

In [1]:

```
import torch

# If there's a GPU available...
if torch.cuda.is_available():

    # Tell PyTorch to use the GPU.
    device = torch.device("cuda")

    print('There are %d GPU(s) available.' % torch.cuda.device_count())

    print('We will use the GPU:', torch.cuda.get_device_name(0))
    !nvidia-smi

# If not...
else:
    print('No GPU available, using the CPU instead.')
    device = torch.device("cpu")
```

There are 1 GPU(s) available.

We will use the GPU: Tesla T4

Fri Jun 10 15:05:30 2022

```
+-----+
+-----+
| NVIDIA-SMI 460.32.03      Driver Version: 460.32.03      CUDA Version:
11.2      |
|-----+-----+-----+-----+
+-----+
| GPU   Name           Persistence-M| Bus-Id        Disp.A | Volatile Un
corr. ECC | Fan  Temp  Perf  Pwr:Usage/Cap|      Memory-Usage | GPU-Util  C
ompute M. |
|
| MIG M. |
|=====+=====+=====+=====
=====|
|    0  Tesla T4             Off  | 00000000:00:04:0  Off  |
0 |
| N/A   38C    P8      9W / 70W |      3MiB / 15109MiB |      0%
Default |
|
| N/A |
+-----+-----+-----+-----+
+-----+
+-----+
| Processes:
|
| GPU   GI    CI          PID    Type    Process name                  G
PU Memory |
|      ID    ID                                   U
sage      |
|=====+=====+=====+=====
=====|
| No running processes found
|
+-----+-----+-----+-----+
+-----+
```

**Importing the libraries needed**

In [2]:

```
!pip install pyarabic  
!pip install emoji  
!pip install pystemmer  
!pip install optuna==2.3.0  
!pip install transformers==4.2.1
```

Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public/simple/>  
Collecting pyarabic  
 Downloading PyArabic-0.6.14-py3-none-any.whl (126 kB)  
 |██| 126 kB 29.3 MB/s  
Requirement already satisfied: six>=1.14.0 in /usr/local/lib/python3.7/dist-packages (from pyarabic) (1.15.0)  
Installing collected packages: pyarabic  
Successfully installed pyarabic-0.6.14  
Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public/simple/>  
Collecting emoji  
 Downloading emoji-1.7.0.tar.gz (175 kB)  
 |██| 175 kB 16.2 MB/s  
Building wheels for collected packages: emoji  
 Building wheel for emoji (setup.py) ... done  
 Created wheel for emoji: filename=emoji-1.7.0-py3-none-any.whl size=171046 sha256=3982bb82c79905d01a1a4e72c967b568b9c255dc757cace82d34e99908218779  
 Stored in directory: /root/.cache/pip/wheels/8a/4e/b6/57b01db010d17ef6ea9b40300af725ef3e210cblacfb7ac8b6  
Successfully built emoji  
Installing collected packages: emoji  
Successfully installed emoji-1.7.0  
Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public/simple/>  
Collecting pystemmer  
 Downloading PyStemmer-2.0.1.tar.gz (559 kB)  
 |██| 559 kB 7.7 MB/s  
Building wheels for collected packages: pystemmer  
 Building wheel for pystemmer (setup.py) ... done  
 Created wheel for pystemmer: filename=PyStemmer-2.0.1-cp37-cp37m-linux\_x86\_64.whl size=425665 sha256=8e37193f072bf2c0ee10f7b388bcfffd218e3891715d6ac60230501008d098c82  
 Stored in directory: /root/.cache/pip/wheels/30/6d/40/0d17a498c5009922dbb3ddaca3d3652387ba94cc96142001f0  
Successfully built pystemmer  
Installing collected packages: pystemmer  
Successfully installed pystemmer-2.0.1  
Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public/simple/>  
Collecting optuna==2.3.0  
 Downloading optuna-2.3.0.tar.gz (258 kB)  
 |██| 258 kB 35.5 MB/s  
Installing build dependencies ... done  
Getting requirements to build wheel ... done  
Preparing wheel metadata ... done  
Requirement already satisfied: joblib in /usr/local/lib/python3.7/dist-packages (from optuna==2.3.0) (1.1.0)  
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.7/dist-packages (from optuna==2.3.0) (21.3)  
Requirement already satisfied: tqdm in /usr/local/lib/python3.7/dist-packages (from optuna==2.3.0) (4.64.0)  
Requirement already satisfied: sqlalchemy>=1.1.0 in /usr/local/lib/python3.7/dist-packages (from optuna==2.3.0) (1.4.36)  
Collecting colorlog  
 Downloading colorlog-6.6.0-py2.py3-none-any.whl (11 kB)  
Collecting cmaes>=0.6.0  
 Downloading cmaes-0.8.2-py3-none-any.whl (15 kB)  
Requirement already satisfied: numpy in /usr/local/lib/python3.7/dist-packages (from optuna==2.3.0) (1.21.6)

