

# Deep Learning Assignment Report

Karthik Babu Nambiar

2320702

## 1 Question 1

Which loss function, out of Cross Entropy and Mean Squared Error, works best with logistic regression because it guarantees a single best answer (no room for confusion)? Explain why this is important and maybe even show how it affects the model's training process.

Logistic function for a single independent variable  $x$  is given as

$$\sigma_{\beta_0, \beta_1}(x) = \frac{1}{1 + \exp(-\beta_0 - \beta_1 x)} \quad (1)$$

The Mean Squared Error (MSE) cost function measures the distances between model and data and is given by

$$h = \sum_{i=1}^N (y_i - \sigma_{\beta_0, \beta_1}(x_i))^2 = \sum_{i=1}^N (y_i - \sigma_i)^2 \quad (2)$$

where  $\sigma_i = \sigma(f(x_i)) = \sigma(\beta_0 + \beta_1 x_i)$ .

Even though MSE cost function is convex for linear regression problem, it is non-convex for logistic regression.

***Proof: Mean squared error is not convex as cost function for logistic regression***

From 2, calculate the first order derivative of summand

$$\frac{\partial (y_i - \sigma_i)^2}{\partial \beta_k} = -2(y_i - \sigma_i) \frac{\partial \sigma_i}{\partial \beta_k} \quad (3)$$

The last factor is according to the chain rule

$$\frac{\partial \sigma(f)}{\partial \beta_k} = \frac{\partial \sigma(f)}{\partial f} \frac{\partial f}{\partial \beta_k} \quad (4)$$

We know that

$$\frac{\partial \sigma(f)}{\partial f} = \sigma(f)(1 - \sigma(f)) \quad (5)$$

Therefore 3 becomes,

$$\frac{\partial (y_i - \sigma_i)^2}{\partial \beta_k} = -2(y_i - \sigma_i) \frac{\partial \sigma}{\partial f} \frac{\partial f}{\partial \beta_k} \quad (6)$$

$$= -2(y_i - \sigma_i) \sigma_i (1 - \sigma_i) \frac{\partial f}{\partial \beta_k} \quad (7)$$

$$= -2(y_i \sigma_i - y_i \sigma_i^2 - \sigma_i^2 + \sigma_i^3) \frac{\partial f}{\partial \beta_k} \quad (8)$$

Now calculate the second order derivative of summand

$$\frac{\partial}{\partial \beta_k} \left( \frac{\partial(y_i - \sigma_i)^2}{\partial \beta_k} \right) = \frac{\partial}{\partial \beta_k} \left( -2(y_i \sigma_i - y_i \sigma_i^2 - \sigma_i^2 + \sigma_i^3) \frac{\partial f}{\partial \beta_k} \right) \quad (9)$$

$$= -2 \left( y_i \frac{\partial \sigma_i}{\partial \beta_k} - y_i \frac{\partial \sigma_i^2}{\partial \beta_k} - \frac{\partial \sigma_i^2}{\partial \beta_k} + \frac{\partial \sigma_i^3}{\partial \beta_k} \right) \frac{\partial f}{\partial \beta_k} \quad (10)$$

$$+ -2(y_i \sigma_i - y_i \sigma_i^2 - \sigma_i^2 + \sigma_i^3) \frac{\partial^2 f}{\partial \beta_k^2} \quad (11)$$

$$= -2(y_i - 2y_i \sigma_i - 2\sigma_i + 3\sigma_i^2) \frac{\partial \sigma_i}{\partial \beta_k} \frac{\partial f}{\partial \beta_k} \quad (12)$$

$$= -2(y_i - 2y_i \sigma_i - 2\sigma_i + 3\sigma_i^2) \frac{\partial \sigma_i}{\partial f} \left( \frac{\partial f}{\partial \beta_k} \right)^2 \quad (13)$$

$$= -2(y_i - 2y_i \sigma_i - 2\sigma_i + 3\sigma_i^2) \sigma_i (1 - \sigma_i) \left( \frac{\partial f}{\partial \beta_k} \right)^2 \quad (14)$$

We know that

$$\sigma_i(1 - \sigma_i) \left( \frac{\partial f}{\partial \beta_k} \right)^2 := g \geq 0 | \sigma_i \in [0, 1] \quad (15)$$

