# Speech Recognition through Mel Frequency Cepstral Coefficients and Dynamic Time Warping

Thomas Smallarz

13 December 2020

## Contents

## Introduction

Speech is a fundamental mode of communication between people, even in this age of non-verbal communication (email, text, etc.), so the benefits of being able to translate human speech to a digital format is obvious.

This report will explain implementing a method of speech recognition in MATLAB R2020b using Mel Frequency Cepstral Coefficients and Dynamic Time Warping with the intention of implementing this method on a low-cost embedded target.

## 1 Method

Most modern automatic speech recognition (ASR) systems have the same structure. It is as follows:

1. Speech is recorded and digitized

2. Features of waveform are extracted

3. Strings of phonemes are recognized based on features

4. Words and strings of words are assembled based on language models
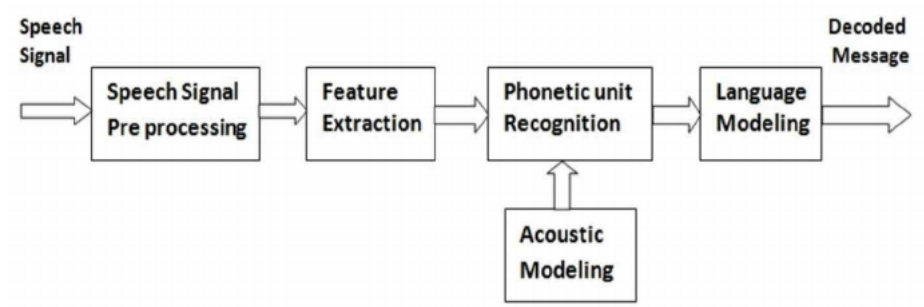


Figure 1: General speech recognition system diagram

As of now, we won't concern ourselves with the first step – recording and digitizing audio. This will be talked about a little in the Implementation section later on. Also, language modeling will not be discussed in this report.

## 1.1 Mel Frequency Cepstral Coefficients

Before using algorithms to determine the likelihood of speech utterances to be specific words, the data is typically manipulated to gain characteristics that describe what is being said.

In this implementation Mel Frequency Cepstral Coefficients (MFCC) are used. There are numerous ways to extract features or characteristics of speech. This is one method that is used.

### 1.1.1 Cepstral

To understand how to derive coefficients of the MFC, it is first important to learn what the "cepstrum" is. The cepstrum of a signal is the result of taking the inverse Fourier transform (IFT) of the logarithm of the Fourier transform (FT) of a signal. Or,

$$|\mathcal{F}^{-1}\{\log|(\mathcal{F}\{x[n]\})^2|\}|^2$$

### 1.1.2 Mel Scale

The Mel scale is a scale based on the perceptions of human hearing. It is a mapping function that separates pitches by what humans consider to be equal distance. A common formula to convert frequency from hertz into mels is:
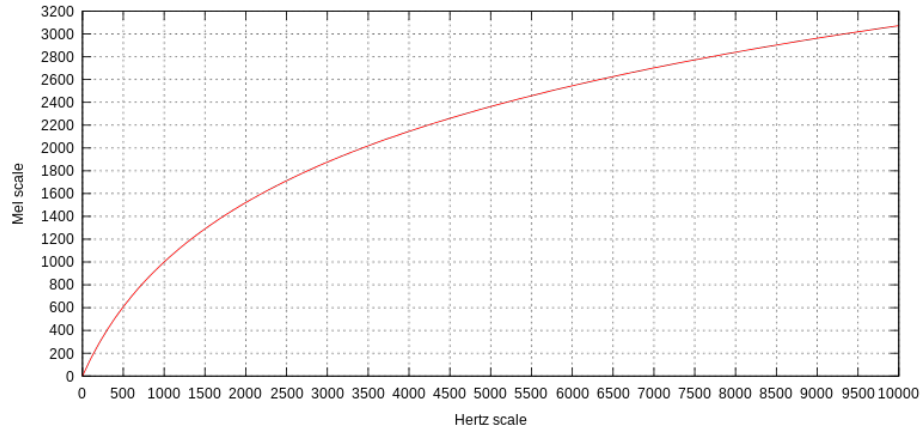
$$m = 2595 log_{10}(1 + f/700)$$



Figure 2: Mel scale vs. hertz

### 1.1.3 Mel Frequency Cepstral

Including both these concepts, the difference between Cepstral and Mel Frequency Cepstral is that MFC has its initial spectrum mapped to the mel scale before taking the logarithm. The common derivation of MFCC's is as follows:

1. Take the Fourier Transform of signal

2. Find the powers of that spectrum

3. Map onto the mel scale (typically with triangle filters)

4. Take the logs at each mel frequency

5. Take the discrete cosine transform (DCT) of the mel log powers

Then, the MFC coefficients are the amplitudes of the resulting spectrum.