Requirement already satisfied: scipy!=1.4.0 in /usr/local/lib/python3.7/dist-packages (from optuna==2.3.0) (1.4.1)

Collecting cliff

Downloading cliff-3.10.1-py3-none-any.whl (81 kB)

|██| 81 kB 5.9 MB/s

Collecting alembic

Downloading alembic-1.8.0-py3-none-any.whl (209 kB)

|██| 209 kB 37.4 MB/s

Requirement already satisfied: pyparsing!=3.0.5,>=2.0.2 in /usr/local/lib/python3.7/dist-packages (from packaging>=20.0->optuna==2.3.0) (3.0.9)

Requirement already satisfied: greenlet!=0.4.17 in /usr/local/lib/python3.7/dist-packages (from sqlalchemy>=1.1.0->optuna==2.3.0) (1.1.2)

Requirement already satisfied: importlib-metadata in /usr/local/lib/python3.7/dist-packages (from sqlalchemy>=1.1.0->optuna==2.3.0) (4.1.1)

Requirement already satisfied: importlib-resources in /usr/local/lib/python3.7/dist-packages (from alembic->optuna==2.3.0) (5.7.1)

Collecting Mako

Downloading Mako-1.2.0-py3-none-any.whl (78 kB)

|██| 78 kB 6.8 MB/s

Requirement already satisfied: PrettyTable>=0.7.2 in /usr/local/lib/python3.7/dist-packages (from cliff->optuna==2.3.0) (3.3.0)

Collecting autopage>=0.4.0

Downloading autopage-0.5.1-py3-none-any.whl (29 kB)

Collecting pbr!=2.1.0,>=2.0.0

Downloading pbr-5.9.0-py2.py3-none-any.whl (112 kB)

|██| 112 kB 75.6 MB/s

Collecting cmd2>=1.0.0

Downloading cmd2-2.4.1-py3-none-any.whl (146 kB)

|██| 146 kB 3.7 MB/s

Requirement already satisfied: PyYAML>=3.12 in /usr/local/lib/python3.7/dist-packages (from cliff->optuna==2.3.0) (3.13)

Collecting stevedore>=2.0.1

Downloading stevedore-3.5.0-py3-none-any.whl (49 kB)

|██| 49 kB 6.6 MB/s

Requirement already satisfied: attrs>=16.3.0 in /usr/local/lib/python3.7/dist-packages (from cmd2>=1.0.0->cliff->optuna==2.3.0) (21.4.0)

Requirement already satisfied: typing-extensions in /usr/local/lib/python3.7/dist-packages (from cmd2>=1.0.0->cliff->optuna==2.3.0) (4.2.0)

Requirement already satisfied: wcwidth>=0.1.7 in /usr/local/lib/python3.7/dist-packages (from cmd2>=1.0.0->cliff->optuna==2.3.0) (0.2.5)

Collecting pyperclip>=1.6

Downloading pyperclip-1.8.2.tar.gz (20 kB)

Requirement already satisfied: zipp>=0.5 in /usr/local/lib/python3.7/dist-packages (from importlib-metadata->sqlalchemy>=1.1.0->optuna==2.3.0) (3.8.0)

Requirement already satisfied: MarkupSafe>=0.9.2 in /usr/local/lib/python3.7/dist-packages (from Mako->alembic->optuna==2.3.0) (2.0.1)

Building wheels for collected packages: optuna, pyperclip

Building wheel for optuna (PEP 517) ... done

Created wheel for optuna: filename=optuna-2.3.0-py3-none-any.whl size=359772 sha256=f7c939d396be3d7f030b944ec7ec12a5c6fc5e134dbb93084b29acbd45671e11