For binary classification problem  $y_i \in \{0, 1\}$

for  $y_i = 0$

$$\frac{\partial^2(y_i - \sigma_i)^2}{\partial \beta_k^2} = -2(-2 + 3\sigma_i) \sigma_i g \quad (16)$$

for  $y_i = 1$

$$\frac{\partial^2(y_i - \sigma_i)^2}{\partial \beta_k^2} = -2(1 - 4\sigma + 3\sigma_i^2)g = -2(1 - \sigma)(1 - 3\sigma)g \quad (17)$$

Hence there is a sign change and  $\frac{\partial^2 h_i^2}{\partial \beta_k^2} \not\geq 0$ .

**Proof: Cross Entropy is convex as cost function for logistic regression**

$$h_i = -y_i \log(\sigma_i) - (1 - y_i) \log(1 - \sigma_i) \quad (18)$$

$$\frac{\partial h_i}{\partial \beta_k} = -y_i \frac{\partial \log(\sigma_i)}{\partial \beta_k} - (1 - y_i) \frac{\partial \log(1 - \sigma_i)}{\partial \beta_k} \quad (19)$$

$$= -y_i \frac{\partial \sigma_i}{\partial \beta_k} \frac{1}{\sigma_i} + (1 - y_i) \frac{\partial \sigma_i}{\partial \beta_k} \frac{1}{1 - \sigma_i} \quad (20)$$

$$\text{with } \frac{\partial \sigma_i}{\partial \beta_k} = \sigma(f)(1 - \sigma(f)) \frac{\partial f}{\partial \beta_k} \quad (21)$$

$$\frac{\partial h_i}{\partial \beta_k} = (-y_i(1 - \sigma_i) + (1 - y_i)\sigma_i) \frac{\partial f}{\partial \beta_k} \quad (22)$$

$$\frac{\partial h_i}{\partial \beta_k} = (\sigma_i - y_i) \frac{\partial f}{\partial \beta_k} \quad (23)$$

$$\frac{\partial^2 h_i}{\partial \beta_k^2} = \frac{\partial \sigma_i}{\partial \beta_k} \frac{\partial f}{\partial \beta_k} + \sigma_i \frac{\partial^2 f}{\partial \beta_k^2} \quad (24)$$

$$\frac{\partial^2 h_i}{\partial \beta_k^2} = \sigma_i(1 - \sigma_i) \left( \frac{\partial f}{\partial \beta_k} \right)^2 \quad (25)$$

$$\frac{\partial^2 h_i}{\partial \beta_k^2} \geq 0 \quad (26)$$

Hence Cross Entropy works best with logistic regression when compared to Mean Squared Error and guarentees a single best answer because of its convexity. This will help during training as it will take less time to converge to the minimum.

## 2 Question 2

For a binary classification task with a deep neural network (containing at least one hidden layer) equipped with linear activation functions, which of the following loss functions guarantees a convex optimization problem? Justify your answer with a formal proof or a clear argument. (a) CE (b) MSE (c) Both (A) and (B) (d) None  
A linear neural network can be denoted as

$$f(x) = W^T x = \beta_0 + \beta_1 x \quad (27)$$

Consider MSE loss function

$$h_i = (y_i - f_i)^2 \quad (28)$$

$$\frac{\partial h_i}{\partial \beta_k} = -2(y_i - f_i) \frac{\partial f_i}{\partial \beta_k} \quad (29)$$

$$\frac{\partial^2 h_i}{\partial \beta_k^2} = -2 \frac{\partial y_i}{\partial \beta_k} \frac{\partial f_i}{\partial \beta_k} + 2 \left( \frac{\partial f_i}{\partial \beta_k} \right)^2 - 2(y_i - f_i) \frac{\partial^2 f_i}{\partial \beta_k^2} \quad (30)$$

$$\frac{\partial^2 h_i}{\partial \beta_k^2} \geq 0 \quad (31)$$

Consider CE loss

$$h_i = -y_i \log(f_i) - (1 - y_i) \log(1 - f_i) \quad (32)$$

$$\frac{\partial h_i}{\partial \beta_k} = -y_i \frac{\partial \log(f_i)}{\partial \beta_k} - (1 - y_i) \frac{\partial \log(1 - f_i)}{\partial \beta_k} \quad (33)$$