## 1.2 Dynamic Time Warping

Dynamic Time Warping (DTW) is an algorithm that measures the similarity between two sequences in time (or other dependent variable) that may vary in speed. It calculates a distance matrix, then the shortest weighted path of that matrix is the "distance" sequence used to map one of the signals onto the other. So, with $A_i$ and $B_i$ sequences

1. Create distance matrix of size i x j

2. Find shortest weighted distance from opposite corners of distance matrix
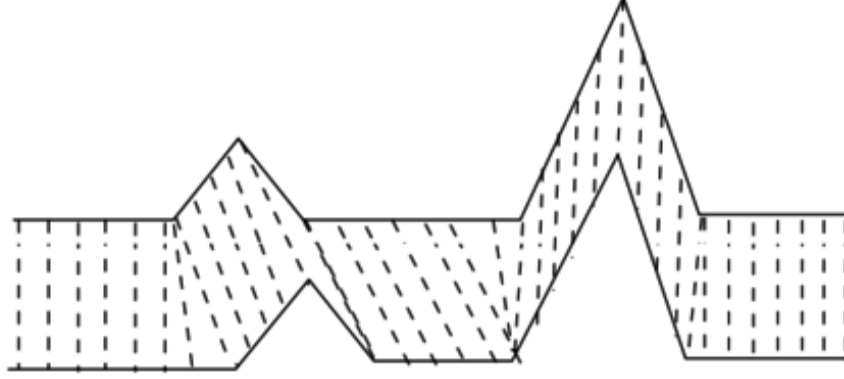
3

Figure 3: Example of mapping one sequence onto another using DTW

### 1.2.1 Creating distance matrix

Creating the distance matrix is done with a simple equation. If the values of $A_i$ and $B_j$ are put on row and column side of the matrix respectively, then the values at $D_{i,j}$ are:

$$D(i,j) = |B_j - A_i| + min(D(i-1,j-1), D(i-1,j), D(i,j-1))$$

This is calculated starting at the first place in the sequence ($D_{0,0}$) and working towards the last place ($D_{i,j}$). Values $D(i-1,j-1)$, $D(i-1,j)$, and $D(i,j-1)$ are used when they have been already computed.

There are more computationally simpler methods of computing this distance matrix, but in short, this is how it is calculated.

### 1.2.2 Finding shortest path

Once a distance matrix is created we need to find the values to warp the two sequences by to to align them. These values are found by starting in the upper corner of the distance matrix ($D(i,j)$) and recursively find the minimum of $D(i-1,j-1)$, $D(i-1,j)$, and $D(i,j-1)$ until reaching the lower corner ($D(i,j)$).

The amplitude values of the distance matrix could be likened to the height of a mountain range, and the path you take to the origin is the valley between the mountain peaks.

## 2 Implementation

Now, armed with the means to find the MFCC's and perform a DTW. We can create some training data (multiple recordings of some word) and have an

input and see if our system can recognized the inputted word.

## 2.1 Digitizing Audio

First is to record the audio. In MATLAB there is a audiorecorder( ) function that uses your default system microphone and has parameters: sampling frequency, bit size to store data in, and number of channels to record.

- Sampling Frequency: A quick Google search brings us to the Voice Frequency Wikipedia page that says a sampling rate of 8,000Hz is sufficient (a 4000Hz bandwidth doubled)

- Bit Size: Let's go with 16-bit Integers for now. Implementing on lower-cost hardware may require less size

- Number of Channels: One channel for one microphone

```matlab
1    clear variables; close all;
2
3    % **************************
4    % *** EDITABLE VARIABLES ***
5    % **************************
6    words = ["Red","Green","Blue"]; % What words you want to record
7    numTrain = 3; % number of training data for each word
8    recTime = 1; % record time in seconds
9    Fs = 8e3; % recording sampling frequency
10   nBits = 16; % number of bits to record audio in
11   NumChannels = 1; % number of channels to record from
12   DistanceMetric = "squared"; % euclidean, absolute, squared as options
13   % ***********
14   % *** END ***
15   % ***********
16
```

Figure 4:

Now, we create an instance of the audiorecorder and record each word "numTrain" times with a one second duration to record.

