



INDRAPRASTHA INSTITUTE OF INFORMATION TECHNOLOGY
NEW DELHI

Department of Computer Science & Engineering

CSE 343/ECE 363 : Machine Learning

Dr. Jainendra Shukla

Assignment - 3: Perceptron / MLP / SVM

Anant Kumar Kaushal (2022067)

SECTION - A

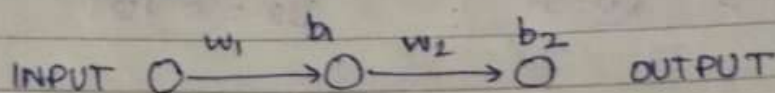
$$X = [1, 2, 3] \quad [\text{Given}]$$

$$Y = [3, 4, 5]$$

Let initial weights and biases be:
 $w_1 = 0.5, w_2 = 0.5$ [since these are normalized]
 $b_1 = 0.5, b_2 = 0.5$

$$\eta = 0.01 \quad (\text{Given})$$

I Forward Pass



(hidden layer)

$$\begin{aligned} \text{Output for the first layer} &= \text{Relu}([X] \times w_1 + b_1) \\ &= \text{Relu}([0.5, 1, 1.5] + 0.5) \\ &= \text{Relu}([1, 1.5, 2]) \\ &= [1, 1.5, 2] \end{aligned}$$

Output for the second layer = $\text{Relu}[h_1 \times w_2 + b_2]$
 (output layer)

$$\begin{aligned} &= [[1, 1.5, 2] \times 0.5 + 0.5] \\ &= [0.5, 0.75, 1] + 0.5 \\ &= [1, 1.25, 1.5] \end{aligned}$$

II Backpropagation

Date			
Page No.			

1. Gradient of loss:

$$L = \frac{1}{2} \sum_{i=1}^2 (y_i - \hat{y}_i)^2$$

$$L = \frac{1}{2} \sum (w^T x_i + b - \hat{y}_i)^2$$

$$\frac{dL}{db} = (w^T x_i + b - \hat{y}_i) \times 1$$

And $w^T x + b =$

Now for output layer

$$w^T x_2 + b = y_2 = [1, 1.25, 1.5]$$

$$\hat{y}_2 = [3, 4, 5]$$

\therefore gradient $b_2 =$

$$\frac{dL}{db} = ([1, 1.25, 1.5] - [3, 4, 5])$$

$$= [-2, -2.75, -3.45]$$

Now for b_2 :

$$\text{grad. } b_2 = \sum \text{gradient}$$

$$= (-2) + (-2.75) + (-3.45)$$

$$= -8.25$$

Since we want to normalize the above calculation so divide the above by Gradient by 3.

$$\therefore \text{grad-}b_2 = \frac{-8.25}{3} = -2.75$$

$$\text{Hly grad-}w_2 = \frac{\partial L}{\partial b} \times \text{Relu}'(\text{input} \cdot w_1 + b_1) \times \text{input}$$

$$= \begin{bmatrix} -2 & -2.75 & -3 \end{bmatrix} \begin{bmatrix} 1 \\ 1.5 \\ 2 \end{bmatrix}$$

$$= (-2) + (-4.125) + (-7)$$

Divide above by 3 to normalize it

$$= \frac{(-2) + (-4.125) + (-7)}{3}$$

$$= -4.375$$

For the first layer:

$$\text{output} = \frac{\partial L}{\partial b} \times w_2$$

$$= [-2, 2.75, 3.5] \times 0.5$$

$$= [-1, 1.375, 1.75]$$

$$\text{Hly Grad-}b_1 = \frac{1}{3} \sum \text{output}$$

$$= \frac{(-1) + (-1.375) + (-1.75)}{3}$$

$$= -1.375$$

$$\text{grad-}w_1 = [-1 \quad -1.375 \quad -1.75] \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

$$\text{grad-}w_1 = \text{output} \times \text{Relu}'(\text{input} \times w_1 + b_1) \times \text{input}$$

$$= [-1 \quad -1.375 \quad -1.75] \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

