# Contents

1. Newton's Method →

gradient → $\hat{y} = f_w(x) = w^T x$

$$J(w) = \frac{1}{2n} \sum_{i=1}^{n} \left( \hat{y}^{(i)} - y^{(i)} \right)^2 \Rightarrow f_{MSE}$$

$$\nabla_w f_{MSE}(y, \hat{y}; w) = \nabla_w \left[ \frac{1}{2n} \sum_{i=1}^{n} \left( x^{(i)^T} w - y^{(i)} \right)^2 \right]$$

$$= \frac{1}{2n} \sum_{i=1}^{n} \nabla_w \left[ \left( x^{(i)^T} w - y^{(i)} \right)^2 \right]$$

$$= \frac{1}{n} \sum_{i=1}^{n} x^{(i)} \left( x^{(i)^T} w - y^{(i)} \right)$$

$$= \frac{1}{n} X \left( X^T w - y \right) \qquad \text{—①}$$

hessian → $\nabla_w \left( \frac{1}{n} X (X^T w - y) \right)$

$$= \boxed{\frac{1}{n} X (X^T)} \qquad \text{—②}$$

Newton's Method →

— let the $2^{nd}$-order Taylor expansion of $f$ around $w^{(k)}$ be:

$$f(w) \approx f(w^{(k)}) + \nabla_w f(w^{(k)})(w - w^{(k)}) + \frac{1}{2}(w - w^{(k)})^T H (w - w^{(k)})$$

→ To minimize $f$, we will find the root of the gradient of $f$'s Taylor expansion:

$$\nabla_w f(w) \approx \nabla_w f(w^{(k)}) + \frac{1}{2} \nabla_w \left( w^T H w - w^T H w^{(k)} - w^{(k)T} H w + w^{(k)T} H w^{(k)} \right)$$

$$= \nabla_w f(w^{(k)}) + H w - \frac{1}{2} H w^{(k)} - \frac{1}{2} H w^{(k)}$$

$$= \nabla_w f(w^{(k)}) + H w - H w^{(k)}$$

equating it to zero →

$$0 = \nabla_w f(w^{(k)}) + H w - H w^{(k)}$$

$$H w = H w^{(k)} - \nabla_w f(w^{(k)})$$

$$w^{(k+1)} = w^{(k)} - H^{-1} \nabla_w f(w^{(k)}) \quad \longleftarrow \quad ③$$

$\rightarrow$ Now, let us substitute eqn ① & ② in eqn ③

$$w^{(k+1)} = w^{(k)} - H^{-1}\left(\nabla_w f\left(w^{(k)}\right)\right)$$

$$= w^{(k)} - \left(\frac{1}{n} X(X^T)\right)^{-1}\left(\frac{1}{n} X\left(X^T w^{(k)} - y\right)\right)$$

$$= w^{(k)} - n\left(XX^T\right)^{-1}\frac{1}{n}\left(XX^T w^{(k)} - XY\right)$$

$$= w^{(k)} - w^{(k)} + \left(XX^T\right)^{-1}XY$$

$$\boxed{w^{(k+1)} = \left(XX^T\right)^{-1}XY}$$

$\rightarrow$ as we can observe that $w^{(k+1)}$ is independent of $w^{(k)}$.
we can conclude that using Newton's method on MSE loss we can reach the optimum solution in just 1 iteration.

$$\boxed{w^* = \left(XX^T\right)^{-1}XY}$$

Q.2

Given :-
$$\hat{y}_k = \frac{e^{z_k}}{\sum_{k'=1}^{c} e^{z_{k'}}}$$

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

Cost $\rightarrow l_{CE}(w,b) = -\frac{1}{n} \sum_{i=1}^{n} \sum_{k=1}^{c} y_k^{(i)} \log \hat{y}_k^{(\ell)}$

Taking Gradient $\Rightarrow \nabla_{w^{(\ell)}} l_{CE}(w,b) = -\frac{1}{n} \sum_{i=1}^{n} \sum_{k=1}^{c} y_k^{(i)} \nabla_{w^{(\ell)}} \log(\hat{y}_k^{(i)})$

$$= -\frac{1}{n} \sum_{i=1}^{n} \sum_{k=1}^{c} y_k^{(i)} \left( \frac{\nabla_{w^{(\ell)}} \hat{y}_k^{(i)}}{\hat{y}_k^{(i)}} \right)$$