Stored in directory: /root/.cache/pip/wheels/38/61/9e/955ab1890f6cab231b1d756db63f36c711968a324296e0b649

Building wheel for pyperclip (setup.py) ... done

Created wheel for pyperclip: filename=pyperclip-1.8.2-py3-none-any.whl size=11137 sha256=f3d67785f928b49457870d502f951a2cde650f399567

8ab4bd26422dcfe0a33c

Stored in directory: /root/.cache/pip/wheels/9f/18/84/8f69f8b08169c7bae2dde6bd7daf0c19fca8c8e500ee620a28

Successfully built optuna pyperclip

Installing collected packages: pyperclip, pbr, stevedore, Mako, cmd2, autopage, colorlog, cmaes, cliff, alembic, optuna

Successfully installed Mako-1.2.0 alembic-1.8.0 autopage-0.5.1 cliff-3.10.1 cmaes-0.8.2 cmd2-2.4.1 colorlog-6.6.0 optuna-2.3.0 pbr-5.9.0 pyperclip-1.8.2 stevedore-3.5.0

Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public/simple/>

Collecting transformers==4.2.1

Downloading transformers-4.2.1-py3-none-any.whl (1.8 MB)

|██| 1.8 MB 29.9 MB/s

Requirement already satisfied: filelock in /usr/local/lib/python3.7/dist-packages (from transformers==4.2.1) (3.7.0)

Requirement already satisfied: importlib-metadata in /usr/local/lib/python3.7/dist-packages (from transformers==4.2.1) (4.11.4)

Requirement already satisfied: tqdm>=4.27 in /usr/local/lib/python3.7/dist-packages (from transformers==4.2.1) (4.64.0)

Collecting tokenizers==0.9.4

Downloading tokenizers-0.9.4-cp37-cp37m-manylinux2010\_x86\_64.whl (2.9 MB)

|██| 2.9 MB 56.0 MB/s

Requirement already satisfied: numpy in /usr/local/lib/python3.7/dist-packages (from transformers==4.2.1) (1.21.6)

Requirement already satisfied: packaging in /usr/local/lib/python3.7/dist-packages (from transformers==4.2.1) (21.3)

Requirement already satisfied: requests in /usr/local/lib/python3.7/dist-packages (from transformers==4.2.1) (2.23.0)

Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.7/dist-packages (from transformers==4.2.1) (2019.12.20)

Collecting sacremoses

Downloading sacremoses-0.0.53.tar.gz (880 kB)

|██| 880 kB 53.8 MB/s

Requirement already satisfied: zipp>=0.5 in /usr/local/lib/python3.7/dist-packages (from importlib-metadata->transformers==4.2.1) (3.8.0)

Requirement already satisfied: typing-extensions>=3.6.4 in /usr/local/lib/python3.7/dist-packages (from importlib-metadata->transformers==4.2.1) (4.2.0)

Requirement already satisfied: pyparsing!=3.0.5,>=2.0.2 in /usr/local/lib/python3.7/dist-packages (from packaging->transformers==4.2.1) (3.0.9)

Requirement already satisfied: urllib3!=1.25.0,!<1.25.1,<1.26,>=1.21.1 in /usr/local/lib/python3.7/dist-packages (from requests->transformers==4.2.1) (1.24.3)

Requirement already satisfied: chardet<4,>=3.0.2 in /usr/local/lib/python3.7/dist-packages (from requests->transformers==4.2.1) (3.0.4)

Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.7/dist-packages (from requests->transformers==4.2.1) (2022.5.18.1)

Requirement already satisfied: idna<3,>=2.5 in /usr/local/lib/python3.7/dist-packages (from requests->transformers==4.2.1) (2.10)

Requirement already satisfied: six in /usr/local/lib/python3.7/dist-packages (from sacremoses->transformers==4.2.1) (1.15.0)

Requirement already satisfied: click in /usr/local/lib/python3.7/dist-packages (from sacremoses->transformers==4.2.1) (7.1.2)

Requirement already satisfied: joblib in /usr/local/lib/python3.7/dist-packages (from sacremoses->transformers==4.2.1) (1.1.0)

