# Lab Assignment 6

Team 08 : EE17B047(Kommineni Aditya) - EE17B035(V Sai Krishna)

March 2020

## 1 Introduction

This lab assignment primarily deals with **text dependent speaker verification** on a closed set data. The dataset contains utterances from 23 speakers and each speaker has about 5 instances out of which 3 of them are being taken as templates and 2 of them as test instances.

The implementation of the verification will involve the use of **Dynamic Time Warping**, a dynamic programming technique which is a robust technique to compare sequences of variable length.

All the utterances were recorded at a sampling frequency of **16kHz** and were **single channel** recordings.

## 2 Cepstrum Calculation

The speech recordings were windowed at regular intervals using the hamming window. Then, the windowed frames were zero padded following which the magnitude spectrum of each frame was computed using the fft function.

Now, since we have the magnitude spectrum, cepstra were calculated using the formula below :
$$cepstra = F[\, log(\, F[w]^2\, )\, ]$$

We must repeat this for all the utterances which have been recorded which facilitates the use of the dynamic time warping algorithm in the next section.

Following is an illustration of a speech frame from an utterance and the corresponding cepstra for the speech frame.
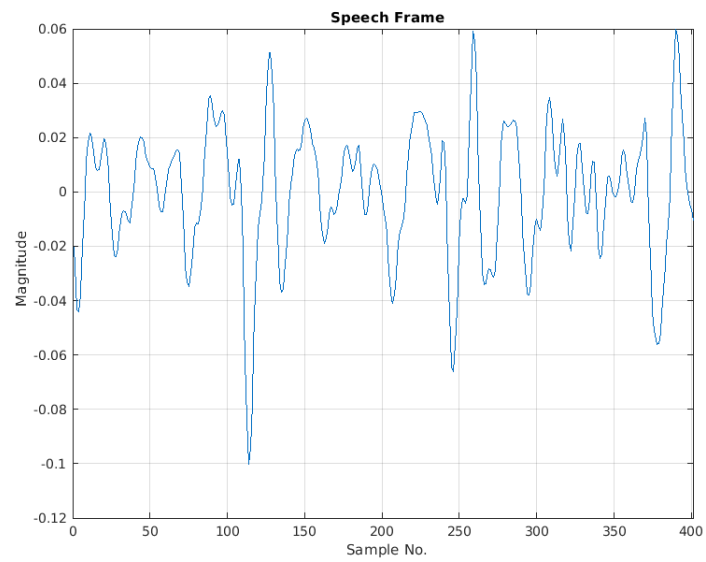
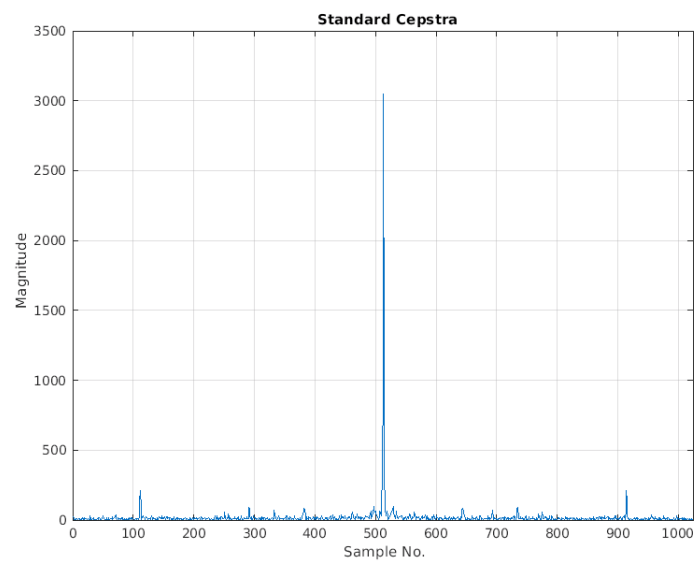Figure 1: A frame of windowed speech



