

Report for Deep Learning (CS541)_Homework3:

Report By:

Swapneel Dhananjay Waghlikar (WPI ID: 257598983)

Anuj Pradeep Pai Raikar (WPI ID: 181758784)

Solution to Question 4:

We have executed grid search on following values of hyperparameters:

Number of epochs = [1, 2, 3, 4]

Learning Rates = [0.000003, 0.000004, 0.000005, 0.000006]

Alphas = [2, 3, 4, 5]

Mini batch sizes = [600, 400, 200, 100]

We found the best sets of hyperparameters as:

[epochs, learning rate, alpha, mini batch size] = [4, 0.0006, 5, 100]

The unregularized fCE on testing dataset is 0.7901629983444174

The accuracy in percentage is 78.28

```
Minimun value of Fce on validation dataset 0.7631654519880522
Fce on Test dataset is: 0.7901629983444174
Best Hyperparameters: epochs (m) 4 Learning rate 6e-06 alpha 5 mini_batch_size 100
Accuracy in percentage 78.28
```

1. Newton's method

2 layer NN: $\hat{y} = f_{\omega}(x) = \omega^T x$

$$J(\omega) = \frac{1}{2n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$$

$$J(\omega) = \frac{1}{2n} \sum_{i=1}^n (\hat{y}^{(i)} - \omega^T x)^2$$

$$H[f](\omega) = \frac{1}{n} X X^T$$

Newton's method:

$$\omega^{(k+1)} = \omega^{(k)} - H^{-1} \nabla_{\omega} J(\omega^{(k)})$$

According to Taylor expansion

$$J(\omega) \approx J(\omega^{(k)}) + \nabla J(\omega^{(k)}) (\omega - \omega^{(k)}) + \frac{1}{2} (\omega - \omega^{(k)}) H (\omega - \omega^{(k)})$$

$$\nabla_{\omega} J(\omega) \approx \nabla_{\omega} J(\omega^{(k)}) + \frac{1}{2} \nabla_{\omega} (\omega H \omega - H \omega)$$

$$\omega^{(k+1)} = \omega^{(k)} - H^{-1} \nabla_{\omega} J(\omega^{(k)})$$

$$H^{-1} = \left(\frac{1}{n} (X X^T) \right)^{-1} \Rightarrow \boxed{H^{-1} = n (X X^T)^{-1}}$$

$$\omega^{(k+1)} - \omega^{(k)} = H^{-1} \nabla_{\omega} J(\omega^{(k)})$$

$$\omega^{(k)} - \omega^{(k+1)} = (H^{-1} \nabla_{\omega} J(\omega^{(k)}))$$

In 1 iteration;

$$\omega^{(1)} = \omega^{(0)} + H^{-1} \nabla_{\omega} J(\omega^{(0)})$$

$$\nabla_{\omega} J(\omega^{(0)}) = \frac{1}{n} X (X^T \omega^{(0)} - Y)$$

$$\omega^{(1)} = \omega^{(0)} + n (X X^T)^{-1} \left(\frac{1}{n} X (X^T \omega^{(0)} - Y) \right)$$

$$= \omega^{(0)} + n (X X^T)^{-1} \left(\frac{1}{n} \right) (X X^T \omega^{(0)} - X Y)$$

$$= \omega^{(0)} - (X X^T)^{-1} (X X^T) (\omega^{(0)}) + (X X^T)^{-1} X Y$$

$$\boxed{\omega^{(i)} = \infty (X X^T)^{-1} X Y} \quad \text{Hence proved}$$

2. Softmax regression gradient updates.

$W = [w^{(0)} \dots w^{(c)}]$ is an $m \times c$ matrix

$c = \text{no. of classes}$

$$\hat{y}_k = \frac{e^{z_k}}{\sum_{k'=1}^c e^{z_{k'}}} \quad \text{and} \quad z_k = X^T w^{(k)} + b_k$$

for each $k = 1, \dots, c$

$$f_{CE}(W, b) = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^c y_i^{(i)} \log \hat{y}_k^{(i)}$$

$$\begin{aligned} \nabla_{w^{(i)}} f_{CE}(W, b) &= -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^c y_i^{(i)} \nabla_{w^{(i)}} \log \hat{y}_k^{(i)} \\ &= -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^c y_i^{(i)} \hat{y}_k^{(i)} \left(\frac{\nabla_{w^{(i)}} \hat{y}_k^{(i)}}{\hat{y}_k^{(i)}} \right) \end{aligned}$$

a. For $l = k$ $\nabla_{w^{(i)}} \hat{y}_k^{(i)} = x^{(i)} \hat{y}_k^{(i)} (1 - \hat{y}_k^{(i)})$

b. For $l \neq k$ $\nabla_{w^{(i)}} \hat{y}_k^{(i)} = -x^{(i)} \hat{y}_k^{(i)} \hat{y}_l^{(i)}$

a. $\nabla_{w^{(i)}} \hat{y}_k^{(i)} = \nabla_{w^{(i)}} \left(\frac{e^{z_k}}{\sum_{k'=1}^c e^{z_{k'}}} \right) = \frac{\nabla_{w^{(i)}} \left[\left(\sum_{k'=1}^c e^{z_{k'}} \right) (\bar{v}_{w^{(i)}} e^{z_k}) - e^{z_k} (e^{z_k})^2 \right]}{\left(\sum_{k'=1}^c e^{z_{k'}} \right)^2}$