Now we have two cases for above expression

① when $\ell = k$

Focussing on $\left( \nabla_w^{(\ell)} \hat{y}_k^{(i)} \right)$ term

$$\nabla_w^{\ell} \hat{y}_k^{(i)} = \nabla_{w^{\ell}} \hat{y}_{\ell}^{(i)} = \nabla_{w^{\ell}} \left( \frac{e^{z_{\ell}}}{\sum_{k'=1}^{c} e^{z_{k'}}} \right)$$

Using Quotient Rule $\Rightarrow \nabla \left( \frac{u}{v} \right) = \frac{u'v - v'u}{v^2}$

$$\nabla_w^{\ell} \hat{y}_k^{(i)} = \frac{e^{z_{\ell}} \times \nabla_{w^{\ell}}(x^T w^{(\ell)} + b_{\ell}) \times \sum_{k'=1}^{c} e^{z_{k'}} - e^{z_{\ell}} \times e^{z_{\ell}} \times \nabla_{w^{\ell}}(x^T w^{\ell} + b_{\ell}}{\left( \sum_{k'=1}^{c} e^{z_{k'}} \right)^2}$$

$$\Rightarrow \frac{x \times e^{z_{\ell}} \times \sum_{k'=1}^{c} e^{z_{k'}} - x \times e^{z_{\ell}})^2 \times 1}{\left( \sum_{k'=1}^{c} e^{z_{k'}} \right)^2} \Rightarrow x \left( \frac{e^{z_{\ell}}}{\sum_{k'=1}^{c} e^{z_{k'}}} - \frac{e^{z_{\ell}} \times e^{z_{\ell}}}{\left( \sum_{k'=1}^{c} e^{z_{k'}} \right)^2} \right)$$

$$\therefore \nabla_{w^{(\ell)}} \hat{y}_k^{(i)} = x_x^{(\ell)} \left( \hat{y}_\ell^{(i)} - \hat{y}_\ell^{(i)2} \right)$$

$$\boxed{\therefore \hat{y}_\ell = \frac{e^{z_\ell}}{\sum\limits_{k=1}^{c} e^{z_{k'}}}}$$

Hence $\boxed{\nabla_{w^{(\ell)}} \hat{y}_k^{(i)} = x^{(\ell)} \hat{y}_\ell^{(i)} \left( 1 - \hat{y}_\ell^{(i)} \right)}$ when $\ell = k$

Now when $\ell \neq k$

$$\nabla_{w^{(\ell)}} \hat{y}_k^{(i)} = \nabla_{w^\ell} \left( \frac{e^{z_k}}{\sum\limits_{k'=1}^{c} e^{z_{k'}}} \right)$$

$$= e^{z_k} \nabla_{w^{(\ell)}} \left( \frac{1}{\sum\limits_{k'=1}^{c} e^{z_{k'}}} \right)$$

$$= e^{z_k} \times \frac{-1}{\left( \sum\limits_{k'=1}^{c} e^{z_{k'}} \right)^2} \times \nabla_{w^{(\ell)}} \left( x^T w^{(\ell)} + b_\ell \right) \times e^{z_\ell}$$

$$= - \frac{e^{z_k} \times e^{z_\ell} \times x}{\left( \sum\limits_{k'=1}^{c} e^{z_{k'}} \right)^2}$$

Hence $\boxed{\nabla_{w^{(\ell)}} \hat{y}_k^{(i)} = -x \times \hat{y}_k^{(i)} \times \hat{y}_\ell^{(i)}}$    $\boxed{\therefore \hat{y}_\ell = \frac{e^{z_\ell}}{\sum\limits_{k'=1}^{c} e^{z_{k'}}}, \hat{y}_k = \frac{e^{z_k}}{\sum\limits_{k'=1}^{c} e^{z_{k'}}}}$

when $\ell \neq k$

Now lets compute the total gradient of $\ell_{CE}$ each wrt $w^{(k)}$

$$\Rightarrow \nabla_{w^{(\ell)}} \ell_{CE}(w,b) = -\frac{1}{n} \sum_{i=1}^{n} \sum_{k=1}^{c} y_k^{(i)} \nabla_{w^{(\ell)}} \log \hat{y}_k^{(i)}$$

