

# Assignment 4: Automatic Speech Recognition

Pranav Sankhe— 150070009

19/11/2018

## 1 Question

Develop an end-pointer using speech/silence detection that enables the automatic segmentation of the individual digit utterances from the continuous audio record. Obtain the pre-emphasised signal corresponding to each utterance.

### 1.1 Answer

Pre-processing of Speech Signal serves various purposes in any speech processing application. It includes Noise Removal, Endpoint Detection, Pre-emphasis, Framing, Windowing, Echo Canceling etc. Out of these, silence/unvoiced portion removal along with endpoint detection is the fundamental step for applications like Speech and Speaker Recognition.

Endpoint detection is used to remove the DC offset value from the signal after silence removal process. Silence removal and Endpoint detection are main part of many applications such as speaker and speech recognition.

Conventional methods using the Zero Crossing Rate (ZCR) and Short Time Energy (STE) have been implemented here. There exist novel methods like the one that uses Probability Density Function (PDF) of the background noise and a Linear Pattern Classifier for classification of Voiced part of a speech from silence/unvoiced part which do indeed perform better than the convention STE method which has been implemented here.

The Short Time Energy (STE) has been calculated and thresholded. We detect speech activity at points where the energy crosses the threshold, and where the energy is below the threshold is labeled as silence. The threshold was decided empirically and the estimation of the threshold by eyeballing was easy since the noise was extremely low.

Later sounds corresponding to different digits were passed through a pre-emphasis filter to take care of lip radiation.

Here is the plot of STE vs time for a particular utterance.

Pre-emphasis filter has been applied on each of the extracted segment. A pre-emphasis filter is useful in several ways:

- balance the frequency spectrum since high frequencies usually have smaller magnitudes compared to lower frequencies,
- avoid numerical problems during the Fourier transform operation and

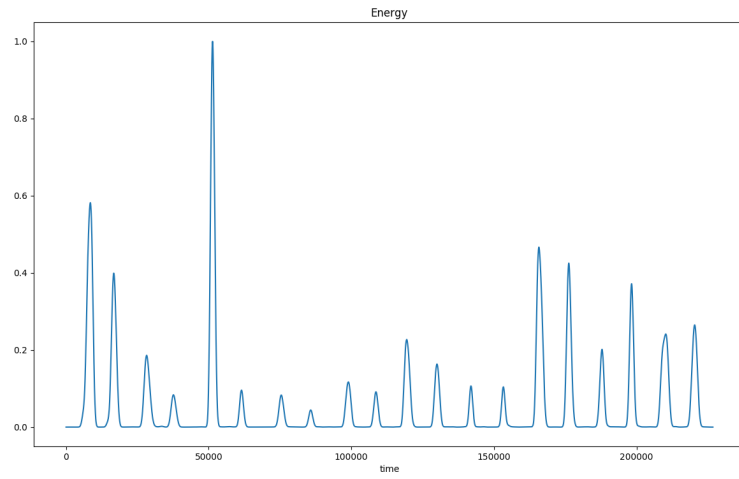
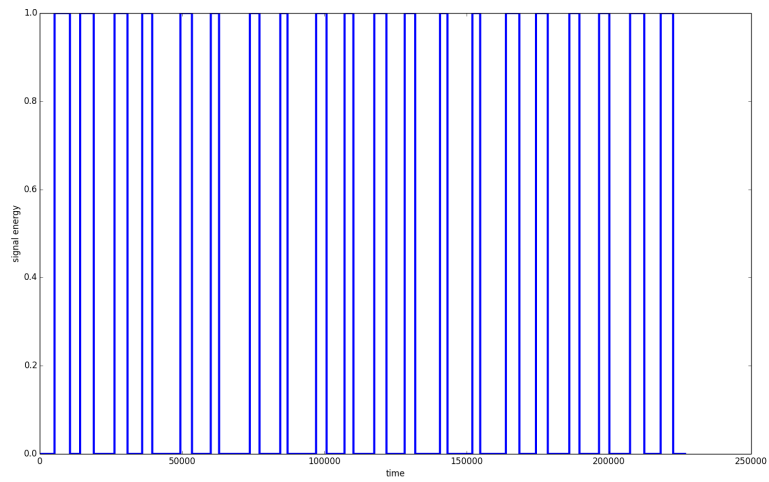


