

10

truth table of XOR :

x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

$$X = \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix}^T$$

$$Y = [y^{(i)}] = [0 \ 1 \ 1 \ 0]^T$$

$$W = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

MSE loss function :-

$$J(\theta) = \frac{1}{4} \sum_{x \in X} (f^*(x) - f(x; \theta))^2$$

where $\theta : w_1, w_2, b$

$$\therefore J(w_1, w_2, b) = \frac{1}{4} \sum_{x \in X} (f^*(x) - (x^T W + b))^2 \quad \text{--- (i)}$$

$$\nabla_{w_1} J(\theta) = \frac{1}{2} \sum_{x \in X} (-x_1) (y - x_1 w_1 - x_2 w_2 - b)$$

$$\nabla_{w_2} J(\theta) = \frac{1}{2} \sum_{x \in X} (-x_2) (y - x_1 w_1 - x_2 w_2 - b)$$

eq. (i) can be rewritten as :-

$$J(w_1, w_2, b) = \frac{1}{4} \left[\begin{bmatrix} 0 & 1 & 1 & 0 \end{bmatrix}_{1 \times 4}^T - \left(\begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}_{4 \times 2} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}_{2 \times 1} + \begin{bmatrix} b \end{bmatrix}_{1 \times 1} \right)^2 \right]$$

$$J(w_1, w_2, h) =$$

$$\frac{1}{4} \left[[0 \ 1 \ 1 \ 0]^T - \left(\begin{bmatrix} 0 \\ w_2 \\ w_1 \\ w_1 + w_2 \end{bmatrix}_{4 \times 1} + h \right) \right]^2$$

$$= \frac{1}{4} \left[[0 \ 1 \ 1 \ 0]^T - \begin{pmatrix} h \\ w_2 + h \\ w_1 + h \\ w_1 + w_2 + h \end{pmatrix}_{4 \times 1} \right]^2$$

$$= \frac{1}{4} \left[\begin{pmatrix} -h \\ 1 - (w_2 + h) \\ 1 - (w_1 + h) \\ -(w_1 + w_2 + h) \end{pmatrix} \right]^2$$

$$= \frac{1}{4} \left[h^2 + (1 - (w_2 + h))^2 + (1 - (w_1 + h))^2 + (w_1 + w_2 + h)^2 \right]$$

$$\therefore J(w_1, w_2, h) = \frac{1}{4} \left[h^2 + (1 - (w_2 + h))^2 + (1 - (w_1 + h))^2 + (w_1 + w_2 + h)^2 \right]$$

(simplifying this form :-)

$$\rightarrow h^2 + 1 + (w_2 + h)^2 - 2(w_2 + h) + 1 + (w_1 + h)^2 + 2(w_1 + h) - 2(w_1 + h) + (w_1 + w_2)^2 + 2(w_1 + w_2)h + h^2$$

$$= h^2 + 2 + w_2^2 + h^2 + 2w_2h - 2w_2 - 2h + 1 + w_1^2 + h^2 + 2w_1h + \cancel{2w_1 + 2h} + w_1^2 + w_2^2 + 2w_1w_2 + 2w_1h + 2w_2h + h^2$$

$$= 4h^2 + 2 + 2(w_1^2 + w_2^2) + 4h(w_1 + w_2) - 2(w_1 + w_2 + 2h) + 2w_1w_2$$

$$J(w_1, w_2, h) = \frac{1}{4} \left(4h^2 + 2 + 2(w_1^2 + w_2^2) + 4h(w_1 + w_2) - 2(w_1 + w_2 + 2h) + 2w_1w_2 \right) \quad \text{--- (ii)}$$