Building wheels for collected packages: sacremoses

```
Building wheel for sacremoses (setup.py) ... done
Created wheel for sacremoses: filename=sacremoses-0.0.53-py3-none-any.whl size=895260 sha256=42d50a10e0287b8d208ade80aa5632fb9f71167ace0c692dac0f9b800b5a512c
Stored in directory: /root/.cache/pip/wheels/87/39/dd/a83eeef36d0bf98e7a4d1933a4ad2d660295a40613079bafc9
Successfully built sacremoses
Installing collected packages: tokenizers, sacremoses, transformers
Successfully installed sacremoses-0.0.53 tokenizers-0.9.4 transformers-4.2.1
```

In [3]:

```
import numpy as np
import pandas as pd
import pyarabic.araby as ar

import re , emoji, Stemmer, functools, operator, string
import torch , optuna, gc, random, os

import matplotlib.pyplot as plt
import seaborn as sns

from tqdm import tqdm_notebook as tqdm
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, accuracy_score, f1_score, confusion_matrix, precision_score , recall_score
from transformers import AutoConfig, AutoModelForSequenceClassification, AutoTokenizer
from transformers.data.processors import SingleSentenceClassificationProcessor
from transformers import Trainer , TrainingArguments
from transformers.trainer_utils import EvaluationStrategy
from transformers.data.processors.utils import InputFeatures
from torch.utils.data import Dataset
from torch.utils.data import DataLoader
from sklearn.utils import resample
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix, classification_report, accuracy_score

import gensim
from gensim.models import KeyedVectors

from keras.preprocessing.text import Tokenizer
from keras.preprocessing.sequence import pad_sequences

import tensorflow as tf
from keras.models import Sequential
from tensorflow.keras.layers import SpatialDropout1D, Conv1D, Bidirectional, LSTM, Dense, Input, Dropout, GlobalMaxPooling1D
from keras.layers.embeddings import Embedding
from tensorflow.keras.callbacks import ModelCheckpoint, ReduceLROnPlateau, EarlyStopping
from tensorflow.keras.optimizers import Adam

import itertools
from numpy import loadtxt
from keras.models import load_model

import logging

logging.basicConfig(level=logging.WARNING)
logger = logging.getLogger(__name__)
```

the function of the preprocessing



In [4]:

```
st = Stemmer.Stemmer('arabic')
def data_cleaning (text):
    text = re.sub(r'^https?:\\/\\.*[\\r\\n]*', '', text, flags=re.MULTILINE)
    text = re.sub(r'^http?:\\/\\.*[\\r\\n]*', '', text, flags=re.MULTILINE)
    text = re.sub(r"http\S+", "", text)
    text = re.sub(r"https\S+", "", text)
    text = re.sub(r'\s+', ' ', text)
    text = re.sub("(\s\d+)", "", text)
    text = re.sub(r"$\d+\W+|\b\d+\b|\W+\d+$", "", text)
    text = re.sub("\d+", " ", text)
    text = ar.strip_tashkeel(text)
    text = ar.strip_tatweel(text)
    text = text.replace("#", " ");
    text = text.replace("@", " ");
    text = text.replace("_", " ");
    translator = str.maketrans('', '', string.punctuation)
    text = text.translate(translator)
    em = text
    em_split_emoji = emoji.get_emoji_regexp().split(em)
    em_split_whitespace = [substr.split() for substr in em_split_emoji]
    em_split = functools.reduce(operator.concat, em_split_whitespace)
    text = " ".join(em_split)
    text = re.sub(r'(\.)\1+', r'\1', text)
    text_stem = " ".join([st.stemWord(i) for i in text.split()])
    text = text + " " + text_stem
    text = text.replace("ا", "آ")
    text = text.replace("إ", "أ")
    text = text.replace("ل", "ل")
    text = text.replace("و", "و")
    text = text.replace("ي", "ي")

    return text
```

## Connecting to google drive

In [5]:

```
from google.colab import drive
drive.mount("/content/gdrive")
```

Mounted at /content/gdrive

## Uploading the dataset

In [6]:

```
Tweets_Ids_Col_Train = "Tweet_ID"
Tweets_Text_Col_Train = "Tweet_text"
Tweets_Sentiment_Col_Train = "Sentiment_label"
Train_Data_File = "/content/gdrive/MyDrive/thesis/modified.csv"

train_data = pd.DataFrame()

train_data = pd.read_csv(Train_Data_File, sep='\t')

train_data.head(3)
```