$$= \frac{-1 - 2.75 - 5.25}{3}$$

$$= -3$$

\therefore

$$w_1 = 0.5 + -\eta(\text{grad-}w_1) = 0.5 + (0.01)(3)$$

$$b_1 = 0.5 - \eta(\text{grad-}b_1) = 0.5 + (0.01)(1.375)$$

$$w_2 = 0.5 - \eta(\text{grad-}w_2) = 0.5 + (0.01)(4.375)$$

$$b_2 = 0.5 - \eta(\text{grad-}b_2) = 0.5 + (0.01)(2.75)$$

$$= 0.53$$

$$= 0.51375$$

$$= 0.54375$$

$$= 0.5275$$

A3

$$w_1 =$$

$$w_2 =$$

$$b$$

$$\|w\|$$

(a)

mar

(b)

(2,3)

(c)

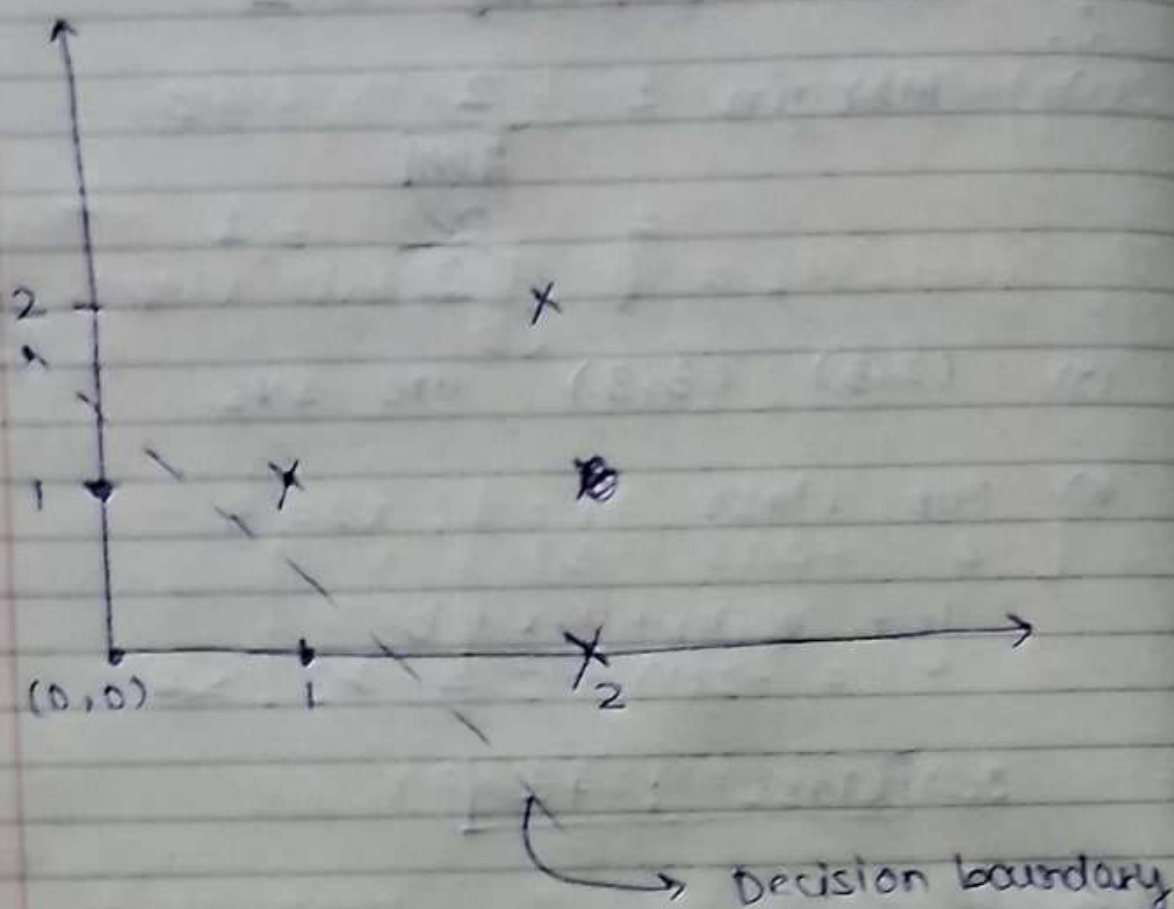
For

$$y =$$

\therefore [C]

A2

(a)



∴ Points are linearly separable as it can be seen from the plot.