Figure 2: Cepstra corresponding to the speech

# 3 Dynamic Time Warping

This algorithm is extremely useful to compare 2 time series of varying length and helps in deducing their similarity. Here, we employ it to compare the similarity between two speech utterances and thereby identify if both of them were spoken by the same speaker or not. An image which depicts the gist of dynamic time warping is as below :
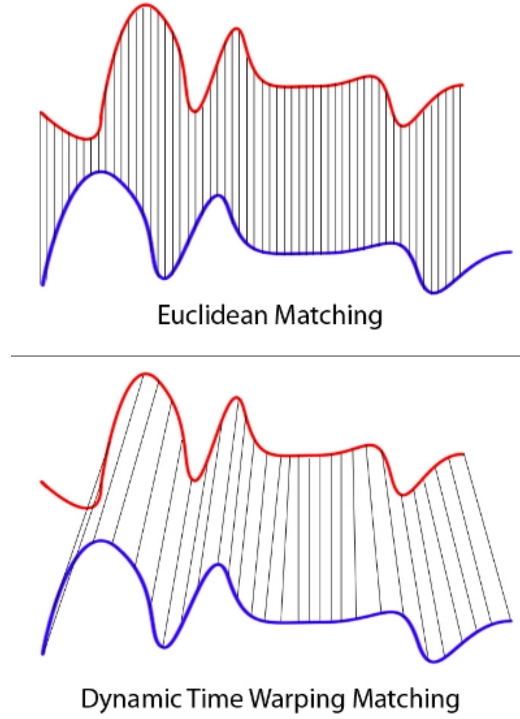


Figure 3: Illustration of Dynamic Time Warping

## 3.1 Implementation

Let the template utterance be a length of M and the test utterance be a length of N. A pseudo code which depicts the implementation of the algorithm is as follows :

```
for i = 1:M
    for j = 1:N
        if (j > 0 and i > 1 )
        {dtw[i,j] = cost(i,j) + min( dtw[i-2,j-1],dtw[i-1,j-1], dtw[i,j-1]) }
        else if (j > 0 and i = 1)
        {dtw[i,j] = cost(i,j) + min(dtw[i-1,j-1],dtw[i,j-1])}
```

```
            else if (j > 0 and i = 1)
            {dtw[i,j] = cost(i,j) + dtw[i,j-1]}
```

The code above makes sure that every test frame is mapped to only one frame in the template and prevents multiple template frames from mapping to one test frame. Also, it allows for the current input to make a jump and ignore one frame of the template if that minimises the cost. Also, the cost function refers to the euclidean distance between between a frame of the test and the template utterance.

## 3.2   Evaluation Metrics

## 3.3   Least Euclidean Distance

In this method, we run the Dynamic Time Warping Algorithm over all the templates for a given test case. Then, the utterance which has the minimum average euclidean distance is taken to be the speaker who spoke the test utterance.

## 3.4   Average Euclidean distance across single speaker

In this evaluation metric, we wanted to exploit the opportunity to utilise all the templates for a speaker. This is done by taking the arithmetic mean of the dtw costs for test utterance across all three templates of a speaker.

However, the accuracy obtained using this method was inferior when compared to using the prior mentioned metric. Therefore, all the accuracy values have been reported for Least Euclidean Distance metric alone.

## 3.5   Observations

### 3.5.1   Accuracy

- The accuracy obtained while using a window of 25 ms is: 82.61%

- The accuracy obtained while using a window of 15 ms is: 91.30%

- The accuracy obtained while using a window of 10 ms is: 84.78%

### 3.5.2   Effect of Window length

As we decrease the window length from very high time instances to considerably lower, we can observe that the accuracy steadily increases upto a certain window length below which the accuracy decreases again. The rough values of the window lengths are as above.

### 3.5.3 Effect of overlapping windows

Using windows that overlap gives a higher accuracy compared to using windows without overlaps while obtaining the cepstra. However, just like in the case of non-overlapping windows, the accuracy increases upto a certain window length below which the accuracy decreases again. We expect it to be better over the non overlapping window case because overlapping windows capture transition feature vectors better. The accuracy obtained are listed below:

- Accuracy while using a window of 25 ms, advanced by 15 ms is: 93.47%

- Accuracy while using a window of 15 ms, advanced by 10 ms is: 89.13%

## 3.6 Confusion Matrices

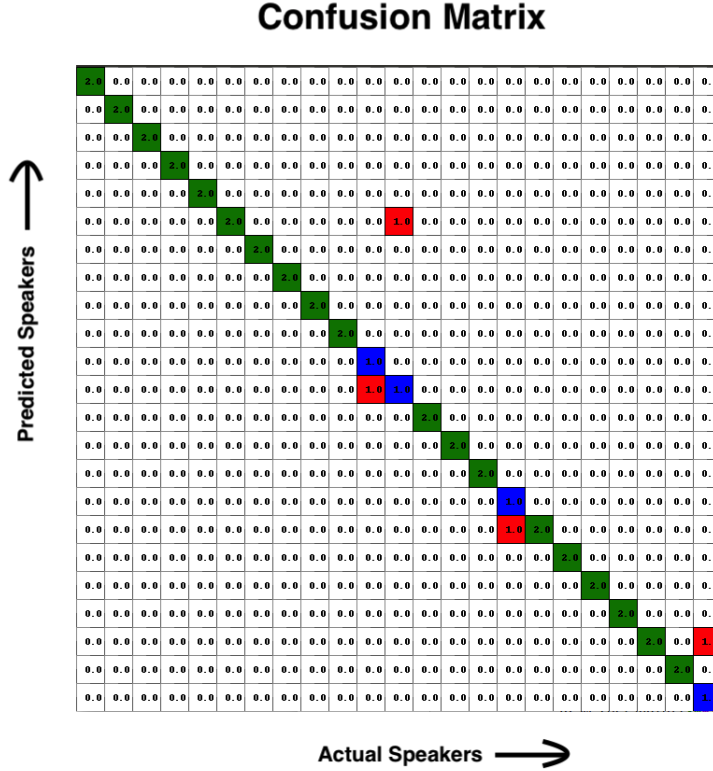Below is the confusion matrix for window length 15ms and the frame is advanced by 15ms each time.



Figure 4: Confusion Matrix for window length 15ms, advanced by 15 ms

Below is the confusion matrix for window length 25ms and the frame is advanced by 25ms each time.
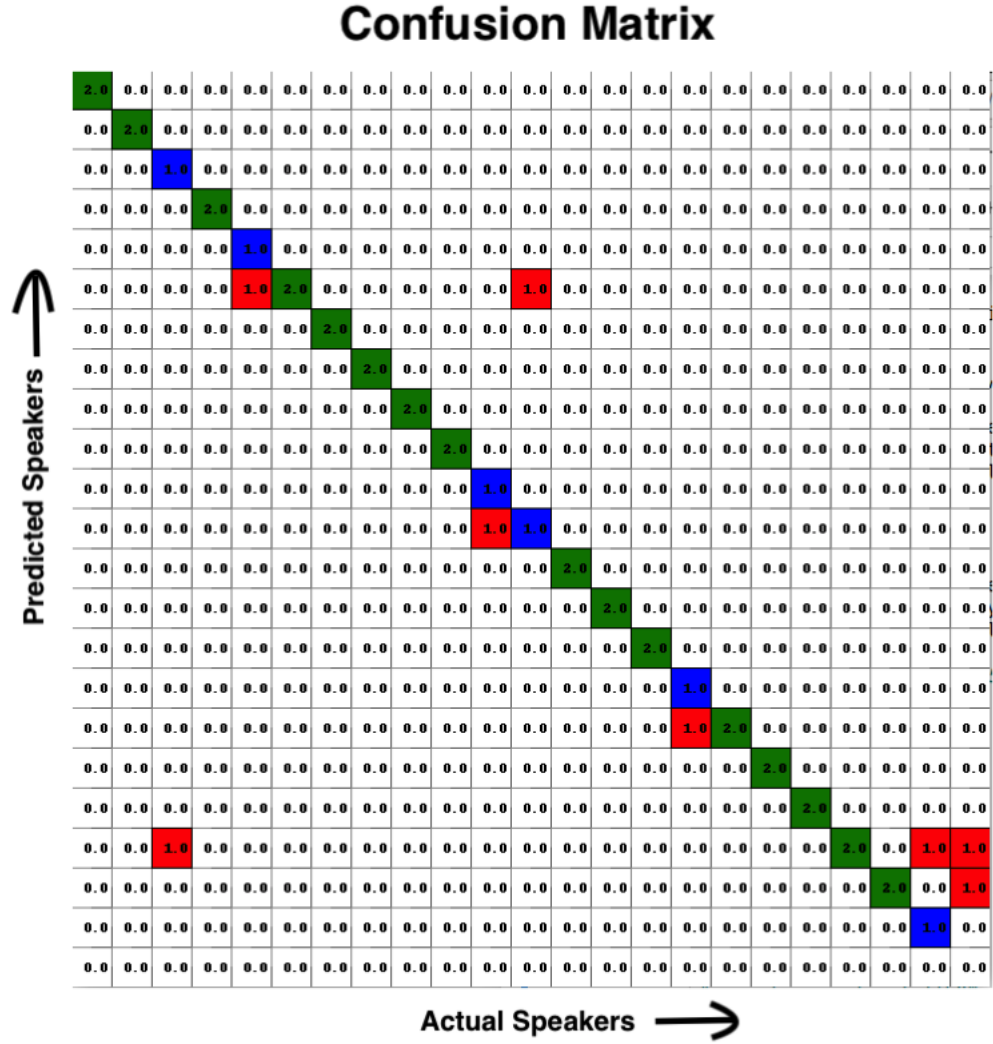


