
Assignment 1 Report

Indian Institute of Technology, Kanpur
Department of Computer Science
CS776 Deep Learning for Computer Vision
Instructor: Priyanka Bagade

Submitted: February 3, 2023 Submitted by :

Chunduri Naga Venkata Hrushikesh
22111280

Gnanendra Sri Phani Sai Channamsetty
22105032

Contents

1	Multilayer Perceptron that has been implemented	3
1.1	Model Description	3
1.2	MLP Diagram	3
1.3	MLP Details	3
1.4	MLP Model Procedure	4
2	Derivations	5
2.1	Activation Functions Derivations	6
2.1.1	ReLU Function	6
2.1.2	Derivative of ReLU Function	6
2.1.3	Softmax Function	6
2.1.4	Derivative of Softmax Function	6
2.2	Cost Function Derivative for initial step of Backpropagation	9
2.3	Back Propagation Derivation	10
2.4	Update Expressions	13
3	Hyperparameter Tuning	14
3.1	Loss vs Epoch Graph for Unaugmented Data Set	14
3.2	Loss vs Epoch Graph for Augmented Data Set	14
4	Evaluation Metrics	15
4.1	Accuracy	15
4.2	Time taken for prediction	15
4.3	Justifications for difference in accuracy and time taken	16
5	References	16

1 Multilayer Perceptron that has been implemented

1.1 Model Description

Implemented MLP, shown in next subsection below, consists of

- 1) Input layer has 512 neurons, whose input is of dimensions (512 x 50000) in case of original data set i.e $m = 50000$ and (512 x 100000) in case of augmented data set i.e $m = 100000$ where m = Number of training examples
- 2) Two hidden layers, each with 64 neurons .
- 3) Output layer has 10 neurons each, representing 10 types of image classes of CIFAR-10 Data Set .

1.2 MLP Diagram

The following figure has been made by us, using LaTeX Code

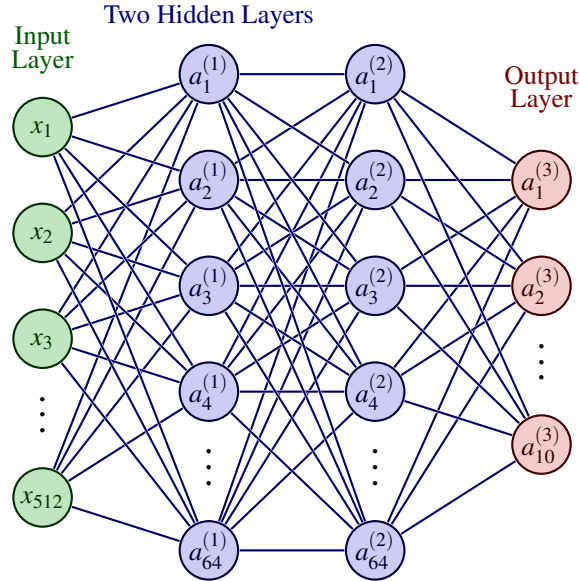


Figure 1: Implemented MLP Model

1.3 MLP Details

Table 1: MLP Details

Layer Name () represents code	No of Neurons $n^{[l]}$	Weight Dims $W^{[l]}$	Bias Dims $b^{[l]}$	Activation Function(Dims) $A^{[l]}(m = \text{No of Training Examples})$
Input Layer (0^{th})	512	-	-	-
Hidden Layer 1(1^{st})	64	64 x 512	64 x 1	RELU(64 x m)
Hidden Layer 2(2^{nd})	64	64 x 64	64 x 1	RELU(64 x m)
Output Layer(3^{rd})	10	10 x 64	10 x 1	Softmax(10 x m)

1.4 MLP Model Procedure

It consists of 5 following steps :

- 1) **Initialization Step** : Weights have been initialized with small random numbers and bias with zeros. This step has been represented in code using "" Function
- 2) **Forward Propagation Step** : Calculate the linear part of a layer's forward propagation step (resulting in $Z^{[l]}$ in code) . Then apply Activation Function RELU on the previous resultant for first two layers and then apply Activation Function Softmax for the last output layer (resulting in $A^{[l]}$ in code) . At every step of forward function, some values are stored in a cache called "". These stored values are useful for computing gradients in step 4 .This step has been represented in code using "" Function representing the following equations.

$$z_{j,i}^{[l]} = \sum_k w_{j,k}^{[l]} a_{k,i}^{[l-1]} + b_j^{[l]}, \quad (1)$$

$$a_{j,i}^{[l]} = g_j^{[l]}(z_{1,i}^{[l]}, \dots, z_{j,i}^{[l]}, \dots, z_{n^{[l]},i}^{[l]}). \quad (2)$$

Using numpy,Above two equations can be vectorized as following

$$Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$$

where $A^{[0]} = X$.

$$A^{[l]} = g(Z^{[l]}) = g(W^{[l]} A^{[l-1]} + b^{[l]})$$

where $g(Z^{[l]})$ could be RELU or Softmax depending on layer number.