(b) For '+' class:

$$\text{Using } w_1 x_1 + w_2 x_2 + b = y$$

$$(1,0) : w_1 + b = 1 \quad \text{①}$$

$$(0,0) : w_2 + b = 1 \quad \text{②}$$

On computing ① and ②, we get $w_1 = w_2$ and $b = 1 - w_1$

For (-) class:

$$(1, 1) : w_1 + w_2 + b = -1 \quad (3)$$

$$(2, 0) : 2w_1 + b = -1$$

~~Using w_1~~

Put values of w_2 and b in (3)

$$w_1 + \cancel{w_1} + 1 - \cancel{w_1} = -1$$

$$w_1 = -2$$

and

$$w_2 = -2$$

nty.

$$b = 3$$

∴ Decision Weight vector : $-2x_1 - 2x_2 = 3$
 $x_1 + x_2 = -\frac{3}{2}$

SVs :
1. $(1, 0)$ & $(2, 0)$
2. $(0, 1)$ & $(1, 1)$

they both
have same
minimum
distance
between them.

Q3

$$W_1 = -2$$

$$W_2 = 0$$

$$b = 5$$

$$\|W\| = \sqrt{W_1^2 + W_2^2} = 2$$

 \therefore

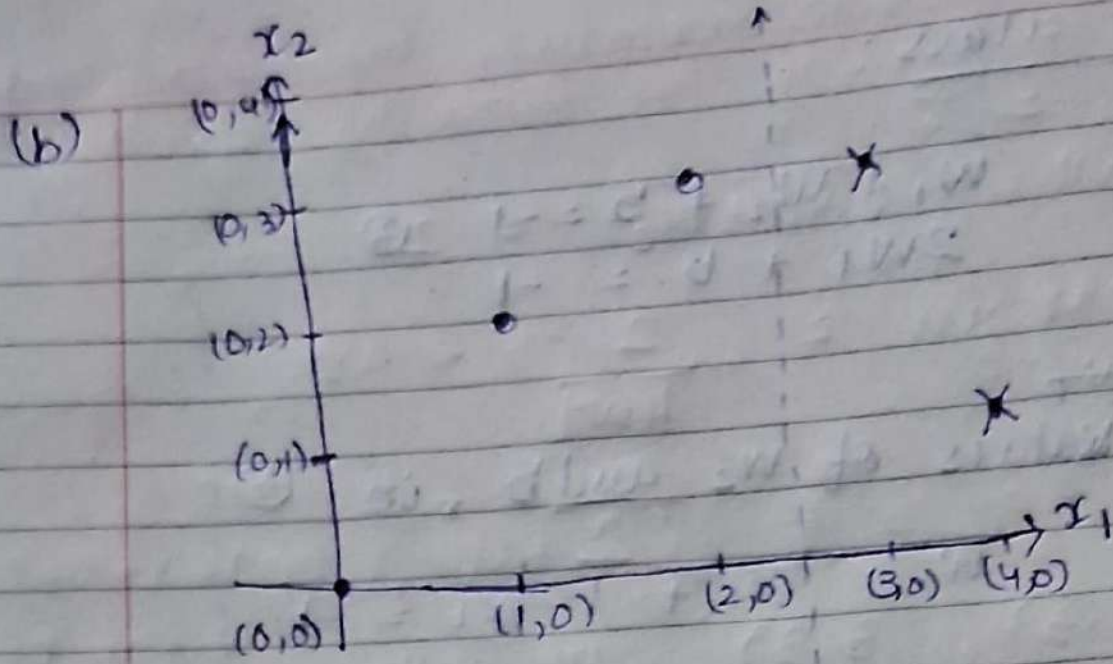
$$\begin{aligned} \text{(a) margin} &= \frac{2}{\|W\|} \\ &= \frac{2}{2} = 1 \end{aligned}$$

~~(b) (2, 3) ; (3, 3) / are svs~~

(c) For class $x_1 = 1, x_2 = 3$

$$\begin{aligned} y &= W_1 x_1 + W_2 x_2 + b \\ &= -2(1) + 5 = +3 \end{aligned}$$

\therefore class : +1



Decision $-2x_1 + 5x_2 = 0$
 Boundary

\therefore Hence $(2,3)$ & $(3,3)$ are SVs as they are closest to the decision boundary can be seen graphically.

