

SPEAKER RECOGNITION

Speaker Recognition is the problem of identifying a speaker from a recording of their speech. It is an important topic in Speech Signal Processing and has a variety of applications, especially in security systems. Voice controlled devices also rely heavily on speaker recognition.

This is already a well-researched problem; my aim was not to come up with a new algorithm for speaker recognition, but to implement some already famous existing methods using Python. My motivation behind doing this independent project was to make a shift from MATLAB to Python for scientific computing. For this I primarily used the NumPy, SciPy and Matplotlib packages that have a huge repository for matrix manipulation, signal processing and plotting.

The main principle behind speaker recognition is extraction of features from speech which are characteristic to a speaker, followed by training on a data set and testing. I have relied heavily on the algorithm suggested in [1], where they extract the Mel-Frequency Cepstral Coefficients from each speaker and train them with Vector Quantization (using the LBG algorithm). I have also tried Vector Quantization by extracting the Linear Prediction Coefficients (LPCs) for training. The data set I have trained and tested on was downloaded from [1], and consists of 8 different female speakers uttering the word 'zero'. This data set is not extensive enough to give conclusive results (with only 8 training and test sets), and the results are far from satisfactory. However, my intention was to learn about and implement algorithms in Python, not carry out accurate tests. I wish to gather more data to extend this work.

1. Speech Signal

Digital speech is a one dimensional time-varying discrete signal as shown in **Figure 1a**. Various mathematical models of speech are available such as the Autoregressive Model and Sinusoidal + Residual model. A popular model of speech production says speech consists of a train of impulses of period equal to its pitch, added with random noise, controlled by a voiced/ unvoiced switch and modulated by the vocal tract which is a time-varying filter.

Speech is quasi-stationary in nature. For short intervals, the signal is stationary but over longer periods, the signal frequency varies. Hence a Short-Time Fourier Transform is needed to visualize its frequency content, as given in **Figure 1b** (also computed in Python with FFT size = Hanning window size = 256 samples and overlap of 50%). The 2D plot for STFT with time and frequency along x and y axis, and log amplitude indicated by intensity of the colour is called a Spectrogram.

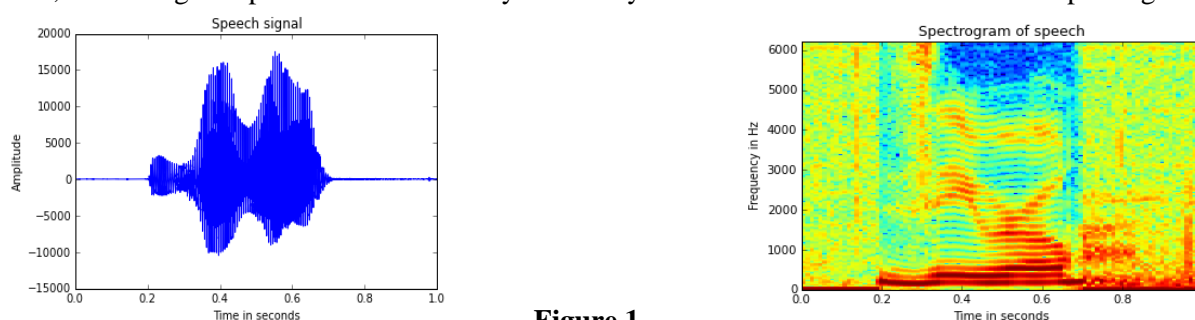


Figure 1

2. Feature Extraction

Choosing which features to extract from speech is the most significant part of speaker recognition. Some popular features are: MFCCs, LPCs, Zero-Crossing Rates etc. In this work, I have concentrated on MFCCs and LPCs. Here is a brief overview of these features.

2.1 Mel-Frequency Cepstral Coefficients

Human hearing is not linear but logarithmic in nature. This implies that our ear acts as a filter. MFCC's are based on the known variation of the human ear's critical bandwidths with frequency. Filters spaced linearly at low frequencies and logarithmically at high frequencies have been used to capture the phonetically important characteristics of speech. This is expressed in the *mel-frequency* scale. The relationship between frequency in Hz and frequency in Mel scale is given by:

$$m = 1125 \ln\left(1 + \frac{f}{700}\right)$$

$$f = 700 \left(e^{\frac{m}{1125}} - 1\right)$$

To calculate MFCCs, the steps are as follows. A very good tutorial is available in [2]. A schematic of this process is given in **Figure 2**.

1. The speech signal is divided into frames of 25ms with an overlap of 10ms. Each frame is multiplied with a Hamming window.
2. The periodogram of each frame of speech is calculated by first doing an FFT of 512 samples on individual frames, then taking the power spectrum as:

$$P(k) = \frac{1}{N} |S(k)|^2$$

Where $P(k)$ refers to power spectral estimate and $S(k)$ refers to Fourier coefficients for the k th frame of speech and N is the length of the analysis window. The last 257 samples of the periodogram are preserved since it is an even function.