It will be the sum over $l = 1, \ldots, C$ terms

so when $k = l$ & $k \neq l$ we know the term

∴ We can split $\sum_k a_k = a_l + \sum_{k \neq l} a_k$

$$\nabla_{w^{(l)}} \Big|_{CE}(w,b) = \frac{-1}{n} \sum_{i=1}^{n} \left( y_l^{(i)} \times \left( \frac{x^{(i)} \times \hat{y}_l^{(i)}(1-\hat{y}_l^{(i)})}{\hat{y}_l^{(i)}} \right) + \sum_{k \neq l}^{c} y_k^{(i)} \times \left( \frac{-x^{(i)} \times \hat{y}_k^{(i)} \hat{y}_l^{(i)}}{\hat{y}_k^{(i)}} \right) \right)$$

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

$$\boxed{\because \sum_{k=1}^{c} y_k^{(i)} = 1 \implies \therefore \sum_{k \neq l}^{c} y_k^{(i)} = \left( 1 - y_l^{(i)} \right)}$$

$$\therefore \implies \frac{-1}{n} \sum_{i=1}^{n} \left( x^{(i)} \times y_l^{(i)} \times (1-\hat{y}_l^{(i)}) - x^{(i)} (1-y_l^{(i)}) \hat{y}_l^{(i)} \right)$$

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

$$\implies \boxed{\frac{-1}{n} \sum_{i=1}^{n} x^{(i)} \left( y_l^{(i)} - \hat{y}_l^{(i)} \right) = \nabla_{w^{(l)}} \Big|_{CE}(w,b)} \quad //$$

Similarly we can prove for $\nabla_b^{(l)} \Big|_{CE}(w,b)$

$$\implies \nabla_b^{(l)} \Big|_{CE}(w,b) = \frac{-1}{n} \sum_{i=1}^{n} \sum_{k=1}^{c} y_k^{(i)} \nabla_{b^{(l)}} \log \hat{y}_k^{(i)} = \frac{-1}{n} \sum_{i=1}^{n} \sum_{k=1}^{c} y_k^{(i)} \left( \frac{\nabla_b \hat{y}_k^{(i)}}{\hat{y}_k^{(i)}} \right)$$

It will be the sum over $l = 1, \ldots, c$ terms

so when $k = l$ & $k \neq l$ we know the term

$\therefore$ we can split $\sum_{k} a_k = a_l + \sum_{k \neq l} a_k$

$$\therefore \nabla_{w^{(l)}} \big|_{CE} (w,b) = \frac{-1}{n} \sum_{i=1}^{n} \left( y_l^{(i)} \times \left( \frac{x^{(i)} \times \hat{y}_l^{(i)} (1 - \hat{y}_l^{(i)})}{\hat{y}_l^{(i)}} \right) + \sum_{k \neq l}^{c} y_k^{(i)} \times \left( \frac{-x^{(i)} \times \hat{y}_k^{(i)} \hat{y}_l^{(i)}}{\hat{y}_k^{(i)}} \right) \right)$$

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

$$\boxed{\therefore \sum_{k=1}^{c} y_k^{(i)} = 1 \implies \therefore \sum_{k \neq l}^{c} y_k^{(i)} = \left( 1 - y_l^{(i)} \right)}$$

$$\therefore \implies \frac{-1}{n} \sum_{i=1}^{n} \left( x^{(i)} \times y_l^{(i)} \times (1 - \hat{y}_l^{(i)}) - x^{(i)} (1 - y_l^{(i)}) \hat{y}_l^{(i)} \right)$$

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

$$\implies \boxed{\frac{-1}{n} \sum_{i=1}^{n} x^{(i)} \left( y_l^{(i)} - \hat{y}_l^{(i)} \right) = \nabla_{w^{(l)}} \big|_{CE} (w,b)} \quad //$$