SECTION - B

1. Class Neural Network with mentioned parameters and functions is implemented as follows:-

```
class NeuralNetwork:
    def __init__(self, N, layer_sizes, lr, activation_func, weight_init,
epochs=100, batch_size=128):
        self.N = N
        self.layer_sizes = layer_sizes
        self.lr = lr
        self.activation_func = activation_func
        self.weight_init = weight_init
        self.epochs = epochs
        self.batch_size = batch_size
        self.weights = self.create_weights()
        self.biases = [np.zeros((1, layer_sizes[i + 1])) for i in range(N - 1)]
def compute_loss(self, X, Y):
    activations = self.forward(X)
    predictions = activations[-1]
    return -np.mean(Y * np.log(predictions + 1e-8))

    def predict_proba(self, X):
        activations= self.forward(X)
        return activations[-1]

    def predict(self, X):
        proba = self.predict_proba(X)
        return np.argmax(proba, axis=1)

    def score(self, X, Y):
        predictions = self.predict(X)
        true_labels = np.argmax(Y, axis=1)
        return np.mean(predictions == true_labels)
```

2.Activation Functions along with their gradients are implemented as follows:-

```
def sigmoid(self, Z):
    return 1 / (1 + np.exp(-Z))

def sigmoid_derivative(self, Z):
    s = self.sigmoid(Z)
    return s * (1 - s)

def tanh(self, Z):
    return np.tanh(Z)

def tanh_derivative(self, Z):
    return 1 - np.tanh(Z) ** 2

def relu(self, Z):
    return np.maximum(0, Z)

def relu_derivative(self, Z):
    return Z > 0

def leaky_relu(self, Z, alpha=0.01):
    return np.where(Z > 0, Z, alpha * Z)

def leaky_relu_derivative(self, Z, alpha=0.1):
    dZ = np.ones_like(Z)
    dZ[Z < 0] = alpha
    return dZ

def softmax(self, Z):
    expZ = np.exp(Z - np.max(Z, axis=1, keepdims=True))
    return expZ / expZ.sum(axis=1, keepdims=True)
```

3.Activation Functions along with their gradients are implemented as follows:-

```
def create_weights(self):
    matrices = []
    for idx in range(1, len(self.layer_sizes)):
        if self.weight_init == 'zero':
            matrix = np.zeros((self.layer_sizes[idx],
self.layer_sizes[idx-1]))
        elif self.weight_init == 'random':
            matrix = np.random.rand(self.layer_sizes[idx],
self.layer_sizes[idx-1]) * 0.01
        elif self.weight_init == 'normal':
            matrix = np.random.randn(self.layer_sizes[idx],
self.layer_sizes[idx-1]) * np.sqrt(2 / self.layer_sizes[idx-1])
        matrices.append(matrix)
    return matrices
```

In this code appropriate scaling factors have been implemented for 'random' and 'normal' initializations . In 'random' initialization weights are scaled by **0.01** to create small random weights, which can help in stabilizing initial gradients. In 'normal' initialization weights are scaled by **np.sqrt(2 / self.layer_sizes[idx-1])**. This is known as He initialization, which helps maintain gradient flow in deeper networks by considering the previous layer's size when scaling weights drawn from a normal distribution.