- 3) **Estimating Overall Cost Function Step** : Cross Entropy has been selected as Cost Function to find deviation of prediction from actual label. This step has been represented in code using "" using the following formula:

$$J = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log(a^{[L](i)}) + (1 - y^{(i)}) \log(1 - a^{[L](i)}))$$

- 4) **Backward Propagation Step** : The following equations have been used in code. Their derivations are in next section.

$$dW^{[l]} = \frac{\partial J}{\partial W^{[l]}} = \frac{1}{m} dZ^{[l]} A^{[l-1]T}$$

$$db^{[l]} = \frac{\partial J}{\partial b^{[l]}} = \frac{1}{m} \sum_{i=1}^m dZ^{[l](i)}$$

$$dA^{[l-1]} = \frac{\partial J}{\partial A^{[l-1]}} = W^{[l]T} dZ^{[l]}$$

$$dZ^{[l]} = dA^{[l]} * g'(Z^{[l]})$$

$$dAL = \frac{\partial J}{\partial a_{j,i}^{[L]}} = \frac{1}{m} \left(\frac{1 - y_{j,i}}{1 - a_{j,i}^{[L]}} - \frac{y_{j,i}}{a_{j,i}^{[L]}} \right),$$

- 5) **Update Expressions** : Gradient Descent algorithm is being used to update expressions using following equations :

$$W^{[l]} = W^{[l]} - \alpha dW^{[l]}$$

$$b^{[l]} = b^{[l]} - \alpha db^{[l]}$$

where α is the learning rate.

2 Derivations

Chain Rule for derivatives can be given by :

$$\frac{\partial y_k}{\partial x_j} = \sum_i \frac{\partial y_k}{\partial u_i} \frac{\partial u_i}{\partial x_j}. \quad (3)$$

Additional Notation :

$$\vec{A}^{[0]} = X_{input} \quad (4)$$

$$\vec{A}^{[L]} = Y_{predicted} \quad (5)$$

Writing equations (1) and (2) in matrix form :

$$\begin{bmatrix} z_{1,i}^{[l]} \\ \vdots \\ z_{j,i}^{[l]} \\ \vdots \\ z_{n^{[l]},i}^{[l]} \end{bmatrix} = \begin{bmatrix} w_{1,1}^{[l]} & \cdots & w_{1,k}^{[l]} & \cdots & w_{1,n^{[l-1]}}^{[l]} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ w_{j,1}^{[l]} & \cdots & w_{j,k}^{[l]} & \cdots & w_{j,n^{[l-1]}}^{[l]} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ w_{n^{[l]},1}^{[l]} & \cdots & w_{n^{[l]},k}^{[l]} & \cdots & w_{n^{[l]},n^{[l-1]}}^{[l]} \end{bmatrix} \begin{bmatrix} a_{1,i}^{[l-1]} \\ \vdots \\ a_{k,i}^{[l-1]} \\ \vdots \\ a_{n^{[l-1]},i}^{[l-1]} \end{bmatrix} + \begin{bmatrix} b_1^{[l]} \\ \vdots \\ b_j^{[l]} \\ \vdots \\ b_{n^{[l]}}^{[l]} \end{bmatrix},$$

$$\begin{bmatrix} a_{1,i}^{[l]} \\ \vdots \\ a_{j,i}^{[l]} \\ \vdots \\ a_{n^{[l]},i}^{[l]} \end{bmatrix} = \begin{bmatrix} g_1^{[l]}(z_{1,i}^{[l]}, \dots, z_{j,i}^{[l]}, \dots, z_{n^{[l]},i}^{[l]}) \\ \vdots \\ g_j^{[l]}(z_{1,i}^{[l]}, \dots, z_{j,i}^{[l]}, \dots, z_{n^{[l]},i}^{[l]}) \\ \vdots \\ g_{n^{[l]}}^{[l]}(z_{1,i}^{[l]}, \dots, z_{j,i}^{[l]}, \dots, z_{n^{[l]},i}^{[l]}) \end{bmatrix},$$

2.1 Activation Functions Derivations

2.1.1 ReLU Function

$$\begin{aligned} a_{j,i}^{[l]} &= g_j^{[l]}(z_{1,i}^{[l]}, \dots, z_{j,i}^{[l]}, \dots, z_{n^{[l]},i}^{[l]}) \\ &= \max(0, z_{j,i}^{[l]}) \\ &= \begin{cases} z_{j,i}^{[l]} & \text{if } z_{j,i}^{[l]} > 0, \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

After vectorization with numpy arrays, RELU function becomes

$$\vec{A}^{[l]} = \max(0, \vec{Z}^{[l]}). \quad (6)$$

2.1.2 Derivative of ReLU Function

$$\begin{aligned} \frac{\partial a_{j,i}^{[l]}}{\partial z_{j,i}^{[l]}} &= \begin{cases} 1 & \text{if } z_{j,i}^{[l]} > 0, \\ 0 & \text{otherwise,} \end{cases} \\ &= I(z_{j,i}^{[l]} > 0), \\ \frac{\partial a_{p,i}^{[l]}}{\partial z_{j,i}^{[l]}} &= 0, \quad \forall p \neq j. \end{aligned}$$

Using above, we can differentiate Cost J w.r.t $z_{j,i}^{[l]}$

$$\begin{aligned} \frac{\partial J}{\partial z_{j,i}^{[l]}} &= \sum_p \frac{\partial J}{\partial a_{p,i}^{[l]}} \frac{\partial a_{p,i}^{[l]}}{\partial z_{j,i}^{[l]}} \\ &= \frac{\partial J}{\partial a_{j,i}^{[l]}} \frac{\partial a_{j,i}^{[l]}}{\partial z_{j,i}^{[l]}} + \sum_{p \neq j} \frac{\partial J}{\partial a_{p,i}^{[l]}} \frac{\partial a_{p,i}^{[l]}}{\partial z_{j,i}^{[l]}} \\ &= \frac{\partial J}{\partial a_{j,i}^{[l]}} I(z_{j,i}^{[l]} > 0), \end{aligned}$$

2.1.3 Softmax Function

$$\begin{aligned} a_{j,i}^{[l]} &= g_j^{[l]}(z_{1,i}^{[l]}, \dots, z_{j,i}^{[l]}, \dots, z_{n^{[l]},i}^{[l]}) \\ &= \frac{\exp(z_{j,i}^{[l]})}{\sum_p \exp(z_{p,i}^{[l]})}. \end{aligned}$$

