

Speech Recognition

Speech Recognition - Reading audio data - Plotting audio data - Transforming audio signals into the frequency domain - Apply Fourier transform signal & plot - Generating audio signals with custom parameters - Generate the time axis - Synthesizing music - Construct audio samples - Amplitude and frequency - Extracting frequency domain features - MFCC & filter bank features - Building HMM - HMM training and prediction - Building a speech recognizer - MFCC features - Case Study.

Introduction

- Speech recognition refers to the process of recognizing & understanding spoken language.
- Input comes in the form of audio data, and the speech recognizer will process this data to extract meaningful information from it.
- This has lot of practical issues such as voice controlled device, transcription of spoken language into words, security systems, and so on.
- Speech signals are very versatile in nature. There are many variations of speech in some language.

→ There are different elements of speech, such as language, emotion, tone, noise, accent and so on.

→ It's difficult to digitally define a set of rules that can constitute speech.

→ famous tasks: automatic speech recognition, has been the focal point of attention for many researchers.

Reading & plotting audio data:
+ Let's see how to read an audio file and visualize the signal.

→ Create a new Python file, import following packages

import numpy as np

import matplotlib.pyplot as plt

from scipy.io import wavfile

2. # Read the .mp3 file

Sampling freq, audio = wavfile.read('input.wav')

3. Parameters of this signal (print)

⇒ print 'In shape:', audio.shape

print 'Data type:', audio.dtype.

(2)

print 'Duration:', sound[0], type(sound[0]) / float(Sampling_freq),
3, 'second'.

4. store audio signal as 16-bit signed integer data.
and normalize the values.

#Normalize the values

$$\text{audio} = \text{audio} / (2^{15})$$

5. Extract the first 30 values to plot, as follows:
Extract the first 30 values for plotting.

$$\text{audio} = \text{audio}[0:30]$$

6. the x-axis is the time axis.

Build axis.

$$x_values = np.arange(0, len(\text{audio}), 1) / float(Sampling_freq)$$

7. convert the units to seconds:

convert to seconds

$$x_values * 1000$$

8. Let's plot this as follows.

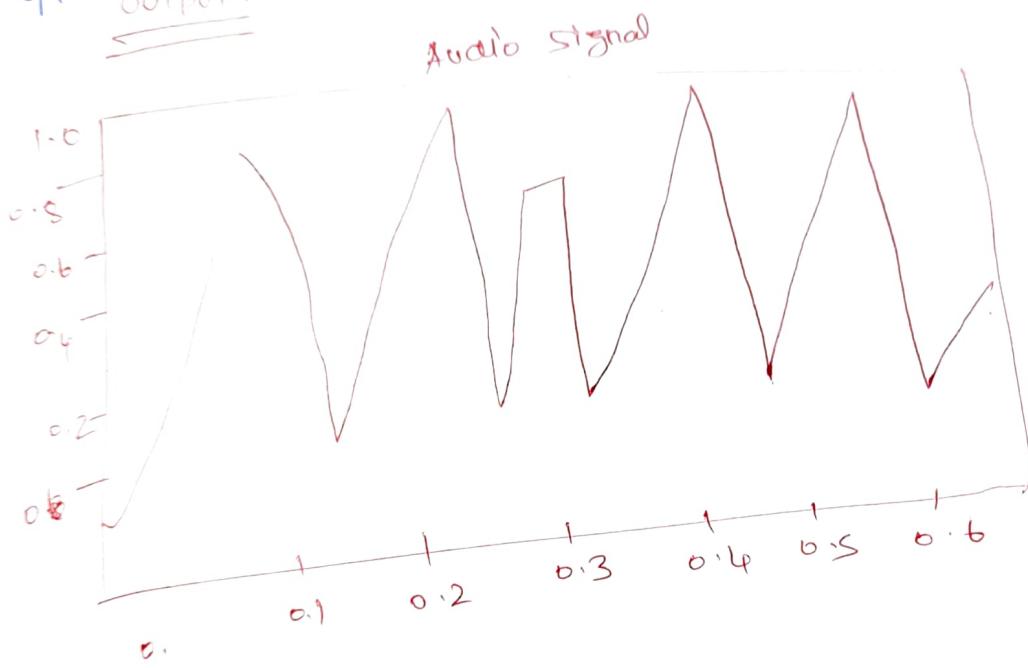
Plotting the chopped audio signal.
plt.plot(x_values, audio, color='black')

```

plt.xlabel('Time (ms)')
plt.ylabel('Amplitude')
plt.title('Audio signal')
plt.show()

```

Q. Output:



Type: (132300,)

datatype: int16.

