

CSE 676: Assignment #1

Name: Koushik Ponugoti
UB ID: kponugot
UBIT#: 50441728

0.1 Softmax

1) given,

$$\text{Softmax}(x)_i \triangleq \frac{e^{x_i}}{\sum_i e^{x_i}}$$

we need to prove,

$$\text{Softmax}(x) = \text{Softmax}(x+c)$$

$$\Rightarrow \text{Softmax}(x+c) = \frac{e^{x_i+c}}{\sum_i e^{x_i+c}}$$

$$= \frac{e^{x_i} * e^c}{e^c \cdot \sum_i e^{x_i}}$$

$$= \frac{e^{x_i}}{\sum_i e^{x_i}}$$

$$= \text{Softmax}(x)$$

$$\therefore \text{Softmax}(x) = \text{Softmax}(x+c)$$

2) given,

$$J = \sum_i \log \text{softmax}(z_i)$$

we need to find $\frac{\partial J}{\partial W}$, $\frac{\partial J}{\partial C}$

To find $\frac{\partial J}{\partial W}$, we can do that by using the chain Rule

$$\text{i.e.; } \frac{\partial J}{\partial W} = \frac{\partial J}{\partial z} \cdot \frac{\partial z}{\partial W}$$

$$\frac{\partial J}{\partial z} = \frac{\partial}{\partial z} (\log(\text{softmax}(z_i))) - \frac{\partial}{\partial z} \sum_{k \neq i} \log(\text{softmax}(z_k))$$

$$\Rightarrow \frac{\partial J}{\partial z} = - \left(\frac{1}{\text{softmax}(z_i)} \cdot \frac{\partial(\text{softmax}(z_i))}{\partial z} \right) - \sum_{k \neq i} y_k \left(\frac{1}{\text{softmax}(z_k)} \cdot \frac{\partial(\text{softmax}(z_k))}{\partial z} \right)$$

Here y_k is one-hot vector representing true label.

$$\frac{\partial J}{\partial z} = - (1 - (\text{softmax}(z_i))) + \sum_{k \neq i} (\text{softmax}(z_k))$$

$$\frac{\partial J}{\partial z} = -1 + (\text{softmax}(z_i)) + \sum_{k \neq i} (\text{softmax}(z_k))$$

$$\frac{\partial J(z)}{\partial z} = -1 + \left(\sum_{k=1}^N 1 \right) (\text{softmax}(z_i))$$

$$\frac{\partial J(z)}{\partial z} = -1 + N (\text{softmax}(z_i))$$

Now, $\frac{\partial J}{\partial w} = (-1 + N(\text{Softmax}(z_i))) \frac{\partial z}{\partial w}$

$$\frac{\partial J}{\partial w} = (-1 + N(\text{Softmax}(z_i))) x$$

we can follow same process for $\frac{\partial J}{\partial c}$

$$\frac{\partial J}{\partial c} = \frac{\partial J}{\partial z} \cdot \frac{\partial z}{\partial c}$$

$$\frac{\partial J}{\partial c} = (-1 + N(\text{Softmax}(z_i)))$$

0.2 - Logistic Regression with Regularization

1) given,

$$p(y_i; x_i, \theta) = \sigma(\theta^T x_i) \triangleq \frac{1}{1 + e^{-\theta^T x_i}}$$

let us suppose the above equation as $g_0(x_i)$ and the Regularization term is

$$R(\theta) = \frac{\lambda}{2} \theta^T \theta$$

we need to come up with our loss function, because if we use MSE as the loss function we will not get the convex function.

the loss function would be

$$\text{Cost}(g_{\theta}(x_i), y_i) = -y_i \log(h_{\theta}(x_i)) - (1-y_i) \log(1-h_{\theta}(x_i))$$

The objective function will be,

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N (\text{Cost}(h_{\theta}(x_i), y_i))$$

To make sure that the model does not overfit, we need to add the regularization term to the objective function.

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N \text{Cost}(g_{\theta}(x_i), y_i) + R(\theta)$$

- ② The gradient descent algorithm helps us to find the minima of the above defined objective function.