Out[6]:

	#Tweet_ID	Tweet_text	Sentiment_label
0	929241870508724224	مصر الجولة الأخيرة x المباراة القادمة #غانا	Positive
1	928942264583376897	هل هذه هي سياسة خارجيه لدوله تحترم نفسها والآخ	Negative
2	928615163250520065	وزير خارجية فرنسا عن منتدى شباب العالم: شعرت ب	Positive

In [7]:

```
print(train_data[Tweets_Sentiment_Col_Train].value_counts())
```

```
Negative    7840
Neutral     7279
Positive    4643
Mixed       1302
Name: Sentiment_label, dtype: int64
```

printing the fiels with missed values

In [8]:

```
train_data.isnull().sum()
```

Out[8]:

```
#Tweet_ID      0
Tweet_text      0
Sentiment_label 0
dtype: int64
```

printing the number of the duplicated rows

In [9]:

```
print("On a {} doublons dans Data.".format(train_data.duplicated().sum()))
```

On a 68 doublons dans Data.

In [10]:

```
train_data.drop_duplicates(inplace = True)
```

In [11]:

```
print("On a {} doublons dans Data.".format(train_data.duplicated().sum()))
```

On a 0 doublons dans Data.

**checking the types of the fiels in the data**

In [12]:

```
train_data.dtypes
```

Out[12]:

```
#Tweet_ID          int64
Tweet_text          object
Sentiment_label     object
dtype: object
```

**function for printing the pie**

In [13]:

```
def pie(data,col):
    labels = data[col].value_counts().keys().tolist()
    n = len(labels)
    if n==2:
        colors = ['#66b3ff', '#fb3999']
    elif n==3:
        colors = ['#66b3ff', '#fb3999', '#ffcc99']
    elif n==4:
        colors = ['#66b3ff', '#fb3999', '#ffcc99', '#66f3ff']
    elif n==5:
        colors = ['#66b3ff', '#fb3999', '#ffcc99', '#66f3ff', '#adcc99']
    elif n==6:
        colors = ['#66b3ff', '#fb3999', '#ffcc99', '#66f3ff', '#adcc99', '#db7f23']

    fig1, f1 = plt.subplots()
    f1.pie(data[col].value_counts(), labels=labels, colors = colors, autopct='%
1.1f%%',shadow=False, startangle=60)
    f1.axis('equal')
    plt.tight_layout()
    plt.show()

def histo(data,col):
    plt.figure(figsize = (10, 8))
    sns.histplot(data=data, x=col, hue = data[col], fill=True)
```

**Counting the % of each classe**

In [14]:

```
train_data.Sentiment_label.value_counts(normalize = True)
```

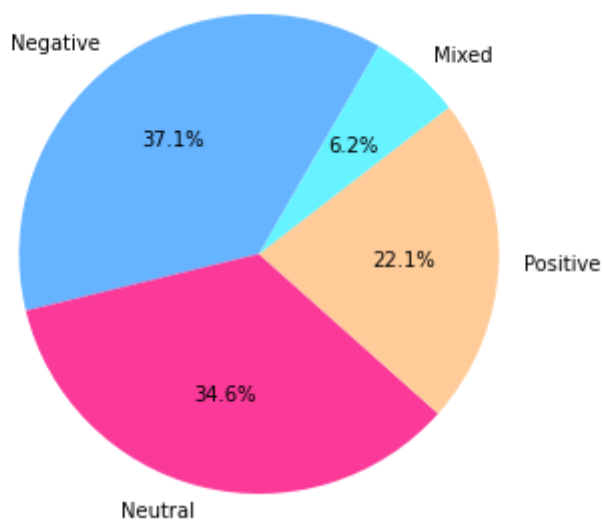
Out[14]:

```
Negative    0.371404  
Neutral     0.346018  
Positive    0.220566  
Mixed       0.062012  
Name: Sentiment_label, dtype: float64
```

**Printing the distribution of the classes**

In [15]:

```
pie(train_data, "Sentiment_label")
```



**preprocessing the reviews and printing the time spent & Deleting unused fields**

In [16]:

```
# Cleaning Training Data
train_data[Tweets_Text_Col_Train] = train_data[Tweets_Text_Col_Train].apply(lambda x: data_cleaning(x))

# Removing un-needed feilds
train_data.drop(['#Tweet_ID'], axis = 1, inplace = True)
train_data.head(3)
```

