

Audio Language Detection

1st Shrutik Pawale - *MS in Data Science*

Roles: Data Pre-processing, EDA , Report

spawale@stevens.edu

2nd Daksh Shah - *ME in Applied Artificial Intelligence*

Roles: Feature Extraction, Model Training, Evaluating Results

dshah94@stevens.edu

May 9, 2023

CONTENTS

1. Introduction	1
2. Dataset and Feature extraction	2
2.1. Data Description	2
2.2. Data Processing	2
2.3. Feature Extraction - MFCC	3
3. Model Implementation	6
4. Results	8
5. Conclusion	8
A. Python Code	A1

LIST OF FIGURES

1.1. Project Outline	1
2.1. Spectrogram of a Sample	2
2.2. DataFrame of Audio Samples Description	3
2.3. DataFrame of Audio Samples Signal and Sample	3
2.4. MFCC Visualization with Spectrogram	4
2.5. MFCC flowchart	5
3.1. CNN Model Structure	7
4.1. Model Evaluation	8

1. INTRODUCTION

Audio language Detection is the method of identifying language from an audio clip by an unidentified speaker, regardless of gender, speaking style, and distinct age speaker. Finding the characteristics that may clearly and effectively differentiate various languages is a challenging task.

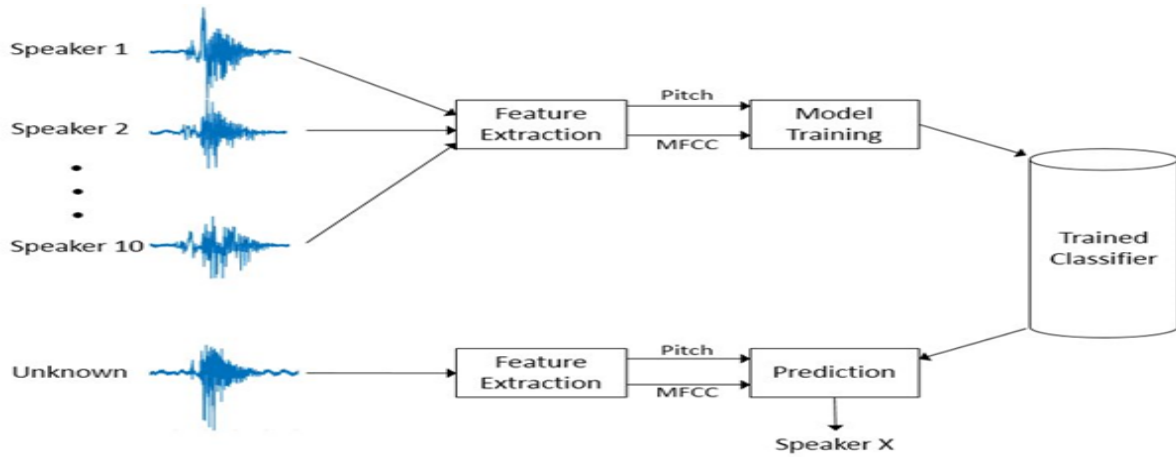


Figure 1.1: Project Outline

The above figure describes the goals of the project. The first step is to read the audio using the Soundfile package. It reads the audio and outputs the signal and sample rate. The next step is to extract features such as MFCC and pitch using the MFCC algorithm. Finally, a convolutional neural network model is trained and the model is evaluated.

2. DATASET AND FEATURE EXTRACTION

2.1. DATA DESCRIPTION

The given dataset contains 10 seconds of speech recorded in English, German, and Spanish languages. Samples are equally balanced between languages, genders, and speakers. The core of the train set is based on 73080 samples after applying several audio transformations (pitch, speed, and noise). No data augmentation has been applied. The number of unique speakers was increased by adjusting pitch (8 different levels) and speed (8 different levels). LibriVox recordings were used to prepare the dataset and particular attention was paid to a big variety of unique speakers since big variance forces the model to concentrate more on language properties than a specific voice.

The spectrogram implies how much the frequency of a signal changes over time. If the frequency matches at that second, then it will show less intensity. That is the reason, why there is a strip of dark line above 4096 Hz. This helps in identifying properties of nonlinear signals and that's why it is helpful in analyzing real-world data with a lot of frequency components and noise.

