# In the name of God

## Voice Recognition Systems

Mansour Saffar Mehrjardi     Ahmad Tavakoli

Voice Recognition systems (VR systems) are widely used in many new technologies. Robust speech recognition systems are used in military applications. For instance in F11 jets these systems recognize pilot's voice for authentication of their instruction. Car production companies are using this technology in their cars so the driver can steer the car only with speaking to the car! Hollywood has also shown the future of speech recognition systems in its movies. We see in the Oscar winner movie "Her" that these systems are used in an interactive super intelligent operating system in order to speak to the computer.

In this project we tried to develop an isolated speaker-dependent voice recognition system using MATLAB software. Isolated voice recognition systems should acquire their data with pause which means every time you want to record voice there must be a pause between every word pronounced. Speaker-dependent voice recognition systems only can recognize one particular's voice whereas voice-independent voice recognition systems can recognize everyone's voice independent of the speaker.

In this project we implemented an isolated speaker-dependent voice recognition system using MATLAB. In Speech signal processing algorithms first we extract information from the signal using several methods and then give these data to a data clustering system which recognizes specific patterns within these data and can categorize these data for further recognition. Among these recognition and data clustering algorithms are neural networks which we used in this project. Artificial Neural networks are inspired by the way biological nervous systems work for processing information. It is composed of a large number of highly interconnected processing elements (neurons) working in unison to solve specific problems. ANNs, like people, learn by example. An ANN is configured for a

specific application, such as pattern recognition or data classification, through a learning process.
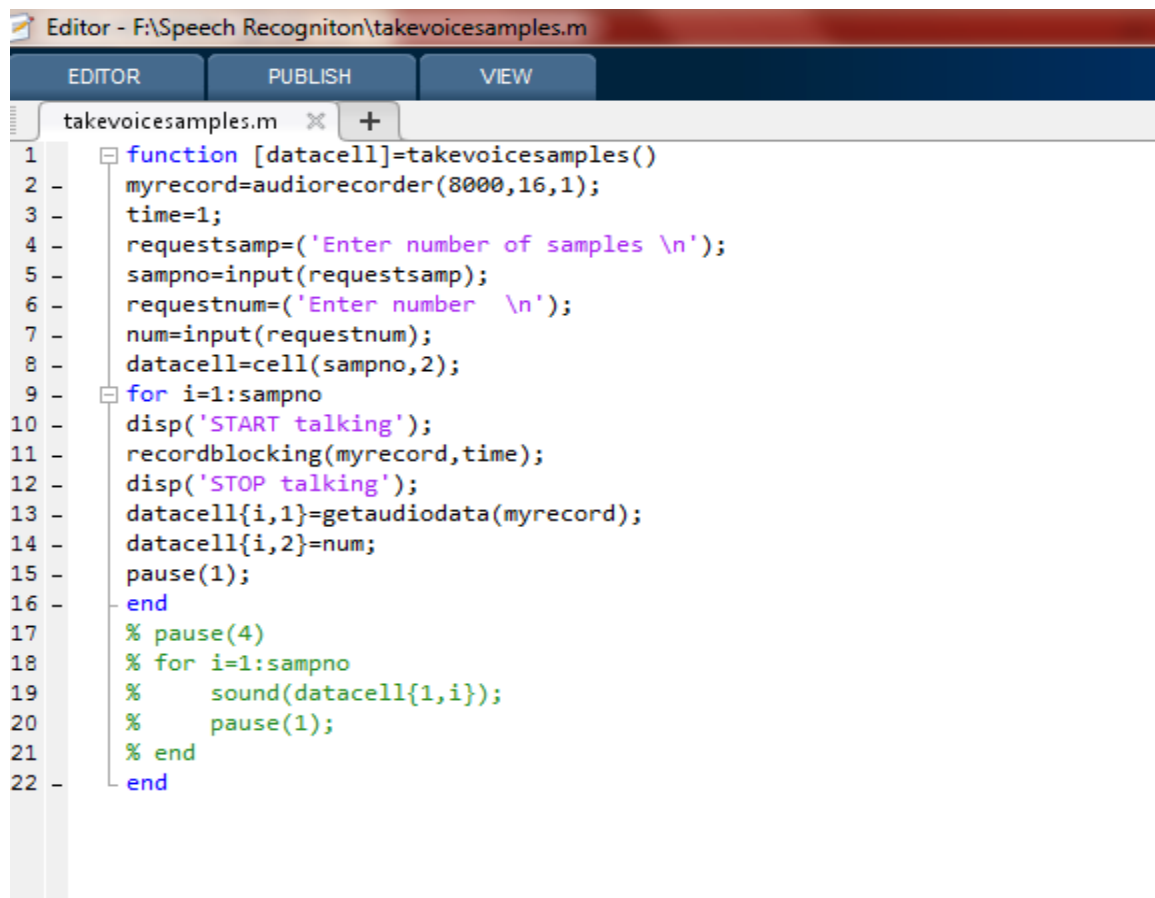
The process of this project is as below:

**Data acquisition**:

I recorded my voice 500 times saying numbers one to ten (each number 50 times). Each recording lasts 1 second and the sampling frequency is 8000Hz. Each recording is the filtered by a high pass IIR filter which is given by the function $H(z) = 1 - 0.98z$. This filter reduces low frequency noises in speech and also amplifies frequencies that is used in.

The magnitude response of this filter is shown in figure 1.1.

The MATLAB function which I used in this project is shown in figure 1.2.

```
Editor - F:\Speech Recogniton\takevoicesamples.m

    EDITOR          PUBLISH          VIEW

    takevoicesamples.m   ✕   +
 1      function [datacell]=takevoicesamples()
 2 -      myrecord=audiorecorder(8000,16,1);
 3 -      time=1;
 4 -      requestsamp=('Enter number of samples \n');
 5 -      sampno=input(requestsamp);
 6 -      requestnum=('Enter number  \n');
 7 -      num=input(requestnum);
 8 -      datacell=cell(sampno,2);
 9 -      for i=1:sampno
10 -      disp('START talking');
11 -      recordblocking(myrecord,time);
12 -      disp('STOP talking');
13 -      datacell{i,1}=getaudiodata(myrecord);
14 -      datacell{i,2}=num;
15 -      pause(1);
16 -      end
17      % pause(4)
18      % for i=1:sampno
19      %      sound(datacell{1,i});
20      %      pause(1);
21      % end
22 -      end
```

Figure 1.2

The function takevoicesamples() stores 1second voice samples in a cell. Samples acquired by this function is the applied to 3 different functions which extract specific information from samples.

**Data extraction**:

MFCC (Mel frequency cepstrum coefficient) is a representation of the short-term power spectrum of a sound, based on a linear cosine transform of a log power spectrum on a nonlinear mel scale of frequency.

MFCCs are commonly derived using these steps:

I. Pre-emphasize the signal using the filter such as figure 1.1
II. Take short-time Fourier transform of the windowed signal. The window is normally hamming window.
III. Compute magnitude cepstrum and then apply the signal to a filter bank. This filter bank is triangular and is equally spaced in mel frequency domain.
IV. Energy of each filter bank part is calculated.
V. These energies are applied to a discrete cosine transform and the MFCCs are the amplitude of the resulting spectrum.

