# Federated Learning Project Report

Katta Sathwik
20mcme16

December 29, 2023

# Contents

**Abstract**

This report presents a detailed analysis of a Log Anomaly Detection using federated learning project focusing on sequence classification using a fine-tuned BERT model. The project involves client-server communication using Flower framework for federated learning. This report outlines the key concepts and tasks involved in the implementation.

# 1    Introduction

Federated learning is a distributed machine learning approach that allows training models across decentralized devices or servers. In this project, I have used the Flower framework for federated learning to implement the solution for sequence classification. The model used is a fine-tuned BERT (Bidirectional Encoder Representations from Transformers) model, a state-of-the-art transformer-based architecture.

# 2    Objective

The main objective of the project is to implement a federated learning system for sequence classification. The client-side code defines a Flower client using NumPyClient. The client utilizes a custom BERT model for sequence classification and a tokenizer for text encoding. The model is fine-tuned on a dataset and saved for later use.

# 3    Explanation of Model

The implementation of federated learning using the Flower framework is explained. Here are the key points:

## 3.1    Starting the Flower Server

1. `fl.server.start_server` is called to start the Flower server. It listens on a specific address and port (localhost and the port specified as a command-line argument).

2. `ServerConfig` is used to configure the server with the number of federated learning rounds (`num_rounds`) and the maximum message length.

## 3.2    Flower Client Implementation

1. Initially, import the necessary libraries for the implementation:

    - `flwr` for federated learning.
    - `tensorflow` and `keras` for machine learning and neural network models.

3

- **sys** to access command-line arguments.
- **numpy** for numerical operations.

2. Define auxiliary methods:

   - **getDist**: Meant for visualizing data class distribution.
   - **getData**: Customize the data distribution.

3. Load and compile the Keras model:

   - Define a neural network model using the Keras Sequential API. It's a simple feedforward model with three layers: input, hidden, and output layers.
   - Compile the model with an optimizer (**"adam"**), a loss function (**"sparse_categorical_crossentropy"**), and a metric to track (**"accuracy"**).

4. Load the dataset:

   - Load the Fashion MNIST dataset using TensorFlow's built-in dataset loader.
   - Normalize the images to values between 0 and 1.
   - Specify a data distribution (**dist**) for Fashion MNIST. This distribution controls how many samples are included from each class (0-9).

5. Define the Flower client:

   - Define a custom Flower client class, **CustomClient**, which extends **fl.client.NumPyClient**.
   - Implement three main methods:
     (a) **get_parameters**: Returns the current model's weights.
     (b) **fit**: Used for model training. Receives model parameters, updates the model, and fits the model to the training data for one epoch. Returns the updated model parameters and the number of training samples used.
     (c) **evaluate**: Evaluates the model on the test data and returns the loss and accuracy.

6. Start the Flower client:

   - **fl.client.start_numpy_client** is used to start the Flower client.
   - It connects to the Flower server at the specified address.
   - Use the custom client class (**CustomClient**), and set the maximum message length to ensure compatibility with the server.

4

## 3.3 Fashion MNIST Dataset

Fashion MNIST is a dataset commonly used in machine learning and computer vision tasks, often serving as a drop-in replacement for the original MNIST dataset. Instead of handwritten digits, Fashion MNIST contains grayscale images of various clothing items and accessories. Here's an explanation of the Fashion MNIST dataset:

### 3.3.1 Images

Fashion MNIST consists of a collection of $28 \times 28$-pixel grayscale images. Each image represents one of ten different fashion categories.

### 3.3.2 Classes

There are ten classes, each corresponding to a specific type of clothing item or accessory. These classes are as follows:

1. T-shirt/top

2. Trouser

3. Pullover

4. Dress

5. Coat

6. Sandal

7. Shirt

8. Sneaker

9. Bag

10. Ankle boot

### 3.3.3 Size

The dataset is relatively small compared to some other image datasets. It includes 60,000 training images and 10,000 test images.

# 4 Basic Implementing a Federated Learning Setup using Terminals as Clients with FashionMNIST Dataset

## 4.1 Starting the Flower Server

To start the Flower server, use the following command in the terminal:

```
1 fl.server.start_server --server_address="127.0.0.1:5012" --
      num_rounds=5 --max_msg_len=1000000
```

<div align="center">Listing 1: Starting the Flower Server</div>

## 4.2 Starting Flower Clients

Open multiple terminals and run the Flower clients using the following command:

```
1 python fashion_mnist_client.py --server_address="
      127.0.0.1:5012"
```

<div align="center">Listing 2: Starting Flower Clients</div>

Ensure that Flower clients are running on different terminals to simulate a federated learning environment.

# 5 Federated client for aggregation

## 5.1 Importing the necessary libraries

These lines import necessary libraries and modules for working with PyTorch, Flower (Federated Learning framework), Hugging Face Transformers (for BERT), and other utilities.

```
1 from dataclasses import dataclass
2 from typing import Any, Dict, List, Union
3 import torch
4 import flwr as fl
5 from transformers import BertForSequenceClassification,
      BertTokenizer, Trainer, TrainingArguments
6 from transformers import AdamW,
      get_linear_schedule_with_warmup
7 import torch.nn.functional as F
8 from torch.utils.data import Dataset, DataLoader
9 import zipfile
10 import pandas as pd
11 import shutil
12 import os
13 from sklearn.model_selection import train_test_split
```