Figure 5: Confusion Matrix for window length 25ms, advanced by 25 ms

Below is the confusion matrix for window length 10ms and the frame is advanced by 10ms each time.
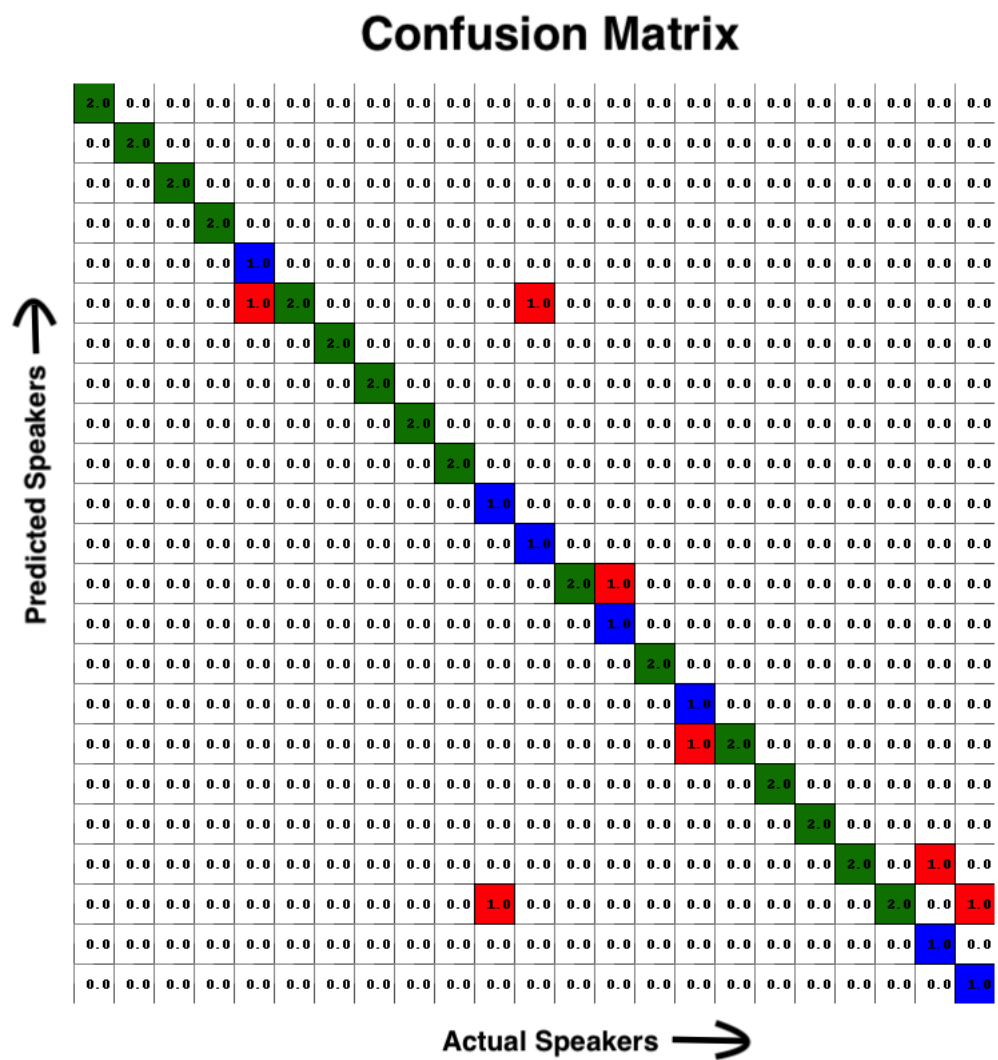