Figure 1: Windowed Energy



S

Figure 2: Endpoints

- may also improve the Signal-to-Noise Ratio (SNR).
- Equation:  $y[t] = x[t] - \alpha x[t - 1]$

Code for segmenting the audio files:

```
1 def pre_emphasis(input_signal):
2     pre_emphasis_alpha = params.pre_emphasis_alpha
3     pre_emphasized_signal = np.append(input_signal[0], input_signal
4     [1:] - pre_emphasis_alpha * input_signal[:-1])
5     return pre_emphasized_signal
6
7 def get_limiting_indices(y):
8     y = y/np.max(y)
9     energy_threshold = params.energy_threshold
10    window_len = 3500
11
12    window = np.hamming(window_len)
13    sig_energy = np.convolve(y**2, window**2, 'same')
14
15    sig_energy = sig_energy/max(sig_energy) #Normalize energy
16    sig_energy_thresh = (sig_energy > energy_threshold).astype('
17    float')
18    import pdb; pdb.set_trace()
19    #convert the bar graph to impulses by subtracting signal from
20    #it's shifted version
21    indices = np.nonzero(abs(sig_energy_thresh[1:] -
22    sig_energy_thresh[0:-1]))[0]
23
24    start_indices = [indices[2*i] for i in range(len(indices)/2)]
25    end_indices = [indices[2*i+1] for i in range(len(indices)/2)]
26
27    return start_indices, end_indices
28
29
30
31 digits = params.digits
32 male_files = np.sort(male_files)
33 male_names = np.sort(male_names)
34 female_files = np.sort(female_files)
35 female_names = np.sort(female_names)
36
37 for i in range(len(male_names)):
38     print('Segmenting audio files of ' + male_names[i])
39     y, sr = librosa.load( male_dir + '/' + male_files[i], sr=None)
40     start_indices, end_indices = get_limiting_indices(y)
41
42     if (not os.path.isdir("male-segmented/" + male_names[i])):
43         os.mkdir("male-segmented/" + male_names[i])
44
45     for p in range(len(end_indices)):
46         sig = y[start_indices[p] : end_indices[p]]
47         digit = pre_emphasis(sig)
48
49         write("male-segmented/" + male_names[i] + '/' + digits[int(
50         np.floor(p/2))] + '-' + str(p%2 + 1) + '.wav', sr, digit)
51
52 for i in range(len(female_names)):
53     print('Segmenting audio files of ' + female_names[i])
54     y, sr = librosa.load( female_dir + '/' + female_files[i], sr=
55     None)
56     start_indices, end_indices = get_limiting_indices(y)
```

```

56     if (not os.path.isdir("female_segmented/" + female_names[i])):
57         os.mkdir("female_segmented/" + female_names[i])
58
59     for p in range(len(end_indices)):
60         sig = y[start_indices[p] : end_indices[p]]
61         digit = pre_emphasis(sig)
62
63         write("female_segmented/" + female_names[i] + '/' + digits[
int(np.floor(p/2))] + '_' + str(p%2 + 1) + '.wav', sr, digit)

```

## 2 Question

Develop a feature extractor that computes an MFCC feature vector for every 10 ms frame of an utterance.

### 2.1 Answer

Mel-Frequency Cepstral Coefficients (MFCCs) were very popular features for a long time for ASR systems. In a nutshell, a signal goes through a pre-emphasis filter; then gets sliced into (overlapping) frames and a window function is applied to each frame; afterwards, we do a Fourier transform on each frame (or more specifically a Short-Time Fourier Transform) and calculate the power spectrum; and subsequently compute the filter banks. To obtain MFCCs, a Discrete Cosine Transform (DCT) is applied to the filter banks retaining a number of the resulting coefficients while the rest are discarded. MFCC features were calculated for every 10ms frame of an utterance for all the digits and all the speakers. MEL filters are uniform in MEL scale which is coherent with the human hearing perception.