3. The entire frequency range is divided into ' n ' Mel filter banks, which is also the number of coefficients we want. 'For ' n ' = 12, the filter bank is shown in **Figure 3** - a number of overlapping triangular filters with increasing bandwidth as the frequency increases.
4. To calculate filter bank energies we multiply each filter bank with the power spectrum, and add up the coefficients. Once this is performed we are left with ' n ' numbers that give us an indication of how much energy was in each filter bank.
5. We take the logarithm of these ' n ' energies and compute its Discrete Cosine Transform to get the final MFCCs.

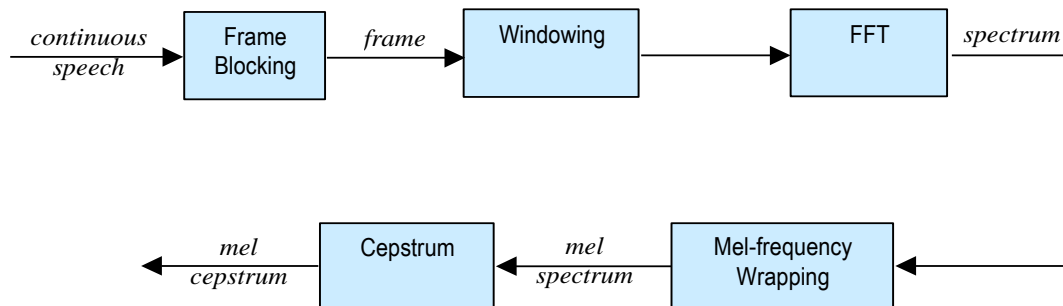


Figure 2 – MFCC Calculation Schematic

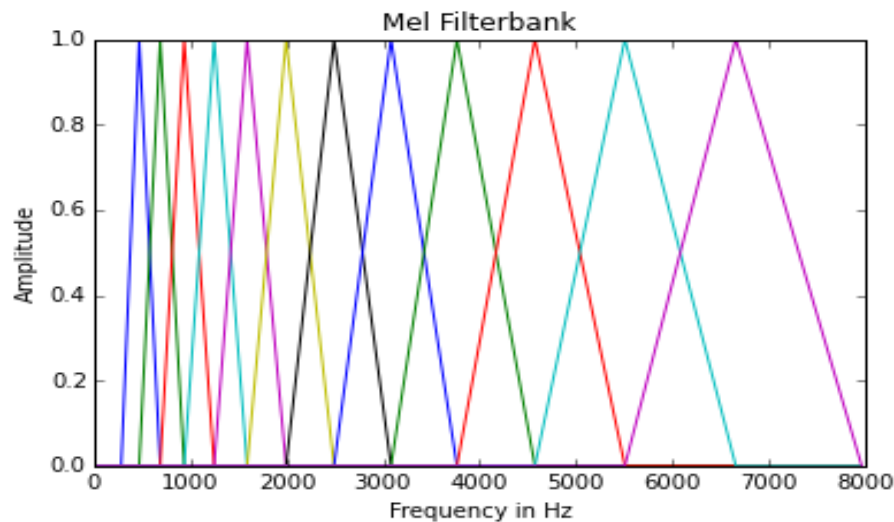


Figure 3 – 12 Mel Filter banks

The Python code for calculating MFCCs from a given speech file (.wav format) is shown in **Listing 1**.

```

1. from __future__ import division
2. from scipy.signal import hamming
3. from scipy.fftpack import fft, fftshift, dct
4. import numpy as np
5. import matplotlib.pyplot as plt
6.
7. def hertz_to_mel(freq):
8.     return 1125*np.log(1 + freq/700)
9.
10. def mel_to_hertz(m):
11.     return 700*(np.exp(m/1125) - 1)
12.
13. #calculate mel frequency filter bank
14. def mel_filterbank(nfft, nfiltbank, fs):
15.
16.     #set limits of mel scale from 300Hz to 8000Hz
17.     lower_mel = hertz_to_mel(300)
18.     upper_mel = hertz_to_mel(8000)
19.     mel = np.linspace(lower_mel, upper_mel, nfiltbank+2)
20.     hertz = [mel_to_hertz(m) for m in mel]
21.     fbins = [int(hz * (nfft/2+1)/fs) for hz in hertz]
22.     fbank = np.empty((nfft/2+1,nfiltbank))
23.     for i in range(1,nfiltbank+1):
24.         for k in range(int(nfft/2 + 1)):
25.             if k < fbins[i-1]:
26.                 fbank[k, i-1] = 0
27.             elif k >= fbins[i-1] and k < fbins[i]:
28.                 fbank[k,i-1] = (k - fbins[i-1])/(fbins[i] - fbins[i-1])
29.             elif k >= fbins[i] and k <= fbins[i+1]:
30.                 fbank[k,i-1] = (fbins[i+1] - k)/(fbins[i+1] - fbins[i])
31.             else:
32.                 fbank[k,i-1] = 0
33.     return fbank
34.
35. def mfcc(s,fs, nfiltbank):
36.
37.     #divide into segments of 25 ms with overlap of 10ms
38.     nSamples = np.int32(0.025*fs)
39.     overlap = np.int32(0.01*fs)
40.     nFrames = np.int32(np.ceil(len(s)/(nSamples-overlap)))
41.

