Pepito, Alyssa Mae & Soleño, Keziah Antonette                          06-21-2022

Group 7

CS171 | BM1

## Module 1: Summative Assessment (Documentation)

I.      Background and Introduction

Speech Emotion Recognition (SER) is a collection of methodologies that process and detect emotions embedded within speech. SER classifies signals from human speech and uses those signals to identify the underlying emotion within them.  The SER system is especially useful in fields where interactive based-assistant or caller-agent conversation analysis is required. Recently, such a system has been primarily utilized in the medical field whenever a patient is being diagnosed over mobile platforms, customer service as conversations between the customer and attendants are analyzed to improve quality of service, and recommender systems as customers' emotions are used to determine which products are most appropriate to recommend to them. (Market Trends, 2020) (Rockikz, 2019)

II.      Data Preparation

Since we have selected emotion recognition through voice (speech emotion recognition), we downloaded the provided data set found in this link: https://www.kaggle.com/datasets/uldisvalainis/audio-emotions?resource=download. The data set contained files from RAVDESS, CREMA-D, SAVEE, and TESS, and .wav recordings sorted into seven (angry, happy, sad, neutral, fearful, disgusted, and surprised) emotions. (*Audio Emotions*, 2020)

## III.    Python Libraries Installation

Before performing the necessary codes for speech emotion recognition, the Librosa Python library and its dependencies had to be installed by entering the following command in the terminal:



Once all the libraries and their dependencies have been installed, these libraries were then imported by typing the following lines of code:

```
import soundfile # to read audio file
import numpy as np
import librosa # to extract speech features
import glob
import os
import pandas as pd
import os
import csv # Preprocessing
import pickle # to save model after training
from sklearn.model_selection import train_test_split # for splitting training and testing
from sklearn.neural_network import MLPClassifier # multi-layer perceptron model
from sklearn.metrics import accuracy_score # to measure how good we are
```
`[1]   ✓ 11.2s`                                                                    `Python`

## IV.    Speech Feature Extraction

Now that the data set, libraries, and dependencies have been downloaded and installed, feature extraction may finally be performed. To prepare the audio files, the kwargs function had to be created to handle the features being extracted. The following are the features supported by kwargs: MFCC, Chroma, MEL Spectogram Frequency, Contrast, and Tonnetz.

```python
def extract_feature(file_name, **kwargs):

    mfcc = kwargs.get("mfcc")
    chroma = kwargs.get("chroma")
    mel = kwargs.get("mel")
    contrast = kwargs.get("contrast")
    tonnetz = kwargs.get("tonnetz")
    with soundfile.SoundFile(file_name) as sound_file:
        X = sound_file.read(dtype="float32")
        sample_rate = sound_file.samplerate
        if chroma or contrast:
            stft = np.abs(librosa.stft(X))
        result = np.array([])
        if mfcc:
            mfccs = np.mean(librosa.feature.mfcc(y=X, sr=sample_rate,
n_mfcc=40).T, axis=0)
            result = np.hstack((result, mfccs))
        if chroma:
            chroma = np.mean(librosa.feature.chroma_stft(S=stft,
sr=sample_rate).T,axis=0)
            result = np.hstack((result, chroma))
        if mel:
            mel = np.mean(librosa.feature.melspectrogram(X,
sr=sample_rate).T,axis=0)
            result = np.hstack((result, mel))
        if contrast:
            contrast = np.mean(librosa.feature.spectral_contrast(S=stft,
sr=sample_rate).T,axis=0)
            result = np.hstack((result, contrast))
        if tonnetz:
```

```
            tonnetz =
np.mean(librosa.feature.tonnetz(y=librosa.effects.harmonic(X),
    sr=sample_rate).T,axis=0)
                result = np.hstack((result, tonnetz))
        return result
```

After applying these features to the .wav files, the speech waveform of each audio will change to a form of parametric representation at a lesser data rate that can be later converted into a .csv file. (Rockikz, 2019)

```python
# Saving the datset into a CSV file

header = 'filename chroma_stft melspectogram spectral_contrast tonnetz'
for i in range(1, 41):
    header += f' mfcc{i}'
header += ' label'
header = header.split()

file = open('emotiondataset.csv', 'w', newline='')
with file:
    writer = csv.writer(file)
    writer.writerow(header)
emotions = 'Angry Disgusted Fearful Happy Neutral Sad Surprised'.split()

for m in emotions:
    for filename in os.listdir(f'Emotions/{m}'):
        audioname = f'Emotions/{m}/{filename}'
        y, sr = librosa.load(audioname, mono=True, duration=30)
        stft = np.abs(librosa.stft(y))
        chroma_stft = librosa.feature.chroma_stft(y=stft, sr=sr)
        mel_spec = librosa.feature.melspectrogram(y=y, sr=sr)
        spec_cont = librosa.feature.spectral_contrast(y=stft, sr=sr)
        tonnetz = librosa.feature.tonnetz(y=librosa.effects.harmonic(y), sr=sr)
        mfcc = librosa.feature.mfcc(y=y, sr=sr, n_mfcc=40)
        to_append = f'{filename} {np.mean(chroma_stft)} {np.mean(mel_spec)}
{np.mean(spec_cont)} {np.mean(tonnetz)}'
        for e in mfcc:
            to_append += f' {np.mean(e)}'
        to_append += f' {m}'
        file = open('emotiondataset.csv', 'a', newline='')
        with file:
            writer = csv.writer(file)
            writer.writerow(to_append.split())
```

