# Machine Learning with Transformers

Vineet Sai Burugu

08/22/2024

## 1  Introduction

Transformers have shown remarkable capabilities in a number of areas, including image classification. In this paper, we present a transformer-based classifier to identify handwritten digits from the MNIST dataset.

## 2  Tasks

### 2.1  Load the MNIST dataset and preprocess the images and labels

To start, we load the MNIST dataset and preprocess the images and labels. The preprocessing includes converting images to RGB format and applying transformations to prepare the data for the transformer model.

```python
# Install necessary packages
!pip install datasets transformers torch evaluate
!pip install --upgrade torch torchvision torchaudio
!pip install scikit-learn

# Import necessary libraries
from huggingface_hub import login
import torch
import numpy as np
from transformers import ViTForImageClassification,
    ViTImageProcessor, TrainingArguments, Trainer,
    TrainerCallback
from datasets import load_dataset, DatasetDict
import matplotlib.pyplot as plt
```

```python
from sklearn.metrics import confusion_matrix
import seaborn as sns

# Login to Hugging Face
login(token='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')

# Set device to GPU if available, else CPU
device = torch.device('cuda' if torch.cuda.is_available()
    else 'cpu')

# Load and split MNIST dataset
dataset = load_dataset("ylecun/mnist")
labels = list(set(dataset['train']['label']))  # Get unique
    labels
splitDataset = dataset['train'].train_test_split(test_size
    =0.2)  # Split into train/eval
# Create dataset dict
ourDataset = DatasetDict({
    'train': splitDataset['train'],
    'validation': splitDataset['test'],
    'test': dataset['test']
})

# Load ViT image processor
processor = ViTImageProcessor.from_pretrained('google/vit-
    base-patch16-224-in21k')

# Preprocess and transform a batch of images
def transformBatch(batch):
    convertedImages = []  # List for converted images
    for image in batch['image']:
        convertedImages.append(image.convert('RGB'))  #
            Convert to RGB
    batch['image'] = convertedImages  # Replace with
        converted images
    inputs = processor(batch['image'], return_tensors='pt')
        # Process images with ViT
    # Assign labels
    inputs['labels'] = []
    for y in batch['label']:
        inputs['labels'].append(y)
    return inputs  # Return processed inputs

processedDataset = ourDataset.with_transform(transformBatch)
    # Apply transformBatch
```

```
50
51  # Collate a batch into a format for training
52  def collateBatch(batch):
53      pixelValuesList = []   # List for pixel values
54      labelsList = []   # List for labels
55      for item in batch:
56          pixelValuesList.append(item['pixel_values'])   # Add
                  pixel values
57          labelsList.append(item['labels'])   # Add labels
58      return {
59          'pixel_values': torch.stack(pixelValuesList),
60          'labels': torch.tensor(labelsList)
61      }
```

## 2.2 Define a transformer model suitable for image classification

Next, we define a Vision Transformer (ViT) model for image classification. This model will be fine-tuned on the MNIST dataset.

```
1   # Load ViT model for image classification
2   model = ViTForImageClassification.from_pretrained(
3       'google/vit-base-patch16-224',
4       num_labels=10,   # Output labels (10 for MNIST)
5       ignore_mismatched_sizes=True   # Ignore size mismatches
6   )
7   model
8   #model returns the text below
9   ViTForImageClassification(
10    (vit): ViTModel(
11      (embeddings): ViTEmbeddings(
12        (patch_embeddings): ViTPatchEmbeddings(
13          (projection): Conv2d(3, 768, kernel_size=(16, 16),
                 stride=(16, 16))
14        )
15        (dropout): Dropout(p=0.0, inplace=False)
16      )
17      (encoder): ViTEncoder(
18        (layer): ModuleList(
19          (0-11): 12 x ViTLayer(
20            (attention): ViTSdpaAttention(
21              (attention): ViTSdpaSelfAttention(
22                (query): Linear(in_features=768, out_features
                     =768, bias=True)
```

