

Vaibhav Nandkumar Kadam
vkadam@wpi.edu

Aadesh Varude
avarude@wpi.edu

Q1 . Solution.

$$J(\theta) = \frac{1}{4} \sum (f^*(x) - f(x, \theta))^2$$

$$X = \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix}$$

$$W = [w_1 \quad w_2] \text{ and we have, } f^*(x) = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

We know $y = W^T X + b$ can be represented as $y = W^T X$ with b augmented in W , Thus we have,

$$X = \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

$$W = [w_1 \quad w_2 \quad b] \text{ and we have, } y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

So we get,

$$W = (X^T X)^{-1} X y$$

$$(X^T X) = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 2 \\ 2 & 1 & 2 \\ 2 & 2 & 4 \end{bmatrix}$$

$$(X^T X)^{-1} = \begin{bmatrix} 1 & 0 & -1/2 \\ 0 & 1 & -1/2 \\ -1/2 & -1/2 & 3/4 \end{bmatrix}$$

$$X y = \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix}$$

$$W = (X^T X)^{-1} X y = \begin{bmatrix} 1 & 0 & -1/2 \\ 0 & 1 & -1/2 \\ -1/2 & -1/2 & 3/4 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1/2 \end{bmatrix}$$

Thus we get $w_1 = 0$, $w_2 = 0$ and $b = 0.5$

For Q 2a and Q 2b please refer last part of the document .

Q2 c. Solution. Training Loss will converge as gradient descent will try minimize log loss as far as the function is differentiable and will satisfies the constraints for optimization.

Q2 d. Solution.

Testing loss convergence significantly depends on whether it is overfitted. Due to overfitting the testing loss would increase. Thus it highly depends on the type and amount of training data its being trained on. Thus the testing loss may or may not converge.

Q3 . Solution.

Deriving $\nabla_{w^{(l)}} z_k$

$$\nabla_{w^{(l)}} z_k = \frac{\partial}{\partial w^{(l)}} [x^T w^k + b_k] \quad (1)$$

$$\nabla_{w^{(l)}} z_k = x^{(i)T} \quad \text{for } l = k \quad (2)$$

$$\nabla_{w^{(l)}} z_k = 0 \quad \text{for } l \neq k \quad (3)$$

Deriving $\nabla_{w^{(l)}} [\sum_{k'=1}^c e^{z_k}]$

$$\nabla_{w^{(l)}} [\sum_{k'=1}^c e^{z_k}] = e^{z_k} \nabla_{w^{(l)}} z_1 \cdot \cdot = e^{z_k} \nabla_{w^{(l)}} z_c \quad (4)$$

$$\begin{aligned} \nabla_{w^{(l)}} [\sum_{k'=1}^c e^{z_k}] &= e^{z_k} x^{(i)} \quad \text{for } k \neq l \\ &= e^{z_l} x^{(i)} \quad \text{for } k = l \end{aligned} \quad (5)$$

Deriving $\nabla_{w^{(l)}} \hat{y}_k$ for $k = l$

$$\nabla_{w^{(l)}} \hat{y}_k = \frac{e^{z_k}}{\sum_{k'=1}^c e^{z_k}} \quad (6)$$

By $\frac{u}{v}$ derivative of division rule.

$$\nabla_{w^{(l)}} \hat{y}_k = \frac{\sum_{k=1}^c e^{z_k} [\nabla_{w^{(l)}} e^{z_k}] - e^{z_k} [\nabla_{w^{(l)}} [\sum_{k=1}^c e^{z_k}]]}{[\sum_{k=1}^c e^{z_k}]^2}$$

From eq. 5

$$\begin{aligned} &= \frac{\sum_{k=1}^c e^{z_k} [e^{z_l} x^{(i)}] - e^{z_l} [e^{z_l} x^{(i)}]}{[\sum_{k=1}^c e^{z_k}]^2} \\ &= \frac{e^{z_l} x^{(i)}}{\sum_{k=1}^c e^{z_k}} - \frac{e^{z_l}}{\sum_{k=1}^c e^{z_k}} \frac{e^{z_l}}{\sum_{k=1}^c e^{z_k}} x^{(i)} \\ &= \hat{y}_l^{(i)} x^{(i)} - (\hat{y}_l^{(i)})^2 x^{(i)} \end{aligned}$$

$$\boxed{\nabla_{w^{(l)}} \hat{y}_k = \hat{y}_l^{(i)} x^{(i)} (1 - \hat{y}_l^{(i)}) \text{ for } l = k} \quad (7)$$

Hence proved.