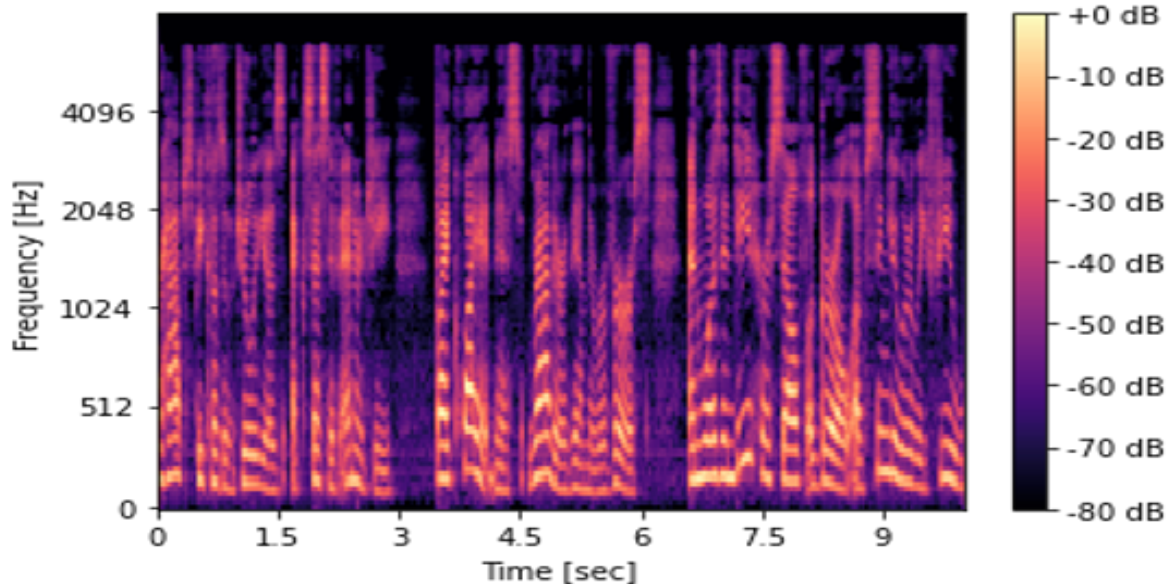


Figure 2.1: Spectrogram of a Sample

2.2. DATA PROCESSING

We can extract the following features from the Audio File Name Format:

```
"(language)_(gender)_(recording_ID).fragment(index)[.(transformation)(index)]"
```

	filename	lang	gender		user_id	fragment	edit
	train\es_f_1d27c6d589eeff17973ffd0b7a77a70a.fr...	es	f	1d27c6d589eeff17973ffd0b7a77a70a	fragment5	speed5	
	train\es_f_53b555eab2b3baada380f7d3ede20b20.fr...	es	f	53b555eab2b3baada380f7d3ede20b20	fragment14	pitch4	
	train\de_f_d94712992f41e3d8d21f22274b3d8fd9.fr...	de	f	d94712992f41e3d8d21f22274b3d8fd9	fragment24	noise6	
	train\en_f_10134f409d9b7b0b95fed6e025febca4.fr...	en	f	10134f409d9b7b0b95fed6e025febca4	fragment25	noise7	
	train\es_m_b8e0e6f56f02e6f8f79cc360958e5982.fr...	es	m	b8e0e6f56f02e6f8f79cc360958e5982	fragment8	noise4	

	train\es_m_d5b91a4ffb1ead826b7968ec19cbfa1c.fr...	es	m	d5b91a4ffb1ead826b7968ec19cbfa1c	fragment3	noise10	
	train\de_m_fc6bd6bb9d66a89bb8d8a8a7efa23e6b.fr...	de	m	fc6bd6bb9d66a89bb8d8a8a7efa23e6b	fragment4	noise6	
	train\de_m_d22535879801cc9c4452d9ed9de5bf61.fr...	de	m	d22535879801cc9c4452d9ed9de5bf61	fragment20	speed4	
	train\de_f_2825fa225d6ca4800f0cf0504b76ca65.fr...	de	f	2825fa225d6ca4800f0cf0504b76ca65	fragment11	pitch6	
	train\es_f_bf4285930fa46f2052e5bdbc37a8a4df.fr...	es	f	bf4285930fa46f2052e5bdbc37a8a4df	fragment21	speed2	

Figure 2.2: DataFrame of Audio Samples Description

After reading the audio file using the sound file python package, we get the sample rate and an array of 220500 numbers. The resulting array contains the audio data as a sequence of samples, where each sample represents the amplitude of the audio signal at a specific point in time.

	0	1	2	3	4	5	6	7	8	9	...	220497	220498	220499	sample_rate
0	-0.020905	-0.031860	-0.028931	-0.020203	-0.002075	0.010193	0.013611	0.004120	-0.009247	-0.012665	...	0.000000	0.000000	0.000000	22050
1	-0.007965	-0.007202	0.003601	0.010895	0.014069	0.005890	-0.005829	-0.019653	-0.028564	-0.035706	...	0.018311	0.018646	0.011688	22050
2	-0.070190	-0.068132	-0.063599	-0.062103	-0.059601	-0.055115	-0.050598	-0.046295	-0.043060	-0.038849	...	0.015778	0.015717	0.017517	22050
3	-0.005951	-0.011993	-0.009888	-0.012848	-0.014374	-0.015961	-0.013062	-0.013824	-0.015961	-0.020050	...	-0.001495	-0.006683	-0.006561	22050
4	0.001556	0.001404	0.001617	0.002106	0.002625	0.002869	0.001709	0.000946	0.001740	0.002380	...	-0.038422	-0.038666	-0.038940	22050
...
495	-0.016266	-0.013611	0.000397	0.020569	0.032318	0.035248	0.035309	0.038361	0.045197	0.046844	...	-0.001801	0.019928	0.016479	22050
496	0.000122	0.001282	0.002045	0.002472	0.002960	0.003479	0.003418	0.002960	0.002747	0.002960	...	0.015900	0.013397	0.016937	22050
497	0.012054	0.015717	0.015717	0.018982	0.018433	0.018677	0.016052	0.018707	0.022522	0.032166	...	-0.009430	-0.018799	-0.018768	22050
498	-0.002380	0.003052	0.003204	-0.004486	-0.004395	0.001587	0.006195	-0.012115	-0.008118	-0.005341	...	0.000000	0.000000	0.000000	22050
499	0.006256	-0.006653	-0.027435	-0.017883	0.003571	0.021820	0.022003	0.009186	-0.001740	-0.008087	...	-0.007996	-0.004730	-0.002106	22050