```

```

42.     #zero padding to make signal length long enough to have nFrames
43.     padding = ((nSamples-overlap)*nFrames) - len(s)
44.     if padding > 0:
45.         signal = np.append(s, np.zeros(padding))
46.     else:
47.         signal = s
48.     segment = np.empty((nSamples, nFrames))
49.     start = 0
50.     for i in range(nFrames):
51.         segment[:,i] = signal[start:start+nSamples]
52.         start = (nSamples-overlap)*i
53.
54.     #compute periodogram
55.     nfft = 512
56.     periodogram = np.empty((nFrames,nfft/2 + 1))
57.     for i in range(nFrames):
58.         x = segment[:,i] * hamming(nSamples)
59.         spectrum = fftshift(fft(x,nfft))
60.         periodogram[i,:] = abs(spectrum[nfft/2-1:])/nSamples
61.
62.     #calculating mfccs
63.     fbank = mel_filterbank(nfft, nfiltbank, fs)
64.     #nfiltbank MFCCs for each frame
65.     mel_coeff = np.empty((nfiltbank,nFrames))
66.     for i in range(nfiltbank):
67.         for k in range(nFrames):
68.             mel_coeff[i,k] = np.sum(periodogram[k,:]*fbank[:,i])
69.
70.     mel_coeff = np.log10(mel_coeff)
71.     mel_coeff = dct(mel_coeff)
72.     #exclude 0th order coefficient (much larger than others)
73.     mel_coeff[0,:] = np.zeros(nFrames)
74.     return mel_coeff
75.

```

Listing 1 – *mel_coefficients.py*

2.2 Linear Prediction Coefficients

LPCs are another popular feature for speaker recognition. To understand LPCs, we must first understand the Autoregressive model of speech. Speech can be modelled as a p^{th} order AR process, where each sample is given by:

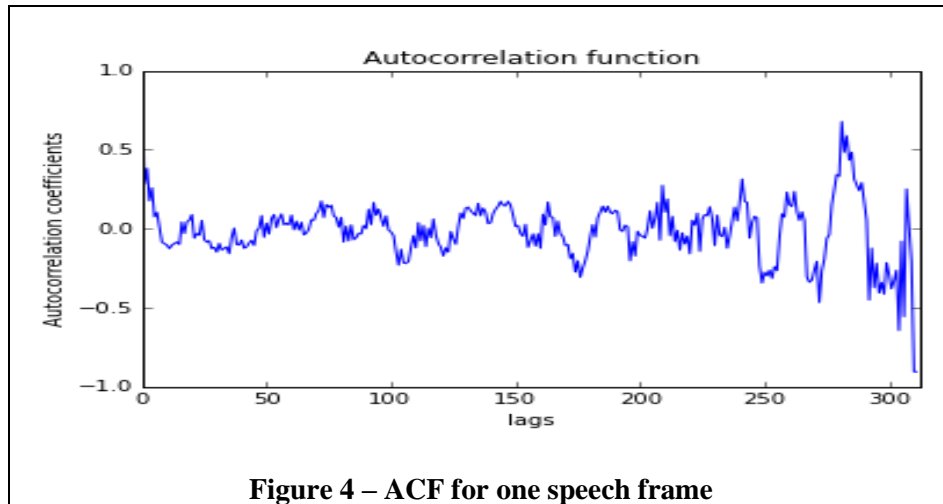
$$x(n) = - \sum_{k=1}^p \alpha_k x(n-k) + u(n)$$

Each sample at the n^{th} instant depends on ' p ' previous samples, added with a Gaussian noise $u(n)$. This model comes from the assumption that a speech signal is produced by a buzzer at the end of a tube (voiced sounds), with occasional added hissing and popping sounds.

LPC coefficients are given by α . To estimate the coefficients, we use the Yule-Walker equations which are explained in [3]. It uses the autocorrelation function R_x . Autocorrelation at lag l is given by:

$$R(l) = \sum_{n=1}^N x(n)x(n-l)$$

While calculating ACF in Python, the Box-Jenkins method is used which scales the correlation at each lag by the sample variance so that the autocorrelation at lag 0 is unity.