By finding the derivative of the ~~cost~~ cost function, we can get the gradient equation.

$$\frac{\partial \text{cost}(g_{\theta}(x_i), y_i)}{\partial \theta_j} = - \left(\frac{y_i}{g_{\theta}(x_i)} - \frac{1-y_i}{1-g_{\theta}(x_i)} \right) \frac{\partial g_{\theta}(x_i)}{\partial \theta_j}$$

$$\Rightarrow - \left(\frac{y_i}{g_{\theta}(x_i)} - \frac{1-y_i}{1-g_{\theta}(x_i)} \right) g_{\theta}(x_i) (1-g_{\theta}(x_i)) \frac{\partial \theta^T x_i}{\partial \theta_j}$$

$$\Rightarrow -(y_i(1-g_{\theta}(x_i)) - (1-y_i)g_{\theta}(x_i)) x_{ij}$$

$$\Rightarrow -(y_i - h_{\theta}(x_i)) x_{ij}$$

Now, using this, we calculate the complete gradient equation.

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \left(\frac{1}{N} \sum_{i=1}^N \text{cost}(g_{\theta}(x_i), y_i) + R(\theta) \right)$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{N} \sum_{i=1}^N \frac{\partial}{\partial \theta_j} \text{cost}(g_{\theta}(x_i), y_i) + \frac{\partial}{\partial \theta_j} R(\theta)$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = -\frac{1}{N} \sum_{i=1}^N (y_i - g_{\theta}(x_i)) x_{ij} + \frac{\partial}{\partial \theta_j} R(\theta)$$

Now, we need to find the derivative of $R(\theta)$

$$\frac{\partial R(\theta)}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \frac{\lambda}{2} \theta^T \theta$$

$$\frac{\partial}{\partial \theta_j} R(\theta) = \frac{\lambda}{2} \frac{\partial}{\partial \theta_j} (\theta^T \theta)$$

$$\frac{\partial}{\partial \theta_j} R(\theta) = \frac{\lambda}{2} (2\theta_j) = \lambda \theta_j$$

Now,

$$\frac{\partial J(\theta)}{\partial \theta_j} = -\frac{1}{N} \sum_{i=1}^N (y_i - g_{\theta}(x_i)) x_{ij} + \lambda \theta_j$$

The equation for updating weights is

$$\theta_j = \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

$$\theta_j = \theta_j - \alpha \left(-\frac{1}{N} \sum_{i=1}^N (y_i - g_{\theta}(x_i)) x_{ij} + \lambda \theta_j \right)$$

$$\Rightarrow \boxed{\theta_j = \theta_j \left(1 - \frac{\lambda}{N} \right) - \alpha \left(\sum_{i=1}^N (g_{\theta}(x_i) - y_i) x_{ij} \right)}$$

0.3 Derivative of the Softmax function.

1.) given, $J(z) = - \sum_{k=1}^K y_k \log \tilde{y}_k$

$$\frac{\partial J(z)}{\partial z} = - \frac{\partial}{\partial z} \sum_{k=1}^K y_k \log \tilde{y}_k$$

This is for class k , we need to consider the derivative for class i also (useful while using chain rule)

$$\frac{\partial J(z)}{\partial z} = - \frac{\partial}{\partial z} (y_i \log(\tilde{y}_i)) - \frac{\partial}{\partial z} \sum_{k \neq i} y_k \log \tilde{y}_k$$

$$\Rightarrow \frac{\partial J(z)}{\partial z} = - \left(y_i \frac{1}{\tilde{y}_i} \frac{\partial \tilde{y}_i}{\partial z} \right) - \sum_{k \neq i} y_k \frac{1}{\tilde{y}_k} \frac{\partial \tilde{y}_k}{\partial z}$$

$$\Rightarrow \frac{\partial J(z)}{\partial z} = - y_i + y_i \tilde{y}_i + \sum_{k \neq i} y_k \tilde{y}_i$$

$$\Rightarrow \frac{\partial J(z)}{\partial z} = - y_i + \left(\sum_{k=1}^N y_k \right) \tilde{y}_i$$

$$\Rightarrow \boxed{\frac{\partial J(z)}{\partial z} = - y_i + \tilde{y}_i}$$

② we are given the cross entropy loss function which is

$$J(z) = - \sum_{k=1}^K y_k \log \tilde{y}_k$$

we know the gradient for the last layer, we found this by computing the derivative of the above function w.r.t output of the final layer.