<div align="center">Listing 3: Hdfs Client for aggregation</div>

## 5.2 Custom Dataset Class:

This class defines a custom dataset for PyTorch. It takes encodings (tokenized input sequences) and labels and provides methods for getting items and the length of the dataset.

```python
1  class CustomDataset(Dataset):
2      def __init__(self, encodings, labels):
3          self.encodings = encodings
4          self.labels = labels
5
6      def __getitem__(self, idx):
7          item = {key: torch.tensor(val[idx]) for key, val in
       self.encodings.items()}
8          item['labels'] = torch.tensor([self.labels[idx], 1 -
        self.labels[idx]]).float()
9          return item
10
11      def __len__(self):
12          return len(self.labels)
```

## 5.3 Initialization and Data Preprocessing:

This section initializes the BERT model and tokenizer and loads a labeled dataset from a CSV file. It selects a subset of the dataset (rows 1 to 200) and preprocesses the data.

```python
1  model_saved = BertForSequenceClassification.from_pretrained(
       "bert-base-uncased")
2  tokenizer_saved = BertTokenizer.from_pretrained("bert-base-
       uncased")
3  model_folder = "fine_tuned_bert_model_for_HDFS"
4  pickle_file = "pickle_fine_tuned_bert_model_for_HDFS.pickle"
5  df = pd.read_csv('Hdfs_labelled_sequence.csv', sep=',',
       quotechar='"', names=["text", "label"])
6  df = df[1:200]
7  X = list(df['text'])
8  y = pd.get_dummies(df['label'], drop_first=True)['Normal'].
       values
9  y = y.astype(int)
```

## 5.4 Model Saving and Loading:

This part saves the initial BERT model's weights and then loads the saved weights from a ZIP file. It ensures that the model and tokenizer are initialized with the same state as before saving.

```python
1  torch.save(model_saved.state_dict(), os.path.join(
       model_folder, "model_weights.pth"))
2  shutil.make_archive(model_folder, 'zip', model_folder)
3  with zipfile.ZipFile(f"{model_folder}.zip", "r") as zip_ref:
4      with zip_ref.open("model_weights.pth", "r") as file:
5          model_weights = torch.load(file)
6  model_saved = BertForSequenceClassification.from_pretrained(
       "bert-base-uncased")
```

```
7 model_saved.load_state_dict(model_weights)
8 tokenizer_saved = BertTokenizer.from_pretrained("bert-base-
      uncased")
```

## 5.5   Training Arguments and Trainer Initialization:

Here, training arguments for the Trainer are defined, specifying various parameters like batch size, evaluation steps, etc. The Trainer is then initialized with the BERT model, training arguments, and training and evaluation datasets.

```
1  training_args = TrainingArguments(
2          output_dir="./bert_base_model",
3          evaluation_strategy="steps",
4          eval_steps=100,
5          per_device_train_batch_size=2,
6          per_device_eval_batch_size=2,
7          save_steps=1000,
8          save_total_limit=2,
9          num_train_epochs=10,
10         logging_dir="./logs",
11     )
12 trainer = Trainer(
13     model=model_saved,
14     args=training_args,
15     train_dataset=dataset_train,
16     eval_dataset=test_dataset,
17     data_collator=None
18 )
```

## 5.6   FlowerClient Class:

This class is a custom client for federated learning using the Flower framework. It has methods for getting parameters, fitting the model, and evaluating the model, among others.

```
1 class FlowerClient(fl.client.NumPyClient):
2     ...
```

## 5.7   Fitting the Model (Training):

loop runs the training process for a specified number of epochs. It iterates through batches of the local dataset, computes the loss, performs backpropagation, and updates the model parameters.

```
1 for epoch in range(3):
2     for batch in DataLoader(test_dataset, batch_size=2,
      shuffle=True):
3          ...
```

## 5.8  Evaluating the Model:

This loop runs the evaluation process on the local dataset. It computes the loss and evaluates the model's performance.

```
1 for batch in DataLoader(test_dataset, batch_size=2, shuffle=
      False):
2      ...
```

## 5.9  Printing Fit History and Evaluation Metrics:

These print statements provide information on the fit history (training loss) and evaluation metrics (evaluation loss and global evaluation accuracy).

```
1 print("Fit history: {'loss':", loss.item(), "}")
2 print("Evaluation metrics: {'eval_loss':", average_loss, "}"
      )
3 print("Now the global Eval accuracy:", 1.0 - average_loss)
```

## 5.10  Flower Client Initialization and Execution:

The script starts the Flower client, connecting it to a server address and using the custom FlowerClient for federated learning.

```
1 if __name__ == '__main__':
2     fl.client.start_numpy_client(
3         server_address="127.0.0.1:5020",
4         client=FlowerClient(),
5     )
```

# 6  HDFS Trained Model Aggregation With Federated Learning

Initially, I have worked on the HDfs trained model in order to aggregate the model using the federated learning, which automatically works for the aggregating three trained models but that was taking too long and more ram even if take a small amount of data for preprocessing.

# 7  Conclusion

This report demonstrates federated learning with BERT using the Flower framework. It includes data preprocessing, model saving/loading, training, and evaluation in a federated learning setup. The Flower client facilitates communication between the server and clients for federated learning. The provided print statements offer insights into the training and evaluation processes.

# 8    Reference

https://github.com/sushanthk-262/Log anomaly detection

Figure 1: Output of the server



Figure 2: Output of client 1



Figure 3: Output of client 2
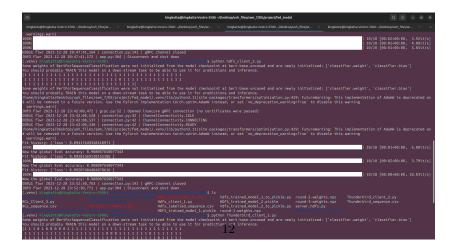
Figure 4: Hdfs Server



Figure 5: Hdfs Client 1



Figure 6: Hdfs Client 1

12

Figure 7: Ram usage while Training