5

```
17      %% Record and save audio
18 -    recObj = audiorecorder(Fs,nBits,NumChannels); % creating record object
19
20 -    fprintf("Get Ready to record!\n");
21 -    pause(1);
22
23 -   ⌐for c = words
24 -         eval(c+"=[];");
25 -    ⌐      for k = 1:numTrain
26 -             fprintf("Click any button to start recording #"+num2str(k)+" for "+c+"\n");
27 -             waitforbuttonpress;
28 -             fprintf("Recording ...\n");
29 -             recordblocking(recObj,recTime);
30 -             fprintf("End recording #"+num2str(k)+" for "+c+"\n");
31 -             pause(0.5);
32 -             fprintf("Playing back recording #"+num2str(k)+" for "+c+"\n");
33 -             play(recObj);
34 -             pause(1);
35 -             if k == 1; eval(c+"=getaudiodata(recObj);");
36 -             else; eval("temp=getaudiodata(recObj);");
37 -                   eval(c+"=["+c+",temp];");
38 -             end
39 -         end
40 -   └end
41
42 -    clear temp;
```

Figure 5:

After this we can record our input audio using the same method.

```
57      %% Record input audio
58 -    recObj = audiorecorder(Fs,nBits,NumChannels); % creating record object
59 -    fprintf("Get Ready to record!\n");
60 -    pause(1);
61 -    fprintf("Click any button to start recording input audio\n");
62 -    waitforbuttonpress;
63 -    fprintf("Recording ...\n");
64 -    recordblocking(recObj,recTime);
65 -    fprintf("End recording\n");
66 -    pause(0.5);
67 -    fprintf("Playing back input audio recording\n");
68 -    play(recObj);
69 -    pause(1);
70 -    inputAudio = getaudiodata(recObj); % save input audio
71
72      % plot input audio
73 -    figure();
74 -    plot(linspace(0,recTime,recTime*Fs),inputAudio,'k');
75 -    title("Recording of Input Audio");
```
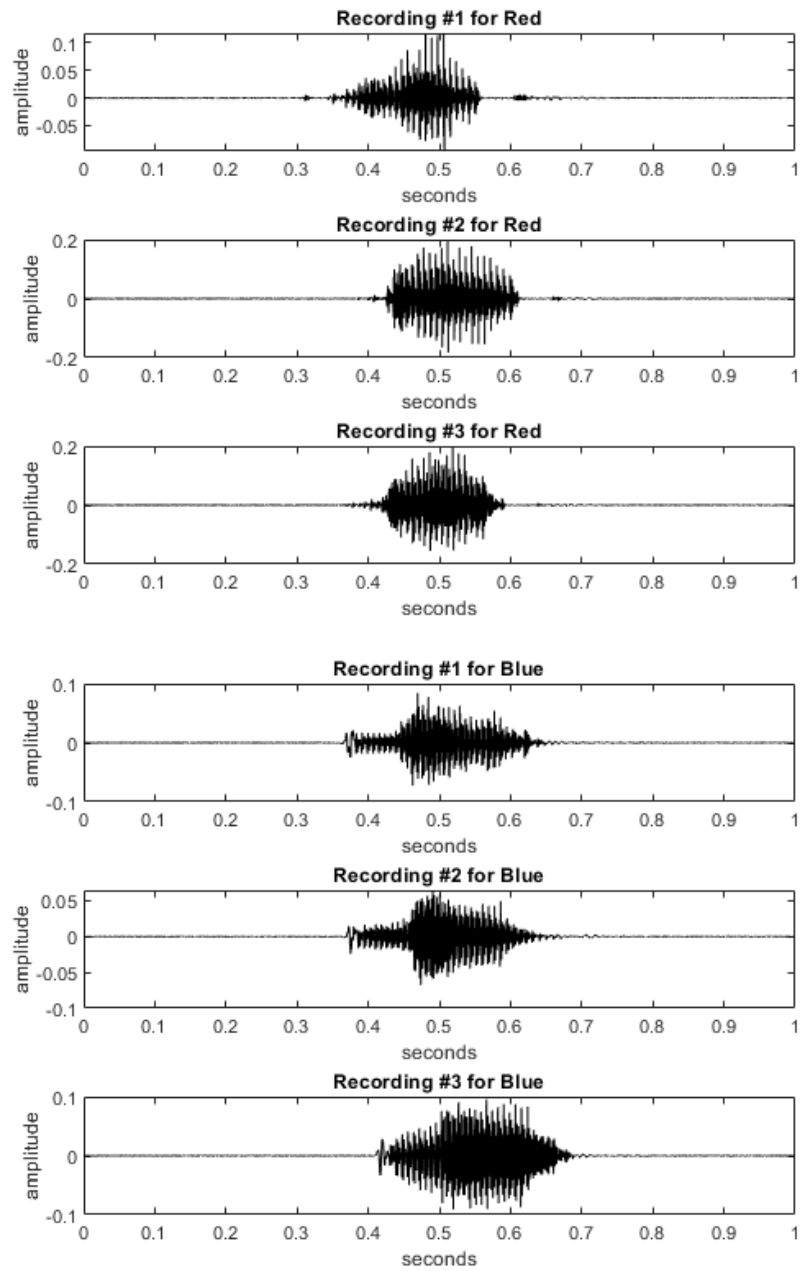
Figure 6:

Plotting recorded audio:

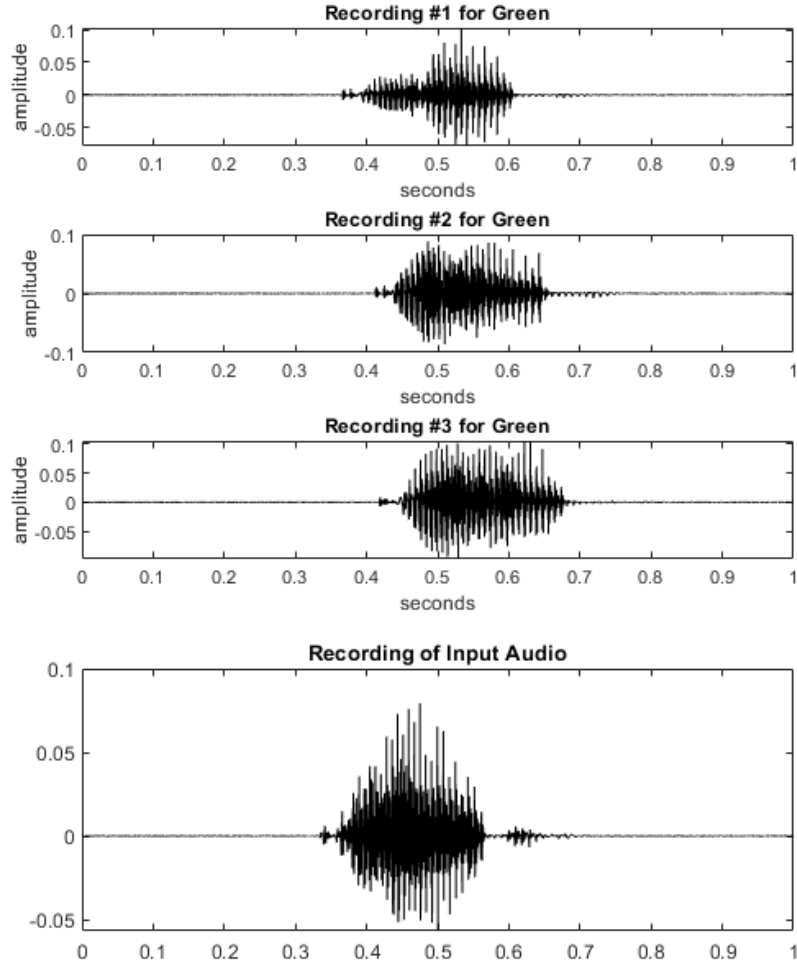Figure 7: "Red" and "Blue" recorded audios

Figure 8: "Green" and input recorded audios

Now that we have all our "training" audio and input audio saved it's time to calculate our MFC coefficients.