Visualize the MFCC series:

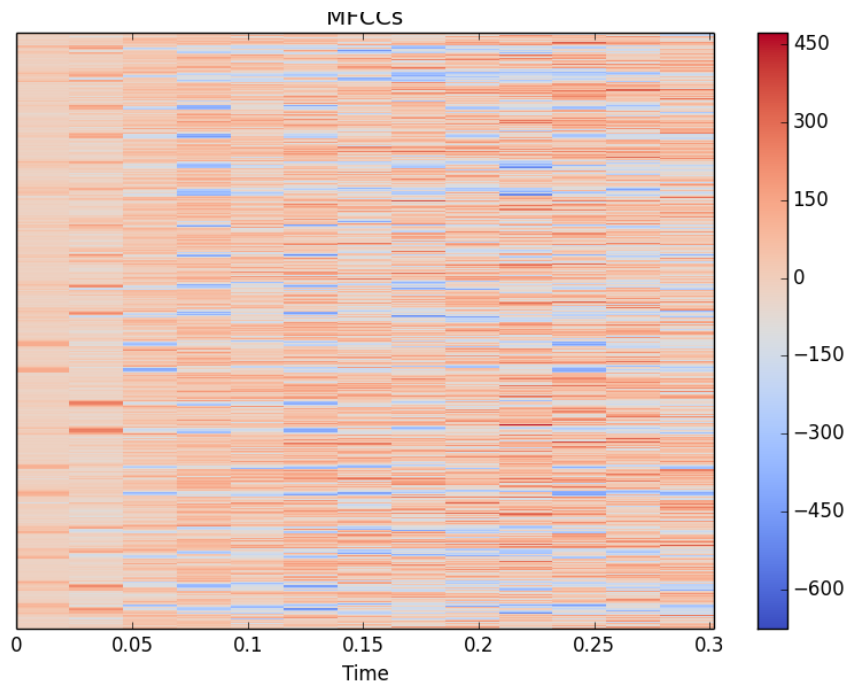


Figure 3: MFCCs

Code for extracting MFCC feature vector:

```
1 import numpy as np
2 import librosa
3 import params
4 import scipy
5 import librosa.display
6 import matplotlib.pyplot as plt
7
8 frame_size = params.frame_size
9 frame_stride = params.frame_stride
10 num_ceps = params.num_ceps
11 def framing(input_signal, sample_rate):
12     frame_length = frame_size * sample_rate
13     frame_step = frame_stride * sample_rate
14     signal_length = len(input_signal)
15     frame_length = int(round(frame_length))
16     frame_step = int(round(frame_step))
17     num_frames = int(np.ceil(float(np.abs(signal_length -
18     frame_length)) / frame_step)) # Make sure that we have at
19     least 1 frame
20
21     pad_signal_length = num_frames * frame_step + frame_length
22     z = np.zeros((pad_signal_length - signal_length))
23     pad_signal = np.append(input_signal, z) # Pad Signal to make
24     sure that all frames have equal number of samples without
25     truncating any samples from the original signal
26
27     indices = np.tile(np.arange(0, frame_length), (num_frames, 1))
28     + np.tile(np.arange(0, num_frames * frame_step, frame_step), (
29     frame_length, 1)).T
30     frames = pad_signal[indices.astype(np.int32, copy=False)]
31     frames = frames * np.hamming(frame_length)
32
33     return frames
34
35 def frame_wise_fft(frames):
36     fft_length = params.fft_length
37     mag_frames = np.abs(np.fft.rfft(frames, fft_length))
38     pow_frames = ((1.0 / fft_length) * ((mag_frames) ** 2)) #
39     frame wise power spectrum
40
41     return pow_frames
42
43 def filter_banks(frames, sample_rate):
44     '''
45     We can convert between Hertz (f) and Mel (m) using the
46     following equation
47     m = 2595 * log10(1 + f * 700)
48     '''
49     fft_length = params.fft_length
50     num_filters = params.num_filters
51
52     low_freq_mel = 0
53     high_freq_mel = (2595 * np.log10(1 + (sample_rate / 2) / 700))
54     # Convert Hz to Mel
55     mel_points = np.linspace(low_freq_mel, high_freq_mel,
56     num_filters + 2) # Equally spaced in Mel scale
57     hz_points = (700 * (10**((mel_points / 2595) - 1))) # Convert
58     Mel to Hz
59     bin = np.floor((fft_length + 1) * hz_points / sample_rate)
```

```

51     fbank = np.zeros((num_filters, int(np.floor(fft_length / 2 + 1)
52                      )))
53     for m in range(1, num_filters + 1):
54         f_m_minus = int(bin[m - 1]) # left
55         f_m = int(bin[m]) # center
56         f_m_plus = int(bin[m + 1]) # right
57
58         for k in range(f_m_minus, f_m):
59             fbank[m - 1, k] = (k - bin[m - 1]) / (bin[m] - bin[m -
60             1])
61         for k in range(f_m, f_m_plus):
62             fbank[m - 1, k] = (bin[m + 1] - k) / (bin[m + 1] - bin[
63             m])
64         filter_banks = np.dot(frames, fbank.T)
65         filter_banks = np.where(filter_banks == 0, np.finfo(float).eps,
66                                filter_banks) # Numerical Stability
67         filter_banks = 20 * np.log10(filter_banks) # dB
68     return filter_banks
69
70 def mfcc(filter_banks, num_ceps):
71     cep_lifter = 22
72     mfcc = scipy.fftpack.dct(filter_banks, type=2, axis=1, norm='
73     ortho')[:, 1 : (num_ceps + 1)] # Keep 2-13
74
75     (nframes, ncoeff) = mfcc.shape
76     n = np.arange(ncoeff)
77     lift = 1 + (cep_lifter / 2) * np.sin(np.pi * n / cep_lifter)
78     mfcc *= lift
79
80     filter_banks -= (np.mean(filter_banks, axis=0) + 1e-8)
81
82     mfcc -= (np.mean(mfcc, axis=0) + 1e-8)
83
84     return mfcc
85
86 def main(filepath):
87     input_signal, sample_rate = librosa.load(filepath, sr=None)
88     frames = framing(input_signal, sample_rate)
89     pow_frames = frame_wise_fft(frames)
90     filter_banks_frames = filter_banks(pow_frames, sample_rate)
91     mfcc_coeffs = mfcc(filter_banks_frames, num_ceps)
92     # librosa.display.specshow(mfcc_coeffs, x_axis='time')
93     # plt.colorbar()
94     # plt.tight_layout()
95     # plt.show()
96
97     return np.transpose(mfcc_coeffs)