Similarly we can prove for $\nabla_{b^{(l)}} \big|_{CE} (w,b)$

$$\implies \nabla_{b^{(l)}} \big|_{CE} (w,b) = \frac{-1}{n} \sum_{i=1}^{n} \sum_{k=1}^{c} y_k^{(i)} \nabla_{b^{(l)}} \log \hat{y}_k^{(i)} = \frac{-1}{n} \sum_{i=1}^{n} \sum_{k=1}^{c} y_k^{(i)} \left( \frac{\nabla_b \hat{y}_k^{(i)}}{\hat{y}_k^{(i)}} \right)$$

We have two cases: ① when $l = k$

then $\nabla_{b^{(l)}} \hat{y}_k^{(i)} = \nabla_{b^{(l)}} \times \left( \dfrac{e^{z_l}}{\sum\limits_{k'=1}^{c} e^{z_{k'}}} \right)$

$\Rightarrow$ Using quotient rule

$\Rightarrow \dfrac{e^{z_l} \times 1 \times \sum\limits_{k'=1}^{c} e^{z_{k'}} - e^{z_l} \times e^{z_l}}{\left( \sum\limits_{k'=1}^{c} e^{z_{k'}} \right)^2}$

$\Rightarrow \dfrac{e^{z_l} \times \sum\limits_{b'=1}^{c} e^{z_{k'}}}{\left( \sum\limits_{b'=1}^{c} e^{z_{k'}} \right)^x} - \dfrac{e^{z_l} \times e^{z_l}}{\left( \sum\limits_{k'=1}^{c} e^{z_{k'}} \right)^2}$

$\boxed{\nabla_{b^{(l)}} \hat{y}_k^{(i)} \Rightarrow \quad e\, \hat{y}_l - \hat{y}_l^2}$

② when $l \neq k$

$\nabla_{b^{(l)}} \hat{y}_k^{(i)} = e^{z_k} \nabla_{b^{(l)}} \left( \dfrac{1}{\sum\limits_{k'=1}^{c} e^{z_{k'}}} \right)$

$= e^{z_k} \times \dfrac{-1}{\left( \sum\limits_{k'=1}^{c} e^{z_{k'}} \right)^2} \times e^{z_l}$

$\boxed{\nabla_{b^{(l)}} \hat{y}_k^{(i)} = - \hat{y}_k\, \hat{y}_l}$

$\therefore \nabla_{b^{(l)}} \mathcal{L}_{CE}(w, b) = \dfrac{-1}{n} \sum\limits_{i=1}^{n} \sum\limits_{k=1}^{c} y_k^{(i)} \times \left( \dfrac{\nabla_{b^{(l)}} \hat{y}_k^{(i)}}{\hat{y}_k} \right)$

$\therefore$ We can split $\boxed{\sum\limits_{k} a_k = a_l + \sum\limits_{k \neq l} a_k}$

$$\nabla_b^{(\ell)} l_{CE}(w,b) = \frac{-1}{n} \sum_{i=1}^{n} \left( \hat{y}_\ell^{(i)} \times \frac{(\hat{y}_\ell - \hat{y}_\ell^2)}{\hat{y}_\ell} + \sum_{k \neq \ell} y_k^{(i)} \times \frac{-\hat{y}_k \times \hat{y}_\ell}{\hat{y}_\ell} \right)$$

$$= \frac{-1}{n} \sum_{i=1}^{n} \left( y_\ell^{(i)} \times (1 - \hat{y}_\ell^{(i)}) + (-(1 - y_\ell^{(i)}) \times \hat{y}_\ell) \right)$$

$$\boxed{\because \sum_{k=1}^{c} y_k = 1 \implies \sum_{k \neq \ell}^{c} y_k = (1 - y_\ell^{(i)})}$$

$$\therefore \implies \frac{-1}{n} \sum_{i=1}^{n} \left( y_\ell^{(i)} - \cancel{y_\ell^{(i)}\hat{y}_\ell} - \hat{y}_\ell + \cancel{y_\ell^{(i)}\hat{y}_\ell^{(i)}} \right)$$

$$\therefore \nabla_b^{(\ell)} l_{CE}(w,b) = \frac{-1}{n} \sum_{i=1}^{n} \left( y_\ell^{(i)} - \hat{y}_\ell^{(i)} \right)$$