Deriving $\nabla_{w^{(l)}} \hat{y}_k$ for $k \neq l$

$$\nabla_{w^{(l)}} \hat{y}_k = \frac{\sum_{k=1}^c e^{z_k} [\nabla_{w^{(l)}} e^{z_k}] - e^{z_k} [\nabla_{w^{(l)}} [\sum_{k=1}^c e^{z_k}]]}{[\sum_{k=1}^c e^{z_k}]^2}$$

From eq. 5

$$\begin{aligned} &= \frac{-e^{z_k} [e^{z_l} x^{(i)}]}{[\sum_{k=1}^c e^{z_k}]^2} \\ &= \frac{-e^{z_k}}{\sum_{k=1}^c e^{z_k}} \frac{e^{z_l}}{\sum_{k=1}^c e^{z_k}} x^{(i)} \\ &= -\hat{y}_k^{(i)} \cdot \hat{y}_l^{(i)} x^{(i)} \end{aligned}$$

$$\boxed{\nabla_{w^{(l)}} \hat{y}_k = -\hat{y}_k^{(i)} \cdot \hat{y}_l^{(i)} x^{(i)} \text{ for } l \neq k}$$

(8)

Hence proved.

Deriving $\nabla_{w^{(l)}} f_{CE}(W, b)$

$$\begin{aligned} \nabla_{w^{(l)}} f_{CE}(W, b) &= \frac{-1}{n} \sum_{i=1}^n \sum_{k=1}^c y_k^{(i)} \nabla_{w^{(l)}} \log \hat{y}_k^{(i)} \\ &= \frac{-1}{n} \sum_{i=1}^n \sum_{k=1}^c y_k^{(i)} \cdot \frac{1}{\hat{y}_k^{(i)}} \cdot \nabla_{w^{(l)}} \hat{y}_k^{(i)} \\ &= \frac{-1}{n} \sum_{i=1}^n \left[y_l^{(i)} \cdot \frac{1}{\hat{y}_l^{(i)}} \cdot \nabla_{w^{(l)}} \hat{y}_l^{(i)} + \sum_{k \neq l} y_k^{(i)} \cdot \frac{1}{\hat{y}_k^{(i)}} \nabla_{w^{(l)}} y_k^{(i)} \right] \\ &= \frac{-1}{n} \sum_{i=1}^n \left[y_l^{(i)} \cdot \frac{1}{\hat{y}_l^{(i)}} \cdot \left(x^{(i)} y_l^{(i)} (1 - \hat{y}_l^{(i)}) \right) + \sum_{k \neq l} y_k^{(i)} \cdot \frac{1}{\hat{y}_k^{(i)}} \cdot \left(-\hat{y}_l^{(i)} \hat{y}_k^{(i)} x^{(i)} \right) \right] \\ &= \frac{-1}{n} \sum_{i=1}^n \left[y_l^{(i)} \left(x^{(i)} - (1 - \hat{y}_l^{(i)}) \right) + \sum_{k \neq l} y_k^{(i)} \left(-\hat{y}_k^{(i)} x^{(i)} \right) \right] \\ &= \frac{-1}{n} \sum_{i=1}^n \left[y_l^{(i)} x^{(i)} - y_l^{(i)} \hat{y}_l^{(i)} x^{(i)} - \sum_{k \neq l} y_k^{(i)} y_k^{(i)} \hat{y}_l^{(i)} x^{(i)} \right] \\ &= \frac{-1}{n} \sum_{i=1}^n \left[y_l^{(i)} x^{(i)} - \sum_{k=1}^c y_k^{(i)} \hat{y}_l^{(i)} x^{(i)} \right] \\ &= \frac{-1}{n} \sum_{i=1}^n \left[y_l^{(i)} x^{(i)} - x^{(i)} \hat{y}_l^{(i)} \sum_{k=1}^c y_k^{(i)} \right] \end{aligned}$$