$\nabla_w z_k = \nabla_w (X^T w^{(k)} + b_k)$
 $= X$

$$\begin{aligned} \nabla_{w^{(i)}} \hat{y}_k^{(i)} &= x^{(i)} \left[e^{z_k} \left(\sum_{k'=1}^c e^{z_{k'}} \right) - (e^{z_k})^2 \right] / \left(\sum_{k'=1}^c e^{z_{k'}} \right)^2 \\ &= x^{(i)} \left(e^{z_k} / \sum_{k'=1}^c e^{z_{k'}} \right) \left(1 - \frac{e^{z_k}}{\sum_{k'=1}^c e^{z_{k'}}} \right) \end{aligned}$$

$$\nabla_{w^{(i)}} \hat{y}_k^{(i)} = x^{(i)} \left(\hat{y}_k^{(i)} \right) \left(1 - \hat{y}_k^{(i)} \right) \Rightarrow \text{check out}$$

For $L = k$

$$\boxed{\nabla_w^{(ii)} \hat{y}_k^{(ii)} = x^{(ii)} y_k^{(ii)} (1 - \hat{y}_k^{(ii)})}$$

$$y = \frac{e^{z_k}}{\sum_{k'=1}^c e^{z_k}} \quad \text{Applying log}$$

$$\log \hat{y} = \log \left(\frac{e^{z_k}}{\sum_{k'=1}^c e^{z_k}} \right) = \log(e^{z_k}) - \log \left(\sum_{k'=1}^c e^{z_k} \right)$$

$$\nabla_w \log \hat{y} = \nabla_w \left(\log e^{z_k} - \log \left(\sum_{k'=1}^c e^{z_k} \right) \right)$$

$$\frac{\nabla_w \hat{y}}{\hat{y}} = \frac{\nabla_w z_k (e^{z_k})}{e^{z_k}} - \frac{\nabla_w \sum e^{z_k}}{\sum e^{z_k}} = \nabla_w z_k - \frac{\nabla_w \sum e^{z_k}}{\sum e^{z_k}}$$

If $k \neq l$

$$\frac{\nabla_w^{(ii)} \hat{y}}{\hat{y}} = \nabla_w z_k - \frac{\nabla_w z_l e^{z_k}}{\sum e^{z_k}} = 0 - \frac{\nabla_w z_l e^{z_k}}{\sum e^{z_k}}$$

$$\text{but, } \frac{e^{z_k}}{\sum e^{z_k}} = \hat{y}_k \quad \text{So, } \frac{\nabla_w \hat{y}}{\hat{y}} = -\nabla_w z_l \hat{y}_l$$

$$\text{If } k \neq l \quad \boxed{\nabla_w \hat{y}^{(ii)} = -x^{(ii)} \hat{y}_k^{(ii)} \hat{y}_l^{(ii)}}$$

$$\begin{aligned} c. \quad \nabla_w f_{CE}(w, b) &= -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^c \hat{y}_k^{(ii)} \nabla_w^{(ii)} \log \hat{y}_k^{(ii)} \\ &= -\frac{1}{n} \sum_{i=1}^n x^{(ii)} (\hat{y}^{(ii)} - \hat{y}^{(ii)}) \end{aligned}$$

$$\begin{aligned} \nabla_w f_{CE}(w, b) &= -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^c \hat{y}_k^{(ii)} \nabla_w^{(ii)} \log \hat{y}_k^{(ii)} \\ &= -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^c \hat{y}_k^{(ii)} \left(\frac{\nabla_w^{(ii)} \hat{y}_k^{(ii)}}{\hat{y}_k^{(ii)}} \right) \end{aligned}$$