3

```
23              (key): Linear(in_features=768, out_features
                    =768, bias=True)
24              (value): Linear(in_features=768, out_features
                    =768, bias=True)
25              (dropout): Dropout(p=0.0, inplace=False)
26            )
27          (output): ViTSelfOutput(
28              (dense): Linear(in_features=768, out_features
                    =768, bias=True)
29              (dropout): Dropout(p=0.0, inplace=False)
30          )
31        )
32        (intermediate): ViTIntermediate(
33          (dense): Linear(in_features=768, out_features
                =3072, bias=True)
34          (intermediate_act_fn): GELUActivation()
35        )
36        (output): ViTOutput(
37          (dense): Linear(in_features=3072, out_features
                =768, bias=True)
38          (dropout): Dropout(p=0.0, inplace=False)
39        )
40        (layernorm_before): LayerNorm((768,), eps=1e-12,
              elementwise_affine=True)
41        (layernorm_after): LayerNorm((768,), eps=1e-12,
              elementwise_affine=True)
42      )
43    )
44    )
45    (layernorm): LayerNorm((768,), eps=1e-12,
          elementwise_affine=True)
46  )
47  (classifier): Linear(in_features=768, out_features=10, bias
        =True)
48 )
49
50 # Compute accuracy metrics
51 def computeMetrics(evalPreds):
52     logits, labels = evalPreds  # Get logits and labels
53     predictions = np.argmax(logits, axis=1)  # Predict labels
            with argmax
54     # Calculate accuracy manually
55     correct_predictions = np.sum(predictions == labels)  #
            Count correct predictions
56     total_predictions = len(labels)  # Total number of labels
```

```
57        accuracy = correct_predictions / total_predictions    #
              Calculate accuracy
58        return {'accuracy': accuracy} # Return accuracy as a
              dictionary
59
60 # Freeze layers
61 for name, param in model.named_parameters():
62     if not name.startswith('classifier'):
63         param.requires_grad = False  # Freeze parameter
```

## 2.3 Train the transformer model on the MNIST training dataset

We set up the training configuration and train the model using the MNIST training dataset.

```
1  # Custom callback to log training metrics
2  class MetricsLogger(TrainerCallback):
3      def __init__(self):
4          super().__init__()
5          self.trainLosses = []   # Store training losses
6          self.valLosses = []   # Store validation losses
7          self.valAccuracies = []   # Store validation
                  accuracies
8          self.currentEpochTrainLosses = []   # Store current
                  epoch's training losses
9
10     def on_log(self, args, state, control, logs=None, **
           kwargs):
11         if logs is not None:
12             if 'loss' in logs:
13                 self.currentEpochTrainLosses.append(logs['
                        loss'])   # Log train loss
14             if 'eval_loss' in logs:
15                 self.valLosses.append(logs['eval_loss'])   #
                        Log val loss
16             if 'eval_accuracy' in logs:
17                 self.valAccuracies.append(logs['eval_accuracy
                        '])   # Log val accuracy
18
19     def on_epoch_end(self, args, state, control, **kwargs):
20         if self.currentEpochTrainLosses:
21             avgTrainLoss = np.mean(self.
                   currentEpochTrainLosses)   # Avg train loss for
```