/usr/local/lib/python3.7/dist-packages/ipykernel\_launcher.py:19: DeprecationWarning: 'emoji.get\_emoji\_regexp()' is deprecated and will be removed in version 2.0.0. If you want to remove emoji from a string, consider the method emoji.replace\_emoji(str, replace=''). To hide this warning, pin/downgrade the package to 'emoji~=1.6.3'

Out[16]:

	Tweet_text	Sentiment_label
0	مصر الجولة الاخيرة من x المباراة القادمة غانا	Positive
1	هل هذه هي سياسة خارجيه لدوله تحترم نفسها والاخ	Negative
2	وزير خارجية فرنسا عن منتدى شباب العالم شعرت با	Positive

### Splitting Data (Train , Evaluation)

In [26]:

```
# First setting the max_len , will be useful later for BERT Model
Extra_Len = 6 # an extra padding in length , found to be useful for increasing F
-score
Max_Len = train_data[Tweets_Text_Col_Train].str.split().str.len().max() + Extra_
Len
print(Max_Len)

#Splitting the Training data
Test_Size = 0.20

Rand_Seed = 42

train_set, evaluation_set = train_test_split( train_data, test_size= Test_Size,
random_state= Rand_Seed)

y=pd.get_dummies(train_data.Sentiment_label)

train_set, X_test, y_train, y_test = train_test_split(train_data, y, test_size =
0.20, random_state = 42)

print("Train set: ")
print(train_set[Tweets_Sentiment_Col_Train].value_counts())
print("-----")
print ("Evaluation set: ")
print (evaluation_set[Tweets_Sentiment_Col_Train].value_counts())
```

```
142
Train set:
Negative      6225
Neutral       5824
Positive      3698
Mixed         1049
Name: Sentiment_label, dtype: int64
-----
Evaluation set:
Negative      1573
Neutral       1441
Positive       933
Mixed         253
Name: Sentiment_label, dtype: int64
```

## Preparing BERTModel Classes

In [18]:

```
Model_Used = "UBC-NLP/MARBERT"
Task_Name = "classification"

class Dataset:
    def __init__(
        self,
        name,
        train,
        test,
        label_list,
    ):
        self.name = name
        self.train = train
        self.test = test
        self.label_list = label_list

class BERTModelDataset(Dataset):
    def __init__(self, text, target, model_name, max_len, label_map):
        super(BERTModelDataset).__init__()
        self.text = text
        self.target = target
        self.tokenizer_name = model_name
        self.tokenizer = AutoTokenizer.from_pretrained(model_name)
        self.max_len = max_len
        self.label_map = label_map

    def __len__(self):
        return len(self.text)

    def __getitem__(self, item):
        text = str(self.text[item])
        text = " ".join(text.split())

        encoded_review = self.tokenizer.encode_plus(
            text,
            max_length=self.max_len,
            add_special_tokens=True,
            return_token_type_ids=False,
            pad_to_max_length=True,
            truncation='longest_first',
            return_attention_mask=True,
            return_tensors='pt'
        )
        input_ids = encoded_review['input_ids'].to(device)
        attention_mask = encoded_review['attention_mask'].to(device)

        return InputFeatures(input_ids=input_ids.flatten(), attention_mask=attention_mask.flatten(), label=self.label_map[self.target[item]])
```

**Defining Needed Methods for training and evaluation**

In [19]:

```
def model_init():
    return AutoModelForSequenceClassification.from_pretrained(Model_Used, return_dict=True, num_labels=len(label_map))

def compute_metrics(p): #p should be of type EvalPrediction
    preds = np.argmax(p.predictions, axis=1)
    assert len(preds) == len(p.label_ids)
    print(classification_report(p.label_ids, preds))
    #print(confusion_matrix(p.label_ids, preds))

    macro_f1_pos_neg = f1_score(p.label_ids, preds, average='macro', labels=[1, 2])
    macro_f1 = f1_score(p.label_ids, preds, average='macro')
    macro_precision = precision_score(p.label_ids, preds, average='macro')
    macro_recall = recall_score(p.label_ids, preds, average='macro')
    acc = accuracy_score(p.label_ids, preds)
    return {
        'macro_f1' : macro_f1,
        'macro_f1_pos_neg' : macro_f1_pos_neg,
        'macro_precision': macro_precision,
        'macro_recall': macro_recall,
        'accuracy': acc
    }

def set_seed(seed):
    torch.manual_seed(seed)
    torch.cuda.manual_seed_all(seed)
    torch.backends.cudnn.deterministic = True
    torch.backends.cudnn.benchmark = False
    np.random.seed(seed)
    random.seed(seed)
    os.environ['PYTHONHASHSEED'] = str(seed)
```

