

Q 2

a) For the bias term:-  
we can find the gradient of the log loss as  $y=1$  as all the training examples are +ve.

$$\therefore \log \text{loss} = -[y \log \hat{y} - (1-y) \log(1-\hat{y})]$$

$$\therefore \nabla_b f_{\text{loss}} = \nabla_b [-1 \log \hat{y} - 0 \cdot \log(1-\hat{y})]$$

$$\therefore \nabla_b f_{\text{loss}} = \nabla_b [-\log \hat{y}]$$

$$= -\nabla_b \log \sigma(x^T w + b)$$

$$= -\frac{\sigma(x^T w + b)(1 - \sigma(x^T w + b))}{\sigma(x^T w + b)}$$

$$\therefore \nabla_b f_{\text{loss}} = (\sigma(x^T w + b) - 1)$$

1) Now, if we assume  $w=0$

$$\nabla_b f_{\text{loss}} = \sigma(b) - 1$$

The bias update is

$$b_{\text{new}} = b_{\text{old}} - [\text{learning rate} \times (\sigma(b) - 1)]$$

$$b_{\text{new}} = b_{\text{old}} - \alpha (\sigma(b) - 1)$$

As  $0 < \sigma(b) < 1 \therefore \sigma(b) - 1 < 0$ ; let  $\sigma(b) - 1 = a$

$$b_{\text{new}} = b_{\text{old}} + \alpha \times a$$

if we assume  $b$  convergence as 10000 using the above equation we won't be able to provide an upper bound for iterations.

2) if  $w \neq 0$

$$\nabla_b f_{\text{loss}} = \sigma(x^T w + b) - 1$$

The bias update is:

$$b_{\text{new}} = b_{\text{old}} - ([\sigma(x^T w + b) - 1] \times \alpha)$$

Now here we can see that the convergence of bias will depend on the values of the training example as well.

$\therefore$  The convergence of the bias will depend and cannot be guaranteed.

2) For the weight vector:-

As 1) we can calculate the gradient of loss function

$$\nabla_w f_{\text{loss}} = x(\hat{y} - y)$$

$$\nabla_w f_{\text{loss}} = x[\hat{y} - y] = x[\hat{y} - 1]$$

$$\nabla_w f_{\text{loss}} = x[\sigma(x^T w + b) - 1]$$

Assume  $b=0$  for simplicity.

$$\nabla_w f_{\text{loss}} = x[\sigma(x^T w) - 1]$$

$$\text{Let } x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}; w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$\nabla_w f_{\text{loss}} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} [\sigma([x_1 \ x_2] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}) - 1]$$

$$\nabla_w f_{\text{loss}} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} [\sigma(w_1 x_1 + w_2 x_2) - 1]$$

Update for  $w$

$$w = w_0 - \alpha \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} [\sigma(x_1 w_1 + x_2 w_2) - 1]$$

Now value of  $\sigma(x_1 w_1 + x_2 w_2) - 1 < 0$

$$w = w_0 - \alpha_{\text{new}} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\alpha_{\text{new}} = \alpha \times a; a = [\sigma(x_1 w_1 + x_2 w_2) - 1]$$

This equation denotes that the convergence of  $w$  is dependent on the exact values of  $x$ .

if  $x = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$  it will never converge.

if  $x = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$  or some suitable value it might converge.

$\therefore$  Convergence of  $w$  depends on the value of training examples.