

HW 8: Gating mechanisms in RNN

David E. Farache, Email ID: dfarache@purdue.edu

April 30, 2023

1 Introduction

For this assignment, we focused on implementing a Vision Transformer (ViT) using a multi-headed self-attention mechanism and transformer architecture for image classification.

2 Theory/Methods

2.1 Scaled Dot Product Attention

Scaled dot product attention is the method utilized to transfer an input sequence into separate relevant vector values to determine importance. These tensors are Q, all query-based vectors into a single data object, K, word-based key vectors, and V, value vectors. The attention is calculated as tensor V being modified by tensors Q and K as seen below:

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}})V \quad (1)$$

In this case, $\sqrt{d_k}$ normalizes the dot product of Q and K^T and softmax turns those normalized values into weights that are then multiplied to V, giving the attention value.

2.2 Multi-Headed Attention

In the case of multi-headed attention, the input tensor is first split along its embedded axis into N slices. Each section is then processed as if it were single-headed, going through scaled dot product attention, with individual Q, K, and V tensors. Each partition, after returning values after applying equation 1, is then concatenated and passed through a fully connected linear layer.

2.3 Transformer

Transformers function by adopting a mechanism of self-attention, then weighting the significance of each part of the input. It does this via an architecture of stacked encoders and decoders. The encoders create an attention map for sentence sequence fed into it while the decoder translates, based on the prior attention map, to a designated sentence.

The encoder structure, in this case, is a 6-layer structure with 3 different sections of 2-layer sections. The first section of the three applies the multi-head attention layer described previously and the section after applies the position-wise fully connected feed-forward network. The final

section is connected with all layers via a residual connection and which passes the information through a normalization layer.

The decoder architecture is the same as the encoder, a 6-layer structure with 3 sub-layers that each have 2 layers. These layers mirror the structure of the encoder. With the first layer being a self-attention layer, the second being a position-wise fully connected layer, and the final layer being a multi-headed attention layer.

This logic, so far, for the transformer has been focused on the assumption that the data being fed into are sentences or sequences. For ViT the issue is that we need to process images which requires transferring image data into sequence data. This is achieved via Patch Embedding, which functions by slicing the image that is then condensed into data vectors. This done via reshaping and positional information is included to enable patch identification. Furthermore, tensors have an input sequence pre-pended to store information from one layer to another and create deeper layers.

3 ViTHelper

For the assignment, the ViTHelper python library from DLStudio was utilized for the development of Transformers and creating the ViT. For the extra credit assignment, the function AttentionHead was copied and updated to add torch.einsum as seen in the code below:

```

1 class AttentionHead_ExtraCredit(nn.Module):
2     def __init__(self, max_seq_length, qkv_size):
3         super().__init__()
4         self.qkv_size = qkv_size
5         self.max_seq_length = max_seq_length
6         self.WQ = nn.Linear(max_seq_length * self.qkv_size,
7                               max_seq_length * self.qkv_size) # (B)
8         self.WK = nn.Linear(max_seq_length * self.qkv_size,
9                               max_seq_length * self.qkv_size) # (C)
10        self.WV = nn.Linear(max_seq_length * self.qkv_size,
11                              max_seq_length * self.qkv_size) # (D)
12        self.softmax = nn.Softmax(dim=1) # (E)
13
14    def forward(self, sentence_portion):
15        Q = self.WQ(sentence_portion.reshape(
16            sentence_portion.shape[0], -1).float()).to(device) # (G)
17        K = self.WK(sentence_portion.reshape(
18            sentence_portion.shape[0], -1).float()).to(device) # (H)
19        V = self.WV(sentence_portion.reshape(
20            sentence_portion.shape[0], -1).float()).to(device) # (I)
21
22        Q = Q.view(sentence_portion.shape[0],
23                    self.max_seq_length, self.qkv_size) # (J)
24        K = K.view(sentence_portion.shape[0],
25                    self.max_seq_length, self.qkv_size) # (K)
26        V = V.view(sentence_portion.shape[0],
27                    self.max_seq_length, self.qkv_size) # (L)
28
29        QK_DotProduct = torch.einsum("bqn, bnk -> bqk", Q, rearrange(K, "b n e
-> b e n"))
30        ScaleCoeff = 1.0 / torch.sqrt(torch.tensor([self.qkv_size]).float()).to(
device)