## Build Train and Evaluation Data Sets



In [20]:

```
label_list = list(train_set[Tweets_Sentiment_Col_Train].unique())

print(label_list)
print(train_set[Tweets_Sentiment_Col_Train].value_counts())

data_set = Dataset( "ArSAS", train_set, evaluation_set, label_list )

label_map = { v:index for index, v in enumerate(label_list) }
print(label_map)

train_dataset = BERTModelDataset(train_set[Tweets_Text_Col_Train].to_list(),
                                  train_set[Tweets_Sentiment_Col_Train].to_list(),
                                  Model_Used,Max_Len,label_map)

evaluation_dataset = BERTModelDataset(evaluation_set[Tweets_Text_Col_Train].to_list(),
                                      evaluation_set[Tweets_Sentiment_Col_Train].to_list(),
                                      Model_Used,Max_Len,label_map)

['Neutral', 'Negative', 'Positive', 'Mixed']
Negative      6225
Neutral       5824
Positive      3698
Mixed         1049
Name: Sentiment_label, dtype: int64
{'Neutral': 0, 'Negative': 1, 'Positive': 2, 'Mixed': 3}
```

## Define Training Arguments

In [21]:

```
#define training arguments
training_args = TrainingArguments("./train")
training_args.lr_scheduler_type = 'cosine'
training_args.evaluate_during_training = True
training_args.adam_epsilon = 1e-8
training_args.learning_rate = 1.78255000000000001e-05 # use this with org data
training_args.fp16 = True
training_args.per_device_train_batch_size = 16
training_args.per_device_eval_batch_size = 128
training_args.gradient_accumulation_steps = 2
training_args.num_train_epochs= 2
training_args.warmup_steps = 0
training_args.evaluation_strategy = EvaluationStrategy.EPOCH
training_args.logging_steps = 200
training_args.save_steps = 100000
training_args.seed = 42
training_args.disable_tqdm = False
```

## Build The Trainer

In [22]:

```
training_args.dataloader_pin_memory = False
gc.collect()
torch.cuda.empty_cache()
set_seed(Rand_Seed)

trainer = Trainer(
    model = model_init(),
    args = training_args,
    train_dataset = train_dataset,
    eval_dataset= evaluation_dataset,
    compute_metrics=compute_metrics
)

print(training_args.seed)
```

Some weights of the model checkpoint at UBC-NLP/MARBERT were not used when initializing BertForSequenceClassification: ['cls.prediction.s.bias', 'cls.predictions.transform.dense.weight', 'cls.predictions.transform.dense.bias', 'cls.predictions.transform.LayerNorm.weight', 'cls.predictions.transform.LayerNorm.bias', 'cls.predictions.decoder.weight', 'cls.seq\_relationship.weight', 'cls.seq\_relationship.bias']

- This IS expected if you are initializing BertForSequenceClassification from the checkpoint of a model trained on another task or with another architecture (e.g. initializing a BertForSequenceClassification model from a BertForPreTraining model).

- This IS NOT expected if you are initializing BertForSequenceClassification from the checkpoint of a model that you expect to be exactly identical (initializing a BertForSequenceClassification model from a BertForSequenceClassification model).