For all cases; i.e., $k=l$ and $k \neq l$

$$\nabla_w f_{CE} = -\frac{1}{n} \sum_{i=1}^n \left[\sum_{k=1}^c \frac{y_k^{(i)}}{\hat{y}_k} g_k^{(i)} (1 - g_k^{(i)}) + \sum_{k \neq l} \frac{y_k^{(i)}}{\hat{y}_k} \frac{y_k^{(i)}}{\hat{y}_k} [-x_i \hat{y}_k \hat{y}_l] \right]$$

$$= -\frac{1}{n} \sum_{i=1}^n \left[\sum_{k=1}^c y_k x^{(i)} (1 - \hat{y}_k) - \sum_{k \neq l} y_k^{(i)} x^{(i)} y_k \hat{y}_l \right]$$

$$= -\frac{1}{n} \sum_{i=1}^n \left[y_i x^{(i)} - \sum_{k=1}^c y_k x^{(i)} y_k - \sum_{k \neq l} y_k x^{(i)} y_k \hat{y}_l \right]$$

$$= -\frac{1}{n} \sum_{i=1}^n \left[y_i x^{(i)} - \sum_{k=1}^c x^{(i)} y_k \hat{y}_k y_l \right]$$

$$= -\frac{1}{n} \sum_{i=1}^n \left[y_i^{(i)} x^{(i)} - x^{(i)} y_i^{(i)} \right]$$

$$\boxed{\nabla_w f_{CE} = -\frac{1}{n} \sum_{i=1}^n x_i^{(i)} [y_i^{(i)} - \hat{y}_i^{(i)}]}$$

$$\text{Now, } d \cdot \nabla_b f_{CE}(w, b) = -\frac{1}{n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})$$

$$\nabla_b \hat{y}_k^{(i)} = \nabla_b \frac{e^{z_k}}{\sum_{j=1}^c e^{z_k}}$$

$$= \frac{e^{z_k} \sum e^{z_k} - e^{z_k} e^{z_k}}{(\sum e^{z_k})^2} = \frac{e^{z_k}}{\sum e^{z_k}} - \left(\frac{e^{z_k}}{\sum e^{z_k}} \right) \left(\frac{e^{z_k}}{\sum e^{z_k}} \right)$$

$$= \hat{y}_k^{(i)} - \hat{y}_k^{(i)} \hat{y}_k^{(i)}$$

$$\text{When } k = l \quad \nabla_b \hat{y}_k^{(i)} = \hat{y}_k^{(i)} (1 - \hat{y}_k^{(i)})$$

When $k \neq l$

$$\nabla_b \hat{y}_k^{(i)} = \nabla_b \left[\frac{e^{z_k}}{\sum e^{z_k}} \right] = \frac{\sum e^{z_k} - e^{z_k} \cdot e^{z_k}}{(\sum e^{z_k})^2}$$
$$= -\frac{e^{z_k}}{\sum e^{z_k}} \cdot \frac{e^{z_k}}{\sum e^{z_k}}$$

$$\therefore \nabla_b \hat{y}_k^{(i)} = -\hat{y}_k^{(i)} \cdot \hat{y}_l^{(i)}$$

and,

$$\nabla_b f_{ce}(w, b) = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^c y_k^{(i)} \nabla_b \log \hat{y}_k^{(i)}$$

$$= -\frac{1}{n} \sum_{i=1}^n \left[\sum_{k \neq l} \frac{y_k \hat{y}_l (1 - \hat{y}_l)}{\hat{y}_k} + \sum_{k \neq l} \frac{y_k}{y_k} (-\hat{y}_k \hat{y}_l) \right]$$

$$= -\frac{1}{n} \sum_{i=1}^n \left[y_i - \sum_{k \neq i} y_k \hat{y}_k + \sum_{k \neq i} y_k \hat{y}_k \right]$$

$$= -\frac{1}{n} \sum_{i=1}^n \left[y_i - \hat{y}_i \sum_{k=1}^c (\hat{y}_k^{(i)}) \right]$$

$$\boxed{\nabla_b f_{ce}(w, b) = -\frac{1}{n} \sum_{i=1}^n [y_i - \hat{y}_i]}$$

3. Softmax c -classes $k=1 \dots c$

$$\hat{y}_k = P(y_k=1 | x, w, b) \quad \forall k \in \{1, \dots, c\}$$

Ground truth label $y = [y_1, \dots, y_c]^T$ y is a one-hot vector

Likelihood of training example,

$$P(y|x, w, b) = P(y_1=1|x, w, b)^{y_1} \times \dots \times P(y_c=1|x, w, b)^{y_c}$$

$$P(y|x, w, b) = \prod_{k=1}^c \hat{y}_k^{y_k}$$

For the dataset,

$$P(D|w, b) = \prod_{i=1}^n \prod_{k=1}^C \hat{y}_k^{(i)} y_k^{(i)}$$

Applying log and multiplying by -1

$$-\log P(D|w, b) = -\log \prod_{i=1}^n \prod_{k=1}^C \hat{y}_k^{(i)} y_k^{(i)}$$

$$= -\prod_{i=1}^n \prod_{k=1}^C \log (\hat{y}_k^{(i)} y_k^{(i)})$$

$$= -\sum_{i=1}^n \sum_{k=1}^C y_k^{(i)} \log \hat{y}_k^{(i)}$$

Hence proved