2.1.4 Derivative of Softmax Function

After constructing a computation graph, for j^{th} activation of present layer

$$\begin{aligned}
u_{-1} &= z_{j,i}^{[l]}, \\
u_{0,p} &= z_{p,i}^{[l]}, & \forall p \neq j, \\
u_1 &= \exp(u_{-1}), \\
u_{2,p} &= \exp(u_{0,p}), & \forall p \neq j, \\
u_3 &= u_1 + \sum_{p \neq j} u_{2,p}, \\
u_4 &= \frac{1}{u_3}, \\
u_5 &= u_1 u_4 = a_{j,i}^{[l]}.
\end{aligned}$$

After applying chain rule,

$$\begin{aligned}
\frac{\partial a_{j,i}^{[l]}}{\partial u_5} &= 1, \\
\frac{\partial a_{j,i}^{[l]}}{\partial u_4} &= \frac{\partial a_{j,i}^{[l]}}{\partial u_5} \frac{\partial u_5}{\partial u_4} = u_1 = \exp(z_{j,i}^{[l]}), \\
\frac{\partial a_{j,i}^{[l]}}{\partial u_3} &= \frac{\partial a_{j,i}^{[l]}}{\partial u_4} \frac{\partial u_4}{\partial u_3} = -u_1 \frac{1}{u_3^2} = -\frac{\exp(z_{j,i}^{[l]})}{(\sum_p \exp(z_{p,i}^{[l]}))^2}, \\
\frac{\partial a_{j,i}^{[l]}}{\partial u_1} &= \frac{\partial a_{j,i}^{[l]}}{\partial u_3} \frac{\partial u_3}{\partial u_1} + \frac{\partial a_{j,i}^{[l]}}{\partial u_5} \frac{\partial u_5}{\partial u_1} \\
&= -u_1 \frac{1}{u_3^2} + u_4 \\
&= -\frac{\exp(z_{j,i}^{[l]})}{(\sum_p \exp(z_{p,i}^{[l]}))^2} + \frac{1}{\sum_p \exp(z_{p,i}^{[l]})}, \\
\frac{\partial a_{j,i}^{[l]}}{\partial u_{-1}} &= \frac{\partial a_{j,i}^{[l]}}{\partial u_1} \frac{\partial u_1}{\partial u_{-1}} \\
&= \left(-u_1 \frac{1}{u_3^2} + u_4 \right) \exp(u_{-1}) \\
&= -\frac{\exp(z_{j,i}^{[l]})^2}{(\sum_p \exp(z_{p,i}^{[l]}))^2} + \frac{\exp(z_{j,i}^{[l]})}{\sum_p \exp(z_{p,i}^{[l]})}.
\end{aligned}$$

But $z_{j,i}^{[l]}$ also affects other activations in the same layer. So, taking derivative for that :

$$\begin{aligned}
u_{-1} &= z_{j,i}^{[l]}, \\
u_{0,p} &= z_{p,i}^{[l]}, & \forall p \neq j, \\
u_1 &= \exp(u_{-1}), \\
u_{2,p} &= \exp(u_{0,p}), & \forall p \neq j, \\
u_3 &= u_1 + \sum_{p \neq j} u_{2,p}, \\
u_4 &= \frac{1}{u_3}, \\
u_5 &= u_{2,p} u_4 = a_{p,i}^{[l]}, & \forall p \neq j.
\end{aligned}$$

Applying backward propagation :

$$\begin{aligned}
\frac{\partial a_{p,i}^{[l]}}{\partial u_5} &= 1, \\
\frac{\partial a_{p,i}^{[l]}}{\partial u_4} &= \frac{\partial a_{p,i}^{[l]}}{\partial u_5} \frac{\partial u_5}{\partial u_4} = u_{2,p} = \exp(z_{p,i}^{[l]}), \\
\frac{\partial a_{p,i}^{[l]}}{\partial u_3} &= \frac{\partial a_{p,i}^{[l]}}{\partial u_4} \frac{\partial u_4}{\partial u_3} = -u_{2,p} \frac{1}{u_3^2} = -\frac{\exp(z_{p,i}^{[l]})}{(\sum_p \exp(z_{p,i}^{[l]}))^2}, \\
\frac{\partial a_{p,i}^{[l]}}{\partial u_1} &= \frac{\partial a_{p,i}^{[l]}}{\partial u_3} \frac{\partial u_3}{\partial u_1} = -u_{2,p} \frac{1}{u_3^2} = -\frac{\exp(z_{p,i}^{[l]})}{(\sum_p \exp(z_{p,i}^{[l]}))^2}, \\
\frac{\partial a_{p,i}^{[l]}}{\partial u_{-1}} &= \frac{\partial a_{p,i}^{[l]}}{\partial u_1} \frac{\partial u_1}{\partial u_{-1}} = -u_{2,p} \frac{1}{u_3^2} \exp(u_{-1}) = -\frac{\exp(z_{p,i}^{[l]}) \exp(z_{j,i}^{[l]})}{(\sum_p \exp(z_{p,i}^{[l]}))^2}.
\end{aligned}$$

From above, we know

$$\begin{aligned}
\frac{\partial a_{j,i}^{[l]}}{\partial z_{j,i}^{[l]}} &= -\frac{\exp(z_{j,i}^{[l]})^2}{(\sum_p \exp(z_{p,i}^{[l]}))^2} + \frac{\exp(z_{j,i}^{[l]})}{\sum_p \exp(z_{p,i}^{[l]})} \\
&= a_{j,i}^{[l]}(1 - a_{j,i}^{[l]}), \\
\frac{\partial a_{p,i}^{[l]}}{\partial z_{j,i}^{[l]}} &= -\frac{\exp(z_{p,i}^{[l]}) \exp(z_{j,i}^{[l]})}{(\sum_p \exp(z_{p,i}^{[l]}))^2} \\
&= -a_{p,i}^{[l]} a_{j,i}^{[l]}, \quad \forall p \neq j.
\end{aligned}$$