Duration: 3.0 seconds

② Transforming audio signals into the frequency domain.

- Audio signals consist of a complex mixture of sine waves of different frequencies, amplitudes, and phases.
- Sine waves are referred to as sinusoids.
- So, there is a lot of information that is hidden in

(3)

go the frequency content of an audio signal.

→ the whole word or speech and music is based on this fact. Before that we need some knowledge on Fourier Transforms.

Implementation:

(1)

```
import numpy as np
from scipy.io import wavfile
import matplotlib.pyplot as plt.
```

(2)

Read the input_freq.wav

Read i/p file.

```
sample_freq, audio = wavfile.read('input_freq.wav')
```

(3)

Normalize the signal, as follows

normalize the val

```
audio = audio / (2.**15)
```

(4)

the audio signal is just a numpy array. So, you can extract the length using the following code.

Extract length

```
len_audio = len(audio).
```

⑤ Apply Farley transform.

$$\text{transformed - signal} = \text{np_fft_fft}(\text{audio})$$

$$\text{half_length} = \text{np.ceil}(\text{len_audio}) / 2.0$$

+ abs(transformed_sigr)

half_length = np.ceil(len(audio))
 transformed_signal = abs(fft.rfft(transformed_signal[:half_length]))

transformed - signal = filtered (den. audio?)

transformed - signal $\star\star = 2$

(b) Extract the design of the signal.

TAKE CARE OF even / odd case
LEN (transformer)

Take care of even/odd cases
 $\rightarrow \text{den_S} = \text{den}(\text{transformed_signal})$.
... den audio '1,2'.

if $\overrightarrow{\text{den_audio}}[1, 2]$
transformed_signal $T1: \text{den_ts})^* = 2$

else
transformed - signal $T[1:n-ts-1]$ * = 2.

1

Extract Power 18 Q.S.

Extract power in dB.
- $10 * \text{np.} \log_{10}(\text{transformed_signal})$

$$\text{Power} = 10 * \text{np. Jpg}$$

68

#Build the time axis:

But
 $x_{\text{values}} = \text{np.arange}(0, \text{hal_length}, 1)^* (\text{competing_freq}/$
 $\text{len_audio}) / 1000.0.$

9

plot the signal as follows :-

(1)

plt.figure()

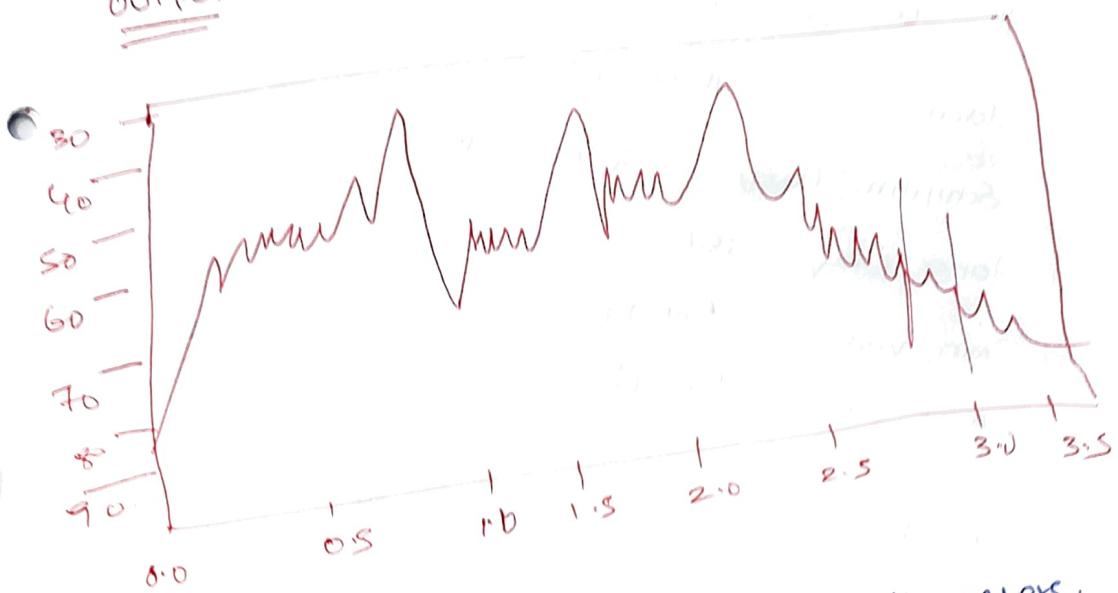
plt.plot(x_normalized, power, color='black')