5

```python
                     epoch
22               self.trainLosses.append(avgTrainLoss)  # Log avg
                     train loss
23               self.currentEpochTrainLosses = []  # Reset for
                     next epoch
24
25       def reset(self):
26           self.trainLosses = []
27           self.valLosses = []
28           self.valAccuracies = []
29           self.currentEpochTrainLosses = []
30
31  # Set up training arguments
32  trainingArgs = TrainingArguments(
33       output_dir="./vit-base-mnist",  # Directory for
             checkpoints
34       per_device_train_batch_size=32,  # Training batch size
35       evaluation_strategy="epoch",  # Evaluate at each epoch
             end
36       save_strategy="epoch",  # Save at each epoch end
37       logging_steps=200,  # Log every 200 steps
38       num_train_epochs=5,  # Number of epochs
39       learning_rate=2e-4,  # Optimizer learning rate
40       save_total_limit=3,  # Max number of checkpoints
41       remove_unused_columns=False,  # Keep all dataset columns
42       report_to='tensorboard',  # Report to TensorBoard
43       load_best_model_at_end=True,  # Load best model at end
44  )
45
46  # Initialize MetricsLogger callback
47  metricsLogger = MetricsLogger()
48  metricsLogger.reset()
49
50  # Create Trainer instance
51  trainer = Trainer(
52       model=model,
53       args=trainingArgs,
54       data_collator=collateBatch,
55       compute_metrics=computeMetrics,
56       train_dataset=processedDataset['train'],
57       eval_dataset=processedDataset['validation'],
58       tokenizer=processor,
59       callbacks=[MetricsLogger()]  # Add custom callback
60  )
61
```

```
62  # Train the model
63  trainer.train()
```

## 2.4 Evaluate the transformer model on the MNIST test dataset

Finally, we evaluate the trained model on the MNIST test dataset and analyze its performance.

```
1   # Plot validation accuracies over epochs
2   plt.plot(metricsLogger.valAccuracies, '-x')  # Plot
        validation accuracies
3   plt.xlabel('Epoch')  # X-axis label
4   plt.ylabel('Validation Accuracy')  # Y-axis label
5   plt.title('Validation Accuracy vs. Epochs')  # Title
6   plt.show()  # Show plot
7
8   # Plot training and validation losses over epochs
9   plt.plot(metricsLogger.valLosses, '-x')  # Plot validation
        losses
10  plt.xlabel('Epoch')  # X-axis label
11  plt.ylabel('Validation Loss')  # Y-axis label
12  plt.title('Validation Loss vs. Epochs')  # Title
13  plt.show()  # Show plot
14
15  # Predict test labels, compute confusion matrix, and plot
        heatmap
16  predictionsOutput = trainer.predict(processedDataset['test'])
         # Predict on test dataset
17  testPreds = np.argmax(predictionsOutput.predictions, axis=1)
         # Get predicted labels
18  testLabels = predictionsOutput.label_ids  # Get true labels
19  confMatrix = confusion_matrix(testLabels, testPreds)  #
        Compute confusion matrix
20  plt.figure(figsize=(10, 7))  # Set figure size
21  sns.heatmap(confMatrix, annot=True, fmt='d', cmap='Blues',
        xticklabels=range(10), yticklabels=range(10))  # Heatmap
        of confusion matrix
22  plt.xlabel('Predicted Label')  # X-axis label
23  plt.ylabel('True Label')  # Y-axis label
24  plt.title('Confusion Matrix')  # Plot title
25  plt.show()  # Show plot
```

# 3 Discussion

## 3.1 Compare the performance of the transformer model with traditional convolutional neural networks (CNNs) typically used for MNIST classification

For the MNIST dataset, I designed a conventional convolutional network. Installing the essential packages and importing the relevant libraries are the first things we do. In the first job, we establish DataLoader objects for training, validation, and testing; we also define transformations for data augmentation, load the MNIST dataset for training and testing, and divide the training dataset into training and validation sets. We define the CNN class with various layers in the following challenge. In the next step, we initialize lists to hold epoch-wise data, define the loss function and optimizer, set the device to GPU if available, otherwise to CPU, and begin the training process. The last process involves testing the model using the test data, calculating the accuracy of the test, plotting the accuracy and loss with time, computing the confusion matrix, and finally plotting the confusion matrix.

| Epoch | Training Loss | Validation Loss | Accuracy |
|-------|---------------|-----------------|----------|
| 1     | 0.3204        | 0.2960          | 0.9220   |
| 2     | 0.2335        | 0.2229          | 0.9385   |
| 3     | 0.1979        | 0.1968          | 0.9461   |
| 4     | 0.1868        | 0.1842          | 0.9489   |
| 5     | 0.1849        | 0.1811          | 0.9503   |