We know $\sum_{k=1}^c y_k^{(i)} = 1$
Thus

$$\boxed{\nabla_{w^{(l)}} f_{CE}(W, b) = \frac{-1}{n} \sum_{i=1}^n x^{(i)} \left[y_l^{(i)} - \hat{y}_l^{(i)} \right]}$$

(9)

Deriving $\nabla_b f_{CE}(W, b)$

$$\begin{aligned}
\nabla_b f_{CE}(W, b) &= \frac{-1}{n} \sum_{i=1}^n \sum_{k=1}^c y_k^{(i)} \nabla_b \log \hat{y}_k^{(i)} \\
&= \frac{-1}{n} \sum_{i=1}^n \sum_{k=1}^c y_k^{(i)} \cdot \frac{1}{\hat{y}_l^{(i)}} \nabla_b \hat{y}_l^{(i)} \\
&= \frac{-1}{n} \sum_{i=1}^n \left[y_l^{(i)} \cdot \frac{1}{\hat{y}_l^{(i)}} \cdot \nabla_{b^{(l)}} \hat{y}_l^{(i)} + \sum_{k \neq l} y_k^{(i)} \cdot \frac{1}{\hat{y}_k^{(i)}} \nabla_{b^{(l)}} y_k^{(i)} \right]
\end{aligned}$$

Similarly we can derive $\nabla_b [\sum_{k'=1}^c e^{z_k}]$ and $\nabla_b y_k^{(i)}$

$$\begin{aligned}
\nabla_b \hat{y}_k &= \frac{\sum_{k=1}^c e^{z_k} [\nabla_b e^{z_k}] - e^{z_k} [\nabla_b [\sum_{k=1}^c e^{z_k}]]}{[\sum_{k=1}^c e^{z_k}]^2} \\
&= \hat{y}_l^{(i)} - (\hat{y}_l^{(i)})^2
\end{aligned}$$

$$\begin{aligned}
\nabla_{w^{(l)}} \hat{y}_k &= \hat{y}_l^{(i)} (1 - \hat{y}_l^{(i)}) \\
&\text{for } l = k
\end{aligned} \tag{10}$$

Therefore further deriving same as $\nabla_{w^{(l)}} f_{CE}(W, b)$ for $\nabla_b f_{CE}(W, b)$ we get

$$\begin{aligned}
\nabla_b f_{CE}(W, b) &= \frac{-1}{n} \sum_{i=1}^n \left[y_l^{(i)} - \hat{y}_l^{(i)} \sum_{k=1}^c y_k^{(i)} \right] \\
\nabla_b f_{CE}(W, b) &= \frac{-1}{n} \sum_{i=1}^n \left[y_l^{(i)} - \hat{y}_l^{(i)} \right] \\
\boxed{\nabla_b f_{CE}(W, b) &= \frac{-1}{n} \sum_{i=1}^n \left[y_l^{(i)} - \hat{y}_l^{(i)} \right]} \tag{11}
\end{aligned}$$

Q4 . Solution.

The unregularized cross-entropy on the test set is 0.47814 loss is and per cent correctly classified examples 83.78% (0.8378).

Q2a and Q2b. Solution.

P.T.O

Q2

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

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

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

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

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

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

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

Now, if we assume $w=0$

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

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 $\alpha(b) - 1 = \alpha$

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

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 $w \neq 0$

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

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.

b) For the weight vector:-
As 1) we can calculate the gradient of loss function

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

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

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

Assume $b = 0$ for simplicity.

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

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

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

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

Update for w

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

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

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

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

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.

Figure 2: Q 2b solution