plt.xlabel('Frequ (in kHz)')

plt.ylabel('Power (in dB)')

plt.show()

OUTPUT:



Generating Audio Signal with custom parameters:

(3)

→ We can use Numpy to generate audio signals. As we discussed earlier, audio signals are complex mixtures of sinusoids.

So keep this in mind when we want to generate our own signal.
~~#~~ required packages

① import numpy as np

import matplotlib.pyplot as plt

from scipy.io.wavfile import write

(2) #File where the output will be saved

output_file = 'output_generated.wav'

(3) Generate a 3-second long signal with a sampling frequency of 44100 and a total frequency of 587 Hz. The values on the time axis will go from -2π to 2π .

Specify audio parameters

duration = 3 # seconds

sampling_freq = 44100 # Hz

tone_freq = 587

min_val = -2π

max_val = 2π

(4) Let's generate the time axis and the audio signal. The audio signal is a simple sinusoid with the previously mentioned parameters:

Generate audio

t = np.linspace(min_val, max_val, duration
 * sampling_freq)

audio = np.sin(2 * np.pi * tone_freq * t)

(5) Add some noise to the signal:

noise = 0.4 + np.random.rand(duration * sampling_freq)

audio += noise

(6) Scale the values to 16-bit integers before we store them:

Scale it to 16-bit integer values

$$\text{scaling_factor} = \text{pow}(2, 15) - 1$$

audio_normalized = audio / np.max(np.abs(audio))

audio_scaled = np.int16(audio_normalized * scaling_factor)

write signal to the op file:

wire (output_file, sampling_freq, audio_scaled)

(7) Plot signals using the first 100 values:-

audio = audio[:100]

Generate the time axis:

Build the time axis

x_values = np.arange(0, len(audio), 1) / float(sampling_freq)

(8) Convert the time axis onto seconds.

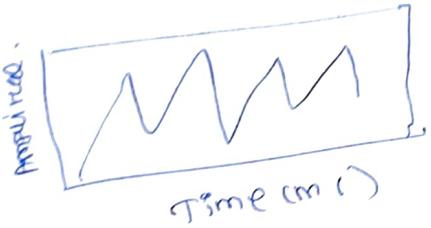
convert to seconds

x_values *= 1000

(9) Plot the signal as follows,

```
plt.plot (y, values, audio, color = 'black')  
plt.xlabel ('time (ms)')  
plt.ylabel ('Amplitude')  
plt.title ('Audio signal')  
plt.show()
```

(2) Output:



(4) Synthesizing music.

→ # Generate simple music.

① → import json
import numpy as np
from scipy.io.wavfile import write
import matplotlib.pyplot as plt

② # Synthesize tone

```
def synthesized (treq, duration, amp = 1.0, sampling_rate = 44100):
```

③ Build the time axis value:

Build the time axis

```
t = np.linspace (0, duration, duration * sampling_rate)
```

6

(d) Define main function:-

import notes:

9 b-name- = :-main -

tone-map-tide = 'tone-treas-map Jason'

(5) Load the file:

Read the frequency map

With open(tone-map-file, "r") as f:

```
with open('tone-map-freq') as f:  
    tone_freq_map = json.loads(f.read())
```

6

Let's assume that we want to generate a G note for a duration of 2 seconds.

set input parameters to generate 'G' tone

$$\text{input_tone} = 161$$

duration = 2 #seconds

$$\text{Amplitude} = 1000$$

$$\text{Sanding - tear} = 44100 \# H2$$

Sompling - today
numbers

P

Call function with `fork()`
Generate the tree

`#Uneren`
`synthesized_tone = synthesized(tone_freq -`
`map_tinput_tone),`
`duration, amplitude, sampling_freq).`