taking gradient of eq (ii) w.r.t w_1, w_2 & h and equating to 0 :-

$$\nabla J(w_1) = [0 + 4w_1 + 4h + 2w_2 - 2] = 0 \quad \text{--- (iii)}$$

$$\nabla J(w_2) = [0 + 4w_2 + 4h + 2w_1 - 2] = 0 \quad \text{--- (iv)}$$

$$\nabla J(h) = [8h + 4(w_1 + w_2) - 4] = 0 \quad \text{--- (v)}$$

eq (iii), (iv) & (v) can be rewritten as :-

$$\begin{aligned} \nabla J(w_1) &\Rightarrow 2w_1 + 2h + w_2 = 1 \quad \text{--- (iii)} \\ \nabla J(w_2) &\Rightarrow 2w_2 + 2h + w_1 = 1 \quad \text{--- (iv)} \\ \nabla J(h) &\Rightarrow 2h + w_1 + w_2 = 1 \quad \text{--- (v)} \end{aligned}$$

subtracting equation eq (iii) & (iv), we will get :-

$$\begin{aligned} w_1 - w_2 &= 0 \\ \Rightarrow w_1 &= w_2 \quad \text{--- (vi)} \end{aligned}$$

equating eq (vi) in eq (v) & (iii).

$$2h + 2w_1 = 1 \quad \text{--- (vii)}$$

$$h + w_1 = 1/2$$

$$3w_1 + 2h = 1 \quad \text{--- (viii)}$$

$$\text{eq. (vii)} \times 3 - \text{eq. (viii)} \times 2$$

$$6h + 6w_1 = 3$$

$$-4h + 6w_1 = 2$$

$$\hline 2h = 1$$

$$\Rightarrow h = 1/2$$

substituting $h = 1/2$ in eq. (vii)

$$2 \times \frac{1}{2} + 2w_1 = 1$$

$$\Rightarrow 2w_1 = 0$$

$$\Rightarrow w_1 = 0$$

from eq. (vi), we can deduce that:

$$w_1 = w_2 = 0$$

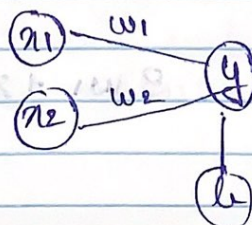
\therefore after equating gradient to 0, we have found the value of w_1, w_2, h as:

$$w_1 = 0$$

$$w_2 = 0$$

$$h = 1/2 = 0.5$$

\therefore hence proved.



2

$$\hat{y} = \sigma(x^T w + b)$$

- (a) ^{for exact,} here, all the training examples are positive. Hence, the 'b' (bias) will converge to a value that will result a positive output for all positive input (i.e. positive training examples), this will not be same for different examples. The training examples should be linearly separable, then 'b' will converge to a unique value, which will minimize log-loss.

- (b) similar to the above scenario, convergence of weight vector 'w' will depend on the exact training examples. For all ^{positive} training examples, weight vector will yield positive outputs for all positive inputs.

Different set of examples will lead to different distribution of weights that minimizes loss log-loss.

- (c) The ^{training} loss will converge. The optimization algorithm will continuously update values of 'w' & 'b' until loss function has reached minima.

- (d) The testing loss may not necessarily converge, even if ~~training~~ training loss converges. This is because, since, training examples contains only positive examples, and testing examples has both positive & negative examples. Hence, when the model predicts on testing set, for positive examples the loss will tend to infinity, however for negative examples, the model will mis-classify.

2

3. Derivation of softmax regression gradient updates

$$W = [w^{(1)} \dots w^{(c)}]$$

$$\hat{y}_k = \frac{\exp z_k}{\sum_{k'=1}^c \exp z_{k'}}$$

$$z_k = x^T w^{(k)} + b_k$$

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

We handle the two cases $l=k$ and $l \neq k$ separately,

(i) Solving first for $l=k$, we get

$$\nabla w^{(k)} \hat{y}_k^{(i)} = \nabla w^{(k)} \left[\frac{\exp(x^{(i)T} w^{(k)} + b_k)}{\sum_{k'=1}^c \exp(x^{(i)T} w^{(k')} + b_{k'})} \right]$$