$$= -y_i \frac{\partial f_i}{\partial \beta_k} \frac{1}{f_i} + (1 - y_i) \frac{\partial f_i}{\partial \beta_k} \frac{1}{1 - f_i} \quad (34)$$

$$= \left( \frac{-y_i}{f_i} + \frac{1 - y_i}{1 - f_i} \right) \frac{\partial f_i}{\partial \beta_k} \quad (35)$$

$$\frac{\partial^2 h_i}{\partial \beta_k^2} = \left( \frac{y_i}{f_i^2} + \frac{1 - y_i}{(1 - f_i)^2} \right) \left( \frac{\partial f_i}{\partial \beta_k} \right)^2 + \left( \frac{-y_i}{f_i} + \frac{1 - y_i}{1 - f_i} \right) \frac{\partial^2 f_i}{\partial \beta_k^2} \quad (36)$$

$$\frac{\partial^2 h_i}{\partial \beta_k^2} \geq 0 | y_i \in \{0, 1\} \quad (37)$$

Therefore both CE and MSE loss functions guarantees a convex optimization problem in this case.

## 3 Question 3

**Dense Neural Network:** Implement a feedforward neural network with dense layers only. Specify the number of hidden layers, neurons per layer, and activation functions. How will you preprocess the input images? Consider hyperparameter tuning strategies.

Importing all necessary libraries.

```
1 import torch
2 import torch.nn as nn
3 import torchvision
4 import torchvision.transforms as transforms
5 from torch.utils.data import DataLoader, random_split
```

Defining the neural network architecture. It has one input layer, one output layer and four hidden layers. The model has the following number of neurons : [784, 256, 256, 128, 128, 10]. The activation function used is ReLU.

```
1 # Define the neural network architecture
2 class NeuralNetwork(nn.Module):
```

```

3     def __init__(self, input_size, hidden_size, num_classes):
4         super(NeuralNetwork, self).__init__()
5         self.fc1 = nn.Linear(input_size, 2*hidden_size)
6         self.relu = nn.ReLU()
7         self.fc2 = nn.Linear(2*hidden_size, 2*hidden_size)
8         self.fc3 = nn.Linear(2*hidden_size, hidden_size)
9         self.fc4 = nn.Linear(hidden_size, hidden_size)
10        self.fc5 = nn.Linear(hidden_size, num_classes)
11
12    def forward(self, x):
13        out = self.fc1(x)
14        out = self.relu(out)
15        out = self.fc2(out)
16        out = self.relu(out)
17        out = self.fc3(out)
18        out = self.relu(out)
19        out = self.fc4(out)
20        out = self.relu(out)
21        out = self.fc5(out)
22        return out
23
24    # Device configuration
25    device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
26
27    # Hyperparameters
28    input_size = 28 * 28 # MNIST image size is 28x28
29    hidden_size = 128 # Number of neurons in the hidden layers
30    num_classes = 10 # Number of output classes (digits 0-9)
31    batch_size = 100
32    num_epochs = 5

```

The original images range from  $[0, 1]$  is transformed into  $[-1, 1]$

```

1 # Load and preprocess the MNIST dataset
2 transform = transforms.Compose([
3     transforms.ToTensor(), # Convert images to tensors
4     transforms.Normalize((0.5,), (0.5,)) # Normalize the pixel values to the range [-1, 1]
5 ])

```

Dataloader for training on MNIST dataset.

```

1 train_dataset = torchvision.datasets.MNIST(root='./data', train=True, transform=transform, download=True)
2 test_dataset = torchvision.datasets.MNIST(root='./data', train=False, transform=transform)
3
4 # Split the dataset into training and validation sets
5 train_size = int(0.8 * len(train_dataset))
6 val_size = len(train_dataset) - train_size
7 train_dataset, val_dataset = random_split(train_dataset, [train_size, val_size])
8
9 train_loader = DataLoader(dataset=train_dataset, batch_size=batch_size, shuffle=True)
10 val_loader = DataLoader(dataset=val_dataset, batch_size=batch_size, shuffle=False)
11 test_loader = DataLoader(dataset=test_dataset, batch_size=batch_size, shuffle=False)