Using above, we can differentiate Cost J w.r.t $z_{j,i}^{[l]}$

$$\begin{aligned}
\frac{\partial J}{\partial z_{j,i}^{[l]}} &= \sum_p \frac{\partial J}{\partial a_{p,i}^{[l]}} \frac{\partial a_{p,i}^{[l]}}{\partial z_{j,i}^{[l]}} \\
&= \frac{\partial J}{\partial a_{j,i}^{[l]}} \frac{\partial a_{j,i}^{[l]}}{\partial z_{j,i}^{[l]}} + \sum_{p \neq j} \frac{\partial J}{\partial a_{p,i}^{[l]}} \frac{\partial a_{p,i}^{[l]}}{\partial z_{j,i}^{[l]}} \\
&= \frac{\partial J}{\partial a_{j,i}^{[l]}} a_{j,i}^{[l]}(1 - a_{j,i}^{[l]}) - \sum_{p \neq j} \frac{\partial J}{\partial a_{p,i}^{[l]}} a_{p,i}^{[l]} a_{j,i}^{[l]} \\
&= a_{j,i}^{[l]} \left(\frac{\partial J}{\partial a_{j,i}^{[l]}} (1 - a_{j,i}^{[l]}) - \sum_{p \neq j} \frac{\partial J}{\partial a_{p,i}^{[l]}} a_{p,i}^{[l]} \right) \\
&= a_{j,i}^{[l]} \left(\frac{\partial J}{\partial a_{j,i}^{[l]}} (1 - a_{j,i}^{[l]}) - \sum_p \frac{\partial J}{\partial a_{p,i}^{[l]}} a_{p,i}^{[l]} + \frac{\partial J}{\partial a_{j,i}^{[l]}} a_{j,i}^{[l]} \right) \\
&= a_{j,i}^{[l]} \left(\frac{\partial J}{\partial a_{j,i}^{[l]}} - \sum_p \frac{\partial J}{\partial a_{p,i}^{[l]}} a_{p,i}^{[l]} \right),
\end{aligned}$$

We can use numpy to vectorize above as

$$\frac{\partial J}{\partial \vec{z}_{:,i}^{[l]}} = \vec{a}_{:,i}^{[l]} \odot \left(\frac{\partial J}{\partial \vec{a}_{:,i}^{[l]}} - \underbrace{\vec{a}_{:,i}^{[l]} \frac{\partial J}{\partial \vec{a}_{:,i}^{[l]}}}_{\text{scalar}} \right).$$

where \odot = Element Wise Multiplication

2.2 Cost Function Derivative for initial step of Backpropagation

Cross Entropy Cost Function "J" is given by

$$\begin{aligned}
 J &= f(\vec{Y}, \vec{Y}) = f(\vec{A}^{[L]}, \vec{Y}) \\
 &= \sum_j \left(-\frac{1}{m} \sum_i (y_{j,i} \log(\hat{y}_{j,i}) + (1 - y_{j,i}) \log(1 - \hat{y}_{j,i})) \right) \\
 &= \sum_j \left(-\frac{1}{m} \sum_i (y_{j,i} \log(a_{j,i}^{[L]}) + (1 - y_{j,i}) \log(1 - a_{j,i}^{[L]})) \right),
 \end{aligned}$$

where $j = 1, \dots, n^{[L]}$

After making computation graph,

$$\begin{aligned}
 u_{0,j,i} &= a_{j,i}^{[L]}, \\
 u_{1,j,i} &= 1 - u_{0,j,i}, \\
 u_{2,j,i} &= \log(u_{0,j,i}), \\
 u_{3,j,i} &= \log(u_{1,j,i}), \\
 u_{4,j,i} &= y_{j,i} u_{2,j,i} + (1 - y_{j,i}) u_{3,j,i}, \\
 u_{5,j} &= -\frac{1}{m} \sum_i u_{4,j,i}, \\
 u_6 &= \sum_j u_{5,j} = J.
 \end{aligned}$$

Using above, finding partial derivatives :

$$\begin{aligned}
 \frac{\partial J}{\partial u_6} &= 1, \\
 \frac{\partial J}{\partial u_{5,j}} &= \frac{\partial J}{\partial u_6} \frac{\partial u_6}{\partial u_{5,j}} = 1, \\
 \frac{\partial J}{\partial u_{4,j,i}} &= \frac{\partial J}{\partial u_{5,j}} \frac{\partial u_{5,j}}{\partial u_{4,j,i}} = -\frac{1}{m}, \\
 \frac{\partial J}{\partial u_{3,j,i}} &= \frac{\partial J}{\partial u_{4,j,i}} \frac{\partial u_{4,j,i}}{\partial u_{3,j,i}} = -\frac{1}{m} (1 - y_{j,i}), \\
 \frac{\partial J}{\partial u_{2,j,i}} &= \frac{\partial J}{\partial u_{4,j,i}} \frac{\partial u_{4,j,i}}{\partial u_{2,j,i}} = -\frac{1}{m} y_{j,i}, \\
 \frac{\partial J}{\partial u_{1,j,i}} &= \frac{\partial J}{\partial u_{3,j,i}} \frac{\partial u_{3,j,i}}{\partial u_{1,j,i}} = -\frac{1}{m} (1 - y_{j,i}) \frac{1}{u_{1,j,i}} = -\frac{1}{m} \frac{1 - y_{j,i}}{1 - a_{j,i}^{[L]}}, \\
 \frac{\partial J}{\partial u_{0,j,i}} &= \frac{\partial J}{\partial u_{1,j,i}} \frac{\partial u_{1,j,i}}{\partial u_{0,j,i}} + \frac{\partial J}{\partial u_{2,j,i}} \frac{\partial u_{2,j,i}}{\partial u_{0,j,i}} \\
 &= \frac{1}{m} (1 - y_{j,i}) \frac{1}{u_{1,j,i}} - \frac{1}{m} y_{j,i} \frac{1}{u_{0,j,i}} \\
 &= \frac{1}{m} \left(\frac{1 - y_{j,i}}{1 - a_{j,i}^{[L]}} - \frac{y_{j,i}}{a_{j,i}^{[L]}} \right).
 \end{aligned}$$