Using quotient rule ~~where~~ $\frac{u'v - v'u}{v^2}$, we get

$$\begin{aligned} &= \cancel{x^{(i)T} \nabla w^{(k)} \exp(x^{(i)T} w^{(k)} + b_k)} - \frac{x^{(i)T} \nabla w^{(k)} \exp(x^{(i)T} w^{(k)} + b_k)}{\left(\sum_{k'=1}^c \exp(x^{(i)T} w^{(k')} + b_{k'}) \right)^2} \cdot \exp(x^{(i)T} w^{(k)} + b_k) \\ &= x^{(i)T} \nabla w^{(k)} \frac{\exp(x^{(i)T} w^{(k)} + b_k)}{\left(\sum_{k'=1}^c \exp(x^{(i)T} w^{(k')} + b_{k'}) \right)^2} \cdot \exp(x^{(i)T} w^{(k)} + b_k) \end{aligned}$$

$$= x^{(i)} \left[\frac{\exp(x^{(i)T} w^{(l)})}{\sum_{k=1}^C \exp(x^{(i)T} w^{(k)})} - \frac{\exp(x^{(i)T} w^{(l)})^2}{\left(\sum_{k=1}^C \exp(x^{(i)T} w^{(k)}) \right)^2} \right]$$

We can see that $y_l^{(i)} = \frac{\exp(x^{(i)T} w^{(l)})}{\sum_{k=1}^C \exp(x^{(i)T} w^{(k)})}$.

Substituting this value above, we get.

$$= x^{(i)} [\hat{y}_l^{(i)} - (\hat{y}_l^{(i)})^2]$$

$$= x^{(i)} \hat{y}_l^{(i)} [1 - \hat{y}_l^{(i)}]$$

ii) Solving now for $l \neq k$, we get

$$\nabla_w^{(l)} \hat{y}_k^{(i)} = \nabla_w^{(l)} \left[\frac{\exp(x^{(i)T} w^{(k)})}{\sum_{k'=1}^C \exp(x^{(i)T} w^{(k')})} \right]$$

$$= -x^{(i)} \frac{\exp(x^{(i)T} w^{(k)}) \cdot \exp(x^{(i)T} w^{(l)})}{\left(\sum_{k'=1}^C \exp(x^{(i)T} w^{(k')}) \right)^2}$$

$$= \frac{-x^{(i)} (\exp(x^{(i)T} w^{(k)})) \exp(x^{(i)T} w^{(l)})}{\left(\sum_{k'=1}^C \exp(x^{(i)T} w^{(k')}) \right)^2}$$

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

(iii) To compute ~~over~~ the total gradient of J_{CE} , wrt each $w^{(k)}$, we have to sum over all examples and over $l=1, \dots, C$.

$$\nabla_{w^{(k)}} J_{CE}(w, b) = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^C y_k^{(i)} \nabla_{w^{(k)}} \log y_k^{(i)}$$

$$\neq \sum_{k=1}^C a_k = a_l + \sum_{k \neq l}^C a_k$$

$$= -\frac{1}{n} \sum_{i=1}^n \left(y_l^{(i)} \times \left(\frac{x^{(i)} y_l^{(i)}}{y_l^{(i)}} (1 - y_l^{(i)}) \right) + \sum_{k \neq l}^C y_k^{(i)} \times \left(\frac{(1 - x^{(i)}) y_k^{(i)}}{y_k^{(i)}} (1 - y_l^{(i)}) \right) \right)$$

$$= -\frac{1}{n} \sum_{i=1}^n \left(y_l^{(i)} x^{(i)} (1 - y_l^{(i)}) - x^{(i)} \left(\sum_{k \neq l}^C y_k^{(i)} \right) y_l^{(i)} \right)$$

$$\sum_{k=1}^C y_k^{(i)} = 1 \quad \therefore \sum_{k \neq l}^C y_k^{(i)} = (1 - y_l^{(i)})$$

$$= -\frac{1}{n} \sum_{i=1}^n x^{(i)} \left(y_l^{(i)} - y_l^{(i)} y_l^{(i)} - y_l^{(i)} + y_l^{(i)} y_l^{(i)} \right)$$