## V. Read Data from the Kaggle Data Set

After declaring a function for feature extraction, we declared two additional functions that will read an audio file by associating the numbers on the .wav files with a corresponding emotion from the data set.

```python
# all emotions on Audio Emotions dataset
int2emotion = {
    "01": "Neutral",
    "02": "Happy",
    "03": "Sad",
    "04": "Angry",
    "05": "Fearful",
    "06": "Disgusted",
    "07": "Surprised"
}


# we allow only these emotions
AVAILABLE_EMOTIONS = {
    "Neutral",
    "Happy",
    "Sad"
}
```

The following code snippet is the main process for extracting features from the data. First, the computer will read the file name and retrieve the emotion label to identify what emotion the .wav

file is expressing. If the retrieved file is not in the list of AVAILABLE_EMOTIONS, the process will continue to the next step of having its speech features extracted. These features will then be added to the data, which will be split into a training and testing set.

```python
def load_data(test_size=0.2):
    X, y = [], []

    for m in emotions:
        for file in glob.glob(f'Emotions/{m}/*.wav'):

            # get the base name of the audio file
            basename = os.path.basename(file)

            # get the emotion label
            emotion = int2emotion[basename.split("-")[2]]

            # we allow only AVAILABLE_EMOTIONS we set
            if emotion not in AVAILABLE_EMOTIONS:
                continue

            # extract speech features
            features = extract_feature(file, mfcc=True, chroma=True, mel=True)

            # add to data
            X.append(features)
            y.append(emotion)

            # split the data to training and testing and return it
            return train_test_split(np.array(X), y, test_size=test_size,
            random_state=7)
```

**VI.** Load and Log Information about the Data Set in 75:25 Proportions

After splitting the data into a training and testing set, the following code was performed to divide them into a 75:25-percent ratio.

```python
# load Audio emotion dataset, 75% training 25% testing
X_train, X_test, y_train, y_test = load_data(test_size=0.25)

# number of samples in training data
print("[+] Number of training samples:", X_train.shape[0])

# number of samples in testing data
print("[+] Number of testing samples:", X_test.shape[0])

# number of features used
```

```python
# this is a vector of features extracted
# using extract_features() function
print("[+] Number of features:", X_train.shape[1])
```

Output:

```
[+] Number of training samples: 23
[+] Number of testing samples: 8
[+] Number of features: 180
```

**VII.** Initialize Model with Given Parameters

The `Multi Layer Perceptron` classifier was initialized using the following parameters.

```python
# best model, determined by a grid search
model_params = {
    'alpha': 0.01,
    'batch_size': 256,
    'epsilon': 1e-08,
    'hidden_layer_sizes': (300,),
    'learning_rate': 'adaptive',
    'max_iter': 500,
}

# initialize Multi Layer Perceptron classifier
# with the parameters declared above
model = MLPClassifier(**model_params)
```

**VIII.** Train Model with Loaded Data Set

After loading the data set, we trained the model by applying the codes below.

```python
# train the model
print("[*] Training the model...")
model.fit(X_train, y_train)
```

Since MLPClassifier was initialized with the given parameters from the previous step, the output below was generated.

```
                         MLPClassifier
MLPClassifier(alpha=0.01, batch_size=256, hidden_layer_sizes=(300,),
              learning_rate='adaptive', max_iter=500)
```
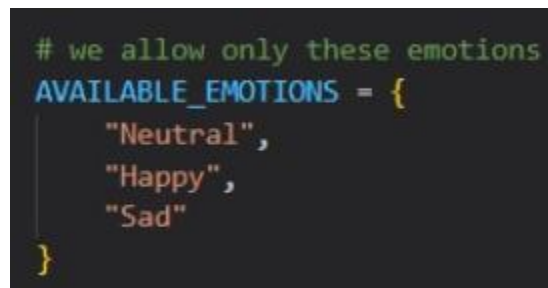
**IX.** Print Accuracy Score