Summarising above as :

$$\frac{\partial J}{\partial a_{j,i}^{[L]}} = \frac{1}{m} \left(\frac{1 - y_{j,i}}{1 - a_{j,i}^{[L]}} - \frac{y_{j,i}}{a_{j,i}^{[L]}} \right),$$

Using numpy, Vectorising previous equations as :

$$\frac{\partial J}{\partial \vec{A}^{[L]}} = \frac{1}{m} \left(\frac{1}{1 - \vec{A}^{[L]}} \odot (1 - \vec{Y}) - \frac{1}{\vec{A}^{[L]}} \odot \vec{Y} \right). \quad (7)$$

Now finding derivative of cost function w.r.t $a_{j,i}^{[L]}$

$$\begin{aligned}
\frac{\partial J}{\partial z_{j,i}^{[L]}} &= \frac{\partial J}{\partial a_{j,i}^{[L]}} a_{j,i}^{[L]} (1 - a_{j,i}^{[L]}) \\
&= \frac{1}{m} \left(\frac{1 - y_{j,i}}{1 - a_{j,i}^{[L]}} - \frac{y_{j,i}}{a_{j,i}^{[L]}} \right) a_{j,i}^{[L]} (1 - a_{j,i}^{[L]}) \\
&= \frac{1}{m} ((1 - y_{j,i}) a_{j,i}^{[L]} - y_{j,i} (1 - a_{j,i}^{[L]})) \\
&= \frac{1}{m} (a_{j,i}^{[L]} - y_{j,i}),
\end{aligned}$$

Vectorizing above equation with numpy as

$$\frac{\partial J}{\partial \vec{Z}^{[L]}} = \frac{1}{m} (\vec{A}^{[L]} - \vec{Y}). \quad (8)$$

2.3 Back Propagation Derivation

Applying Chain Rule (3)

$$\frac{\partial J}{\partial w_{j,k}^{[l]}} = \sum_i \frac{\partial J}{\partial z_{j,i}^{[l]}} \frac{\partial z_{j,i}^{[l]}}{\partial w_{j,k}^{[l]}} = \sum_i \frac{\partial J}{\partial z_{j,i}^{[l]}} a_{k,i}^{[l-1]}, \quad (9)$$

$$\frac{\partial J}{\partial b_j^{[l]}} = \sum_i \frac{\partial J}{\partial z_{j,i}^{[l]}} \frac{\partial z_{j,i}^{[l]}}{\partial b_j^{[l]}} = \sum_i \frac{\partial J}{\partial z_{j,i}^{[l]}}. \quad (10)$$

Vectorization of above two equations results in :

$$\begin{aligned}
& \begin{bmatrix} Jw_{1,1}^{[l]} & \dots & Jw_{1,k}^{[l]} & \dots & Jw_{1,n^{[l]-1}}^{[l]} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ Jw_{j,1}^{[l]} & \dots & Jw_{j,k}^{[l]} & \dots & Jw_{j,n^{[l]-1}}^{[l]} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ Jw_{n^{[l]},1}^{[l]} & \dots & Jw_{n^{[l]},k}^{[l]} & \dots & Jw_{n^{[l]},n^{[l]-1}}^{[l]} \end{bmatrix} \\
&= \begin{bmatrix} Jz_{1,1}^{[l]} & \dots & Jz_{1,i}^{[l]} & \dots & Jz_{1,m}^{[l]} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ Jz_{j,1}^{[l]} & \dots & Jz_{j,i}^{[l]} & \dots & Jz_{j,m}^{[l]} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ Jz_{n^{[l]},1}^{[l]} & \dots & Jz_{n^{[l]},i}^{[l]} & \dots & Jz_{n^{[l]},m}^{[l]} \end{bmatrix} \\
&\cdot \begin{bmatrix} a_{1,1}^{[l-1]} & \dots & a_{k,1}^{[l-1]} & \dots & a_{n^{[l]-1},1}^{[l-1]} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ a_{1,i}^{[l-1]} & \dots & a_{k,i}^{[l-1]} & \dots & a_{n^{[l]-1},i}^{[l-1]} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ a_{1,m}^{[l-1]} & \dots & a_{k,m}^{[l-1]} & \dots & a_{n^{[l]-1},m}^{[l-1]} \end{bmatrix}, \\
&\begin{bmatrix} Jb_1^{[l]} \\ \vdots \\ Jb_j^{[l]} \\ \vdots \\ Jb_{n^{[l]}}^{[l]} \end{bmatrix} = \begin{bmatrix} Jz_{1,1}^{[l]} \\ \vdots \\ Jz_{j,1}^{[l]} \\ \vdots \\ Jz_{n^{[l]},1}^{[l]} \end{bmatrix} + \dots + \begin{bmatrix} Jz_{1,i}^{[l]} \\ \vdots \\ Jz_{j,i}^{[l]} \\ \vdots \\ Jz_{n^{[l]},i}^{[l]} \end{bmatrix} + \dots + \begin{bmatrix} Jz_{1,m}^{[l]} \\ \vdots \\ Jz_{j,m}^{[l]} \\ \vdots \\ Jz_{n^{[l]},m}^{[l]} \end{bmatrix},
\end{aligned}$$