## 2.2 Feature Extraction

Using the MFCC method we talked about earlier, we can use the built-in mfcc( ) MATLAB function. This function has input parameters for the audio amplitude values, and sampling frequency, and returns the mel frequency cepstral coefficients. The algorithm from MATLAB's documentation for this function is:
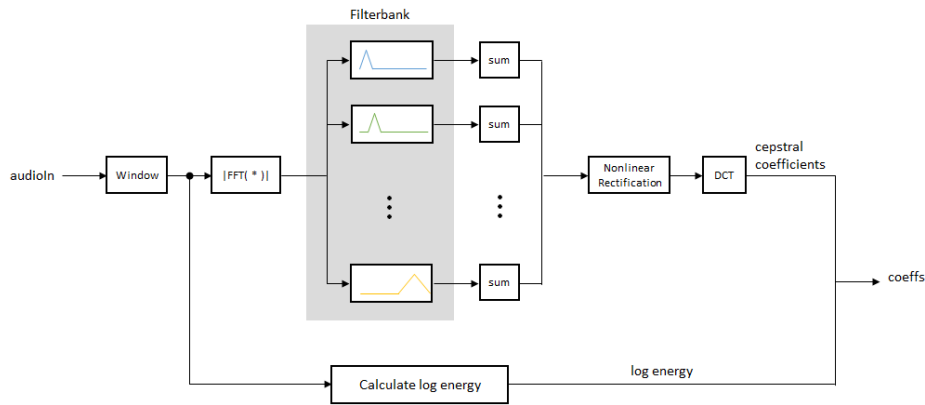
Figure 9:

This is implemented in MATLAB as:

```
78          %% Calculate MFCC Coefficients
79
80 -    ┌for c = words
81 -    │    for k = 1:numTrain
82 -    │        K = num2str(k);
83 -    │        C = c+K;
84 -    │        eval("coef"+C+"=mfcc("+c+"(:,k),Fs);");
85 -    │    end
86 -    └end
87
88 -    coefInputAudio = mfcc(inputAudio,Fs);
```

Figure 10:

## 2.3 Feature Recognition

Now that we've calculated the MFC coefficients for each audio sequence we can compare the inputted audio against the templates using dynamic time warping to see which template has the lowest distance when compared to our inputted audio.

For this, MATLAB has a built-in dtw( ) function that has input parameters of an input vector "x" and an input vector "y". The order of inputting "x" and "y" does not matter. This function outputs: the total calculated euclidean distance of "x" and "y", the warping path for "x", and the warping path for "y".
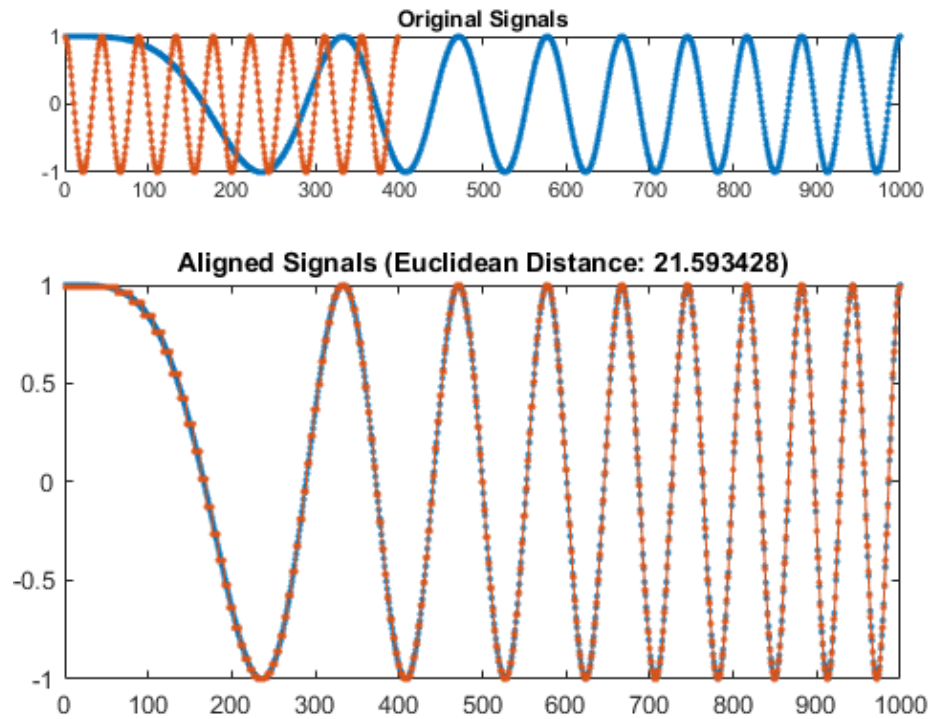
9