```

```

31         return ScaleCoeff * torch.einsum("ban, bne -> bae", self.softmax(
            QK_DotProduct), V)

```

Listing 1: *ViT Extra Credit*

4 Setup Code

```

1  # Import Libraries
2  import numpy as np
3  import torch
4  import torchvision.transforms as tvt
5  import torch.utils.data
6  import torch.nn as nn
7  import torch.nn.functional as F
8  import matplotlib.pyplot as plt
9  from PIL import Image
10 import os
11 from pprint import pprint
12 import seaborn as sns
13 import cv2
14 from ViTHelper import MasterEncoder
15 from einops import repeat
16 from einops.layers.torch import Rearrange, Reduce
17 import time
18 import datetime
19
20 device = "cuda"
21 device = torch.device("cuda")
22
23 # GLOBAL VARIABLES
24 categories = ["airplane", "bus", "cat", "dog", "pizza"]
25 num_classes = len(categories)
26 class_encoding = {0: "airplane", 1: "bus", 2: "cat", 3: "dog", 4: "pizza"}
27
28 patch_size = 16 # pixels
29 C_in = 3
30 embedded_size = 64
31 max_seq_length = patch_size + 1 # class token
32 image_size = 64 # h x w
33
34 images_Train = '/scratch/gilbreth/dfarache/ece60146/David/HW9/trainingData'
35 images_Val = '/scratch/gilbreth/dfarache/ece60146/David/HW9/valData'
36
37 #DataLoader
38 def ImageProcessing(directory, cat):
39     dir = os.path.join(directory, cat)
40
41     imageFile = [image for image in os.listdir(dir)]
42     imagesPIL = [Image.open(os.path.join(dir, image)).convert("RGB") for image
43                  in imageFile]
44
45     toTensor = [tvt.ToTensor()(image) for image in imagesPIL]
46     toNormalize = [tvt.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5))(tensor) for
47                    tensor in toTensor]
48     return toNormalize

```

```

47
48 class MyDataset(torch.utils.data.Dataset):
49
50     def __init__(self, root, categories):
51         super().__init__()
52         self.categories = categories
53         self.root = root
54         self.data = []
55
56         for i, cat in enumerate(self.categories):
57             images = ImageProcessing(self.root, cat)
58             for image in images:
59                 self.data.append([image, i])
60
61     def __len__(self):
62         return len(self.data)
63
64     def __getitem__(self, i):
65         image = self.data[i][0]
66         label = torch.tensor(self.data[i][1])
67         return image, label

```

Listing 2: Setup Code

4.1 Task 1: Network

For this work, the batch size is 16, and the parameters of the network are a learning rate of $1e-4$ and trained for 20 epochs. The difference between the given assignment and extra credit is using the prior code of AttentionHead_ExtraCredit vs AttentionHead_Assignment.

```

1 # Network
2 # Based on https://towardsdatascience.com/implementing-visualtransformer-in-
  pytorch-184f9f16f632
3 # Based on https://medium.com/mllearning-ai/vision-transformers-from-scratch-
  pytorch-a-step-by-step-guide-96c3313c2e0c
4
5 class PatchEmbedding(nn.Module):
6     def __init__(self, patch_size, embedded_size, image_size, in_channels=3):
7         super(PatchEmbedding, self).__init__()
8
9         self.in_channels = in_channels
10        self.patch_size = patch_size
11        self.embedded_size = embedded_size
12        self.image_size = image_size
13
14        self.sequential = nn.Sequential(
15            # batch_size, embedded_size, patch_height, patch_width
16            nn.Conv2d(in_channels, embedded_size, kernel_size=patch_size, stride
17                    =patch_size),
18
19            # prior to batch_size, patch_height * patch_width, embedded_size
20            Rearrange('b e (h) (w) -> b (h w) e')
21        )
22
23        self.class_token = nn.Parameter(torch.rand(1, 1, self.embedded_size))