The final form of the Yule-Walker equations is:

$$\sum_{k=1}^p \alpha_k R(l-k) = -R(l)$$

$$\begin{bmatrix} R(0) & R(1) & \dots & R(p-1) \\ R(1) & R(0) & \ddots & R(p-2) \\ \vdots & \vdots & \ddots & \vdots \\ R(p-1) & R(p-2) & \dots & R(0) \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_p \end{bmatrix} = - \begin{bmatrix} R(1) \\ R(2) \\ \vdots \\ R(p) \end{bmatrix}$$

The solution for α is given by:

$$\alpha = -R^{-1}r$$

In this case, I have normalised the LPC coefficients estimated so that they lie between $[-1,1]$. This was seen to give more accurate results. The Python code for calculating LPCs is given in **Listing 2**. We first divide speech into frames of 25ms with 10ms overlap, then calculate ' p ' LPCs for each frame.

```

1. #calculate LPC coefficients from sound file
2. from __future__ import division
3. import numpy as np
4.
5. def autocorr(x):
6.     n = len(x)
7.     variance = np.var(x)
8.     x = x - np.mean(x)
9.     #n numbers from last index
10.    r = np.correlate(x, x, mode = 'full')[-n:]
11.    result = r/(variance*(np.arange(n, 0, -1)))
12.    return result
13.
14. def createSymmetricMatrix(acf,p):
15.     R = np.empty((p,p))
16.     for i in range(p):
17.         for j in range(p):
18.             R[i,j] = acf[np.abs(i-j)]

```

```

19.     return R
20.
21. def lpc(s,fs,p):
22.
23.     #divide into segments of 25 ms with overlap of 10ms
24.     nSamples = np.int32(0.025*fs)
25.     overlap = np.int32(0.01*fs)
26.     nFrames = np.int32(np.ceil(len(s)/(nSamples-overlap)))
27.
28.     #zero padding to make signal length long enough to have nFrames
29.     padding = ((nSamples-overlap)*nFrames) - len(s)
30.     if padding > 0:
31.         signal = np.append(s, np.zeros(padding))
32.     else:
33.         signal = s
34.     segment = np.empty((nSamples, nFrames))
35.     start = 0
36.     for i in range(nFrames):
37.         segment[:,i] = signal[start:start+nSamples]
38.         start = (nSamples-overlap)*i
39.
40.     #calculate LPC with Yule-Walker
41.     lpc_coeffs = np.empty((p, nFrames))
42.     for i in range(nFrames):
43.         acf = autocorr(segment[:,i])
44.         r = -acf[1:p+1].T
45.         R = createSymmetricMatrix(acf,p)
46.         lpc_coeffs[:,i] = np.dot(np.linalg.inv(R),r)
47.         lpc_coeffs[:,i] = lpc_coeffs[:,i]/np.max(np.abs(lpc_coeffs[:,i]))
48.
49.     return lpc_coeffs

```

Listing 2 – LPC.py

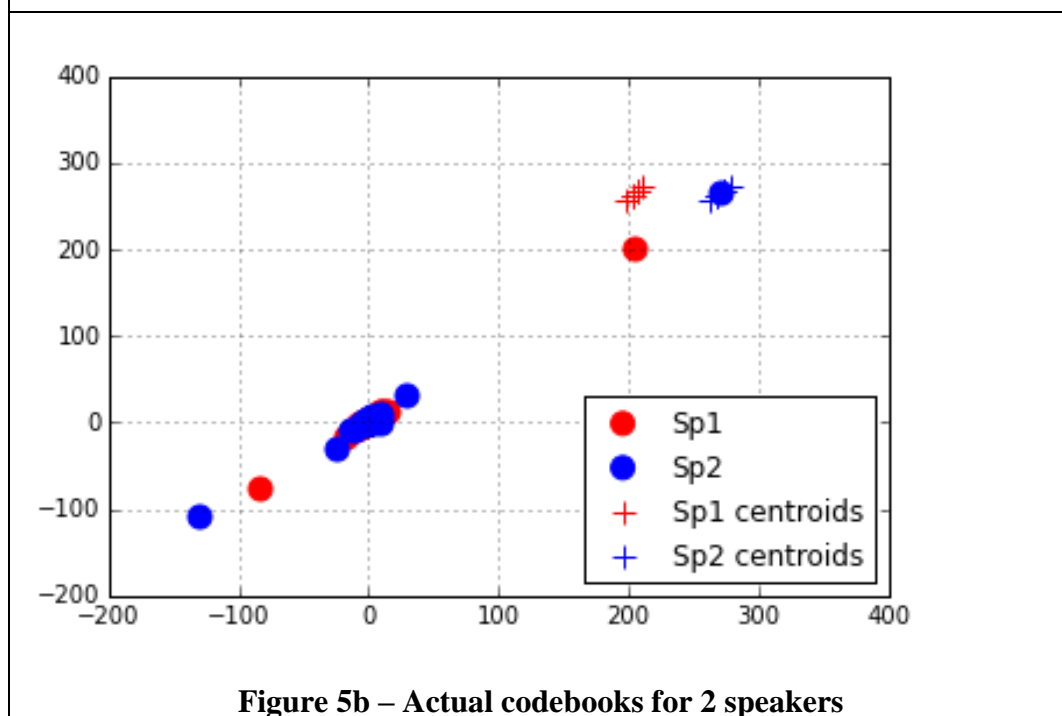
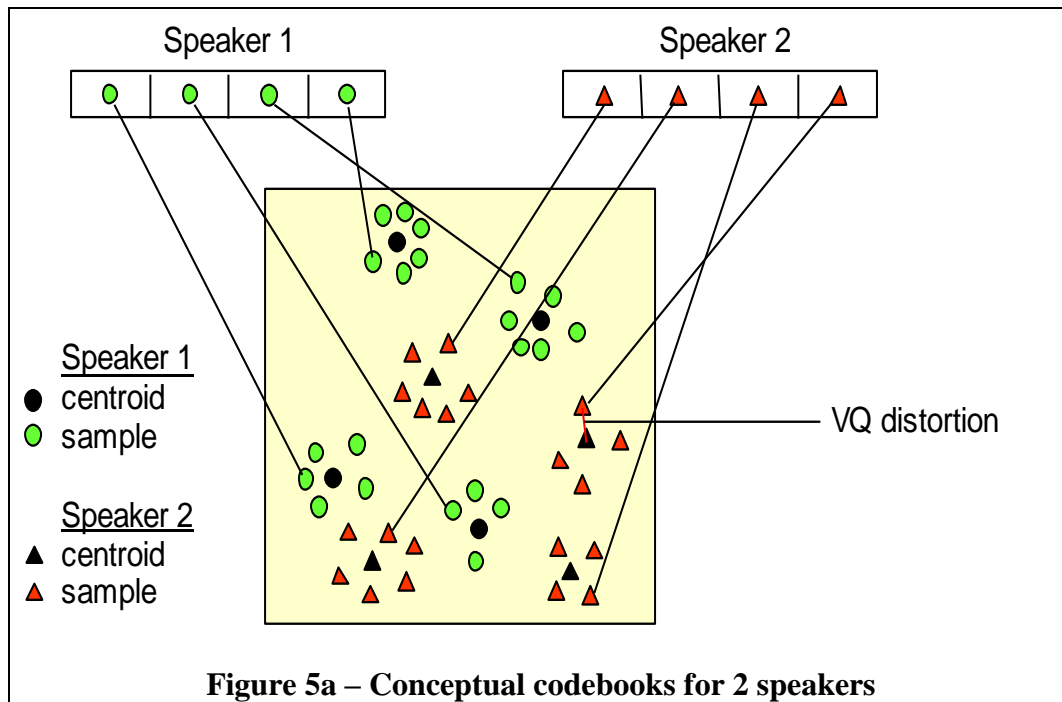
3. Feature Matching