⑧ wire the generated signal into the output tide:

wire the o/p tide.

wire ('output_tone.wav', sampling_freq,
synthesized_tone)

⑨ open this file in a media player and listen to it. that's the G note! let's do something more interesting. let's generate some notes in sequence to give it a musical feel. define a note sequence along with their duration in seconds

Tone-duration sequence

tone_seq = [('D', 0.3), ('G', 0.6), ('C', 0.5), ('A', 0.3),
('Fsharp', 0.7)]

⑩ iterate through the list and call the synthesizer function for each of them:

construct the audio signal based on the chord

def

output = np.array([t])

for item in tone_seq:

input_tone = item[0]

duration = item[1]

synthesized_tone

= synthesizes(tone_freq_map[input_tone],
duration, amplitude, sampling_freq)

output = np.append(output, synthesized_tone, axis=0).

~~7~~

7

wire the signal to o/p file:

Twice to the o/p file.

wire ('output_tone-seq.wav', 'mp3', 'output')

⑤ extracting frequency domain features.

→ After converting a signal into the frequency domain, you need to convert it into a usable form. MFCC is a good way to do this. MFCC takes the power spectrum of a signal & then uses a combination of filter banks, and discrete cosine transformation to extract features.

code:-

import numpy as np

import matplotlib.pyplot as plt

from scipy.io import wavfile

from features import mfcc, logfbnk.

→ Read the input_tone.wav input file that is already provided:

Sampling-trea, audio = wavfile.read('input_tone.wav')

→ # extract MFCC & filter bank features

mfcc_features = mfcc(audio, sampling_trea).

filter_bank_features = logfbnk(audio, sampling_trea.)

→ Check the number of windows generated

print parameters

print 'InMFCC: In Number of windows = ', mfc_features.shape[0]

print 'Length of each feature = ', mfc_features.shape[1]

print 'In Filter bank: In Number of windows = ',

filterbank_features.shape[0]

print 'Length of each feature = ', filterbank_features.shape[1]

→ visualize MFCC features.

transform the matrix so that the time domain is horizontal:

plot the features

mfc_features = mfc_features.T

plt.imshow(mfc_features)

plt.title('MFCC')

→ visualize the filter bank features

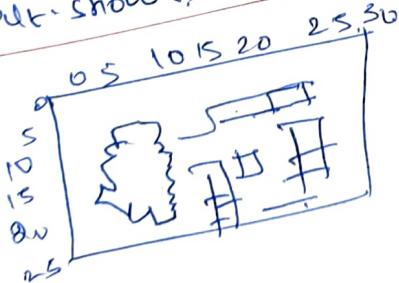
filterbank_features = filterbank_features.T

plt.imshow(filterbank_features)

plt.title('Filter bank!')

plt.show()

→ O/P:



MFCC:

Number of windows = 40

Length of each feature = 13

Filter bank:

Number of windows = 40

Length of each feature = 26.

(5) Building Hidden markov model

→ HMM is something like an NFA. An HMM is a model that represents probability distributions over sequences of observations. We assume that the outputs are generated by hidden states. Goal: Find the hidden states so that we can model the signal.

Code:

- ① A class to handle all HMM related processing.
works HMM training (object).
- ② Initialize the class. Gauss are used to model the data. The n-components parameter defines the number of hidden states. The cov-type defines the type of covariance matrix. In our transition matrix, and (n_iter) indicates the number of iterations it will go through before it stops training.
- ```
def __init__(self, model_name='GaussianHMM')
 n_components=4, cov_type='diag', n_iter=100)
```

## (3) Initialize the variables,

```
self.model_name = model_name
self.n_components = n_components
self.cov_type = cov_type
self.n_iter = n_iter
self.models = t]
```

- ④ Define the model with training parameters  
if self.model\_name == 'GaussianHMM':

```
self.model = hmm.Gaussian(n_components = self.n_comp,
 concentrations, types = self.cov_type,
 n_iters = self.n_iters)
```

```
else:
 raise TypeError('Invalid model type!')
```

5. the input data is a Numpy array, where each element is a feature vector consisting of k-dimensions.