Figure 11: Example of how dtw( ) "maps" one sequence to another

Implementing this in MATLAB:

```
90      %% Calculate Dynamic Time Warping distance for MFCC coefficients
91 -    for i = 1:width(coefInputAudio)
92 -        for c = words
93 -            for k = 1:numTrain
94 -                K = num2str(k);
95 -                C = c+K;
96 -                Ci = C+"_"+num2str(i);
97 -                eval("[dist"+Ci+",x"+Ci+",y"+Ci+"]=dtw(coef"+C+"(:,i),coefInputAudio(:,i),'"+DistanceMetric+"');");
98 -            end
99 -        end
100 -   end
```

Figure 12:

Then plotting results gives us:

```
102          %% Plot new warped signals
103 —    ┌ for c = words
104 —    │       figure();
105 —    │ ┌    for k = 1:numTrain
106 —    │ │        subplot(numTrain,1,k);
107 —    │ │        K = num2str(k);
108 —    │ │        C = c+K;
109 —    │ │ ┌      for i = width(coefInputAudio)
110 —    │ │ │          Ci = C+"_"+num2str(i);
111 —    │ │ │          hold on;
112 —    │ │ │          eval("plot(coef"+C+"(x"+Ci+","+num2str(i)+"),'k')");
113 —    │ │ │          hold on;
114 —    │ │ │          eval("plot(coefInputAudio(y"+Ci+",i),'b');");
115 —    │ │ │          hold on;
116 —    │ │ └      end
117 —    │ │        title("DTW of Input Audio vs. " + C);
118 —    │ └    end
119 —    └ end
```
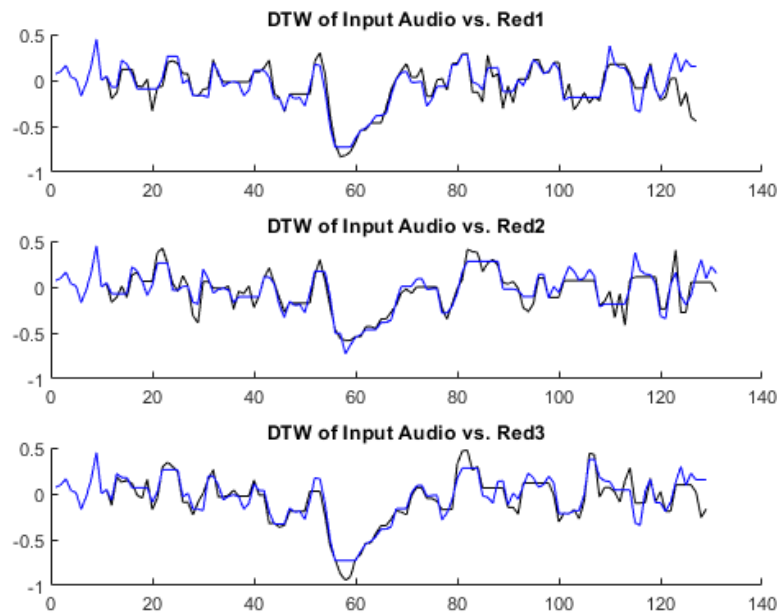
Figure 13:



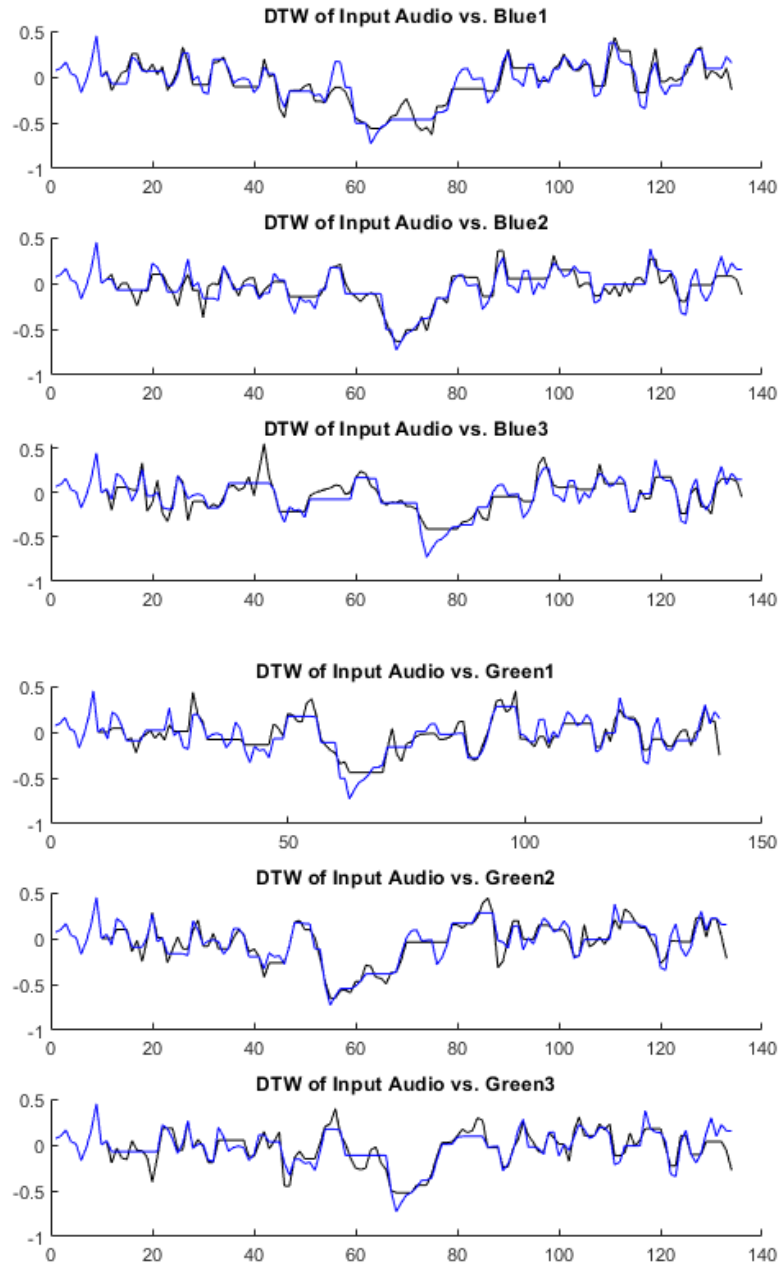Figure 14: DTW of Input vs. Red MFC Coefficients

Figure 15: DTW of Input vs. Blue/Green MFC Coefficients

Then, summing all the distances from the Red, Blue, and Green coefficient matrices gives us:

```
121      %% Calculate sum of distances
122
123 -  ┌for c = words
124 -  ├    for k = 1:numTrain
125 -  │        C = c+num2str(k);
126 -  │        eval(C+"sum=0;");
127 -  ├        for i = 1:width(coefInputAudio)
128 -  │            eval(C+"sum="+C+"sum + dist"+C+"_"+num2str(i)+";");
129 -  │        end
130 -  │    end
131 -  └end
132
133 -  ┌for c = words
134 -  ├    for k = 1:numTrain
135 -  │        C = c+num2str(k);
136 -  │        eval(C+"sum");
137 -  │    end
138 -  └end
```

Figure 16:

Summing each gives us:

| Color | Set | Distance |
|-------|-----|----------|
|       | 1   | 95.78    |
| Red   | 2   | 180.09   |
|       | 3   | 139.43   |
|       | 1   | 179.83   |
| Blue  | 2   | 172.95   |
|       | 3   | 292.08   |
|       | 1   | 208.80   |
| Green | 2   | 241.63   |
|       | 3   | 201.73   |

From this we can see that training sets from the color red have the lowest two distances. This is a good thing because in this case I did say "Red" for my input recording.

# A    MATLAB

```
1   clear variables; close all;
2
3   % ************************
4   % *** EDITABLE VARIABLES ***
5   % ************************
6   words = ["Red","Green","Blue"]; % What words you want to
        record
```