The most popular feature matching algorithms for speaker recognition are Dynamic Time Warping (DTW), Hidden Markov Model (HMM) and Vector Quantization (VQ). Here, I have used Vector Quantization as suggested in [1].

VQ is a process of mapping vectors from a large vector space to a finite number of regions in that space. Each region is called a *cluster* and can be represented by its center called a *codeword*. The collection of all codewords is called a *codebook*.

Figure 5a shows a conceptual diagram to illustrate this recognition process. Only 2 feature dimensions are shown on a 2D plane for two different speakers. We can identify clusters of vectors in the plane. During training, a codeword is chosen for each cluster by minimizing the distortion between each vector in the cluster and the codeword. The collection of codewords forms a speaker specific codebook unique to each speaker. The codebook for each speaker is determined by the **LBG** algorithm, which is described later. To identify a speaker, the distance (or distortion) of the speaker's features from the all the trained codebooks is calculated. The codebook that has minimum distortion with the speaker's features is identified.

Figure 5b shows the actual 2D diagram of 5th and 6th features of two speakers, and their respective codebooks.



3.1 LBG Algorithm

The LBG algorithm [Linde, Buzo and Gray], is used for clustering a set of L training vectors into a set of M codebook vectors. The algorithm is formally implemented by the following recursive procedure:

1. Design a 1-vector codebook; this is the centroid of the entire set of training vectors (hence, no iteration is required here).
2. Double the size of the codebook by splitting each current codebook \mathbf{y}_n according to the rule

$$\mathbf{y}_n^+ = \mathbf{y}_n(1 + \varepsilon)$$

$$\mathbf{y}_n^- = \mathbf{y}_n(1 - \varepsilon)$$

where n varies from 1 to the current size of the codebook, and ε is a splitting parameter (we choose $\varepsilon=0.01$).

3. Nearest-Neighbor Search: for each training vector, find the codeword in the current codebook that is closest (in terms of similarity measurement), and assign that vector to the corresponding cell (associated with the closest codeword).
4. Centroid Update: update the codeword in each cell using the centroid of the training vectors assigned to that cell.
5. Iteration 1: repeat steps 3 and 4 until vector distortion for current iteration falls below a fraction of the previous iteration's distortion. This is to ensure that the process has converged.
6. Iteration 2: repeat steps 2, 3 and 4 until a codebook size of M is designed.

Intuitively, the LBG algorithm designs an M -vector codebook in stages. It starts first by designing a 1-vector codebook, then uses a splitting technique on the codewords to initialize the search for a 2-vector codebook, and continues the splitting process until the desired M -vector codebook is obtained.