Figure 6: Confusion Matrix for window length 10ms, advanced by 10ms

Below is the confusion matrix for window length 25ms and the frame is advanced by 15ms each time.
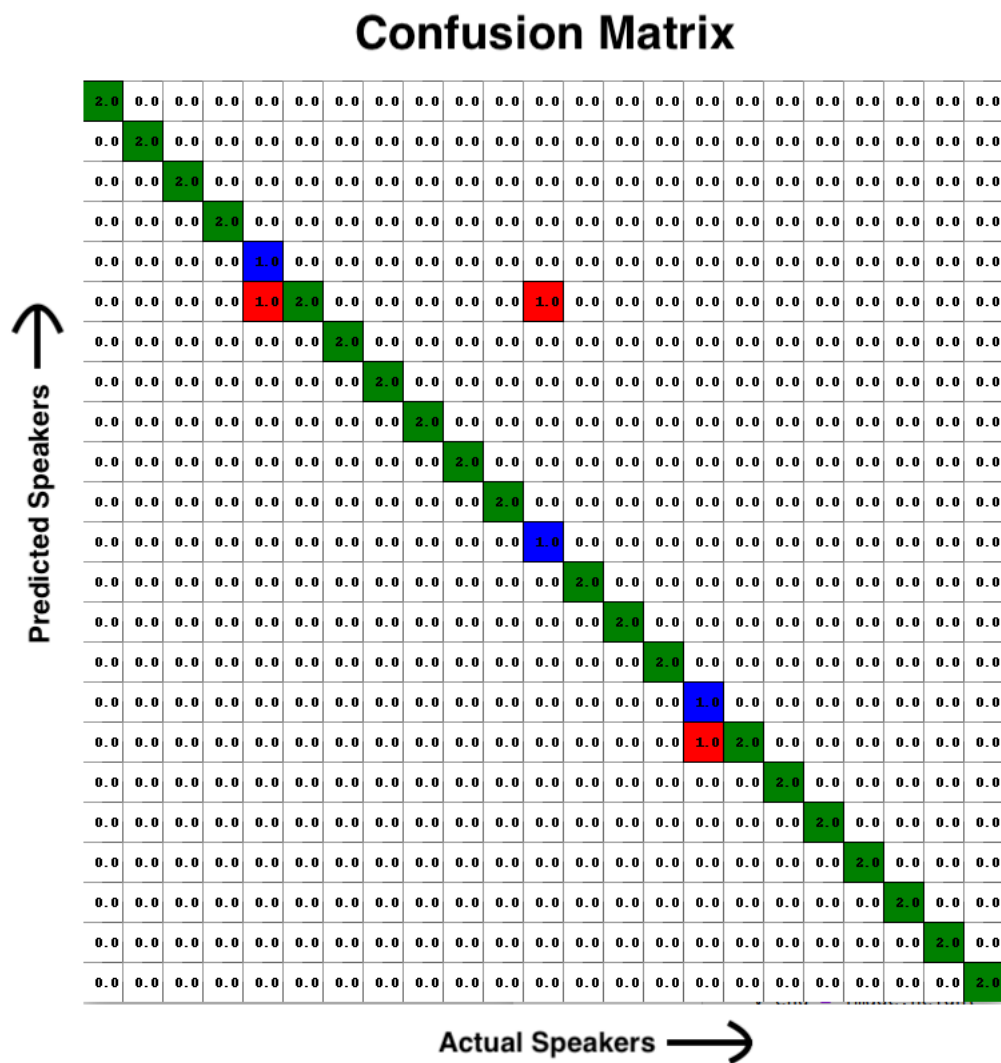
## Confusion Matrix