```

Hyperparameter tuning of learning rate using Grid Search strategy.

```

1 # Hyperparameter tuning - Grid Search for learning rate
2 learning_rates = [0.001, 0.01, 0.1]
3 best_accuracy = 0
4 best_learning_rate = None
5
6 for lr in learning_rates:
7     # Initialize the model
8     model = NeuralNetwork(input_size, hidden_size, num_classes).to(device)

```

```

9
10 # Loss and optimizer
11 criterion = nn.CrossEntropyLoss()
12 optimizer = torch.optim.Adam(model.parameters(), lr=lr)
13
14 # Training the model
15 for epoch in range(num_epochs):
16     model.train()
17     for i, (images, labels) in enumerate(train_loader):
18         images = images.reshape(-1, 28 * 28).to(device)
19         labels = labels.to(device)
20
21         # Forward pass
22         outputs = model(images)
23         loss = criterion(outputs, labels)
24
25         # Backward pass and optimization
26         optimizer.zero_grad()
27         loss.backward()
28         optimizer.step()
29
30 # Validation
31 model.eval()
32 with torch.no_grad():
33     correct = 0
34     total = 0
35     for images, labels in val_loader:
36         images = images.reshape(-1, 28 * 28).to(device)
37         labels = labels.to(device)
38         outputs = model(images)
39         _, predicted = torch.max(outputs.data, 1)
40         total += labels.size(0)
41         correct += (predicted == labels).sum().item()
42
43     accuracy = correct / total
44     print(f'Learning Rate: {lr}, Validation Accuracy: {accuracy}')
45
46     # Update best learning rate if current accuracy is higher
47     if accuracy > best_accuracy:
48         best_accuracy = accuracy
49         best_learning_rate = lr
50
51 print(f'Best Learning Rate: {best_learning_rate}, Best Validation Accuracy: {best_accuracy}')

```

Training the model with the best learning rate from the hyperparameter tuning.

```

1
2 # Testing with best learning rate
3 # Initialize the model with the best learning rate
4 model = NeuralNetwork(input_size, hidden_size, num_classes).to(device)
5 optimizer = torch.optim.Adam(model.parameters(), lr=best_learning_rate)
6
7 # Training the model
8 for epoch in range(num_epochs):
9     model.train()
10    for i, (images, labels) in enumerate(train_loader):
11        images = images.reshape(-1, 28 * 28).to(device)
12        labels = labels.to(device)
13
14    # Forward pass

```

```

15     outputs = model(images)
16     loss = criterion(outputs, labels)
17
18     # Backward pass and optimization
19     optimizer.zero_grad()
20     loss.backward()
21     optimizer.step()
22
23 # Testing
24 model.eval()
25 with torch.no_grad():
26     correct = 0
27     total = 0
28     for images, labels in test_loader:
29         images = images.reshape(-1, 28 * 28).to(device)
30         labels = labels.to(device)
31         outputs = model(images)
32         _, predicted = torch.max(outputs.data, 1)
33         total += labels.size(0)
34         correct += (predicted == labels).sum().item()
35
36     print(f'Test Accuracy with Best Learning Rate: {100 * correct / total}%')

```

Result

```

1 Learning Rate: 0.001, Validation Accuracy: 0.9658333333333333
2 Learning Rate: 0.01, Validation Accuracy: 0.9268333333333333
3 Learning Rate: 0.1, Validation Accuracy: 0.1158333333333333
4 Best Learning Rate: 0.001, Best Validation Accuracy: 0.9658333333333333
5 Test Accuracy with Best Learning Rate: 97.05%

```

## 4 Question 4

Build a classifier for Street View House Numbers (SVHN) (Dataset) using pretrained model weights from PyTorch. Try multiple models like LeNet-5, AlexNet, VGG, or ResNet(18, 50, 101). Compare performance comment why a particular model is well suited for SVHN dataset. (You can use a subset of dataset (25%) in case you do not have enough compute.)