```

```

23         self.positions = nn.Parameter(torch.randn((self.image_size // self.
patch_size)**2 + 1, self.embedded_size))
24
25     def forward(self, x):
26         x = self.sequential(x)
27         # Repeat batch
28         class_tokens = repeat(self.class_token, '() n e -> b n e', b=batch_size)
29         x = torch.cat([class_tokens, x], dim=1) # prepend token to match
sequence len
30         x = x + self.positions
31         return x
32
33 # Based on https://towardsdatascience.com/implementing-visualttransformer-in-
pytorch-184f9f16f632
34 class ClassificationHead(nn.Sequential):
35     def __init__(self, embedded_size, num_classes):
36         super().__init__(
37             Reduce('b n e -> b e', reduction='mean'),
38             nn.Linear(embedded_size, num_classes))
39
40 # Assignment Network
41 # Based on https://towardsdatascience.com/implementing-visualttransformer-in-
pytorch-184f9f16f632
42 class ViT_Assignment(nn.Sequential):
43     def __init__(self, how_many_basic_encoders=2, num_attention_heads=2,
in_channels=3):
44         super(ViT_Assignment, self).__init__(
45             PatchEmbedding(patch_size, embedded_size, image_size, in_channels),
46             MasterEncoder(max_seq_length, embedded_size, how_many_basic_encoders
, num_attention_heads),
47             ClassificationHead(embedded_size, num_classes)
48         )
49
50 # Extra Credit Network
51 # Based on https://towardsdatascience.com/implementing-visualttransformer-in-
pytorch-184f9f16f632
52 class ViT_ExtraCredit(nn.Sequential):
53     def __init__(self, how_many_basic_encoders=2, num_attention_heads=2,
in_channels=3):
54         super(ViT_ExtraCredit, self).__init__(
55             PatchEmbedding(patch_size, embedded_size, image_size, in_channels),
56             MasterEncoder(max_seq_length, embedded_size, how_many_basic_encoders
, num_attention_heads, method=2),
57             ClassificationHead(embedded_size, num_classes)
58         )

```

Listing 3: *Network*

4.2 Training, Testing, and Plotting

```

1 # Training
2 def train(net, epochs, lr, dataloader):
3     net = net.to(device)
4     criterion = torch.nn.CrossEntropyLoss()
5     optimizer = torch.optim.Adam(net.parameters(), lr=lr, betas=(0.9, 0.999))
6     lossRun = []

```

```

7
8     for epoch in range(1, epochs+1):
9         running_loss = 0.0
10        for i, data in enumerate(dataloader):
11
12            inputs, labels = data
13            inputs = inputs.to(device)
14            labels = labels.to(device)
15            optimizer.zero_grad()
16            outputs = net(inputs)
17            loss = criterion(outputs, labels)
18            loss.backward()
19            optimizer.step()
20            running_loss += loss.item()
21
22            if((i + 1) % 100 == 0):
23                print("[epoch: %d, batch: %5d] loss: %.3f" % (epoch, i + 1,
running_loss / 100))
24                lossRun.append(running_loss / 100)
25                running_loss = 0.0
26
27    return lossRun

```

Listing 4: *Training*

```

1 # Testing
2 def test(net, dataloader, numCat):
3     net = net.to(device)
4     confusion_matrix = np.zeros((numCat, numCat))
5
6     with torch.no_grad():
7         for inputs, labels in dataloader:
8             inputs = inputs.to(device)
9             labels = labels.to(device)
10            outputs = net(inputs)
11            _, predicted = torch.max(outputs, dim=1)
12            for label, prediction in zip(labels, predicted):
13                confusion_matrix[label][prediction] += 1
14
15    accuracy = np.trace(confusion_matrix) / np.sum(confusion_matrix)
16    return confusion_matrix, accuracy