Figure 7: Confusion Matrix for window length 25ms shifted by 15ms

Below is the confusion matrix for window length 15ms and the frame is advanced by 10ms each time.
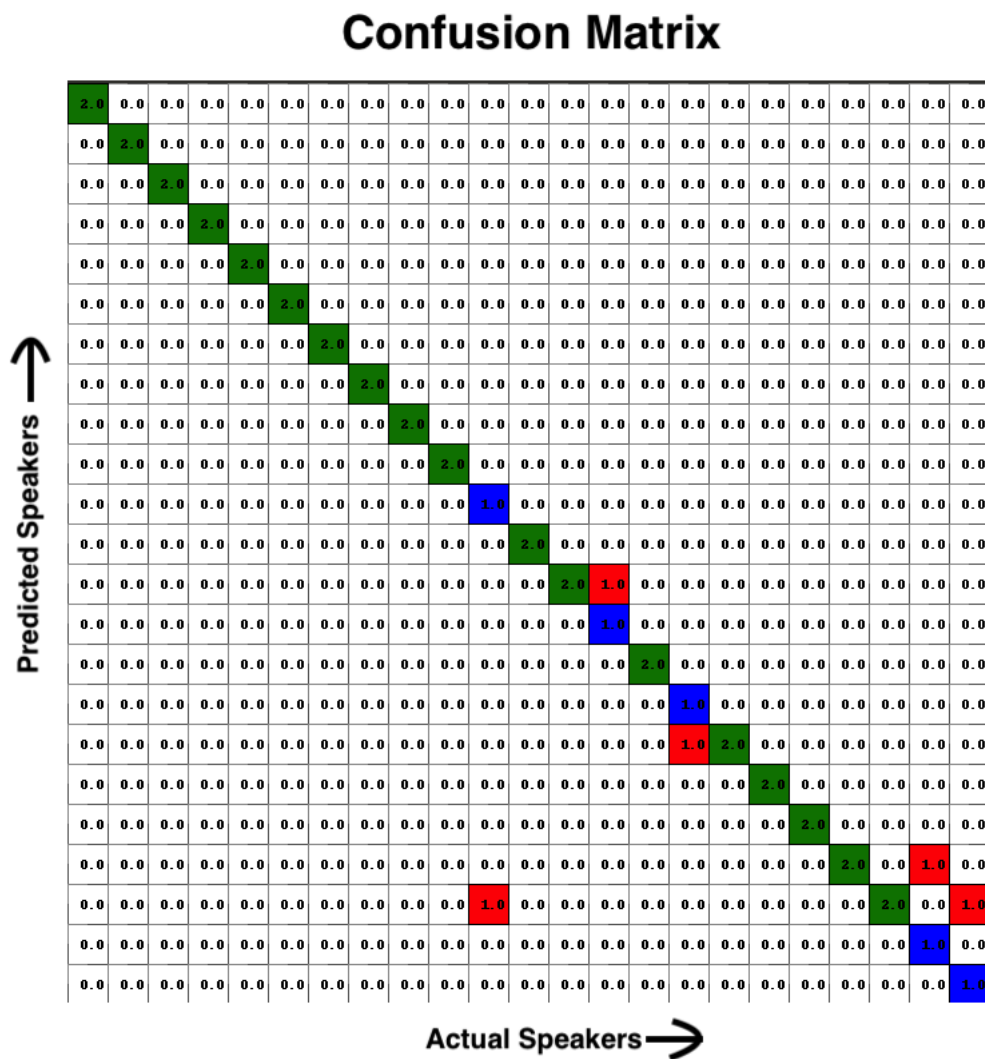


Figure 8: Confusion Matrix for window length 15ms shifted by 10ms

## 3.7 Inference

- We can conclude from the above observations that the text dependent speech recognition systems are prone to errors whereby the utterance can be mis-classified owing to two key phrases being similar.

- Although it is text dependent, we can see from the above confusion matrices that for speakers with identical key phrases the model is able to disambiguate the identity of the speaker to a good extent. This could be attributed to property of the cepstra capturing speaker dependent features.

- The window size which one chooses while calculating the cepstral coefficients will play an important role in the accuracy of the model.

# 4   Codes used for the assignment

Please click on the following link to access the codes:
https://drive.google.com/drive/folders/134ExedMVua-dFOQFvmrUDKACZlwCIPum?usp=sharing