4. Pre processing and Plots Of The Respective Models :-

```
def load_images(file_path):
    with open(file_path, 'rb') as file:
        cool, num_images, height, width = struct.unpack(">IIII", file.read(16))
        images = np.fromfile(file, dtype=np.uint8).reshape(num_images, height *
width) / 255.0
    return images

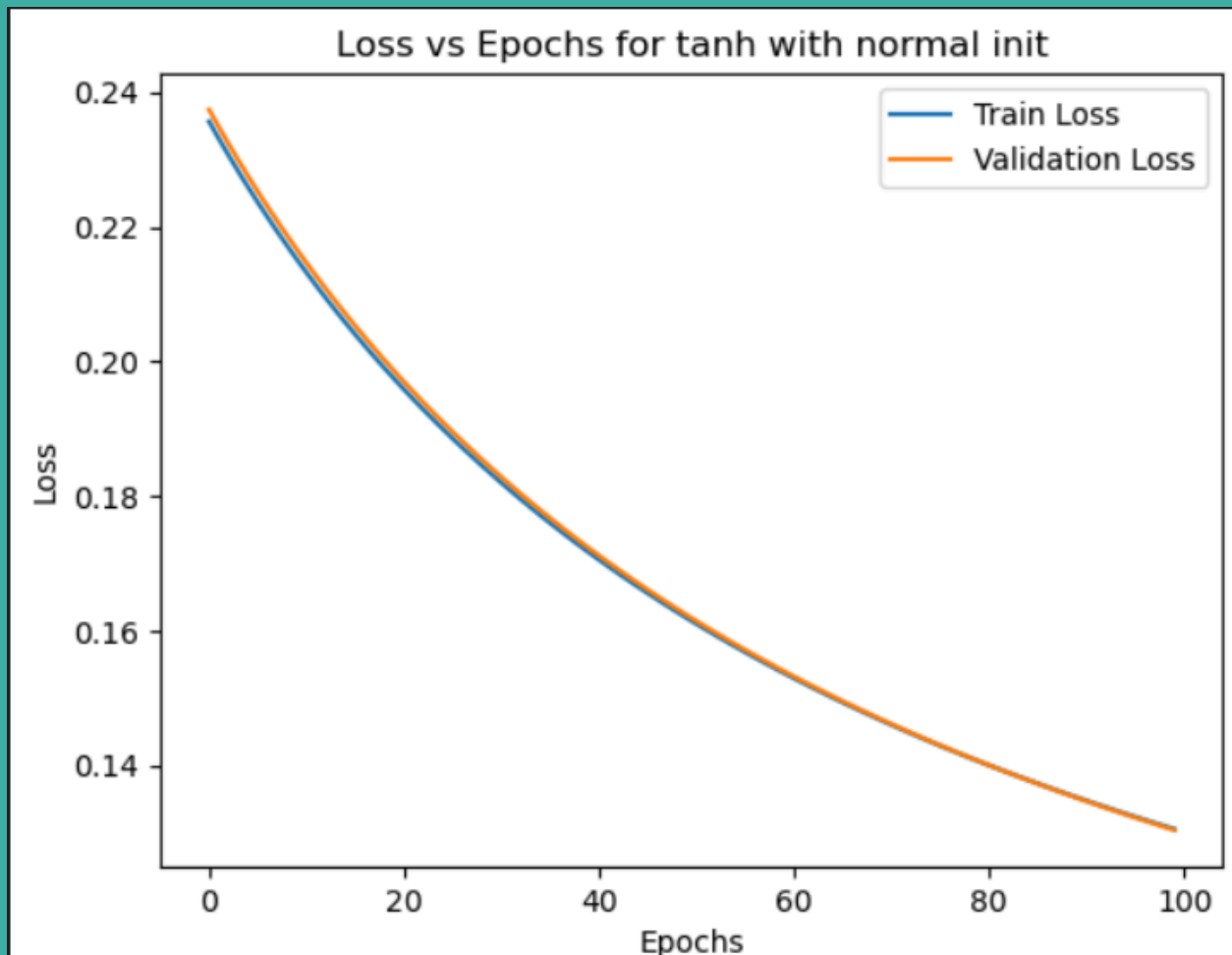
def load_labels(file_path):
    with open(file_path, 'rb') as file:
        cool, num_labels = struct.unpack(">II", file.read(8))
        labels = np.fromfile(file, dtype=np.uint8)
    return labels

def convert_to_one_hot(labels, classes=10):
    return np.eye(classes)[labels]
```

```
model = NeuralNetwork(
    N=6,
    layer_sizes=[784, 256, 128, 64, 32, 10],
    lr=2e-5,
    activation_func=activation,
    weight_init=init,
    epochs=100,
    batch_size=128,
)
```

The ***tanh normal initialization*** function stands out as the best-performing method due to its high accuracy on the training set, reaching approximately 63%, which is the highest among other initialization-function combinations tested. This