13

```matlab
7   numTrain = 3; % number of training data for each word
8   recTime = 1; % record time in seconds
9   Fs = 8e3; % recording sampling frequency
10  nBits = 16; % number of bits to record audio in
11  NumChannels = 1; % number of channels to record from
12  DistanceMetric = "squared"; % euclidean, absolute,
        squared as options
13  % ***********
14  % *** END ***
15  % ***********
16
17  %% Record and save audio
18  recObj = audiorecorder(Fs, nBits, NumChannels); % creating
        record object
19
20  fprintf("Get Ready to record!\n");
21  pause(1);
22
23  for c = words
24      eval(c+"=[];");
25      for k = 1:numTrain
26          fprintf("Click any button to start recording #"+
                num2str(k)+" for "+c+"\n");
27          waitforbuttonpress;
28          fprintf("Recording ...\n");
29          recordblocking(recObj, recTime);
30          fprintf("End recording #"+num2str(k)+" for "+c+"\n
                ");
31          pause(0.5);
32          fprintf("Playing back recording #"+num2str(k)+"
                for "+c+"\n");
33          play(recObj);
34          pause(1);
35          if k == 1; eval(c+"=getaudiodata(recObj);");
36          else; eval("temp=getaudiodata(recObj);");
37                  eval(c+"=["+c+",temp];");
38          end
39      end
40  end
41
42  clear temp;
43
44  %% Plot saved audio
45  close all;
46
47  for c = words
```

```matlab
48        figure();
49        for k = 1:numTrain
50            subplot(numTrain,1,k);
51            eval("plot(linspace(0,"+num2str(recTime)+","+
                 num2str(recTime*Fs)+"),"+c+"(:,k),'k');");
52            title("Recording #"+num2str(k)+" for "+c);
53            xlabel("seconds"); ylabel("amplitude");
54        end
55  end
56
57  %% Record input audio
58  recObj = audiorecorder(Fs,nBits,NumChannels); % creating
        record object
59  fprintf("Get Ready to record!\n");
60  pause(1);
61  fprintf("Click any button to start recording input audio\
        n");
62  waitforbuttonpress;
63  fprintf("Recording ...\n");
64  recordblocking(recObj,recTime);
65  fprintf("End recording\n");
66  pause(0.5);
67  fprintf("Playing back input audio recording\n");
68  play(recObj);
69  pause(1);
70  inputAudio = getaudiodata(recObj); % save input audio
71
72  % plot input audio
73  figure();
74  plot(linspace(0,recTime,recTime*Fs),inputAudio,'k');
75  title("Recording of Input Audio");
76
77
78  %% Calculate MFCC Coefficients
79
80  for c = words
81      for k = 1:numTrain
82          K = num2str(k);
83          C = c+K;
84          eval("coef"+C+"=mfcc("+c+"(:,k),Fs);");
85      end
86  end
87
88  coefInputAudio = mfcc(inputAudio,Fs);
89
90  %% Calculate Dynamic Time Warping distance for MFCC
```

```matlab
        coefficients
91  for i = 1:width(coefInputAudio)
92      for c = words
93          for k = 1:numTrain
94              K = num2str(k);
95              C = c+K;
96              Ci = C+"_"+num2str(i);
97              eval("[dist"+Ci+",x"+Ci+",y"+Ci+"]=dtw(coef"+C
                    +"(:,i),coefInputAudio(:,i),'"+
                    DistanceMetric+"');");
98          end
99      end
100 end
101
102 %% Plot new warped signals
103 for c = words
104     figure();
105     for k = 1:numTrain
106         subplot(numTrain,1,k);
107         K = num2str(k);
108         C = c+K;
109         for i = width(coefInputAudio)
110             Ci = C+"_"+num2str(i);
111             hold on;
112             eval("plot(coef"+C+"(x"+Ci+","+num2str(i)+"),
                    'k')");
113             hold on;
114             eval("plot(coefInputAudio(y"+Ci+",i),'b');");
115             hold on;
116         end
117         title("DTW of Input Audio vs. " + C);
118     end
119 end
120
121 %% Calculate sum of distances
122
123 for c = words
124     for k = 1:numTrain
125         C = c+num2str(k);
126         eval(C+"sum=0;");
127         for i = 1:width(coefInputAudio)
128             eval(C+"sum="+C+"sum + dist"+C+"_"+num2str(i)
                    +";");
129         end
130     end
131 end
```

```matlab
132
133    for  c  =  words
134         for  k  =  1:numTrain
135              C  =  c+num2str(k);
136              eval(C+"sum");
137         end
138    end
```

# References

1. Bhadragiri, J.M., and Ramesh, B.N., *Speech Recognition using MFCC and DTW*, VIT University, Vellore, India, 2014

2. Kare, C.B., and Navale, V.S., *Speech recognition by Dynamic Time Warping*, IOSR Journal of Electronics and Communication Engineering, Pune, 2015

3. Kesarkar, M.P., *Feature Extraction for Speech Recognition*, IIT Bombay, Mumbai, India, 2003

4. Narang, S., and Gupta, D., *Speech Feature Extraction Techniques: A Review*, Amity University, Noida, India, 2015

5. Rao, K.S., and Manjunath K.E., *Speech Recognition Using Articulatory and Excitation Source Features*, 2017