Table 1: Performance of Vision Transformer (VIT) model

| Epoch | Training Loss | Validation Loss | Accuracy |
|-------|---------------|-----------------|----------|
| 1     | 0.8500        | 0.2990          | 0.8979   |
| 2     | 0.3229        | 0.1593          | 0.9522   |
| 3     | 0.2344        | 0.1165          | 0.9630   |
| 4     | 0.1884        | 0.0938          | 0.9701   |
| 5     | 0.1621        | 0.0798          | 0.9735   |

Table 2: Performance of CNN model

**Comparison between models**

1. The Vision Transformer (VIT) achieves an accuracy of 0.95, whereas the Convolutional Neural Network (CNN) reaches a higher accuracy of 0.97. This indicates that the CNN is more effective at classification than the VIT.

2. The training loss for the VIT decreases from 0.32 to 0.18, while the CNN's training loss drops from 0.85 to 0.16. This suggests that the CNN learns more efficiently than the VIT.

3. For validation loss, the VIT shows a reduction from 0.30 to 0.18, whereas the CNN's validation loss decreases from 0.30 to 0.08. This demonstrates that the CNN generalizes better compared to the VIT.



Figure 1: VIT Validation Accuracy



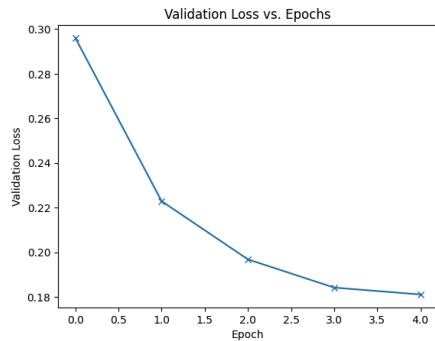Figure 2: CNN Validation Accuracy
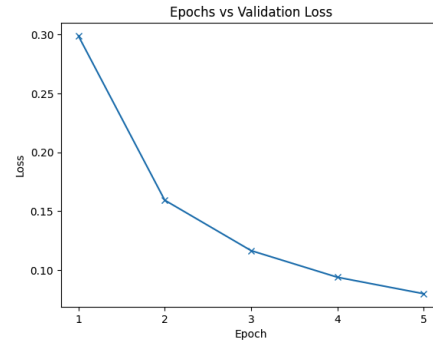


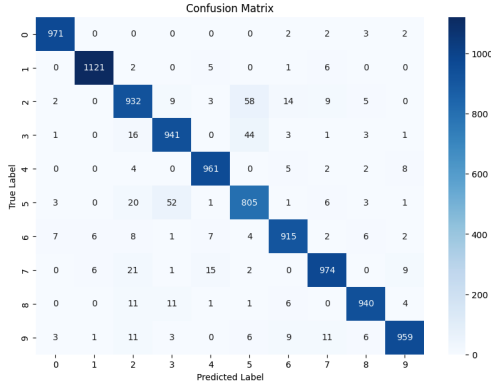Figure 3: VIT Validation Loss



Figure 4: CNN Validation Loss
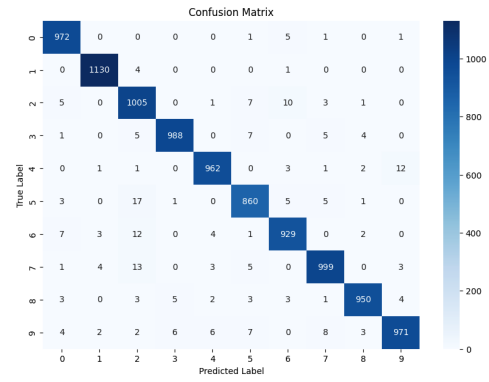
9

Figure 5: VIT Confusion Matrix



Figure 6: CNN Confusion Matrix

## 3.2 Highlight the strengths and weaknesses of using a transformer model for this type of task.

**Strengths**

1. Due to their ability to interpret images in their entirety, Visual Transformers capture context between different portions of an image, allowing the model to effectively identify subtle relationships.