```

Listing 5: *Testing*

```

1 # Plotting
2 def trainingPlot(loss, epochs):
3     plt.plot(range(len(loss)), loss)
4
5     plt.xlabel("Iterations")
6     plt.ylabel("Loss")
7     plt.legend(loc="lower right")
8     plt.show()
9
10    def confusionMatrix(conf, categories, acc):
11        sns.heatmap(conf, xticklabels=categories, yticklabels=categories, annot=True
)
12        plt.xlabel(f"True Label \n Accuracy: {acc}")
13        plt.ylabel("Predict Label")

```

```
14 plt.show()
```

Listing 6: *Plotting*

5 Results

```
1 # Main
2 ## DataLoaders
3
4 batch_size = 16
5
6 # Create DataLoaders
7 trainDataset = MyDataset(images_Train, categories)
8 trainDataloader = torch.utils.data.DataLoader(trainDataset, batch_size=
    batch_size, shuffle=True, num_workers=2, drop_last=True)
9
10 testDataset = MyDataset(images_Val, categories)
11 testDataloader = torch.utils.data.DataLoader(testDataset, batch_size=batch_size,
    shuffle=True, num_workers=2, drop_last=True)
12
13 ## Parameters for Training
14 lr = 1e-4 # Learning Rate
15 epochs = 20 # Epochs
16
17 ## Assignment
18 net1 = ViT_Assignment()
19 training_loss = train(net1, epochs, lr, trainDataloader)
20 trainingPlot(training_loss, epochs)
21
22 confusion_matrix, accuracy = test(net1, testDataloader, num_classes)
23 confusionMatrix(confusion_matrix, categories, accuracy)
24
25 ## Extra Credit
26 net2 = ViT_ExtraCredit()
27 training_loss = train(net2, epochs, lr, trainDataloader)
28 trainingPlot(training_loss, epochs)
29
30 confusion_matrix, accuracy = test(net2, testDataloader, num_classes)
31 confusionMatrix(confusion_matrix, categories, accuracy)
```

Listing 7: *Main To Run Code*

As seen in Figures 1 and 2, the loss per iteration for both the original assignment and extra credit have the same loss value decrease and plateau. In terms of accuracy, the extra credit slightly outperformed the original assignment code with 54% beginning greater than 53%. Both performed less than the homework 4 assignment with that being 79% greater than 54% and 53%. Overall the DCNN performed better than the ViT model created in this case for both extra credit and assignment-based tasks. This could be due to low embedding size which had to be used due to lack of GPU space.