Listing 3 gives the details of implementing the LBG algorithm in Python. The print statements should be used for debugging, especially to see if distortion is converging.

```

1. #speaker specific Vector Quantization codebook using LBG algorithm
2. from __future__ import division
3. import numpy as np
4.
5. #calculate Euclidean distance between two matrices
6. def EUDistance(d,c):
7.
8.     # np.shape(d)[0] = np.shape(c)[0]
9.     n = np.shape(d)[1]
10.    p = np.shape(c)[1]
11.    distance = np.empty((n,p))
12.
13.    if n<p:
14.        for i in range(n):
15.            copies = np.transpose(np.tile(d[:,i], (p,1)))
16.            distance[i,:] = np.sum((copies - c)**2,0)
17.    else:
18.        for i in range(p):
19.            copies = np.transpose(np.tile(c[:,i], (n,1)))
20.            distance[:,i] = np.sum((d - copies)**2,0)
21.
22.    distance = np.sqrt(distance)
23.    return distance
24.
25. def lbg(features, M):
26.     eps = 0.01
27.     codebook = np.mean(features, 1)
28.     distortion = 1
29.     nCentroid = 1
30.     while nCentroid < M:
31.
32.         #double the size of codebook
33.         new_codebook = np.empty((len(codebook), nCentroid*2))
34.         if nCentroid == 1:

```



```

35.         new_codebook[:,0] = codebook*(1+eps)
36.         new_codebook[:,1] = codebook*(1-eps)
37.     else:
38.         for i in range(nCentroid):
39.             new_codebook[:,2*i] = codebook[:,i] * (1+eps)
40.             new_codebook[:,2*i+1] = codebook[:,i] * (1-eps)
41.
42.         codebook = new_codebook
43.         nCentroid = np.shape(codebook)[1]
44.         D = EUDistance(features, codebook)
45.
46.         while np.abs(distortion) > eps:
47.             #nearest neighbour search
48.             prev_distance = np.mean(D)
49.             nearest_codebook = np.argmin(D,axis = 1)
50.
51.             #cluster vectors and find new centroid
52.             for i in range(nCentroid):
53.                 #add along 3rd dimension
54.                 codebook[:,i] = np.mean(features[:,np.where(nearest_codebook == i)], 2).T
55.
56.             #replace all NaN values with 0
57.             codebook = np.nan_to_num(codebook)
58.             D = EUDistance(features, codebook)
59.             distortion = (prev_distance - np.mean(D))/prev_distance
60.             #print 'distortion' , distortion
61.
62.         #print 'final codebook', codebook, np.shape(codebook)
63.         return codebook
64.
65.

```

Listing 3 – LBG.py

4. Feature Training

The main algorithms needed for speaker recognition have been implemented. Now, everything needs to be brought together to train our dataset and derive codebooks for each speaker using VQ. The Python code is given in Listing 4. I have hard-coded the name of the directory where the speech files are stored and the .wav filenames, but that can be easily changed by giving them as parameters to the **training()** function. The number of speakers is **nSpeaker = 8**. As mentioned before, speech recordings of 8 female speakers uttering the word ‘zero’ has been taken for training and testing. Each codebook should have 16 codewords, hence **nCentroid = 16** (it is highly recommended to keep this number a power of 2).

Codebooks for both MFCC features and LPC features are plotted for all 8 speakers. One of them is shown in Figure 6. Lines 44 to 61 may be commented out as they are only used to plot the 5th and 6th dimension MFCC features for the first two speakers on a 2D plane.

```

1. from __future__ import division
2. import numpy as np
3. from scipy.io.wavfile import read
4. from LBG import lbg
5. from mel_coefficients import mfcc
6. from LPC import lpc
7. import matplotlib.pyplot as plt
8.
9. def training(nfiltbank, orderLPC):
10.     nSpeaker = 8
11.     nCentroid = 16