2. They use patch embeddings that allocate weights to various visual segments, or "patches," enabling the model to understand the essence of each portion of the image and quickly find relationships.

3. Transformers utilize self-attention, which assigns similarity scores between different portions of the image using the patch embeddings.

**Weaknesses**

1. Due to the fact that transformers look at an image as a whole, they are very demanding on the device they train on. CNN look at localized regions, so less computational resources are required.

2. By finding unique peculiarities in the training dataset, transformers run the risk of over fitting their data, meaning they fit well on the training data, but not for the test data.

## 3.3 Reflect on the challenges you faced during the implementation and training of the transformer model

**Finding the Right Training Arguments**  Because the program takes a long time to run (almost an hour or two), making modifications to the training arguments was often difficult. These training arguments include factors like batch size, logging steps, and learning rate. Initially, I had set the batch size to 64, which was the same as in my CNN model, to ensure a fair comparison between the CNN and ViT models. However, I noticed that the program's performance wasn't as good. I experimented with different batch sizes—both increasing and decreasing it. Decreasing the batch size actually resulted in a more accurate model.

**Metrics Logger for the Graphs**  Originally, there was no way to keep track of metrics like validation loss and training loss because the HuggingFace interface didn't allow me to store that data in a list or anything similar. The only place these metrics were displayed was in the table that showed up during training. To resolve this, I created a MetricsLogger callback to store these variables in a list. I then used this data to plot validation graphs.

**Problems Plotting the Confusion Matrix**  I initially struggled with plotting the confusion matrix and wasn't sure how to approach it. In addition, I wasn't sure how to make predictions from the model. After some research on Seaborn, I learned that it can be used to draw a heatmap for the confusion matrix, and to get all the predictions, I simply called the trainer to predict based on the test dataset.

## 3.4 Consider the potential improvements or alternative approaches that could enhance the performance of the transformer classifier on the MNIST dataset

**Data Augmentation**  As of this moment, my code doesn't use data augmentation although it could be beneficial. In data augmentation, flipping, rotation, scaling, and other transformations are applied to data (not just images). Transformers require vast amounts of data for high accuracy, and

augmentation helps to combat this problem by artificially creating additional data. More data generally correlates with better performance for many models. It also helps with overfitting, which occurs when our model performs very well on the training data but not on the test data. This indicates that the model is not generalized enough for all data, which can hinder model performance.

**Additional Metrics**   The metric analyzed in this model is accuracy, which allows us to see how "correct" our model is, in other words, how true our predictions were. Despite this, there are other metrics that allow for a better assessment of our model. Precision provides the portion of all positive classifications that turn out to actually be positive. Recall gives the portion of all actual positive classifications that are classified as actual positives. The F1-score calculates the average of both precision and recall.

**Alternative Models**   The model in this current project is a pre-trained Vision Transformer that has been adjusted for the MNIST database. If we wanted better accuracy and speed, other transformer models that could have been used include the Swin Transformer, DeiT (Data Efficient Image Transformer), CvT (Convolutional Vision Transformer), and more. The paper CvT: Swin Transformer: Hierarchical Vision Transformer Using Shifted Windows introduces the Swin Transformer, which has a hierarchical vision Transformer architecture using shifted windows, enabling efficient computation and strong performance across various vision tasks. Additionally, we could have created a Vision Transformer from scratch with just PyTorch. This approach allows for more customization but would likely result in lower accuracy compared to both the CNN and pre-trained Vision Transformer, and the training speed would be considerably slower, potentially taking many hours.

# 4   Code

In this study, we have provided detailed analyses and implementations for both the Vision Transformer (ViT) and the Convolutional Neural Network (CNN) models. The corresponding Jupyter Notebooks can be accessed via the following links:

**Vision Transformer (ViT)** ViT Notebook (if you decide to run the code, here is token you can use: "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")

**Convolutional Neural Network (CNN)** CNN Notebook