```

### 3 Question

Develop a digit recognizer based on the “bag of frames” approach with a codebook for each digit created out of training set speakers’ data. Provide the achieved word error rate (WER) in terms of % words incorrectly detected in the N-fold CV testing using a VQ codebook for each digit obtained via K-means clustering. Provide the achieved WER for with different numbers of clusters (e.g. 4, 8, 16, 64). Observe the common confusions, and comment on your results.

#### 3.1 Answer

The MFCC feature vectors extracted for frames in the train set are used as a bag of frames for the recognition. We run a nearest neighbor classifier for an incoming frame MFCC vector. The output of the classifier for a given word is given by the majority vote of frame classification results.

The general N-fold cross validation technique is used to evaluate performance i.e. model is trained on N-1 speakers and tested on 1 left out speaker, which is iterated over for all N speakers. For every frame in the test pattern, minimum distortions from all reference vectors (codebooks of digits) are calculated and the digit which minimizes this distortion is predicted as an output.

If we include all the vectors of a given digit uttered by all the speakers, the task of predicting the correct digit becomes computationally inefficient. Hence we will employ vector quantization which reduces all the utterances by fewer vectors which we call as ‘representative’ vectors. We use Kmeans to get these ‘representative’ vectors. In the terminology of Kmeans, these representative vectors are the centres/centroids of the clusters.

Here’s a brief overview of k-means algorithm:

- Generate random k points from data as initial estimate of centroids
- Assign each vector to the cluster whose centroid yields the least distance
- Update the centroids of the clusters
- Repeat step 2 and step 3 until convergence

The accuracy is affected by number of clusters (which will be discussed in results section) We get the following WER for different number of clusters considered:

- WER for 4 clusters = 0.29
- WER for 8 clusters = 0.23
- WER for 16 clusters = 0.17
- WER for 64 clusters = 0.128

It is evident that as the number of clusters increases, the WER decreases which is something we can indeed expect since more number of clusters will be able to explain data better and hence less WER. There is, a trade-off between the WER and the number of clusters we can permit. The computational complexity



of the system while testing increases with the number of clusters, but the word error rate decreases. We cannot expect very good WERs from either BOF or VQ directly, because they do not capture phoneme-level information. Two words may have the same phoneme, and a classifier like the one we built will get confused between the two.