$$\text{i.e., } \frac{\partial J(z)}{\partial z} = -y_i + \tilde{y}_i \text{ as } z = W^T h + b.$$

In order to update W_{ij} which is the value of each weights, we have to find the value of
$$\frac{\partial J(z)}{\partial W_{ij}}$$

So, To find this, we need to propagate the error that ~~backwards~~ we get at the last layer and then we can use the chain rule in the calculus to propagate the error backwards.

Now, $\frac{\partial J(z)}{\partial W_{ij}}$ can be calculated for each weight.

when we calculate the gradient equation, based on the incorrectness of the classes, the incorrect classes also get affected which is not the property of the loss function.

③

we have $\frac{\partial J(z)}{\partial z} = -y_i + \tilde{y}_i$

we need to find $\frac{\partial J}{\partial w}$ and $\frac{\partial J}{\partial b}$
we can find $\frac{\partial J}{\partial w}$ by using chain rule.

$$\frac{\partial J(z)}{\partial w} = \frac{\partial J(z)}{\partial z} \cdot \frac{\partial z}{\partial w}$$

$$\frac{\partial J(z)}{\partial w} = (-y_i + \tilde{y}_i) \frac{\partial z}{\partial w}$$

$$\Rightarrow \frac{\partial J(z)}{\partial w} = (-y_i + \tilde{y}_i) \frac{\partial (w^T h + b)}{\partial w}$$

$$\Rightarrow \boxed{\frac{\partial J(z)}{\partial w} = (-y_i + \tilde{y}_i) h}$$

Similarly,

$$\frac{\partial J(z)}{\partial b} = (-y_i + \tilde{y}_i) \frac{\partial (w^T h + b)}{\partial b}$$

$$\boxed{\frac{\partial J(z)}{\partial b} = (-y_i + \tilde{y}_i)}$$

0.4 MNIST with FNN [30 points]

- 1) [10 points] Design an FNN for MNIST classification. Draw the computational graph of your model.
- 2) [20 points] Implement the model and plot two curves in one figure: i) training loss vs. training iterations; ii) test loss vs. training iterations.

The Github link for the above code is: <https://github.com/Koushik-24/CSE676-DeepLearning>

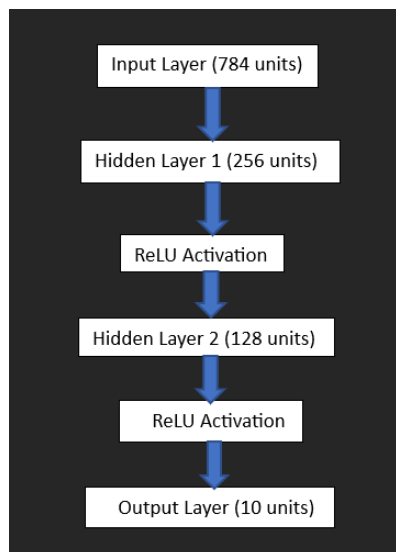


fig1 Computational Graph

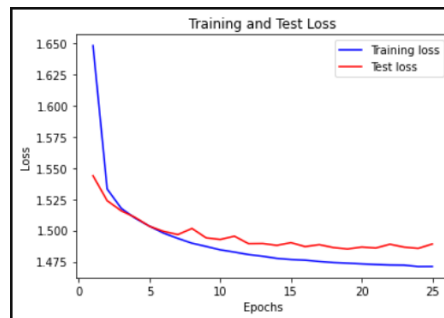


fig2 Training and Testing loss vs Epochs

References

1. <https://www.youtube.com/watch?v=PcU2RCWuDSI>
2. <https://www.kaggle.com/code/poojandabhi/fnn-mnist>
3. <https://www.baeldung.com/cs/training-validation-loss-deep-learning>