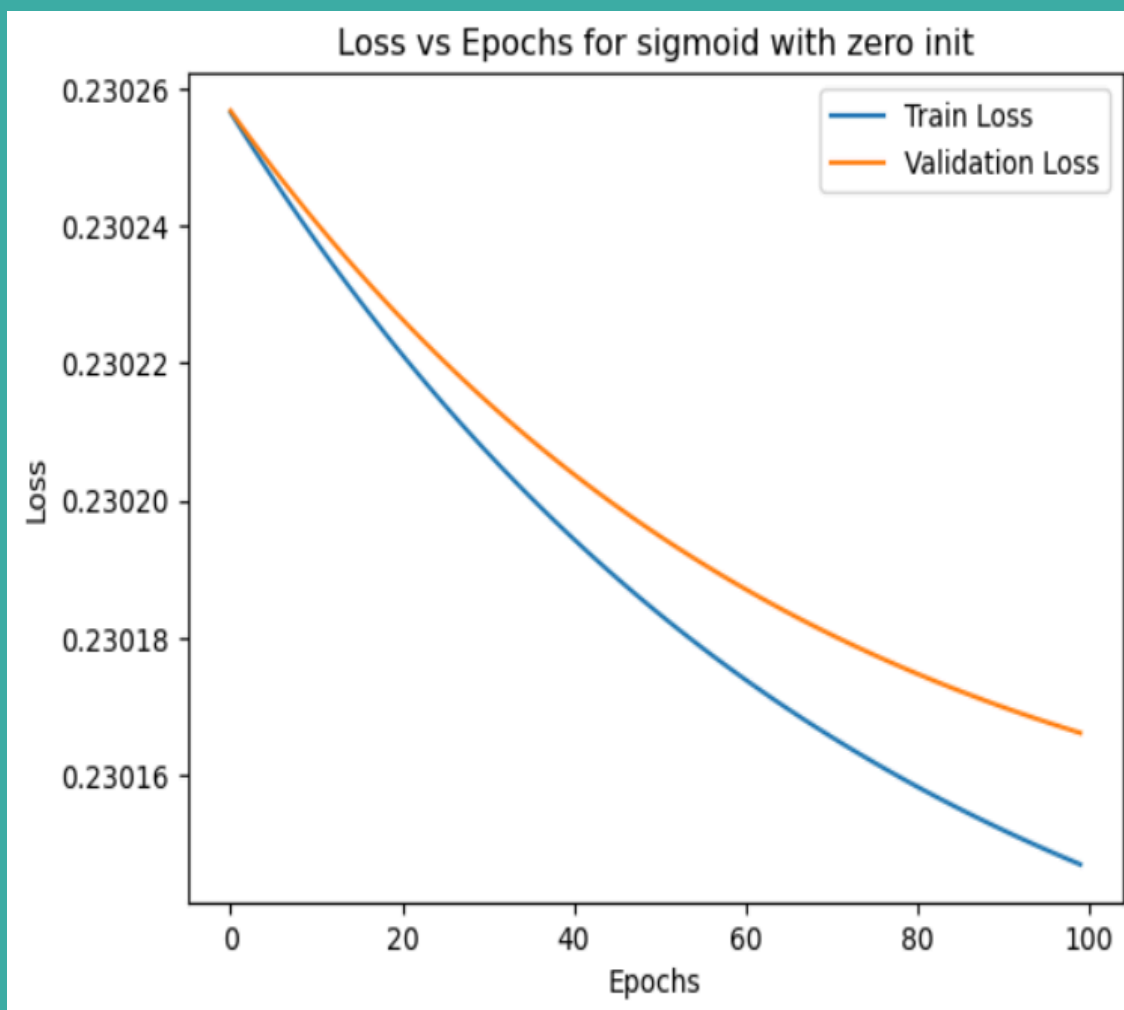
performance indicates that using the hyperbolic tangent (tanh) activation function combined with a normal distribution for weight initialization effectively supports the network's learning process. The training and validation losses for this setup were also closely monitored to ensure the model's generalization, indicating consistent improvement across the epochs without significant overfitting or underfitting.