Figure 2.3: DataFrame of Audio Samples Signal and Sample

2.3. FEATURE EXTRACTION - MFCC

MFCC (Mel-Frequency Cepstral Coefficients) is a technique used in signal processing to analyze and represent the sound of a human voice or other sound signals. The sound is first broken down into many tiny pieces called frames, and then for each frame, the MFCC algorithm measures the power of different frequency bands within that frame. The frequency bands are spaced out in a way that is more like how the human ear perceives sound, which is why it's called Mel-frequency. Next, the algorithm applies some math operations to these frequency band measurements to reduce the dimensionality and capture the most

important features of the sound. The resulting features are called cepstral coefficients and they can be used as inputs to machine learning models for tasks like speech recognition or music genre classification.

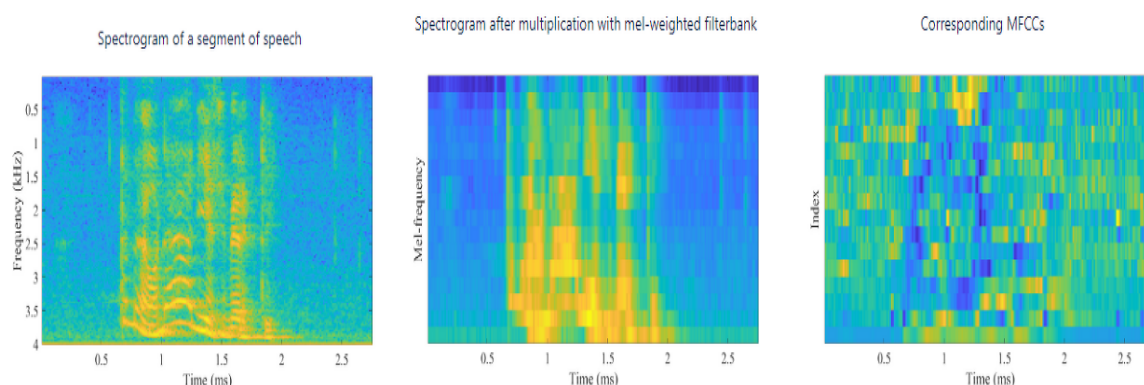


Figure 2.4: MFCC Visualization with Spectrogram

Here is a step-wise illustration of how MFCC works for audio detection:

- Select an Audio Sample and apply Pre-emphasis, which amplifies higher frequencies to balance the audio spectrum.
- Perform Frame blocking and windowing to divide the audio into small frames and apply a window function to each frame.
- Apply Fast Fourier Transformation (FFT) to convert the audio from the time domain to the frequency domain.
- Construct the Mel-Scale Filter Bank, which is a set of triangular filters spaced according to the Mel scale, to capture relevant frequency bands.
- Compute the log and modulus of the filtered signals to enhance the representation of the audio's spectral characteristics.
- Apply Discrete Cosine Transformation (DCT) to the logarithmic filter bank energies to obtain the MFCC coefficients.
- Obtain the MFCC Output, which is a feature vector representing the audio sample's spectral properties.

By following these steps, the MFCC algorithm transforms the audio signal into a compact representation suitable for further analysis and classification in audio detection tasks.