```

```

12. codebooks_mfcc = np.empty((nSpeaker,nfiltbank,nCentroid))
13. codebooks_lpc = np.empty((nSpeaker, orderLPC, nCentroid))
14. directory = 'C:/Users/ORCHISAMA/Documents/Audio Signal Processing/Speaker recognition/test'
15. fname = str()
16.
17. for i in range(nSpeaker):
18.     fname = '/s' + str(i+1) + '.wav'
19.     print 'Now speaker ', str(i+1), 'features are being trained'
20.     (fs,s) = read(directory + fname)
21.     mel_coeff = mfcc(s, fs, nfiltbank)
22.     lpc_coeff = lpc(s, fs, orderLPC)
23.     codebooks_mfcc[i,:,:] = lbg(mel_coeff, nCentroid)
24.     codebooks_lpc[i,:,:] = lbg(lpc_coeff, nCentroid)
25.
26.     plt.figure(i)
27.     plt.title('Codebook for speaker ' + str(i+1) + ' with ' + str(nCentroid) + ' centroids')
28.
29.     for j in range(nCentroid):
30.         plt.subplot(211)
31.         plt.stem(codebooks_mfcc[i,:,j])
32.         plt.ylabel('MFCC')
33.         plt.subplot(212)
34.         markerline, stemlines, baseline = plt.stem(codebooks_lpc[i,:,j])
35.         plt.setp(markerline,'markerfacecolor','r')
36.         plt.setp(baseline,'color', 'k')
37.         plt.ylabel('LPC')
38.         plt.axis(ymin = -1, ymax = 1)
39.         plt.xlabel('Number of features')
40.
41. plt.show()
42. print 'Training complete'
43.
44. #plotting 5th and 6th dimension MFCC features on a 2D plane
45. codebooks = np.empty((2, nfiltbank, nCentroid))
46. mel_coeff = np.empty((2, nfiltbank, 68))
47.
48. for i in range(2):
49.     fname = '/s' + str(i+1) + '.wav'
50.     (fs,s) = read(directory + fname)
51.     mel_coeff[i,:,:] = mfcc(s, fs, nfiltbank)[:,:68]
52.     codebooks[i,:,:] = lbg(mel_coeff[i,:,:], nCentroid)
53.
54. plt.figure(nSpeaker + 1)
55. s1 = plt.scatter(mel_coeff[0,4,:], mel_coeff[0,5,:], s = 100, color = 'r', marker = 'o')
56. c1 = plt.scatter(codebooks[0,4,:], codebooks[1,5,:], s = 100, color = 'r', marker = '+')
57. s2 = plt.scatter(mel_coeff[1,4,:], mel_coeff[1,5,:], s = 100, color = 'b', marker = 'o')
58. c2 = plt.scatter(codebooks[1,4,:], codebooks[1,5,:], s = 100, color = 'b', marker = '+')
59. plt.grid()
60. plt.legend((s1, s2, c1, c2), ('Sp1','Sp2','Sp1 centroids', 'Sp2 centroids'),scatterpoints=1,
61. loc = 'lower right')
62. plt.show()
63.
64. return (codebooks_mfcc, codebooks_lpc)
65.

```

Listing 4 – train.py

5. Testing

It is finally time to test our speaker recognition algorithm. Listing 5 contains the code for identifying the speaker by comparing their feature vector to the codebooks of all trained speakers and computing the minimum distance between them. Heads up, the results are not as accurate as I thought they'd be,

yielding an accuracy of **50% with MFCC** and **37.5% with LPC**. Reasons for this low accuracy can be the fact that there wasn't enough data to train on. Other complex classification algorithms such as ANNs and SVMs should yield better results. I observed that training with **12 MFCC features** and **LPC coefficients of order 15** gives the best results. There are other parameters that can be varied, such as number of codewords in a codebook and FFT size while computing MFCCs. It is possible that a different combination of these will give higher accuracy.

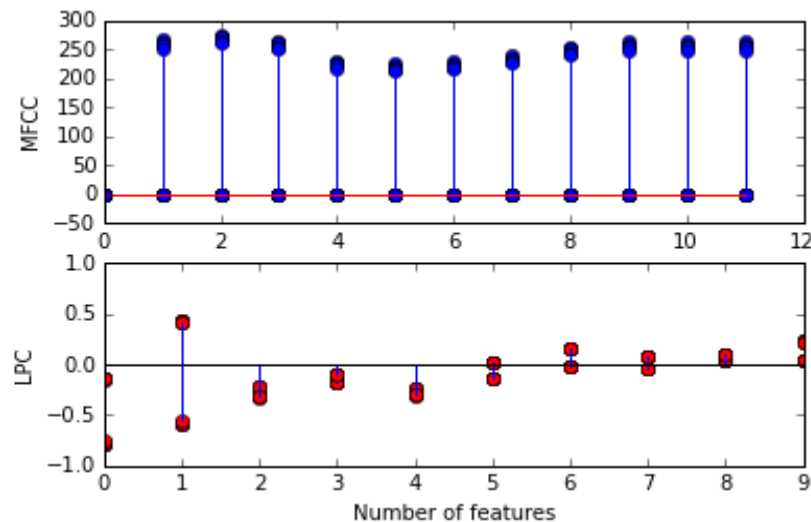


Figure 6 – Codebooks for a) MFCC features and b) LPC features

```

1. from __future__ import division
2. import numpy as np
3. from scipy.io.wavfile import read
4. from LBG import EUDistance
5. from mel_coefficients import mfcc
6. from LPC import lpc
7. from train import training
8.
9. nSpeaker = 8
10. nfilbank = 12
11. orderLPC = 15
12. (codebooks_mfcc, codebooks_lpc) = training(nfilbank, orderLPC)
13. directory = 'C:/Users/ORCHISAMA/Documents/Audio Signal Processing/Speaker recognition/train'
14. fname = str()
15. nCorrect_MFCC = 0
16. nCorrect_LPC = 0
17.
18. def minDistance(features, codebooks):
19.     speaker = 0
20.     distmin = np.inf
21.     for k in range(np.shape(codebooks)[0]):
22.         D = EUDistance(features, codebooks[k, :, :])
23.         dist = np.sum(np.min(D, axis = 1))/(np.shape(D)[0])
24.         if dist < distmin:
25.             distmin = dist
26.             speaker = k
27.     return speaker
28.
29. for i in range(nSpeaker):
30.     fname = '/s' + str(i+1) + '.wav'