The common confusions are represented in a confusion matrix, whose image is shown below:

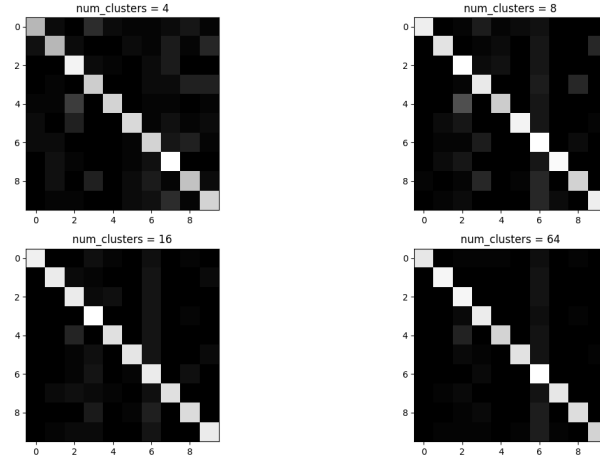


Figure 4: MFCCs

The row index in the confusion matrix indicates the ground truth, that is, the actual identity of the test utterance. The column index indicates the classification output.

1 is added to confusion matrix index  $(x+1, y+1)$  every time digit  $x$  is classified as  $y$ , thus the common classification outputs are brighter than the rest. Each row is scaled such that the maximum in the row is white. This indicates that the most common confusions were

- 3 confused as 8
- 4 confused as 2
- 8 confused as 6

Code for bag of frames approach:

```
1 import os
2 import scipy
3 import numpy as np
4 from scipy import signal
5 import matplotlib.pyplot as plt
6 from scipy import spatial
7 from scipy.cluster.vq import vq, kmeans
8 import params
9 import feature_extractor
10
11 sample_rate = params.sample_rate
12
13 digits = ['zero', 'one', 'two', 'three', 'four', 'five', 'six', 'seven', 'eight', 'nine']
14
15 seg_male_dir = params.seg_male_dir
16 seg_female_dir = params.seg_female_dir
17
18 male_speakers = os.listdir(seg_male_dir)
19 female_speakers = os.listdir(seg_female_dir)
20
21 all_speakers = male_speakers + female_speakers
22 all_speakers.remove('.DS_Store')
23 codeBook = {}
24 n_frame_dict = {}
25
26 for digit in digits:
27     print("Preparing codebook for digit", digit)
28     codeBook[digit] = {}
29     n_frame_dict[digit] = {}
30
31     for speaker in all_speakers:
32         print("speaker = ", speaker)
33         parent_dir = ''
34         if speaker in male_speakers:
35             parent_dir = seg_male_dir
36
37         if speaker in female_speakers:
38             parent_dir = seg_female_dir
39         codeBook[digit][speaker] = []
40         n_frame_dict[digit][speaker] = []
41         for iteration in range(1,3):
42             file = parent_dir + '/' + speaker + '/' + str(digit) +
43             '_' + str(iteration) + '.wav'
44             feature_mat = feature_extractor.main(file)
45             n_frame_dict[digit][speaker].append(feature_mat.shape
46             [1])
47             for i in range(feature_mat.shape[1]):
48                 codeBook[digit][speaker].append(feature_mat[:, i])
49
50     print('Codebook created')
51     confusion_matrix = np.zeros([10,10])
52
53     VQCodeBook = {}
54     centroids = {}
55     n_clusters = params.n_clusters
56
57     for test_speaker in all_speakers:
58         train_speakers = all_speakers
```

```

59 train_speakers.remove(test_speaker)
60
61 for digit in digits:
62     VQCodeBook[digit] = []
63     centroids[digit] = []
64     for speaker in train_speakers:
65         mat = np.asarray(codeBook[digit][speaker])
66         if (VQCodeBook[digit] == []):
67             VQCodeBook[digit] = mat
68         else:
69             VQCodeBook[digit] = np.concatenate([VQCodeBook[
digit],mat],0)
70
71     centroids[digit] = (kmeans(VQCodeBook[digit],n_clusters))
[0]
72
73 for test_digit in digits:
74     for utterance in range(1,3):
75
76         n_frames = n_frame_dict[test_digit][test_speaker]
77         test_mat = codeBook[test_digit][test_speaker][sum(
n_frames[0:(utterance-1)]):sum(n_frames[0:utterance])]
78
79         sum_dist = np.zeros(10)
80         for digit in digits:
81             for l in range(len(test_mat)):
82                 test_vec = np.asarray(test_mat)[l,:]
83                 temp_list = np.asarray(centroids[digit])
84                 min_dist, index = spatial.KDTree(temp_list).
query(test_vec)
85                 sum_dist[digits.index(digit)]+=min_dist
86
87                 pred_digit = np.argmin(sum_dist)
88                 print( "For ", test_speaker, " predicted digit = ",
pred_digit, " ground truth = ", test_digit)
89                 confusion_matrix[digits.index(test_digit),pred_digit]
+= 1.0
90
91
92 np.save('VQ_confusion_matrix_n_cluster'+str(n_clusters),
confusion_matrix)
93
94
95 wer = 1 - np.trace(confusion_matrix)/(len(all_speakers)*20)
96 print wer