Preceding matrix equations can be compressed into :

$$\frac{\partial J}{\partial \vec{W}^{[l]}} = \sum_i \frac{\partial J}{\partial \vec{z}_{:,i}^{[l]}} \vec{a}_{:,i}^{[l-1]} = \frac{\partial J}{\partial \vec{Z}^{[l]}} \vec{A}^{[l-1]}, \quad (11)$$

$$\frac{\partial J}{\partial \vec{b}^{[l]}} = \sum_i \frac{\partial J}{\partial \vec{z}_{:,i}^{[l]}} = \underbrace{\sum_{\text{axis}=1} \frac{\partial J}{\partial \vec{Z}^{[l]}}}_{\text{column vector}}, \quad (12)$$

where dimensions are :

$$\frac{\partial J}{\partial \vec{z}_{:,i}^{[l]}} \in n^{[l]}, \frac{\partial J}{\partial \vec{Z}^{[l]}} \in n^{[l]} \times m, \frac{\partial J}{\partial \vec{W}^{[l]}} \in n^{[l]} \times n^{[l-1]} \text{ and } \frac{\partial J}{\partial \vec{b}^{[l]}} \in R^{n^{[l]}}$$

The only unknown partial derivative can be found through Chain Rule (3) :

$$\frac{\partial J}{\partial z_{j,i}^{[l]}} = \sum_p \frac{\partial J}{\partial a_{p,i}^{[l]}} \frac{\partial a_{p,i}^{[l]}}{\partial z_{j,i}^{[l]}}, \quad (13)$$

where $p = 1, \dots, n^{[l]}$

After vectorising previous equation :

$$\begin{bmatrix} Jz_{1,i}^{[l]} \\ \vdots \\ Jz_{j,i}^{[l]} \\ \vdots \\ Jz_{n^{[l]},i}^{[l]} \end{bmatrix} = \begin{bmatrix} a_{1,i}^{[l]} z_{1,i}^{[l]} & \dots & a_{j,i}^{[l]} z_{1,i}^{[l]} & \dots & a_{n^{[l]},i}^{[l]} z_{1,i}^{[l]} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ a_{1,i}^{[l]} z_{j,i}^{[l]} & \dots & a_{j,i}^{[l]} z_{j,i}^{[l]} & \dots & a_{n^{[l]},i}^{[l]} z_{j,i}^{[l]} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ a_{1,i}^{[l]} z_{n^{[l]},i}^{[l]} & \dots & a_{j,i}^{[l]} z_{n^{[l]},i}^{[l]} & \dots & a_{n^{[l]},i}^{[l]} z_{n^{[l]},i}^{[l]} \end{bmatrix} \begin{bmatrix} Ja_{1,i}^{[l]} \\ \vdots \\ Ja_{j,i}^{[l]} \\ \vdots \\ Ja_{n^{[l]},i}^{[l]} \end{bmatrix},$$

which compresses into

$$\frac{\partial J}{\partial \vec{z}_{:,i}^{[l]}} = \frac{\partial \vec{a}_{:,i}^{[l]}}{\partial \vec{z}_{:,i}^{[l]}} \frac{\partial J}{\partial \vec{a}_{:,i}^{[l]}}, \quad (14)$$

where $\frac{\partial J}{\partial \vec{a}_{:,i}^{[l]}} \in \mathbb{R}^{n^{[l]}}$ and $\frac{\partial \vec{a}_{:,i}^{[l]}}{\partial \vec{z}_{:,i}^{[l]}} \in \mathbb{R}^{n^{[l]} \times n^{[l]}}$

We already know from previous equations that :

$$\frac{\partial J}{\partial \vec{Z}^{[l]}} = \begin{bmatrix} J\vec{z}_{:,1}^{[l]} & \dots & J\vec{z}_{:,i}^{[l]} & \dots & J\vec{z}_{:,m}^{[l]} \end{bmatrix}, \quad (15)$$

$$\frac{\partial J}{\partial \vec{A}^{[l]}} = \begin{bmatrix} J\vec{a}_{:,1}^{[l]} & \dots & J\vec{a}_{:,i}^{[l]} & \dots & J\vec{a}_{:,m}^{[l]} \end{bmatrix}, \quad (16)$$

where : $\frac{\partial J}{\partial \vec{A}^{[l]}} \in \mathbb{R}^{n^{[l]} \times m}$

$\frac{\partial a_{j,i}^{[l]}}{\partial z_{j,i}^{[l]}}$ has already been derived previously in page 8.