$\Rightarrow$ The gradient of cross-entropy loss wrt bias of class $\ell'$ is as shown above

$\Rightarrow$ Now collecting all the gradient wrt bias of classes $l = 1, \ldots, k$ we can represent this in vector form as :-

$$\nabla_b^{1} l_{CE}(w,b) = \frac{-1}{n} \sum_{i=1}^{n} \left( y_1^{(i)} - \hat{y}_1^{(i)} \right)$$
$$\vdots$$
$$\nabla_b^{c} l_{CE}(w,b) = \frac{-1}{n} \sum_{i=1}^{n} \left( y_c^{(i)} - \hat{y}_c^{(i)} \right)$$

$$\Rightarrow \quad \boxed{\nabla_b \, l_{CE}(w,b) = \frac{-1}{n} \sum_{i=1}^{n} \left( y^{(i)} - \hat{y}^{(i)} \right)}$$

## Problem 3: Logistic Sigmoid Identity

③ Derivation of Cross-Entropy as Negative Log-likelihood-

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^{n}$$

Now,

$$P(D|w,b) = P\left(\{(x^{(i)}, y^{(i)})\}_{i=1}^{n} \mid w, b\right)$$

$$= P(y^1|x^1, w, b) \cdot \prod_{i=2}^{n} P(y^i|x^i, w, b, y_2, y_3 \cdots y_n)$$

— assuming Conditional Independence

$$P(D|w,b) = \prod_{i=1}^{n} P(y_i|w,b)$$

$$P(D|w,b) = \prod_{i=1}^{n} P(y_i|w,b)$$

→ from the given eqn in the question the above eqn can be re-written as —

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

taking negative log on both sides.

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

$$= -\sum_{i=1}^{m} \left( \log \prod_{k=1}^{c} \hat{y}_k^{(i)}{}^{y_k^{(i)}} \right)$$

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

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

$$\boxed{-\log\left(P(D|w,b)\right) = f_{CE}(D; W, b)}$$

## Problem 4: Implementation of SoftMax Regression

- We developed a 2-layer SoftMax neural network using the Fashion MNIST dataset. The file is saved as **homework3_590146168_249812226.py.**
- The model is trained via stochastic gradient descent method to minimize cross-entropy loss.
- There are 4 Hyperparameters that can be altered:
  - Mini-Batch Size o Learning Rate o Number of Epochs o L$_2$ Regularization Strength
- The training data is split into Training (80%) and Validation (20%) dataset.
- We have selected a set of 8 values for each hyperparameter i.e., $8^4$ =4096 combinations.
- The weights and bias are initialized to zero.
- Once training is done, we calculate the cost on validation set. If the cost is minimum, the best hyperparameter is updated. This process is done till all hyperparameters are checked.
- It took approximately 30 hours to check all $8^4$ combinations.
- At the end we obtain the best set of hyperparameters. Using these hyperparameters we train the model on validation + training dataset.
- Now this trained model is used to calculate the cost (unregularized MSE) and percent correctly classified examples for the test dataset and the and the result is reported below:

```
(base) saammmy@saammmy-xps15:~/projects/3$ /usr/bin/python3 /home/saammmy/projects/3/neural_network.py
------------------------------------------------------------------
Performing Grid Search for 4096 combinations of Hyperparameters:
------------------------------------------------------------------
 Grid Search Completed
 Results after performing Grid Search:
 Best Hyperparameters:
   Epochs=  35
   Alpha=  1
   Learning Rate=  5e-07
   Mini Batch Size=  32
 Cost on Validation Set with Best Hyperparameters=  0.433649244826

------------------------------------------------------------------
Training on Training + Validation Dataset:
------------------------------------------------------------------
 Training Completed

------------------------------------------------------------------
Performance Evaluation
------------------------------------------------------------------
 Cost on Test Dataset=  0.455827235063
 Accuracy on Test Dataset= 84.26 %
```

## *End Of Assignment*