Importing the necessary libraries

```

1 import torch
2 import torchvision
3 from torchvision import transforms, datasets
4 import torch.nn as nn
5 import torch.optim as optim
6 from torch.utils.data import DataLoader
7 from sklearn.metrics import accuracy_score

```

Dataloader and data transformation of SVHN dataset

```

1 # Step 1: Load and preprocess the SVHN dataset
2 transform = transforms.Compose([
3     transforms.Resize((224, 224)),
4     transforms.ToTensor(),
5     transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5))
6 ])
7
8 train_data = datasets.SVHN(root='./data', split='train', download=True, transform=transform)
9 test_data = datasets.SVHN(root='./data', split='test', download=True, transform=transform)
10

```

```

11 # Split the dataset into train and validation sets
12 train_size = int(0.75 * len(train_data))
13 val_size = len(train_data) - train_size
14 train_dataset, val_dataset = torch.utils.data.random_split(train_data, [train_size, val_size])
15
16 # Create data loaders
17 train_loader = DataLoader(train_dataset, batch_size=64, shuffle=True)
18 val_loader = DataLoader(val_dataset, batch_size=64, shuffle=False)
19 test_loader = DataLoader(test_data, batch_size=64, shuffle=False)

```

training loop

```

1 # Step 3: Fine-tune the pretrained models
2 def train_model(model, criterion, optimizer, train_loader, val_loader, num_epochs=5):
3     for epoch in range(num_epochs):
4         model.train()
5         running_loss = 0.0
6         for inputs, labels in train_loader:
7             inputs, labels = inputs.cuda(), labels.cuda()
8             optimizer.zero_grad()
9             outputs = model(inputs)
10            loss = criterion(outputs, labels)
11            loss.backward()
12            optimizer.step()
13            running_loss += loss.item() * inputs.size(0)
14
15        epoch_loss = running_loss / len(train_loader.dataset)
16        print(f"Epoch {epoch+1}/{num_epochs}, Loss: {epoch_loss:.4f}")
17
18        # Validate the model
19        model.eval()
20        correct = 0
21        total = 0
22        with torch.no_grad():
23            for inputs, labels in val_loader:
24                inputs, labels = inputs.cuda(), labels.cuda()
25                outputs = model(inputs)
26                _, predicted = torch.max(outputs, 1)
27                total += labels.size(0)
28                correct += (predicted == labels).sum().item()
29
30        val_accuracy = correct / total
31        print(f"Validation Accuracy: {val_accuracy:.4f}")

```

Comparing result of alexnet, vgg, resnet-18, resnet-50, resnet-101. Only the last FC layer is being fine tuned

```

1 # AlexNet
2 alexnet = torchvision.models.alexnet(pretrained=True)
3 for param in alexnet.parameters():
4     param.requires_grad = False
5 alexnet.classifier[6] = nn.Linear(4096, 10) # Modify the last fully connected layer for 10 classes
6 alexnet.classifier[6].requires_grad = True # Set requires_grad to True for the fully connected layer
7 alexnet = alexnet.cuda()
8 criterion = nn.CrossEntropyLoss()
9 optimizer = optim.Adam(alexnet.classifier[6].parameters(), lr=0.001) # Only optimize parameters of the fully
    connected layer
10 train_model(alexnet, criterion, optimizer, train_loader, val_loader)
11
12 # VGG
13 vgg = torchvision.models.vgg16(pretrained=True)
14 for param in vgg.parameters():