```

```

31.     print 'Now speaker ', str(i+1), 'features are being tested'
32.     (fs,s) = read(directory + fname)
33.     mel_coefs = mfcc(s,fs,nfiltbank)
34.     lpc_coefs = lpc(s, fs, orderLPC)
35.     sp_mfcc = minDistance(mel_coefs, codebooks_mfcc)
36.     sp_lpc = minDistance(lpc_coefs, codebooks_lpc)
37.
38.     print 'Speaker', (i+1), ' in test matches with speaker ', (sp_mfcc+1), 'in train for training
39.     with MFCC'
40.     print 'Speaker', (i+1), ' in test matches with speaker ', (sp_lpc+1), 'in train for training
41.     with LPC'
42.
43.     if i == sp_mfcc:
44.         nCorrect_MFCC += 1
45.     if i == sp_lpc:
46.         nCorrect_LPC += 1
47.
48.
49. percentageCorrect_MFCC = (nCorrect_MFCC/nSpeaker)*100
50. print 'Accuracy of result for training with MFCC is ', percentageCorrect_MFCC, '%'
51. percentageCorrect_LPC = (nCorrect_LPC/nSpeaker)*100
52. print 'Accuracy of result for training with LPC is ', percentageCorrect_LPC, '%'
53.
54.

```

Listing 5 – test.py

6. Results

The following table gives the identification results for each of the 8 speakers with MFCC and LPC coefficients and Vector Quantization with LBG algorithm for classification.

Serial Number	True speaker	Identified as (MFCC)	Identified as(LPC)
1.	S1	S1	S3
2.	S2	S6	S2
3.	S3	S6	S3
4.	S4	S8	S7
5.	S5	S3	S1
6.	S6	S6	S7
7.	S7	S7	S7
8.	S8	S8	S7
		Accuracy = 50%	Accuracy = 37.5%

7. Wrapping up

I had a lot of fun doing this small project, understanding algorithms and implementing them in Python. It was a good learning experience, and the resources available online helped a lot. I would like to extend it by experimenting with more features and more classification algorithms. This was not an original research work (I don't want to be accused of plagiarism ☺), but rather an exploration and implementation of existing methods (the code is mine). It might be particularly helpful for those who want to get started with Python for signal processing, or those who want to explore the vast topic of Speech Processing.

Supplement

I would like to include the Python program for plotting the spectrogram of a signal by doing an STFT. A spectrogram carries fundamental information about speech signals, and it is a basic tool that is used in all audio analysis. The default window type has been taken as a Hanning window, and the window length is equal to the number of FFT points.

```

1. #calculate short time fourier transform and plot spectrogram
2. from __future__ import division
3. import matplotlib.pyplot as plt
4. import numpy as np
5. from scipy.fftpack import fft, fftshift
6. from scipy.signal import hann
7.
8. def nearestPow2(inp):
9.     power = np.ceil(np.log2(inp))
10.    return 2**power
11.
12. def stft(signal, fs, nfft, overlap):
13.     #plotting time domain signal
14.     plt.figure(1)
15.     t = np.arange(0,len(signal)/fs, 1/fs)
16.     plt.plot(t,signal)
17.     plt.axis(xmax = 1)
18.     plt.xlabel('Time in seconds')
19.     plt.ylabel('Amplitude')
20.     plt.title('Speech signal')
21.
22.     if not np.log2(nfft).is_integer():
23.         nfft = nearestPow2(nfft)
24.     slength = len(signal)
25.     hop_size = np.int32(overlap * nfft)
26.     nFrames = int(np.round(len(signal)/(nfft-hop_size)))
27.     #zero padding to make signal length long enough to have nFrames
28.     signal = np.append(signal, np.zeros(nfft))
29.     STFT = np.empty((nfft, nFrames))
30.     segment = np.zeros(nfft)
31.     start = 0
32.     for n in range(nFrames):
33.         segment = signal[start:start+nfft] * hann(nfft)
34.         padded_seg = np.append(segment,np.zeros(nfft))
35.         spec = fftshift(fft(padded_seg))
36.         spec = spec[len(spec)/2:]
37.         spec = abs(spec)/max(abs(spec))
38.         powerspec = 20*np.log10(spec)
39.         STFT[:,n] = powerspec
40.         start = start + nfft - hop_size
41.
42.     #plot spectrogram
43.     plt.figure(2)
44.     freq = (fs/(2*nfft)) * np.arange(0,nfft,1)
45.     time = np.arange(0,nFrames)*(slength/(fs*nFrames))
46.     plt.imshow(STFT, extent = [0,max(time),0,max(freq)], origin='lower',
47.         cmap='jet', interpolation='nearest', aspect='auto')
48.     plt.ylabel('Frequency in Hz')
49.     plt.xlabel('Time in seconds')
50.     plt.axis([0,max(time),0,np.max(freq)])
51.     plt.title('Spectrogram of speech')
52.     plt.show()
53.     return (STFT, time, freq)

```

Listing 6 – spectrogram.py