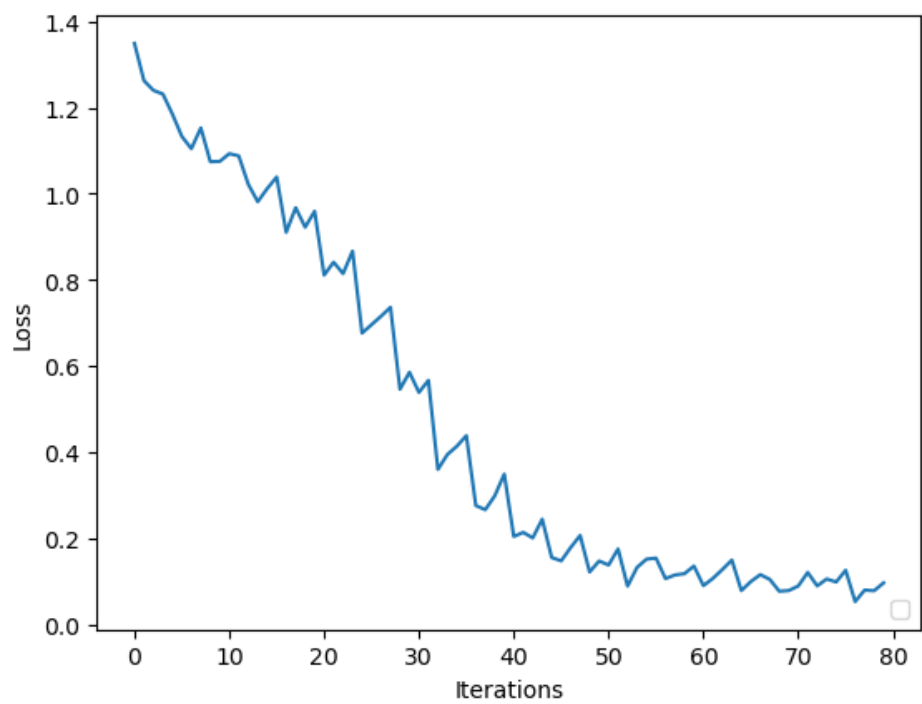


Figure 1: *Assignment Loss Plot*

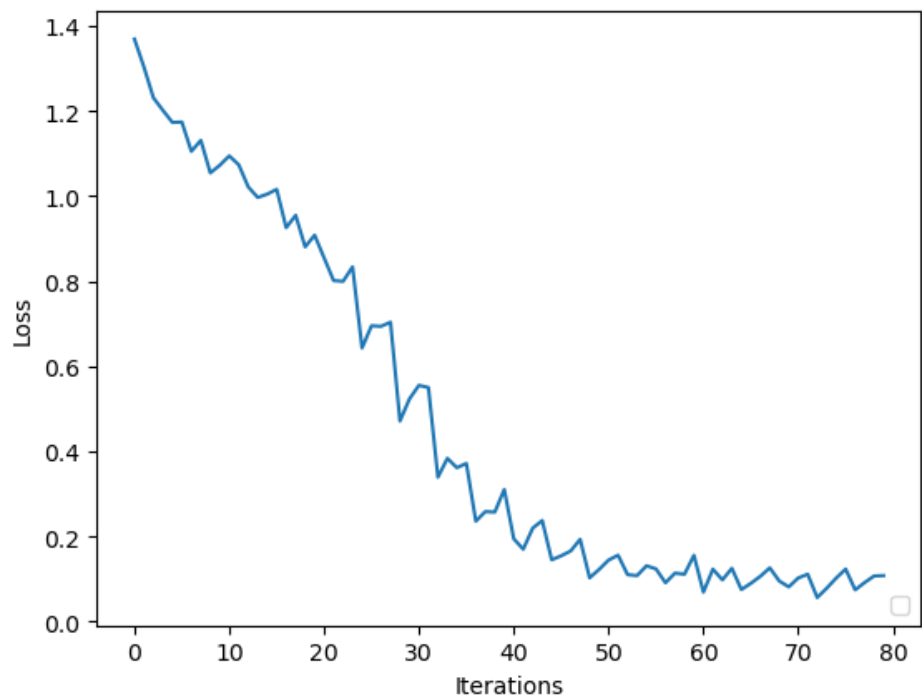


Figure 2: *Extra Credit Loss Plot*

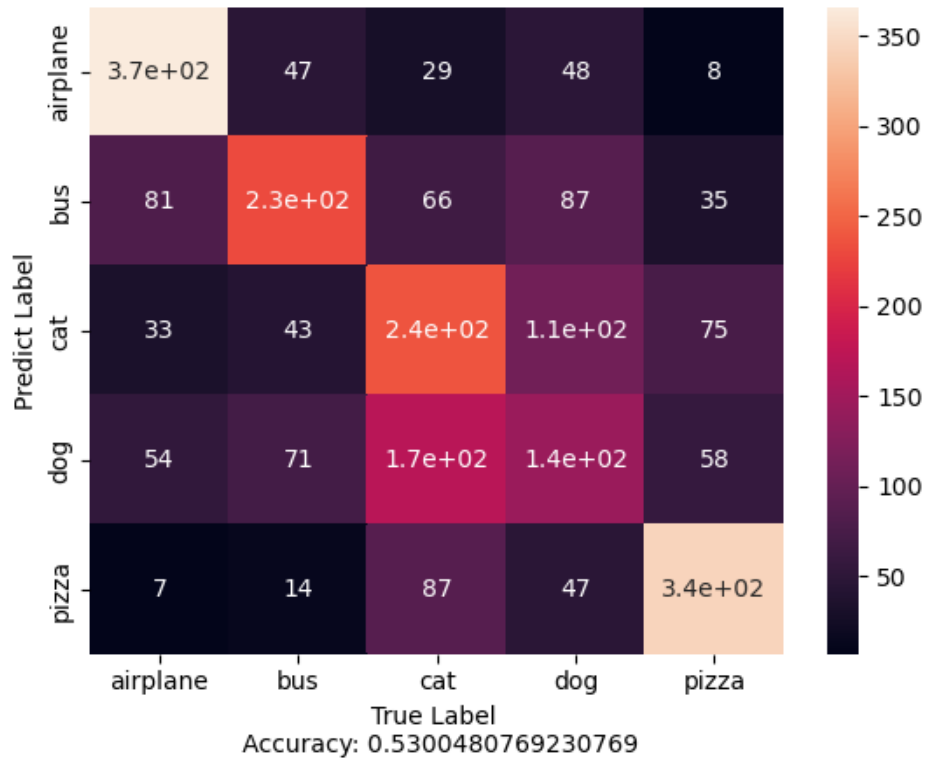


Figure 3: *Assignment Confusion Matrix*