```

```

15     param.requires_grad = False
16 vgg.classifier[6] = nn.Linear(4096, 10) # Modify the last fully connected layer for 10 classes
17 vgg.classifier[6].requires_grad = True # Set requires_grad to True for the fully connected layer
18 vgg = vgg.cuda()
19 criterion = nn.CrossEntropyLoss()
20 optimizer = optim.Adam(vgg.classifier[6].parameters(), lr=0.001) # Only optimize parameters of the fully
    connected layer
21 train_model(vgg, criterion, optimizer, train_loader, val_loader)
22
23 # ResNet-18
24 resnet18 = torchvision.models.resnet18(pretrained=True)
25 for param in resnet18.parameters():
26     param.requires_grad = False
27 resnet18.fc = nn.Linear(resnet18.fc.in_features, 10) # Modify the last fully connected layer for 10 classes
28 resnet18.fc.requires_grad = True # Set requires_grad to True for the fully connected layer
29 resnet18 = resnet18.cuda()
30 criterion = nn.CrossEntropyLoss()
31 optimizer = optim.Adam(resnet18.fc.parameters(), lr=0.001) # Only optimize parameters of the fully connected
    layer
32 train_model(resnet18, criterion, optimizer, train_loader, val_loader)
33
34 # ResNet-50
35 resnet50 = torchvision.models.resnet50(pretrained=True)
36 for param in resnet50.parameters():
37     param.requires_grad = False
38 resnet50.fc = nn.Linear(resnet50.fc.in_features, 10) # Modify the last fully connected layer for 10 classes
39 resnet50.fc.requires_grad = True # Set requires_grad to True for the fully connected layer
40 resnet50 = resnet50.cuda()
41 criterion = nn.CrossEntropyLoss()
42 optimizer = optim.Adam(resnet50.fc.parameters(), lr=0.001) # Only optimize parameters of the fully connected
    layer
43 train_model(resnet50, criterion, optimizer, train_loader, val_loader)
44
45 # ResNet-101
46 resnet101 = torchvision.models.resnet101(pretrained=True)
47 for param in resnet101.parameters():
48     param.requires_grad = False
49 resnet101.fc = nn.Linear(resnet101.fc.in_features, 10) # Modify the last fully connected layer for 10
    classes
50 resnet101.fc.requires_grad = True # Set requires_grad to True for the fully connected layer
51 resnet101 = resnet101.cuda()
52 criterion = nn.CrossEntropyLoss()
53 optimizer = optim.Adam(resnet101.fc.parameters(), lr=0.001) # Only optimize parameters of the fully
    connected layer
54 train_model(resnet101, criterion, optimizer, train_loader, val_loader)

```

## Results

### Alexnet

```

1 Epoch 1/5, Loss: 1.7448
2 Validation Accuracy: 0.5084
3 Epoch 2/5, Loss: 1.6341
4 Validation Accuracy: 0.5263
5 Epoch 3/5, Loss: 1.6068
6 Validation Accuracy: 0.5163
7 Epoch 4/5, Loss: 1.6011
8 Validation Accuracy: 0.5335
9 Epoch 5/5, Loss: 1.5893
10 Validation Accuracy: 0.5411

```

### VGG



```
1 Epoch 1/5, Loss: 1.7812
2 Validation Accuracy: 0.4805
3 Epoch 2/5, Loss: 1.6833
4 Validation Accuracy: 0.5086
5 Epoch 3/5, Loss: 1.6733
6 Validation Accuracy: 0.5195
7 Epoch 4/5, Loss: 1.6636
8 Validation Accuracy: 0.5343
9 Epoch 5/5, Loss: 1.6614
10 Validation Accuracy: 0.5080
```

#### Resnet-18

```
1 Epoch 1/5, Loss: 1.8437
2 Validation Accuracy: 0.4131
3 Epoch 2/5, Loss: 1.6680
4 Validation Accuracy: 0.4353
5 Epoch 3/5, Loss: 1.6341
6 Validation Accuracy: 0.4480
7 Epoch 4/5, Loss: 1.6121
8 Validation Accuracy: 0.4530
9 Epoch 5/5, Loss: 1.6033
10 Validation Accuracy: 0.4631
```

#### Resnet-50

```
1 Epoch 1/5, Loss: 1.8692
2 Validation Accuracy: 0.3840
3 Epoch 2/5, Loss: 1.7183
4 Validation Accuracy: 0.3951
5 Epoch 3/5, Loss: 1.6795
6 Validation Accuracy: 0.4355
7 Epoch 4/5, Loss: 1.6445
8 Validation Accuracy: 0.4414
9 Epoch 5/5, Loss: 1.6166
10 Validation Accuracy: 0.4413
```

#### Resnet-101

```
1 Epoch 1/5, Loss: 1.8666
2 Validation Accuracy: 0.3762
3 Epoch 2/5, Loss: 1.7301
4 Validation Accuracy: 0.4121
5 Epoch 3/5, Loss: 1.6848
6 Validation Accuracy: 0.4319
7 Epoch 4/5, Loss: 1.6553
8 Validation Accuracy: 0.4110
9 Epoch 5/5, Loss: 1.6310
10 Validation Accuracy: 0.4245
```

Typically, deeper models like VGG, ResNet-50, and ResNet-101 tend to perform better on complex datasets like SVHN due to their deeper architectures which allow them to learn more intricate features. These models have shown great success in various image classification tasks and are well-suited for datasets like SVHN.