvantages

At the very first time step t=0, the first and second pulse terms m0 and v0 are set to zero. After the first update of the second momentum v1, this term is still very close to zero. When we update the weight parameters in the last expression in GL. 5, we divide by a very small second momentum term v1. This leads to a very large first update step.

To address the problem of large update steps happening at the beginning of training, ADAM includes a correction clause:



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 respectively.

These correction 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.

9a)AdaGrad optimization strategy

Another optimization strategy I would like to introduce is called 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 , which is **the sum of squared gradients** up to the time .

SGD mit Impuls

AdaGrad

 $v_{t+1} \leftarrow \rho v_t + \nabla_{\theta} \mathcal{L}(\theta) \qquad g_0 = 0$ $\theta_j \leftarrow \theta_j - \epsilon v_{t+1} \qquad g_{t+1} \leftarrow g_t + \nabla_{\theta} \mathcal{L}(\theta)^2$ $\theta_j \leftarrow \theta_j - \epsilon \frac{\nabla_{\theta} \mathcal{L}}{\sqrt{g_{t+1}} + 1e^{-5}}$

When we optimize a weights θ_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.

(45 a) A local max/min occurs at x=c when f(c)>f(x) for >c values around A global max/min occurs at x = c if $f(cc) \leq f(c)$ for all values of x in the It is possible to have several global max/min if the function reaches its feak value at more than one faint. domain. 956) The Hessian matrix is a sequence materix that represents the 2nd order partial derivatives of a scalar valued function. It provides information about the local Intervature of the function and is used in curvature of the function of functions to oftenesation and analysis of functions to determine and in determine concernity or convexity. For a function f: R ->R that is twice tor a function $f: \mathcal{K} \longrightarrow \mathcal{K}$ that is twice continuously differentiable, the Hessian matrix is defined as 3f(3x; 3x) $H(f) = \begin{bmatrix} 3f(3x; 3x) \\ 3f($

The the relative extraction of the function

$$f(x, y) = \frac{1}{3}x^3 + xy^2 - 8xy + 3$$

$$f_x = x^2 + y^2 - 8y \quad f_y = 2xy - 8x$$

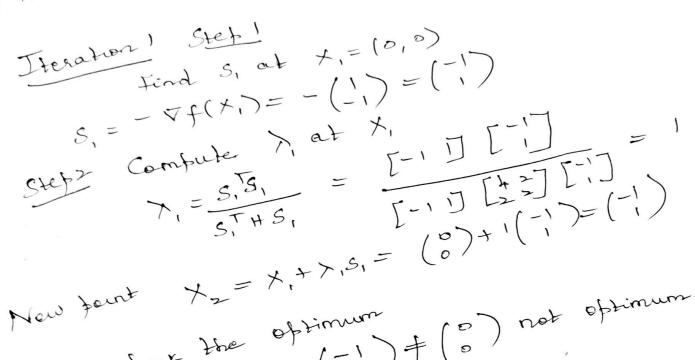
$$f_x = 0 \quad f_y = 0 \quad 2xy - 8x = 0 - 0$$

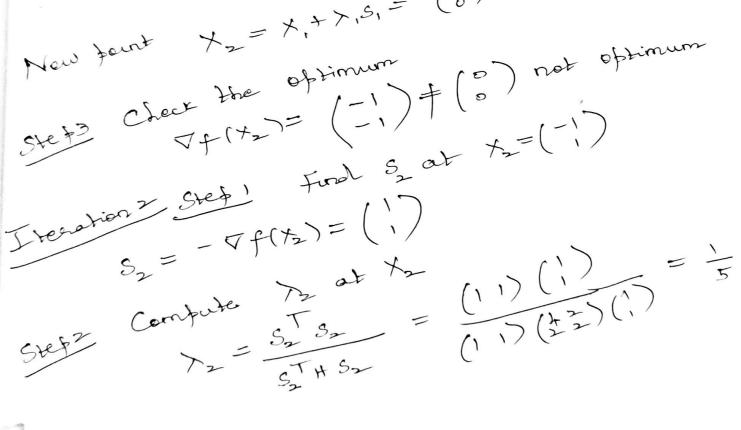
$$x^2 + y^2 - 8y = 0 - 0 \quad 2x(y - k) = 0$$

$$= x^2 + y^2 - 8y = 0 - 0 \quad 2x(y - k) = 0$$

$$= x^2 + y^2 - 8y = 0 - 0 \quad 2x(y - k) = 0$$

$$\begin{array}{l} \underbrace{344} (1) & \underline{x} = 0 \\ (1) & \underline{y}^2 - 8 & \underline{y} = 0 \Rightarrow y & (\underline{y} - 8) = 0 \Rightarrow y = 0, 8 \\ \hline (1) & \underline{y}^2 - 8 & \underline{y} = 0 \Rightarrow x = 1 & \underline{h} \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 & \underline{y} = 1 & \underline{y} = 0 \Rightarrow x = 1 & \underline{h} \\ \hline (1) & \underline{x} = 2 & \underline{y} = 1 & \underline{y} = 1 & \underline{y} = 0 \Rightarrow x = 1 & \underline{h} \\ \hline (1) & \underline{x} = 2 & \underline{y} = 1 & \underline{y} = 1 & \underline{y} = 0 \Rightarrow x = 1 & \underline{h} \\ \hline (1) & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y} = 1 & \underline{y} = 1 \\ \hline (1) & \underline{y}$$





Hence the new faint

$$X_{5} = X_{2} + \lambda_{2} S_{2} = \begin{pmatrix} -1 \\ -1 \end{pmatrix} + \frac{1}{5} \begin{pmatrix} 1 \\ -1 \end{pmatrix} = \begin{pmatrix} -0.8 \\ -0.2 \end{pmatrix}$$

$$G(e) = X_{2} + \lambda_{2} S_{2} = \begin{pmatrix} -0.2 \\ -0.2 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \end{pmatrix} \text{ nat optimum}$$

$$\nabla f(X_{5}) = \begin{pmatrix} +0.2 \\ -0.2 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \end{pmatrix} \text{ nat optimum}$$

$$X_{5} \text{ is net optimum}$$

$$X_{5} \text{ is net optimum}$$

$$X_{5} = -\nabla f(X_{5}) = -\begin{pmatrix} 0.2 \\ -0.2 \end{pmatrix} = \begin{pmatrix} -0.2 \\ 0.2 \end{pmatrix}$$

$$S_{3} = -\nabla f(X_{5}) = -\begin{pmatrix} 0.2 \\ -0.2 \end{pmatrix} = \begin{pmatrix} -0.2 \\ 0.2 \end{pmatrix}$$

$$S(e) = \frac{5}{3} \frac{5_{3}}{3} = \begin{pmatrix} -0.2 \\ -0.2 \end{pmatrix} = \begin{pmatrix} -0.2 \\ 0.2 \end{pmatrix}$$

$$(-0.2 \ 0.2 \end{pmatrix} \begin{pmatrix} -0.2 \\ 0.2 \end{pmatrix}$$
Hence the new fourt is
$$= 1$$

$$Hence \text{ the new fourt is } = \begin{pmatrix} -0.8 \\ -0.2 \end{pmatrix} + 1 \begin{pmatrix} -0.2 \\ 0.2 \end{pmatrix}$$

$$= (1.2)^{2} \begin{pmatrix} -0.2 \\ 0.2 \end{pmatrix}$$

Hence the new point is
$$X_5 = X_4 + X_5 S_4$$

$$= \begin{pmatrix} -1 \\ 1.4 \end{pmatrix} + \frac{1}{5} \begin{pmatrix} 0.2 \\ 0.2 \end{pmatrix} = \begin{pmatrix} -0.96 \\ 1.4 \end{pmatrix}$$
Stef 3 $\nabla f(X_5) = \begin{pmatrix} 1+4(-0.96)+2(1.44) \\ -1+2(-0.96)+2(1.44) \end{pmatrix} \land \begin{pmatrix} 0 \\ 0 \end{pmatrix}$
 $X_5 = 4$ or optimum

÷

$$f(\infty) = \infty - \tan x \quad \text{frag a simple}$$

soot near 4.5

$$f'(\infty) = 1 - \sec^2 x$$

$$f'(\infty) = -\frac{f(\infty)}{2}$$

$$x_0 = h.5 \quad x_1 = x_0 - \frac{f(\infty)}{f'(\infty)}$$

$$f(\mathbf{x}, 5)$$

$$x_{1} = 4.5 - \frac{1}{f'(4.5)}$$

$$f'(4.5) = -21.5048$$

$$f(4.5) = -21.5048$$

$$x_{1} = 4.5 - \frac{-0.1373}{-21.5048} = 4.4936$$

$$\begin{aligned} x_{1} = x_{0} - \frac{f(x_{0})}{f'(x_{0})} \\ = 0 \cdot 1 - \frac{f(0 \cdot 1)}{f'(0 \cdot 1)} = 0 \cdot 1 - \frac{(0 \cdot 10033 - 0 \cdot 1)}{gec^{2}(0 \cdot 1) - 1} \\ & 0 \cdot 1 - \frac{(0 \cdot 10033) - (0 \cdot 1)}{(0 \cdot 99025 - 1)} \\ = 0 \cdot 1 - \frac{(0 \cdot 00033)}{(-0009915)} \\ & 0 \cdot 1 + 0 \cdot 033 \\ & = \frac{0 \cdot 1033}{1} \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & \\ & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & &$$