$$= -\frac{1}{n} \sum_{i=1}^n x^{(i)} (y_l^{(i)} - y_l^{(i)})$$

$$= \nabla_{w^{(l)}} J_{CE}(w, b)$$

iv) The gradient of b will follow the same steps

For $l = k$

$$\nabla_{b_l} \hat{y}_k^{(l)} = \frac{\nabla_{b_l} \exp(z_l)}{\sum_{k'=1}^C \exp(z_{k'})}$$

$$= \frac{\exp(x^{(l)T} w^{(l)})}{\sum_{k'=1}^C \exp(x^{(l)T} w^{(k')})} - \frac{\exp(x^{(l)T} w^{(l)})}{\left(\sum_{k'=1}^C \exp(x^{(l)T} w^{(k')})\right)^2} \exp(x^{(l)T} w^{(l)})$$

$$= \left[\frac{\exp(x^{(l)T} w^{(l)})}{\sum_{k'=1}^C \exp(x^{(l)T} w^{(k')})} - \frac{(\exp(x^{(l)T} w^{(l)})^2)}{\left(\sum_{k'=1}^C \exp(x^{(l)T} w^{(k')})\right)^2} \right]$$

$$= \hat{y}_l^{(l)} - (\hat{y}_l^{(l)})^2$$

$$= \hat{y}_l^{(l)} (1 - \hat{y}_l^{(l)})$$

For $l \neq k$

$$\nabla_{b^{(l)}} \hat{y}_k^{(l)} = \nabla_{b^{(l)}} \left[\frac{\exp(x^{(l)T} w^{(k)})}{\sum_{k'=1}^C \exp(x^{(l)T} w^{(k')})} \right]$$

$$= \frac{\exp(x^{(l)T} w^{(k)})}{\left(\sum_{k'=1}^C \exp(x^{(l)T} w^{(k')})\right)^2} \exp(x^{(l)T} w^{(l)})$$

$$= \frac{\exp(x^{(l)T} w^{(k)}) (\exp(x^{(l)T} w^{(l)}))}{\left(\sum_{k'=1}^C \exp(x^{(l)T} w^{(k')})\right)^2}$$

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

To compute the total gradient of J_{CE} wrt each $b^{(L)}$,

$$\nabla_{b^{(L)}} J_{CE}(w, b) = -\frac{1}{n} \left[y_L^{(i)} \hat{y}_L^{(i)} (1 - \hat{y}_L^{(i)}) - \sum_{k \neq 1} \frac{y_k^{(i)} \hat{y}_k^{(i)} \hat{y}_L^{(i)}}{\hat{y}_k^{(i)}} \right]$$

$$= -\frac{1}{n} \left[\sum_{i=1}^n [y_L^{(i)} (1 - \hat{y}_L^{(i)})] - \sum_{k \neq L} y_k^{(i)} \hat{y}_L^{(i)} \right]$$

$$= -\frac{1}{n} \left[\sum_{i=1}^n [y_L^{(i)} (y_L^{(i)} + \sum_{k \neq L} y_k^{(i)}) - y_L^{(i)}] \right]$$

$$\sum_k y_k = 1 \quad \text{and} \quad y_L + \sum_{k \neq L} y_k = 1.$$

$$\therefore \nabla_{b^{(L)}} J_{CE}(w, b) = -\frac{1}{n} \sum_{i=1}^n (y_L^{(i)} - \hat{y}_L^{(i)})$$

Hence Proved.


```
[Running] python -u "c:\Users\Sumeet\OneDrive - Worcester Polytechnic Institute (wpi.edu)\WPI Sem 2\Deep Learning\Submission\HW3\homework3_SSHANBHAG_UMAHANTI.py"
Finding best hyperparameters from 3*3*3*3 combinations
Best Hyperparameters:

Epochs = 35
Alpha = 0.001
Learning Rate = 1e-05
Mini Batch Size = 32
Cost function value using above hyperparameter values 1.4192011617995408

Now combining Training and validation Sets

Cost & Accuracy values

Cost = 1.8416920764703457
Accuracy = 76.51 %

[Done] exited with code=0 in 19.675 seconds
```

```
# Below values are tuned values for the same for highest accuracy
Mb = [32]
Epoch = [35]
Alpha = [0.001]
Eps = [0.00001]
```