```

## 4 Question

Develop a template-matching digit recognizer based on DTW alignment and distance computation. Provide the achieved WER in N-fold CV evaluation. Observe the common confusions, and comment on your results.

### 4.1 Answer

In this approach, we align frames after doing a non-linear time warp. This takes into account some kind of sequence information about the utterances, and hence is expected to perform better than the previous methods and indeed it does. With this approach we get a WER of 6.25%. The common confusions are represented in a confusion matrix, whose image is shown below:

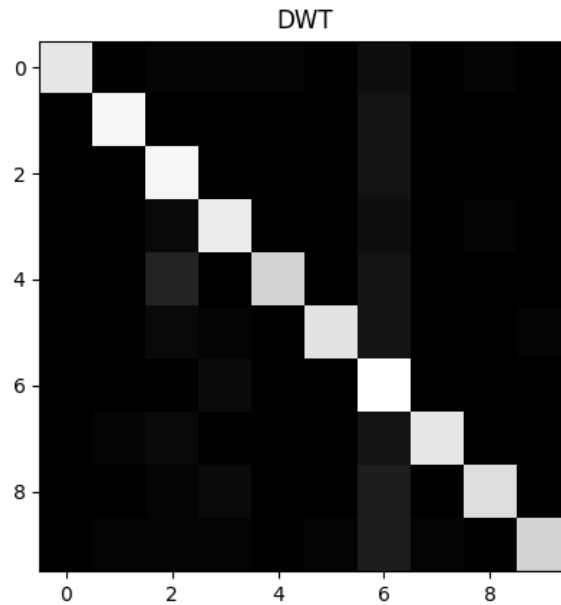


Figure 5: MFCCs

As earlier, the row index in the confusion matrix indicates the ground truth, that is, the actual identity of the test utterance. The column index indicates the classification output. 1 is added to confusion matrix index  $(x+1, y+1)$  every time digit  $x$  is classified as  $y$ , thus the common classification outputs are brighter than the rest. Each row is scaled such that the maximum in the row is white. This indicates that the most common confusions were

- 2 confused as 4
- 8 confused as 6

We notice that the confusion matrix in this case is much cleaner than the previous approaches. The WER is also much better. This is because this

approach is more principled in capturing phoneme timing information. The previous two methods get confused between similar phonemes, but this method uses sequence information to disambiguate word identity.