Some weights of BertForSequenceClassification were not initialized from the model checkpoint at UBC-NLP/MARBERT and are newly initialized: ['classifier.weight', 'classifier.bias']

You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

42

**Train**

In [23]:

```
all_results = []

print(Max_Len)
print(training_args.learning_rate)
print(training_args.adam_epsilon)
print(training_args.warmup_steps)
trainer.train()

results = trainer.evaluate()
all_results.append(results)
print(results)
```

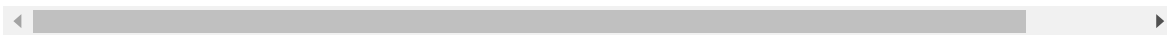
142  
1.78255e-05  
1e-08  
0

/usr/local/lib/python3.7/dist-packages/transformers/tokenization\_utils\_base.py:2143: FutureWarning: The `pad\_to\_max\_length` argument is deprecated and will be removed in a future version, use `padding=True` or `padding='longest'` to pad to the longest sequence in the batch, or use `padding='max\_length'` to pad to a max length. In this case, you can give a specific length with `max\_length` (e.g. `max\_length=45`) or leave max\_length to None to pad to the maximal input size of the model (e.g. 512 for Bert).

FutureWarning,

[1050/1050 05:58, Epoch 2/2]

Epoch	Training Loss	Validation Loss	Macro F1	Macro F1 Pos Neg	Macro Precision	Macro Recall	Accuracy	Runtime
1	0.607700	0.569811	0.613973	0.786606	0.694312	0.625396	0.789286	11.586500
2	0.400900	0.597260	0.653482	0.791730	0.668300	0.650245	0.788810	11.723100



	precision	recall	f1-score	support
0	0.81	0.87	0.84	1441
1	0.80	0.87	0.83	1573
2	0.74	0.74	0.74	933
3	0.43	0.02	0.04	253
accuracy			0.79	4200
macro avg	0.69	0.63	0.61	4200
weighted avg	0.77	0.79	0.77	4200

/usr/local/lib/python3.7/dist-packages/transformers/tokenization\_utils\_base.py:2143: FutureWarning: The `pad\_to\_max\_length` argument is deprecated and will be removed in a future version, use `padding=True` or `padding='longest'` to pad to the longest sequence in the batch, or use `padding='max\_length'` to pad to a max length. In this case, you can give a specific length with `max\_length` (e.g. `max\_length=45`) or leave max\_length to None to pad to the maximal input size of the model (e.g. 512 for Bert).

FutureWarning,

	precision	recall	f1-score	support
0	0.81	0.86	0.84	1441
1	0.83	0.85	0.84	1573
2	0.74	0.74	0.74	933
3	0.28	0.15	0.19	253
accuracy			0.79	4200
macro avg	0.67	0.65	0.65	4200
weighted avg	0.77	0.79	0.78	4200

```
/usr/local/lib/python3.7/dist-packages/transformers/tokenization_utils_base.py:2143: FutureWarning: The `pad_to_max_length` argument is deprecated and will be removed in a future version, use `padding=True` or `padding='longest'` to pad to the longest sequence in the batch, or use `padding='max_length'` to pad to a max length. In this case, you can give a specific length with `max_length` (e.g. `max_length=45`) or leave max_length to None to pad to the maximal input size of the model (e.g. 512 for Bert).
```

```
FutureWarning,
```

```
[33/33 00:11]
```

	precision	recall	f1-score	support
0	0.81	0.86	0.84	1441
1	0.83	0.85	0.84	1573
2	0.74	0.74	0.74	933
3	0.28	0.15	0.19	253
accuracy			0.79	4200
macro avg	0.67	0.65	0.65	4200
weighted avg	0.77	0.79	0.78	4200

```
{'eval_loss': 0.5972602367401123, 'eval_macro_f1': 0.6534821491758084, 'eval_macro_f1_pos_neg': 0.7917304439371178, 'eval_macro_precision': 0.6682995144970398, 'eval_macro_recall': 0.6502452795740927, 'eval_accuracy': 0.7888095238095238, 'eval_runtime': 11.7431, 'eval_samples_per_second': 357.655, 'epoch': 2.0}
```

## Results

```
In [24]:
```

```
all_results
```

```
Out[24]:
```

```
[{'epoch': 2.0, 'eval_accuracy': 0.7888095238095238, 'eval_loss': 0.5972602367401123, 'eval_macro_f1': 0.6534821491758084, 'eval_macro_f1_pos_neg': 0.7917304439371178, 'eval_macro_precision': 0.6682995144970398, 'eval_macro_recall': 0.6502452795740927, 'eval_runtime': 11.7431, 'eval_samples_per_second': 357.655}]
```

```
In [25]:
```

```
from statistics import mean
mean([x['eval_macro_f1'] for x in all_results])
```

```
Out[25]:
```

```
0.6534821491758084
```