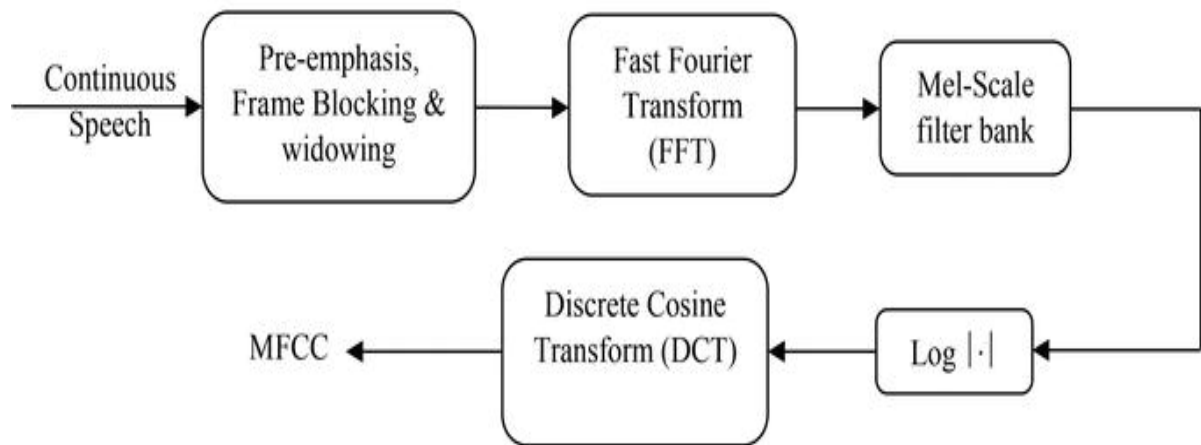


Figure 2.5: MFCC flowchart

3. MODEL IMPLEMENTATION

Since the output from MFCC algorithm is in a 2D array, we need to train it on a compatible model like Convolutional Neural Networks (CNNs). CNNs are particularly effective in capturing spatial dependencies and extracting relevant features from structured data, making them suitable for tasks such as image recognition, object detection, and audio analysis.

Why we used CNN for our audio language detection model:

- **Spatial Relationships:** CNNs excel at capturing spatial relationships in data. In audio language detection, the arrangement of audio frames or spectrogram slices is crucial for accurate classification. CNNs leverage convolutional layers, which apply filters to small regions of the input data, allowing them to detect local patterns and capture spatial dependencies within audio signals.
- **Feature Extraction:** CNNs are capable of automatically learning hierarchical representations of data. In the context of audio language detection, CNNs can automatically extract discriminative features from audio spectrograms or other time-frequency representations. This ability to learn and extract relevant features makes CNNs well-suited for distinguishing language-specific patterns and acoustic cues.
- **Invariance to Local Variations:** CNNs exhibit a degree of translation invariance, meaning they can recognize patterns regardless of their position in the input data. This property is advantageous for audio language detection, as the specific location of language-related features within an audio sample may vary. CNNs can effectively capture invariant features, enabling robust classification even with slight variations in temporal alignment.
- **Parameter Sharing and Efficiency:** CNNs employ parameter sharing, where the same set of weights is used across different regions of the input data. This sharing reduces the number of learnable parameters in the network, making CNNs computationally efficient and less prone to overfitting, especially in cases with limited training data.
- **Previous Success in Audio Analysis:** CNNs have demonstrated remarkable success in various audio-related tasks, such as speech recognition, music classification, and sound event detection. The effectiveness of CNNs in these areas highlights their potential for audio language detection, which also relies on analyzing audio signals for language identification.

Considering these advantages, we selected CNNs as the go-to model for our audio language detection project. The model's structure, as depicted in the figure below, consists of 5 convolutional and max pooling layers, followed by 2 dense layers with batch normalization applied at each step.

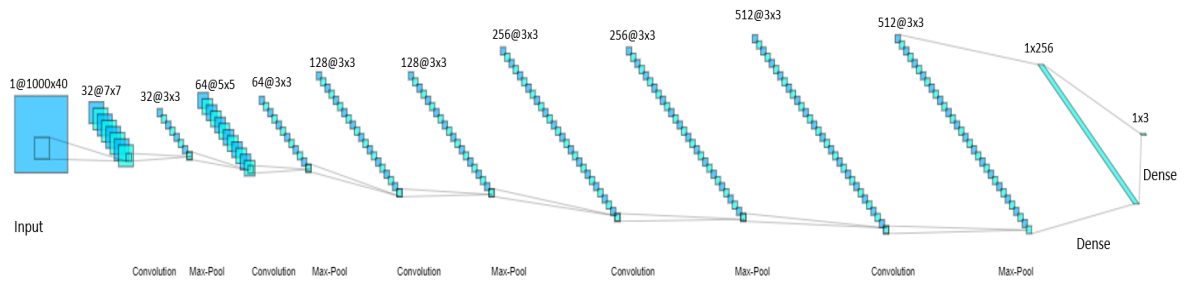


Figure 3.1: CNN Model Structure

The design of the CNN architecture is driven by specific considerations. Firstly, the number of kernels (parameters) doubles in size after each convolutional layer. This progressive incrementation allows deeper layers to capture more complex hidden features such as syllables or even words, while the base layers focus on extracting fundamental features like pitch and frequency. Additionally, the kernel size decreases after each layer, enabling the network to effectively learn and capture finer details in the audio data.

By leveraging this CNN structure, our model can effectively extract and utilize hierarchical representations of audio features, enhancing its ability to classify languages accurately. The combination of convolutional and pooling layers, along with dense layers and batch normalization, enables the model to capture spatial relationships, extract discriminative features, and improve generalization.

4. RESULTS

The trained model utilized 13,200 audio samples for training, achieving an impressive training accuracy of nearly 100%. To evaluate the model's performance, a separate test set consisting of 3,200 samples was used. The evaluation resulted in an accuracy exceeding 99

The effectiveness of the model is further supported by the analysis of the confusion matrix for both the training and test datasets. It reveals minimal mispredictions, with the number of incorrect classifications in the single-digit range and deemed negligible. This indicates the model's ability to accurately recognize and classify audio samples in English, German, and Spanish languages.

Train Set

In [72]:

1

confusion_matrix(y_train1,y_pred1)

Out[72]:

array([[4262, 4, 8],

[0, 4256, 0],

[3, 0, 4267]], dtype=int64)

In [73]:

1

print(classification_report(y_train1,y_pred1))

precision

recall

f1-score

support

de

1.00

1.00

1.00

4274

en

1.00

1.00

1.00

4256

es

1.00

1.00

1.00

4270

accuracy

1.00

12800

macro avg

1.00

1.00

1.00

12800

weighted avg

1.00

1.00

1.00

12800

Test Set

In [75]:

1

confusion_matrix(y_test1,y_pred2)

Out[75]:

array([[1054, 5, 9],

[5, 1055, 4],

[2, 1, 1065]], dtype=int64)

In [76]:

1

print(classification_report(y_test1,y_pred2))

precision

recall

f1-score

support

de

0.99

0.99

0.99

1068

en

0.99

0.99

0.99

1064

es

0.99

1.00

0.99

1068

accuracy

0.99

3200

macro avg

0.99

0.99

0.99

3200

weighted avg

0.99

0.99

0.99

3200

Figure 4.1: Model Evaluation

5. CONCLUSION

Based on these results, we can confidently conclude that the developed model has effectively learned the distinguishing features of different languages, demonstrating its capability to identify and classify audio samples accurately. The model's high training accuracy, coupled with its exceptional performance on the test set, signifies its proficiency in language recognition. This model holds significant potential in various practical applications, such as speech processing, voice assistants, and language identification systems. Its ability to recognize languages with a high degree of accuracy can facilitate multilingual support in various domains, contributing to enhanced user experiences and improved language-based services.

Further research and development may involve expanding the model to handle additional languages or exploring techniques to handle variations in dialects or accents within a specific language. Additionally, fine-tuning the model on larger and more diverse datasets could potentially enhance its performance and robustness. In summary, the developed model has proven its effectiveness in audio language detection, achieving high accuracy and showcasing its proficiency in recognizing English, German, and Spanish languages.

REFERENCES

- [1] Singh, Gundeep and Sharma, Sahil and Chahar, Vijay and Kaur, Manjit and Baz, Mohammed and Masud, Mehedi, *Spoken Language Identification Using Deep Learning*, (2021)
- [2] G. Montavon, *Deep learning for spoken language identification*, (2009)
- [3] P. Kumar, A. Biswas, A. N. Mishra, and M. Chandra, *Spoken language identification using hybrid feature extraction methods*, (2010)

A. PYTHON CODE

```
1 import os
2 import time
3 import glob
4 import pandas as pd
5 import soundfile as sf
6 import scipy.signal as signal
7 import matplotlib.pyplot as plt
8 import gc
9 import IPython.display as ipd
10 import pickle
11 from sklearn.preprocessing import StandardScaler
12 import numpy as np
13 import scipy.signal
14 import warnings
15 warnings.filterwarnings('ignore')
16
17 from sklearn import preprocessing
18 from sklearn.metrics import classification_report, confusion_matrix
19 from sklearn.model_selection import train_test_split
20
21 from keras.models import Model, load_model, Sequential
22 from keras.layers import Conv2D, MaxPooling2D, AveragePooling2D, Dense,
    Flatten
23 from keras.layers import Dropout, Input, Activation
24 from keras.optimizers import Nadam, SGD, Adam
25 from keras.preprocessing.image import ImageDataGenerator
26 from keras.utils import np_utils
27 from keras.callbacks import EarlyStopping, TensorBoard, ModelCheckpoint
28 from keras.models import load_model
29 from keras.layers import BatchNormalization
30 from keras import regularizers
31 import math
32 from keras.callbacks import LearningRateScheduler
33
34 def generate_fb_and_mfcc(signal, sample_rate):
35
36     # Pre-Emphasis
37     pre_emphasis = 0.97
38     emphasized_signal = np.append(
39         signal[0],
40         signal[1:] - pre_emphasis * signal[:-1])
41
42     # Framing
43     frame_size = 0.025
44     frame_stride = 0.01
45
46     # Convert from seconds to samples
47     frame_length, frame_step = (
48         frame_size * sample_rate,
49         frame_stride * sample_rate)
```

```

50 signal_length = len(emphasized_signal)
51 frame_length = int(round(frame_length))
52 frame_step = int(round(frame_step))
53
54 # Make sure that we have at least 1 frame
55 num_frames = int(
56     np.ceil(float(np.abs(signal_length - frame_length)) /
57 frame_step))
58
59 pad_signal_length = num_frames * frame_step + frame_length
60 z = np.zeros((pad_signal_length - signal_length))
61
62 # Pad Signal to make sure that all frames have equal
63 # number of samples without truncating any samples
64 # from the original signal
65 pad_signal = np.append(emphasized_signal, z)
66
67 indices = (
68     np.tile(np.arange(0, frame_length), (num_frames, 1)) +
69     np.tile(
70         np.arange(0, num_frames * frame_step, frame_step),
71         (frame_length, 1)
72     ).T
73 )
74 frames = pad_signal[indices.astype(np.int32, copy=False)]
75
76 # Window
77 frames *= np.hamming(frame_length)
78
79 # Fourier-Transform and Power Spectrum
80 NFFT = 512
81
82 # Magnitude of the FFT
83 mag_frames = np.absolute(np.fft.rfft(frames, NFFT))
84
85 # Power Spectrum
86 pow_frames = ((1.0 / NFFT) * ((mag_frames) ** 2))
87
88 # Filter Banks
89 nfilt = 40
90
91 low_freq_mel = 0
92
93 # Convert Hz to Mel
94 high_freq_mel = (2595 * np.log10(1 + (sample_rate / 2) / 700))
95
96 # Equally spaced in Mel scale
97 mel_points = np.linspace(low_freq_mel, high_freq_mel, nfilt + 2)
98
99 # Convert Mel to Hz
100 hz_points = (700 * (10**(mel_points / 2595) - 1))
101 bin = np.floor((NFFT + 1) * hz_points / sample_rate)

```

```

101
102 fbank = np.zeros((nfilt, int(np.floor(NFFT / 2 + 1))))
103 for m in range(1, nfilt + 1):
104     f_m_minus = int(bin[m - 1])    # left
105     f_m = int(bin[m])              # center
106     f_m_plus = int(bin[m + 1])    # right
107
108     for k in range(f_m_minus, f_m):
109         fbank[m - 1, k] = (k - bin[m - 1]) / (bin[m] - bin[m - 1])
110     for k in range(f_m, f_m_plus):
111         fbank[m - 1, k] = (bin[m + 1] - k) / (bin[m + 1] - bin[m])
112 filter_banks = np.dot(pow_frames, fbank.T)
113
114 # DCT
115 filter_banks = np.where(
116     filter_banks == 0,
117     np.finfo(float).eps,
118     filter_banks)
119
120 # dB
121 filter_banks = 20 * np.log10(filter_banks)
122
123 return filter_banks
124
125 list1=list(pd.DataFrame(os.listdir('train'))[0].apply(lambda x: "train
126 \\"+x).values)
127 df1=pd.DataFrame(list1)
128 df1['lang']=df1[0].apply(lambda x: x.split('_')[0][-2:])
129 df1['gender']=df1[0].apply(lambda x: x.split('_')[1])
130 df1['user_id']=df1[0].apply(lambda x: x.split('_')[-1].split('.')[0])
131 df1['fragment']=df1[0].apply(lambda x: x.split('_')[-1].split('.')[1])
132 df1['edit']=df1[0].apply(lambda x: x.split('_')[-1].split('.')[2])
133
134 german_m = []
135 german_f = []
136 spanish_m = []
137 spanish_f = []
138 english_m = []
139 english_f = []
140 for i in range(len(df1)):
141     if df1['lang'][i]=='de' and df1['gender'][i]=='f':
142         german_f.append(df1[0][i])
143     if df1['lang'][i]=='de' and df1['gender'][i]=='m':
144         german_m.append(df1[0][i])
145     if df1['lang'][i]=='en' and df1['gender'][i]=='f':
146         english_f.append(df1[0][i])
147     if df1['lang'][i]=='en' and df1['gender'][i]=='m':
148         english_m.append(df1[0][i])
149     if df1['lang'][i]=='es' and df1['gender'][i]=='f':
150         spanish_f.append(df1[0][i])
151     if df1['lang'][i]=='es' and df1['gender'][i]=='m':
152         spanish_m.append(df1[0][i])

```

```

152
153 sig_spanish_m, sr_spanish_m = sf.read(file+spanish_m[0])
154 sig_spanish_f, sr_spanish_f = sf.read(file+spanish_f[0])
155 sig_english_m, sr_english_m = sf.read(file+english_m[0])
156 sig_english_f, sr_english_f = sf.read(file+english_f[0])
157 sig_german_m, sr_german_m = sf.read(file+german_m[0])
158 sig_german_f, sr_german_f = sf.read(file+german_f[0])
159
160 (f, S)= scipy.signal.welch(sig_german_m, sr_german_m, nperseg=1024)
161
162 plt.semilogy(f, S)
163 plt.xlabel('frequency [Hz]')
164 plt.ylabel('PSD [V**2/Hz]')
165 plt.show()
166
167 # using padding
168 plt.subplots_adjust(left=0,
169                     bottom=0,
170                     right=2.0,
171                     top=2.0,
172                     wspace=0.4,
173                     hspace=1)
174 #plt.show()
175 plt.subplot(6,2,1)
176 plt.title('Spanish Male')
177 Pxx, freqs, bins, im = plt.specgram(sig_spanish_m, Fs=sr_spanish_m)
178
179 # add axis labels
180 plt.ylabel('Freq')
181 plt.xlabel('Time in samples')
182 plt.subplot(6,2,2)
183 plt.title('Spanish Female')
184 Pxx, freqs, bins, im = plt.specgram(sig_spanish_f, Fs=sr_spanish_f)
185 plt.ylabel('Freq')
186 plt.xlabel('Time in samples')
187
188 plt.subplot(6,2,3)
189 plt.plot(sig_spanish_m)
190 plt.ylabel('Amplitude')
191 plt.xlabel('sample')
192
193 plt.subplot(6,2,4)
194 plt.plot(sig_spanish_f)
195 plt.ylabel('Amplitude')
196 plt.xlabel('sample')
197
198 # add axis labels
199 # plt.ylabel('Freq')
200 # plt.xlabel('Time in samples')
201 plt.subplot(6,2,5)
202 plt.title('English Male')
203 Pxx, freqs, bins, im = plt.specgram(sig_english_m, Fs=sr_english_m)

```

```

204 # add axis labels
205 plt.ylabel('Freq')
206 plt.xlabel('Time in samples')
207 plt.subplot(6,2,6)
208 plt.title('English Female')
209 Pxx, freqs, bins, im = plt.specgram(sig_english_f ,Fs=sr_english_f)
210 plt.ylabel('Freq')
211 plt.xlabel('Time in samples')
212
213 plt.subplot(6,2,7)
214 plt.plot(sig_english_m)
215 plt.ylabel('Amplitude')
216 plt.xlabel('sample')
217
218 plt.subplot(6,2,8)
219 plt.plot(sig_english_f)
220 plt.ylabel('Amplitude')
221 plt.xlabel('sample')
222
223 # add axis labels
224 # plt.ylabel('Freq')
225 # plt.xlabel('Time in samples')
226 plt.subplot(6,2,9)
227 plt.title('German Male')
228 Pxx, freqs, bins, im = plt.specgram(sig_german_m ,Fs=sr_german_m)
229 # add axis labels
230 plt.ylabel('Freq')
231 plt.xlabel('Time in samples')
232 plt.subplot(6,2,10)
233 plt.ylabel('Freq')
234 plt.xlabel('Time in samples')
235 plt.title('German Female')
236 Pxx, freqs, bins, im = plt.specgram(sig_german_f ,Fs=sr_german_f)
237
238 plt.subplot(6,2,11)
239 plt.plot(sig_german_m)
240 plt.ylabel('Amplitude')
241 plt.xlabel('sample')
242
243 plt.subplot(6,2,12)
244 plt.plot(sig_german_f)
245 plt.ylabel('Amplitude')
246 plt.xlabel('sample')
247
248 X_train,X_test,y_train,y_test = train_test_split(df1,df1['lang'],
249         stratify = df1[['lang','gender','user_id','edit']],test_size =
250         0.78,random_state = 420)
251
252 list1=X_train[0].values
253
254 for i in range(round(len(list1)/500)):
255     sigdf=pd.DataFrame()

```



```

254     start = time.time()
255     for j in range(500):
256         signal, sample_rate = sf.read(list1[i*500+j])
257         df=pd.DataFrame(signal.reshape(-1, len(signal)))
258         df['sample_rate']=sample_rate
259         df['filename']=list1[i*500+j]
260         sigdf=pd.concat([sigdf,df])
261         sigdf=sigdf.reset_index(drop=True)
262         sigdf['lang']=sigdf['filename'].apply(lambda x: x.split('_')[0][-2:])
263         sigdf['gender']=sigdf['filename'].apply(lambda x: x.split('_')[1])
264         sigdf['user_id']=sigdf['filename'].apply(lambda x: x.split('_')[-1].split('.')[0])
265         sigdf['fragment']=sigdf['filename'].apply(lambda x: x.split('_')[-1].split('.')[1])
266         sigdf['edit']=sigdf['filename'].apply(lambda x: x.split('_')[-1].split('.')[2])
267         sigdf.drop([220500],axis=1)
268         filehandler = open(r'ex_data\singal_df_'+str(i*500+1)+'_'+str(i*500+j+1)+'.pkl','wb")
269         pickle.dump(sigdf,filehandler)
270         filehandler.close()
271         hours, rem = divmod(time.time()-start, 3600)
272         minutes, seconds = divmod(rem, 60)
273         print('ex_data\singal_df_'+str(i*500+1)+'_'+str(i*500+j+1)+'.csv')
274         print("{:0>2}:{:0>2}:{:05.2f}".format(int(hours),int(minutes),seconds))
275
276 language_dummies=pd.DataFrame()
277 MFCC_array = []
278 sc = StandardScaler()
279 start=time.time()
280 for i in range(round(len(list1)/500)):
281     file = open(r'ex_data\singal_df_'+str(i*500+1)+'_'+str(i*500+500)+'.pkl','rb')
282     sigdf = pickle.load(file)
283     singal_values=np.array(sigdf.iloc[:,220500])
284     language_dummies = pd.concat([language_dummies,pd.get_dummies(sigdf['filename'].apply(lambda x: x.split('_')[0][-2:]))])
285     for i in range(0,len(singal_values)):
286         MFCC = generate_fb_and_mfcc(singal_values[i], sigdf['sample_rate'][i])
287         MFCC_sc = sc.fit_transform(MFCC)
288         MFCC_array.append(MFCC_sc)
289
290 MFCC_array = np.array(MFCC_array)
291 filehandler = open(r'ex_data\MFCC_arrays.pkl','wb")
292 pickle.dump(MFCC_array,filehandler)
293 filehandler.close()
294
295 filehandler = open(r'ex_data\language_dummies.pkl','wb")
296 pickle.dump(language_dummies,filehandler)

```

```

297 filehandler.close()
298
299 hours, rem = divmod(time.time()-start, 3600)
300 minutes, seconds = divmod(rem, 60)
301
302 print("{:0>2}:{:0>2}:{:05.2f}".format(int(hours),int(minutes),seconds))
303
304 start=time.time()
305 input_shape = (1000,40,1)
306
307 model = Sequential()
308 model.add(Conv2D(32,(7, 7), activation='relu', padding='valid',
    input_shape=input_shape))
309 model.add(BatchNormalization())
310 model.add(MaxPooling2D(pool_size=(3,3), strides=2, padding='same'))
311 model.add(Conv2D(64,(5,5), activation='relu', padding='same'))
312 model.add(BatchNormalization())
313 model.add(MaxPooling2D(pool_size=(3,3), strides=2, padding='same'))
314 model.add(Conv2D(128,(3,3), activation='relu', padding='same'))
315 model.add(BatchNormalization())
316 model.add(MaxPooling2D(pool_size=(3,3), strides=2, padding='same'))
317 model.add(Conv2D(256,(3,3), activation='relu', padding='same'))
318 model.add(BatchNormalization())
319 model.add(MaxPooling2D(pool_size=(3,3), strides=2, padding='same'))
320 model.add(Conv2D(512,(3,3), activation='relu', padding='same'))
321 model.add(BatchNormalization())
322 model.add(MaxPooling2D(pool_size=(3,3), strides=2, padding='same'))
323 model.add(Flatten())
324 model.add(BatchNormalization())
325 model.add(Dense(256, activation='relu'))
326 model.add(BatchNormalization())
327 model.add(Dropout(0.5))
328 model.add(Dense(3, activation='softmax'))
329
330
331 adam = Adam()
332 def step_decay(epoch):
333     initial_lrate = 0.00158
334     drop = 0.9
335     epochs_drop = 1
336     lrate = initial_lrate * math.pow(drop, math.floor((1+epoch)/
    epochs_drop))
337     return lrate
338
339
340 model.compile(loss='categorical_crossentropy',optimizer=adam,metrics=['
    accuracy'])
341 checkpoint = ModelCheckpoint(
342     'model.h5',
343     monitor='val_acc',
344     verbose=0,
345     save_best_only=True,

```

```

346         mode='max'
347     )
348
349     lr_scheduler = LearningRateScheduler(step_decay)
350     model.fit(X_train_MFCC,
351             y_train_MFCC,
352             epochs=9,
353             callbacks=[checkpoint, lr_scheduler],
354             verbose=2,
355             validation_data=(X_test_MFCC, y_test_MFCC),
356             batch_size=32)
357
358     hours, rem = divmod(time.time()-start, 3600)
359     minutes, seconds = divmod(rem, 60)
360
361     print("{:0>2}:{:0>2}:{:05.2f}".format(int(hours),int(minutes),seconds))
362
363     filehandler = open(r'ex_data\model.pkl','wb')
364     pickle.dump(model,filehandler)
365     filehandler.close()
366
367     y_pred = model.predict(X_train_MFCC)
368     y_train1 = []
369     label={0:'de',1:'en',2:'es'}
370     for i in range(0,len(y_train_MFCC)):
371         argmax = label[np.argmax(y_train_MFCC.iloc[i,:])]
372         y_train1.append(argmax)
373     y_pred1 = []
374     for i in range(0,len(y_train_MFCC)):
375         argmax = label[np.argmax(y_pred[i,:])]
376         y_pred1.append(argmax)
377
378     print(classification_report(y_train1,y_pred1))
379     print(classification_report(y_train1,y_pred1))
380
381     y_pred_test = model.predict(X_test_MFCC)
382     y_test1 = []
383     label={0:'de',1:'en',2:'es'}
384     for i in range(0,len(y_test_MFCC)):
385         argmax = label[np.argmax(y_test_MFCC.iloc[i,:])]
386         y_test1.append(argmax)
387     y_pred2 = []
388     for i in range(0,len(y_test_MFCC)):
389         argmax = label[np.argmax(y_pred_test[i,:])]
390         y_pred2.append(argmax)
391
392     confusion_matrix(y_test1,y_pred2)
393     print(classification_report(y_test1,y_pred2))

```