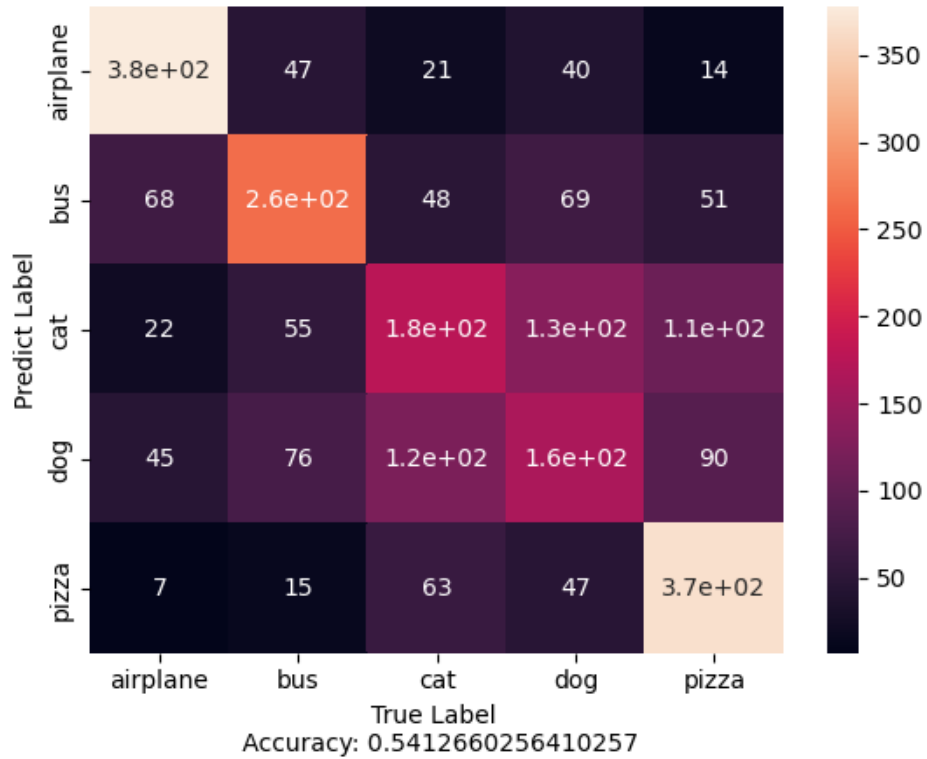


Figure 4: *Extra Credit Confusion Matrix*

6 Lessons learned

What has been learned is the utilization of a transformer and the functions of attention mechanisms.