# x is a 2D numpy array where each row is

```
13D def train(self, x):
```

```
 np.seterr(all='ignore').
```

```
 self.models.append(self.model.fit(x)).
```

6. define a method to extract the score, based on the

model.

# run the model on input data

```
def get_score(self, input_data):
```

```
 return self.model.score(input_data).
```

7. we build a class to handle HMM training and prediction, but we need some data to see it in action. we will use it in the next recipe to build a speech recognition.

→ building a speech recognizer:

→ A database containing 7 different words where each word has 15 audio files associated with it. This is a small dataset, but this is sufficient to understand how to build a speech recognizer that can recognize seven different words. We need to build an HMM model for each class.

→ Import os

```
import argparse
```

```
import numpy as np
```

```
from scipy.io import wavfile
```

```
from hmmlearn import hmm
```

```
from feature import mfc
```

(2) Parse the input arguments in the command line:

# Function to parse command line arguments

```
def build_arg_parser():
```

```
 parser = argparse.ArgumentParser
```

```
(description='Trains the HMM models')
```

```
 parser.add_argument('--input-folders',
```

```
 dest='input_folders', required=True,
```

```
 help='Input folder containing the
```

```
 audio files in subfolders')
```

```
 return parser
```

(3) Define the main function, & parse the command line arguments;

```
it_name = '_main_'
args = build_arg_parser()
parse_args()
input_folders = args.input_folders
```

(4) Initiate the variable that will hold all the HMM models: hmm\_models = []

(5) Parse the GIP directory for dirname in os.listdir

(input\_folders):

→

(6) Extract the name of the subfolders:

# Get the name of the

subfolders

```
subfolders = os.path.join
```

(input\_folders, dirname)

```
if not os.path.isdir
```

(subfolders):

continue

⑦ The name of the subfolder  
is the label of this class.  
# Extract the label.

label = subfolders[i].subfolders[  
j].name('!') + ']'

⑧ Initialize the variables for  
training:

# Initialize Variables

x = np.array([])

y\_words = []

⑨ Iterate through the list  
of audiotracks in each subfolder.

# Iterate through the audiotracks

for filename in tx:  
os.listdir(subfolders)

if x.endswith('.wav')]  
x[-1]:

⑩ Read the each audiotrack,  
as follows:

# Read the I/P file.

full path = os.path.join  
(subfolders, filename)

sampling - freq, audio  
= wavfile.read(fullpath)

⑪ Extract the MFCC features.  
mfcc\_features = mfcc(audio)  
sampling - freq

⑫ Keep appending this to the  
X variable:

# Append to the variable X.

if len(x) == 0:

x = mfcc\_features

else:

x = np.append(x, mfcc\_features,  
axis=0).

⑬ Append the corresponding label too.

# Append the label

y\_words.append(label)

⑭ # Train and save HMM model.

hmm\_trainer = HMM.train()

hmm\_trainer.train(x)

hmm\_models.append((hmm\_trainer,  
label))

hmm\_trainers = None.

⑮ Over or file to test.

# Test file

input\_file = 'claro...data'

(16) parse the inputs

# Read input data  
for input\_true in input\_files

(17) Read in each audio file:

# Read input file

sampling\_freq, audio = wavfile.read(input\_file)

(18) Read in each audio file:

sampling\_freq, audio = wavfile.read(input\_file)

(19)

Extract the MFCC features:

# Extract MFCC features

mfcc\_features = mfcc(audio, sampling\_freq)

(20)

# max\_score

max\_score = None

output\_label = None

(21) # iterate through all HMM models and pick  
the one with highest score

for item in hmm\_models:

hmm\_model, label = item

(22)

Extract the score > score max\_score

score = hmm\_model.get\_score(mfcc\_features)

if score > max\_score:

MAX\_score = score

Output\_label = label.

(23)

print true & predicted labels:

# print output

print "In True:",  
input\_bill [input\_bill: find ('/ /') + 1 : input\_bill  
find ('/ /')]

print "predicted": output\_label

(24)

OUTPUT

true: pineapple

predicted: pineapple

true: orange

predicted: orange

true: apple

predicted: apple