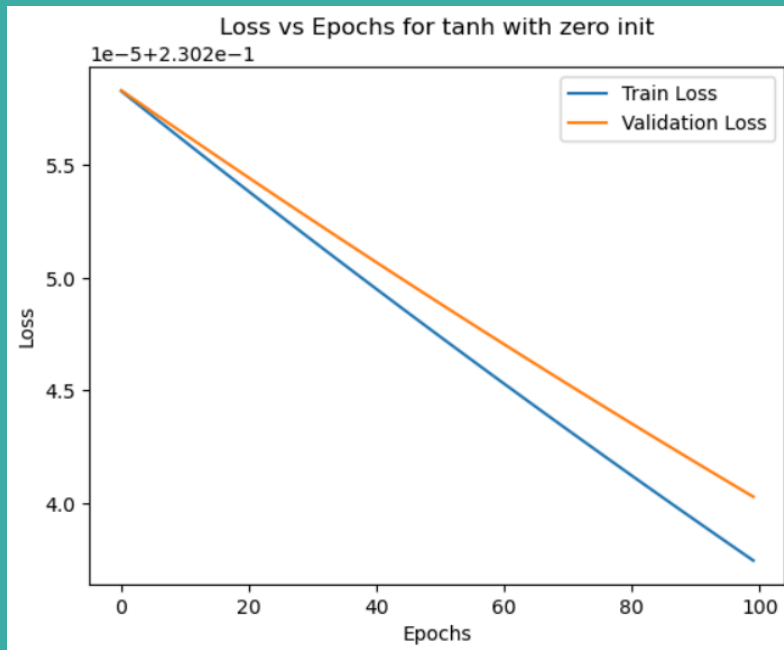
```
Epoch 100/100 - Train Loss: 0.130473  
Accuracy on train set: 0.63798828125  
Validation Loss: 0.130331
```

All three *zero initialization* functions are suboptimal for training neural networks because it initializes all weights to the same value, causing neurons to produce identical outputs and gradients during the forward and backward passes. This symmetry prevents the network from learning effectively, as each neuron in a layer receives the same updates and fails to diversify in function. Plots for the same are shown below :

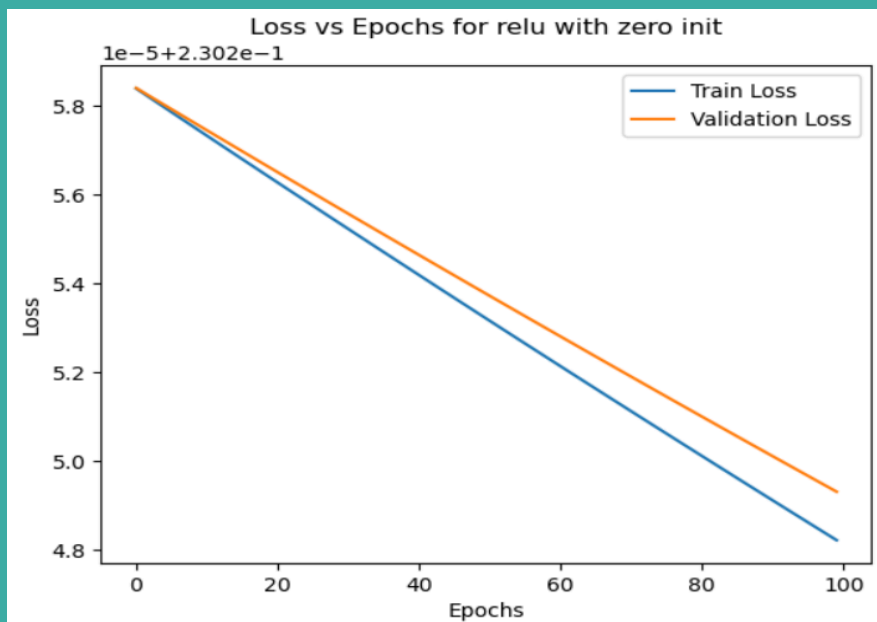
SIGMOID WITH ZERO INIT:



TANH WITH ZERO INIT:



LEAKY WITH ZERO INIT:



RELU LEAKY WITH ZERO INIT:

