

10/11/24

Neural Networks

Gokhulath ①

Thirumaran

675086474

1) a) Gradient Descent

Given:

$$R(w) = 13w_1^2 - 10w_1w_2 + 4w_1 + 2w_2^2 - 2w_2 + 1$$

Partial Derivatives of R with respect to w_1 & w_2

$$\frac{\partial R}{\partial w_1} = 26w_1 - 10w_2 + 4 = 0 \quad \text{--- (1)}$$

$$\frac{\partial R}{\partial w_2} = -10w_1 + 4w_2 - 2 = 0 \quad \text{--- (2)}$$

To find the optimal solution, equate (1) & (2) to zero

$$26w_1 - 10w_2 + 4 = 0 \quad \text{--- (3)}$$

$$-10w_1 + 4w_2 - 2 = 0 \quad \text{--- (4)}$$

From (3),

$$26w_1 - 10w_2 + 4 = 0$$

$$26w_1 = 10w_2 - 4$$

$$w_1 = \frac{10w_2 - 4}{26}$$

$$w_1 = \frac{5w_2 - 2}{13} \rightarrow \text{--- (5)}$$

Substitute (5) in (4)

$$-10\left(\frac{5w_2 - 2}{13}\right) + 4w_2 - 2 = 0$$

$$\frac{-50w_2 + 20}{13} + 4w_2 - 2 = 0$$

$$-50w_2 + 20 + 52w_2 - 26 = 0$$

$$-50w_2 + 52w_2 + 20 - 26 = 0$$

$$2w_2 - 6 = 0$$

$$w_2 - 3 = 0$$

$$\boxed{w_2 = 3} \rightarrow (6)$$

sub (6) in (5)

$$w_1 = \frac{5w_2 - 2}{13}$$

$$= \frac{5(3) - 2}{13}$$

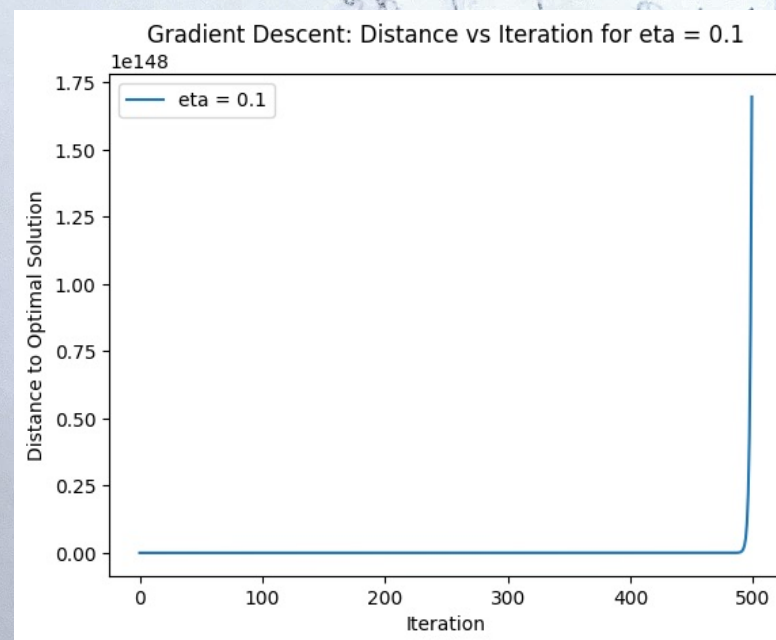
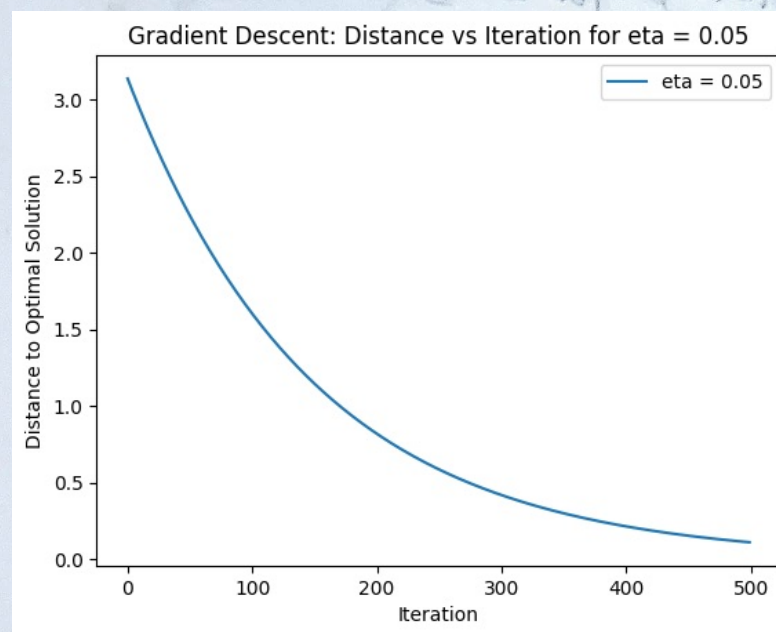
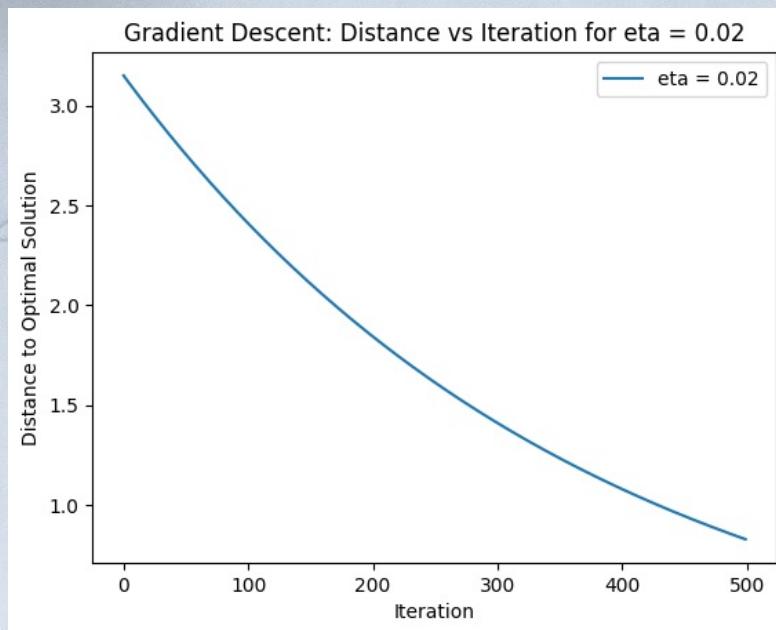
$$= \frac{15 - 2}{13} = \frac{13}{13} = 1$$

$$\boxed{w_1 = 1}$$

Thus, the optimal solution is $w_1 = 1$

$$w_2 = \underline{\underline{3}}$$

b)



c) If η is large, the algorithm may overshoot the optimal point, which might cause divergence. ④

If η is small, the convergence may be slow.

Thus choosing an optimal η is important.

2) a) Backpropagation.

The Sigmoid function ϕ for x is

$$\phi(x) = \frac{1}{1 + e^{-5x}}$$

Derivative of $\phi(x)$

$$\phi'(x) = (1 + e^{-ax})^{-1} \quad \text{where } a = +5$$

$$\begin{aligned} \phi'(x) &= (-1) (1 + e^{-ax})^{-2} \cdot (e^{-ax} \cdot -a) \\ &= (1 + e^{-ax})^{-2} (ae^{-ax}) \end{aligned}$$

$$= \frac{ae^{-ax}}{(1 + e^{-ax})^2}$$

$$= a \cdot \frac{1}{(1 + e^{-ax})} \cdot \frac{e^{-ax}}{(1 + e^{-ax})}$$

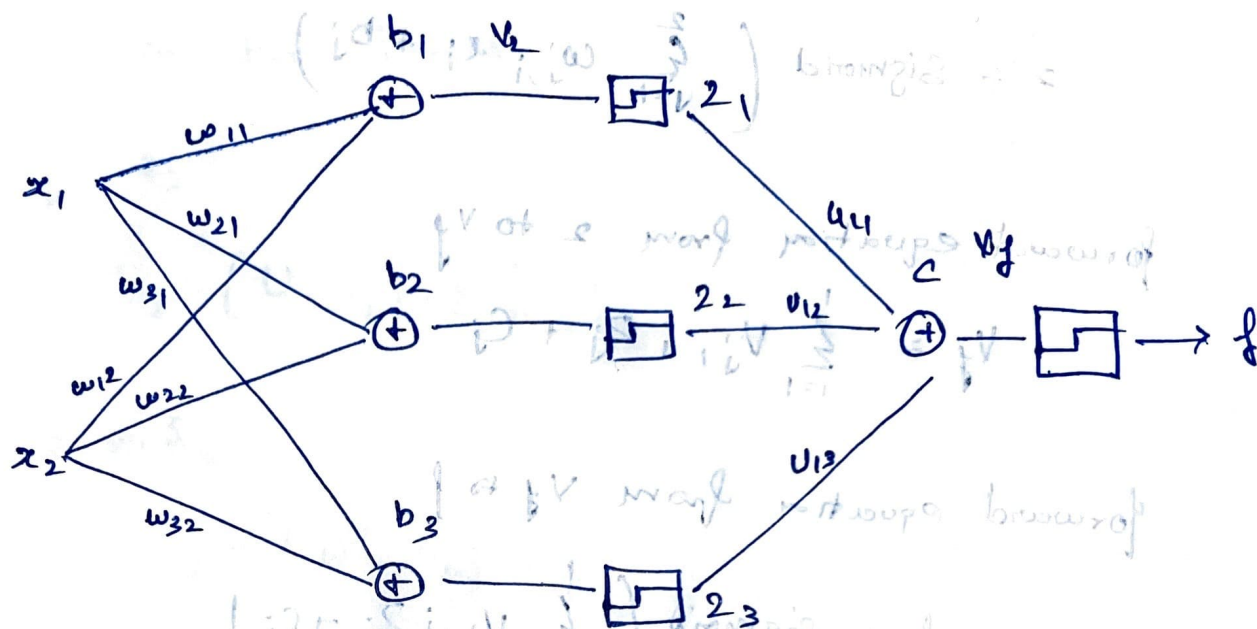
$$\text{here } \frac{e^{-ax}}{1 + e^{-ax}} = 1 - \phi(x)$$

Thus

$$\phi'(x) = a \cdot \phi(x)(1 - \phi(x))$$

$$\phi'(x) = 5 \cdot \phi(x)(1 - \phi(x))$$

2) b)



input vector $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_{2 \times 1}$

weight matrix $w = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix}_{3 \times 2}$

bias vector $b = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}_{3 \times 1}$ hidden layer o/p vector $z = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix}_{3 \times 1}$

weight matrix $w = [w_{11}, w_{12}, w_{13}]_{1 \times 3}$ bias $b = [b]$ 1×1 output $f = [f]$ 1×1

weight matrix

Bias

output

⑥

forward equation from x to v_2

$$v_2 = \sum_{i=1}^2 \omega_{j,i} x_i + b_j$$

forward equation from v_2 to z

$$z = \text{sigmoid} \left(\sum_{i=1}^2 \omega_{j,i} v_2 + b_j \right)$$

forward equation from z to v_f

$$v_f = \sum_{i=1}^1 v_{j,i} z_i + c_j$$

forward equation from v_f to f

$$f = \text{sigmoid} \left(\sum_{i=1}^1 v_{j,i} z_i + c_j \right)$$

2) c) Given,

Squared loss function,

$$L(y, f) = (f - y)^2$$

also we know that

Sigmoid function

$$\phi(v) = \frac{1}{1 + e^{-v}}$$

derivation sigmoid fu

$$\phi'(v) = \phi(v)(1 - \phi(v))$$

$$\nabla_f L = \frac{\partial L}{\partial f} = 2(f - y)$$

$$\delta_f = \nabla_f L \cdot \phi(v_f) = 2(f - y) \phi'(v_f)$$

backward equation from

i) δ_f to δ_2

$$\delta_2 = (v \cdot \delta_f) \cdot \phi'(v_2)$$

ii) δ_2 to δ_x

$$\delta_x = (w^T \cdot \delta_2)$$

Gradients for

$$\nabla_w L = 2 \cdot \delta_f$$

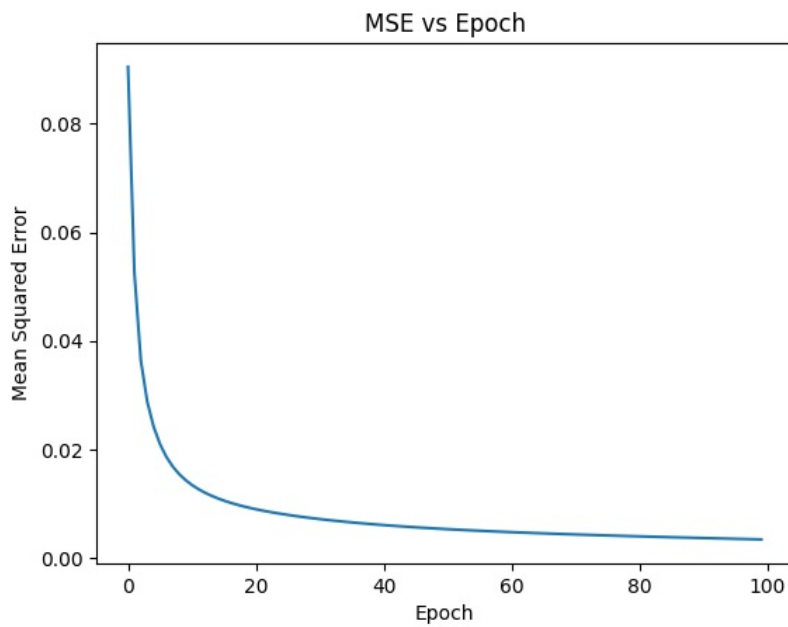
$$\nabla_b L = \delta_f$$

$$\nabla_w L = x \cdot \delta_2$$

$$\nabla_b L = \delta_2$$

d)

8



$$f(x) = \frac{1}{1 + e^{-x}}$$

$$\frac{df}{dx} = f(1-f)$$

Derivative of sigmoid function

$$\frac{df}{dx} = f(1-f)$$

$$\frac{df}{dx} = f(1-f)$$

$$\frac{df}{dx} = f(1-f)$$

$$= f(1-f)$$

Derivative of sigmoid function

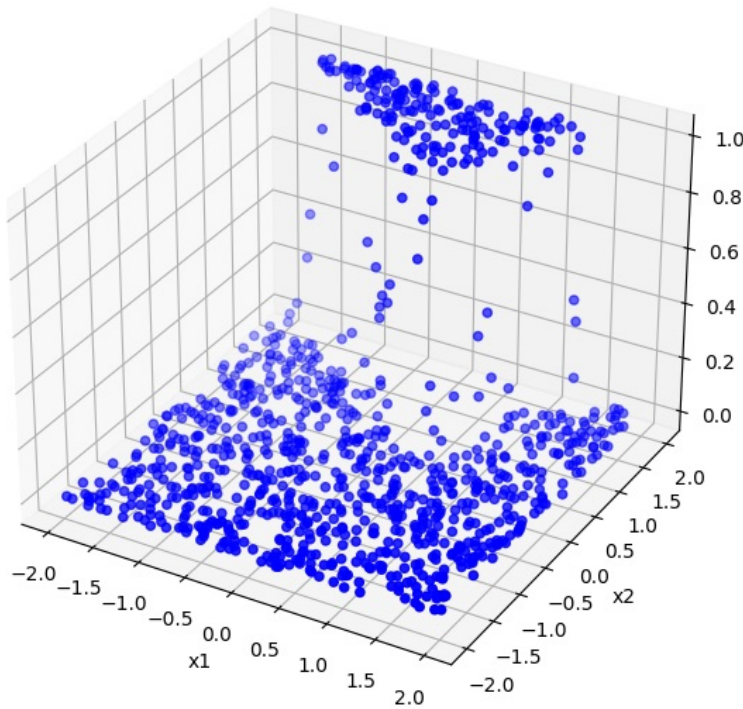
$$\Delta f = f(1-f)$$

$$\Delta f = f(1-f)$$

$$\Delta f = f(1-f)$$

$$\Delta f = f(1-f)$$

e)



f) we have chosen

Learning rate change

Mini batch gradient Descent

Start η high but multiply by 0.9 if MSE \uparrow

Results

Starting with a high learning rate can accelerate convergence early in the process. But as we approach a minimum, smaller learning rate stabilize it. Mini Batches can help reduce variance in the gradient, thus leading to smoother convergence.