Code for bag of DWT approach:

```
1 import os
2 import scipy
3 import numpy as np
4 from scipy import signal
5 import matplotlib.pyplot as plt
6 from scipy.cluster.vq import vq, kmeans
7 import params
8 import feature_extractor
9
10
11 def find_dtw_distance(test_pattern, ref_pattern):
12     n = test_pattern.shape[1]
13     m = ref_pattern.shape[1]
14
15     distMat = np.zeros([n, m])
16
17     for i in range(n):
18         for j in range(m):
19
20             distMat[i, j] = np.linalg.norm(np.subtract(test_pattern
21 [:, i], ref_pattern[:, j]))
22
23     DIW = np.zeros([n+1, m+1])
24
25     for i in range(1, n+1):
26         DIW[i, 0] = float('Inf')
27
28     for i in range(1, m+1):
29         DIW[0, i] = float('Inf')
30
31     DIW[0, 0] = 0
32
33     for i in range(1, n+1):
34         for j in range(1, m+1):
35             cost = distMat[i-1, j-1]
36             DIW[i, j] = cost + np.min([DIW[i-1, j], np.min([DIW[i
37 -1, j-1], DIW[i, j-1]])])
38
39     return DIW[n, m]
40
41 sample_rate = params.sample_rate
42
43 digits = ['zero', 'one', 'two', 'three', 'four', 'five', 'six', 'seven', '
44 eight', 'nine']
45
46 seg_male_dir = params.seg_male_dir
47 seg_female_dir = params.seg_female_dir
48
49 male_speakers = os.listdir(seg_male_dir)
50 female_speakers = os.listdir(seg_female_dir)
51 all_speakers = male_speakers + female_speakers
52 all_speakers.remove('.DS_Store')
53 codeBook = {}
54 n_frame_dict = {}
55
56
57 for digit in digits:
58     print("Preparing codebook for digit", digit)
```

```

59     codeBook[digit] = {}
60     n_frame_dict[digit] = {}
61
62     for speaker in all_speakers:
63         print("speaker = ", speaker)
64         parent_dir = ''
65         if speaker in male_speakers:
66             parent_dir = seg_male_dir
67
68         if speaker in female_speakers:
69             parent_dir = seg_female_dir
70         codeBook[digit][speaker] = []
71         n_frame_dict[digit][speaker] = []
72         for iteration in range(1,3):
73             file = parent_dir + '/' + speaker + '/' + str(digit) +
74             '_' + str(iteration) + '.wav'
75             feature_mat = feature_extractor.main(file)
76             n_frame_dict[digit][speaker].append(feature_mat.shape
77             [1])
78             for i in range(feature_mat.shape[1]):
79                 codeBook[digit][speaker].append(feature_mat[:,i])
80
81     print('Codebook created')
82
83     confusion_matrix = np.zeros([10,10])
84
85     for test_speaker in all_speakers:
86         train_speakers = all_speakers
87         train_speakers.remove(test_speaker)
88
89         for test_digit in digits:
90             for utterance in range(1,3):
91
92                 n_frames = n_frame_dict[test_digit][test_speaker]
93                 test_mat = codeBook[test_digit][test_speaker][sum(
94                 n_frames[0:(utterance-1)]):sum(n_frames[0:utterance])]
95
96                 sum_dist = np.zeros(10)
97                 min_dist = float('Inf')
98                 for digit in digits:
99
100                     for speaker in train_speakers:
101                         for ref_utterance in range(1,3):
102                             n_ref_frames = n_frame_dict[digit][speaker]
103                             ref_mat = codeBook[digit][speaker][sum(
104                             n_ref_frames[0:(ref_utterance-1)]):sum(n_ref_frames[0:
105                             ref_utterance])]
106
107                             test_pattern = np.transpose(np.asarray(
108                             test_mat))
109                             ref_pattern = np.transpose(np.asarray(
110                             ref_mat))
111
112                             if np.sum(ref_pattern.shape) == 0:
113                                 continue
114                             curr_dist = find_dtw_distance(test_pattern,
115                             ref_pattern)
116
117                             if(curr_dist < min_dist):
118                                 min_dist = curr_dist
119                                 pred_digit = digits.index(digit)
120
121                     print("For ", test_speaker, " predicted digit = ",
122                     pred_digit, " ground truth = ", test_digit)

```

```

112         confusion_matrix[digits.index(test_digit), pred_digit]
113         += 1
114     np.save('DTW_confusion_matrix', confusion_matrix)
115
116 wer = 1 - np.trace(confusion_matrix)/(len(all_speakers)*20)
117 print wer

```