Remaining $\frac{\partial J}{\partial z_{j,i}^{[l]}}$ is dependant on $\frac{\partial J}{\partial a_{j,i}^{[l]}}$ which is computed as follows :

$$\frac{\partial J}{\partial a_{k,i}^{[l-1]}} = \sum_j \frac{\partial J}{\partial z_{j,i}^{[l]}} \frac{\partial z_{j,i}^{[l]}}{\partial a_{k,i}^{[l-1]}} = \sum_j \frac{\partial J}{\partial z_{j,i}^{[l]}} w_{j,k}^{[l]}. \quad (17)$$

Vectorizing equation (17) as :

$$\begin{aligned}
& \begin{bmatrix} Ja_{1,1}^{[l-1]} & \dots & Ja_{1,i}^{[l-1]} & \dots & Ja_{1,m}^{[l-1]} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ Ja_{k,1}^{[l-1]} & \dots & Ja_{k,i}^{[l-1]} & \dots & Ja_{k,m}^{[l-1]} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ Ja_{n^{[l-1]},1}^{[l-1]} & \dots & Ja_{n^{[l-1]},i}^{[l-1]} & \dots & Ja_{n^{[l-1]},m}^{[l-1]} \end{bmatrix} \\
&= \begin{bmatrix} w_{1,1}^{[l]} & \dots & w_{j,1}^{[l]} & \dots & w_{n^{[l]},1}^{[l]} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ w_{1,k}^{[l]} & \dots & w_{j,k}^{[l]} & \dots & w_{n^{[l]},k}^{[l]} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ w_{1,n^{[l-1]}}^{[l]} & \dots & w_{j,n^{[l-1]}}^{[l]} & \dots & w_{n^{[l]},n^{[l-1]}}^{[l]} \end{bmatrix} \\
&\cdot \begin{bmatrix} Jz_{1,1}^{[l]} & \dots & Jz_{1,i}^{[l]} & \dots & Jz_{1,m}^{[l]} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ Jz_{j,1}^{[l]} & \dots & Jz_{j,i}^{[l]} & \dots & Jz_{j,m}^{[l]} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ Jz_{n^{[l]},1}^{[l]} & \dots & Jz_{n^{[l]},i}^{[l]} & \dots & Jz_{n^{[l]},m}^{[l]} \end{bmatrix},
\end{aligned}$$

Preceding matrix equations can be re written as :

$$dA^{[l-1]} = \frac{\partial J}{\partial A^{[l-1]}} = W^{[l]T} dZ^{[l]} \quad (18)$$

where $\frac{\partial J}{\partial A^{[l-1]}} \in R^{n^{[l-1]} \times m}$

2.4 Update Expressions

Gradient Descent algorithm is being used to update expressions using following equations :

$$\begin{aligned}
W^{[l]} &= W^{[l]} - \alpha dW^{[l]} \\
b^{[l]} &= b^{[l]} - \alpha db^{[l]}
\end{aligned}$$

where α is the learning rate.

3 Hyperparameter Tuning

Selected learning Rate of $\alpha = 0.01$ and $epochs = 20$ has been found through linear search. Also the reason to choose 20 epochs is loss is not converging beyond it.

3.1 Loss vs Epoch Graph for Unaugmented Data Set

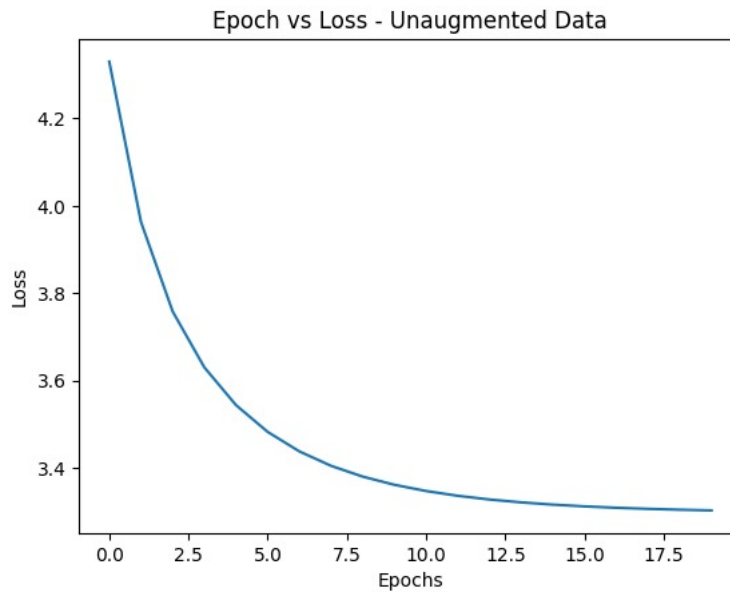


Figure 2: Loss vs Epoch Graph for Original Data Set

3.2 Loss vs Epoch Graph for Augmented Data Set

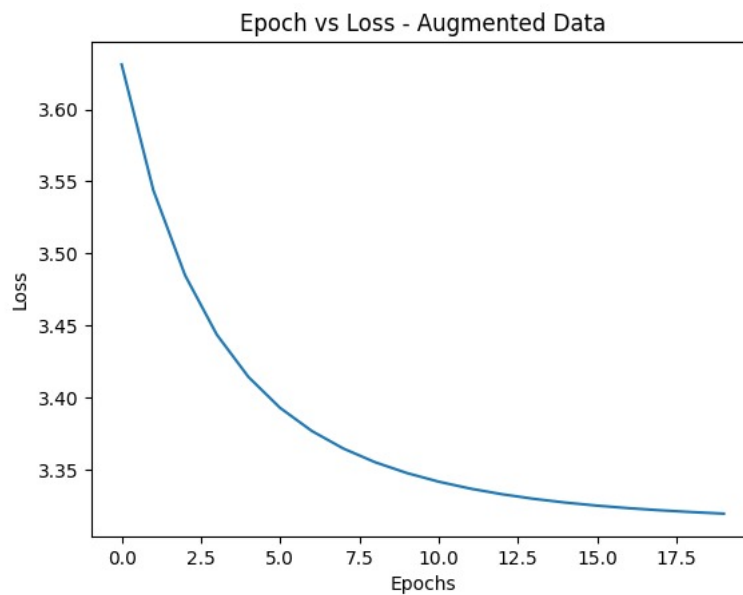


Figure 3: Loss vs Epoch Graph for Augmented Data Set

4 Evaluation Metrics

Implemented MLP Model achieved 90 % Train accuracies for both augmented and unaugmented data set , The test accuracies for unaugmented and augmented test data is 82.1 % and 82 % respectively.