After training our model with our dataset, we can now calculate the accuracy score and print it to measure how good we are.

```python
# predict 25% of data to measure how good we are
y_pred = model.predict(X_test)

# calculate the accuracy
accuracy = accuracy_score(y_true=y_test, y_pred=y_pred)

print("Accuracy: {:.2f}%".format(accuracy*100))
```
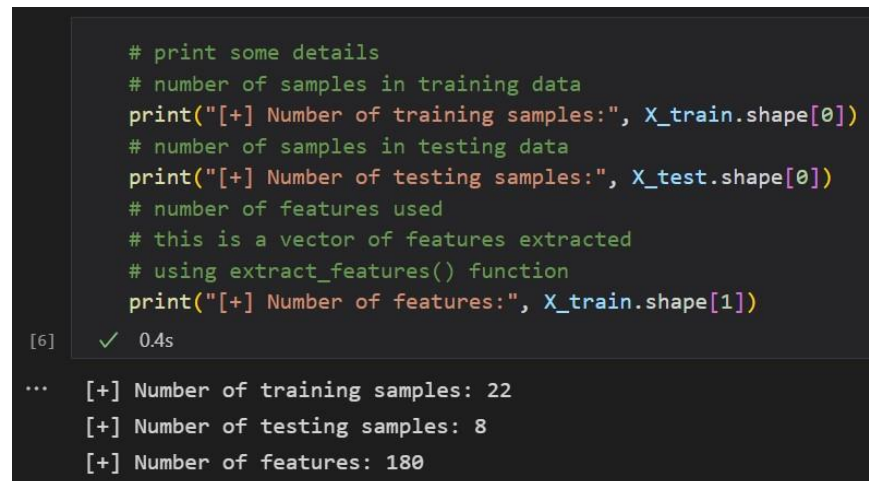
The accuracy score depends on how many emotions we want to classify. For example, let's say we only want to classify neutral, happy, and sad emotions.

```python
# we allow only these emotions
AVAILABLE_EMOTIONS = {
    "Neutral",
    "Happy",
    "Sad"
}
```

Since we have a total of 10 audio files for each emotion, we have 30 audio files all in all. It is then divided into training and testing samples. We notice that we only have a small sample size to work with.

```python
# print some details
# number of samples in training data
print("[+] Number of training samples:", X_train.shape[0])
# number of samples in testing data
print("[+] Number of testing samples:", X_test.shape[0])
# number of features used
# this is a vector of features extracted
# using extract_features() function
print("[+] Number of features:", X_train.shape[1])
```
```
[6]   ✓  0.4s

...   [+] Number of training samples: 22
      [+] Number of testing samples: 8
      [+] Number of features: 180
```

However, we get a 62.50% accuracy for classifying neutral, happy, and sad emotions. This is quite a good score since the accuracy rate is more than 50%.

```
    # predict 25% of data to measure how good we are
    y_pred = model.predict(X_test)

    # calculate the accuracy
    accuracy = accuracy_score(y_true=y_test, y_pred=y_pred)

    print("Accuracy: {:.2f}%".format(accuracy*100))
[10]  ✓  0.1s

...   Accuracy: 62.50%
```

If we classify all emotions, it will indicate an inaccurate model due to the sample size. We are aware that this could cause a problem in dealing with more than two classes.

```
# we allow only these emotions
AVAILABLE_EMOTIONS = {
    "Neutral",
    "Angry",
    "Happy",
    "Sad",
    "Fearful",
    "Disgusted",
    "Surprised"
}
```

As we predicted, it becomes more inaccurate if the model tries to classify more than two.

```
    # predict 25% of data to measure how good we are
    y_pred = model.predict(X_test)

    # calculate the accuracy
    accuracy = accuracy_score(y_true=y_test, y_pred=y_pred)

    print("Accuracy: {:.2f}%".format(accuracy*100))
[19]  ✓  0.6s

...   Accuracy: 33.33%
```

In another situation, if we only classify two emotions, such as Neutral and Angry.

```
# we allow only these emotions
AVAILABLE_EMOTIONS = {
    "Neutral",
    "Angry"
}
```

As there are only two, fewer samples were being used to train the model with.

```
# print some details
# number of samples in training data
print("[+] Number of training samples:", X_train.shape[0])
# number of samples in testing data
print("[+] Number of testing samples:", X_test.shape[0])
# number of features used
# this is a vector of features extracted
# using extract_features() function
print("[+] Number of features:", X_train.shape[1])
✓ 0.1s

[+] Number of training samples: 15
[+] Number of testing samples: 5
[+] Number of features: 180
```

We believe that the model finds it easier to classify two emotions due to the sample size and how it will only differentiate between two emotions.

```
# predict 25% of data to measure how good we are
y_pred = model.predict(X_test)

# calculate the accuracy
accuracy = accuracy_score(y_true=y_test, y_pred=y_pred)

print("Accuracy: {:.2f}%".format(accuracy*100))
✓ 0.6s

Accuracy: 100.00%
```

**X.**    Save the Model
We created a result directory named "result," where we saved our model
`mlp_classifier.model`.

```
if not os.path.isdir("result"):
    os.mkdir("result")


pickle.dump(model, open("result/mlp_classifier.model", "wb"))
```

# References

*Audio emotions*. (2020, June 9). [Dataset]. https://www.kaggle.com/datasets/uldisvalainis/audio-

    emotions?resource=download

M.T. (2020, July 25). *Speech emotion recognition (SER) through machine learning*. Analytics

    Insight. Retrieved June 21, 2022, from https://www.analyticsinsight.net/speech-emotion-

    recognition-ser-through-machine-learning/

Rockikz, A. (2019, July 29). *How to make a speech emotion recognizer using python and scikit-

    learn - python code*. PythonCode. Retrieved June 21, 2022, from

    https://www.thepythoncode.com/article/building-a-speech-emotion-recognizer-using-

    sklearn