4.1 Accuracy

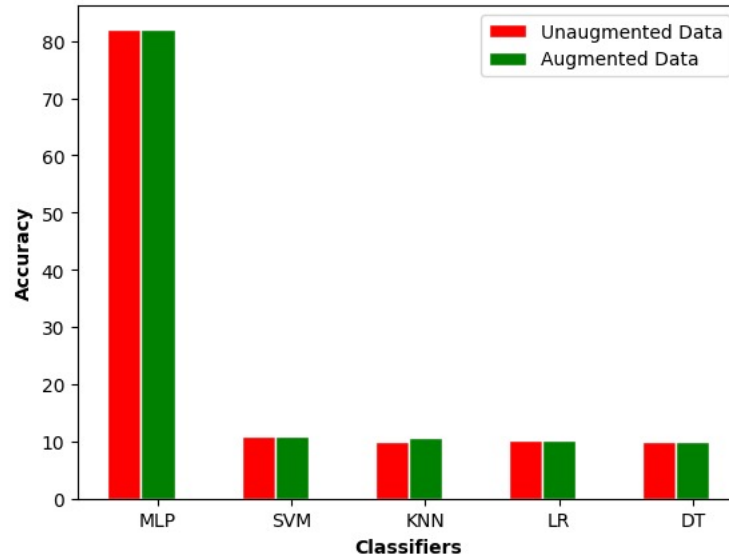


Figure 4: Accuracy Graph for both Original and Augmented Data set

4.2 Time taken for prediction

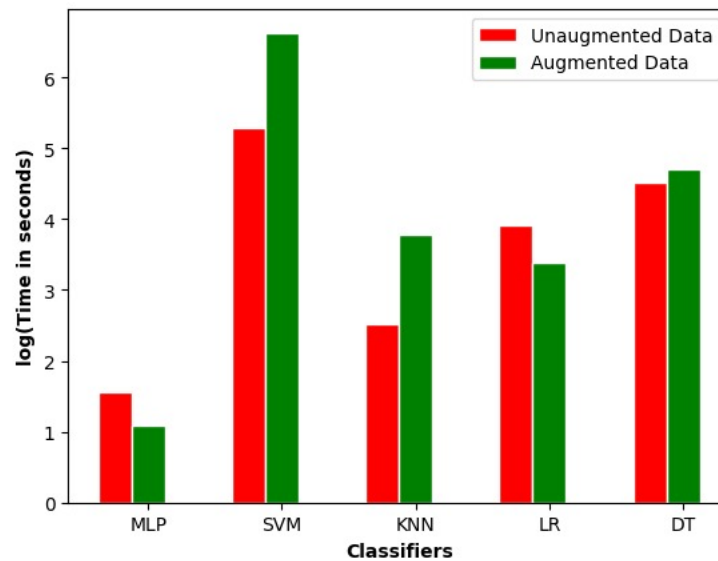


Figure 5: Prediction Time Graph for both Original and Augmented Data set

4.3 Justifications for difference in accuracy and time taken

- SVM, LR , KNN and DT Classifiers perform relatively poorly on non-linear data sets as compared to MLP Model.
- Computation time is less for MLP Model due to vectorization of forward computation process using numpy and parallel computing.

5 References

- [1] <https://neptune.ai/blog/performance-metrics-in-machine-learning-complete-guide>
- [2] https://gautamnagrawal.medium.com/rotating-image-by-any-angle-shear-transformation-using-only-numpy-d23Aeq%3Achain_rule
- [3] https://jonaslalin.com/2021/12/10/feedforward-neural-networks-part-1/#mjax-eqn%3Aeq%3Achain_rule
- [4] <https://jonaslalin.com/2021/12/22/feedforward-neural-networks-part-3/>
- [5] <https://jonaslalin.com/2021/12/21/feedforward-neural-networks-part-2/>
- [6] https://www.cs.umd.edu/class/fall2018/cmsc470/slides/slides_8.pdf
- [7] <https://www.youtube.com/watch?v=7HPwo4wnJeA>
- [8] <https://www.binarystudy.com/2021/09/how-to-load-preprocess-visualize-CIFAR-10-and-CIFAR-100.html#cifar10-load>
- [9] <https://www.geeksforgeeks.org/rotate-image-without-cutting-off-sides-using-python-opencv/>
- [10] <https://towardsdatascience.com/data-augmentation-compilation-with-python-and-opencv-b76b1cd500e0>
- [11] <https://towardsdatascience.com/a-beginners-guide-to-image-augmentations-in-machine-learning-22c48a2fb>
- [12] <https://stackoverflow.com/questions/49643907/clipping-input-data-to-the-valid-range-for-imshow-with-r>
- [13] <https://towardsdatascience.com/a-beginners-guide-to-image-augmentations-in-machine-learning-22c48a2fb>
- [14] <https://stackoverflow.com/questions/43391205/add-padding-to-images-to-get-them-into-the-same-shape>
- [15] <https://towardsdatascience.com/a-beginners-guide-to-image-augmentations-in-machine-learning-22c48a2fb>
- [16] https://www.youtube.com/watch?v=YP6APKLRf58&ab_channel=AdrianDolinay
- [17] <https://towardsdatascience.com/how-to-cluster-images-based-on-visual-similarity-cd6e7209fe34>
- [18] <https://www.kaggle.com/mtax687/single-layer-neural-network-using-numpy>
- [19] Various forums such as StackOverflow and StackExchange for debugging the code