**STFT (short-time Fourier transform)**

Short time Fourier transform gives us signal's information both in time domain and also the frequency domain. Speech signal are non-stationary signals which means the amplitude of a certain frequency can vary with time, so we need to know which time does a frequency appears and when it vanishes. In STFT first we window the original signal, using hamming window which smooths signal fluctuations. Then we take FFT of the windowed signal. Then we go further in signal and according to hop size of SFTF which is actually the value of pace we take when going further in the signal and then windowing that part of signal.
If we choose the length of the window to be large then we can discern adjacent frequencies in the signal but processing of signal will take much more time. In contrast to what said before, if we choose the length of window to be short it reduces processing time but also reduces resolution of resulting STFT. Note that

hop size or simply the pace we choose, can also affect processing time and resolution of resulting STFT. To demonstrate what I said, figure 1.3 shows STFT of 5 seconds excerpt of a music sampled with frequency of 8KHz. Length of original signal is 40000. In figure 1.3 window length is 50 and hop size is also 50, so we get a high resolution STFT.
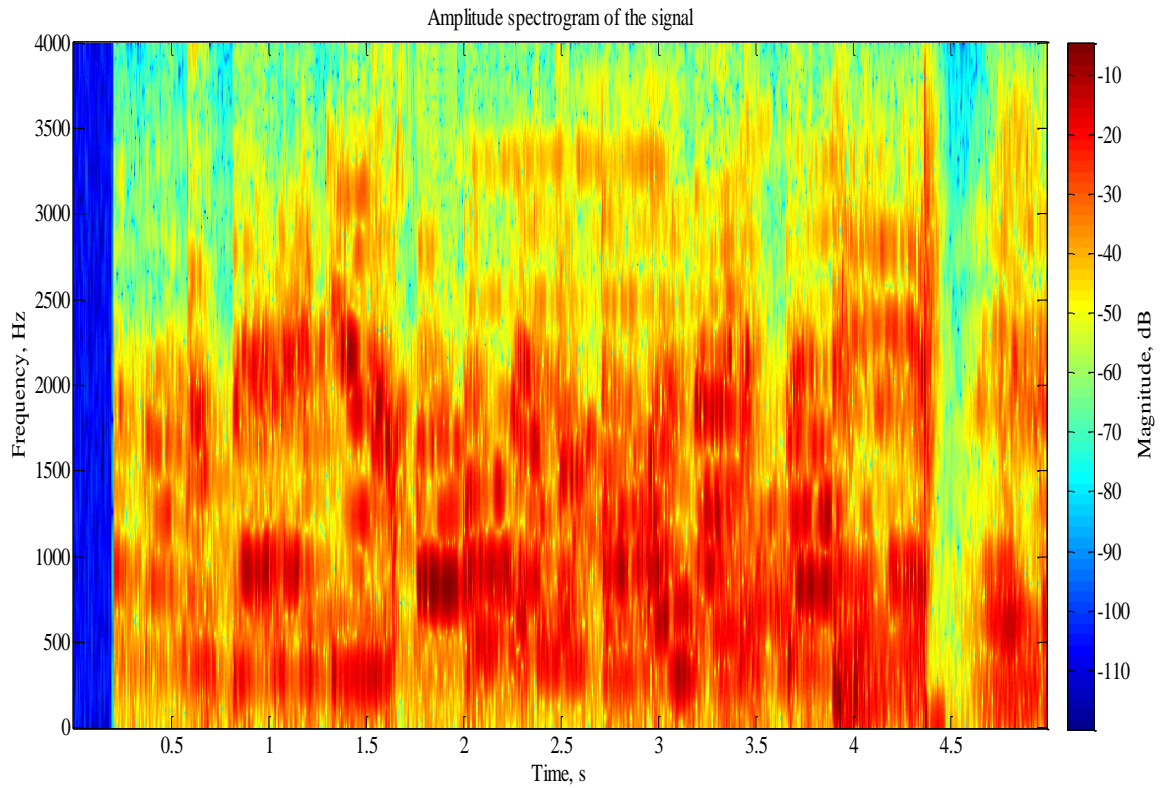


Figure 1.3

Note that horizontal axis is time (in seconds) and vertical axis is frequency (in Hertz). The color bar right to the signal shows the magnitude of signal in dB.

Figure 1.4 shows STFT of the previous signal in which window length is 200 and hop size is also 200. As you can see it has considerably lower resolution than Figure 1.3 but it takes less time to achieve this plot.
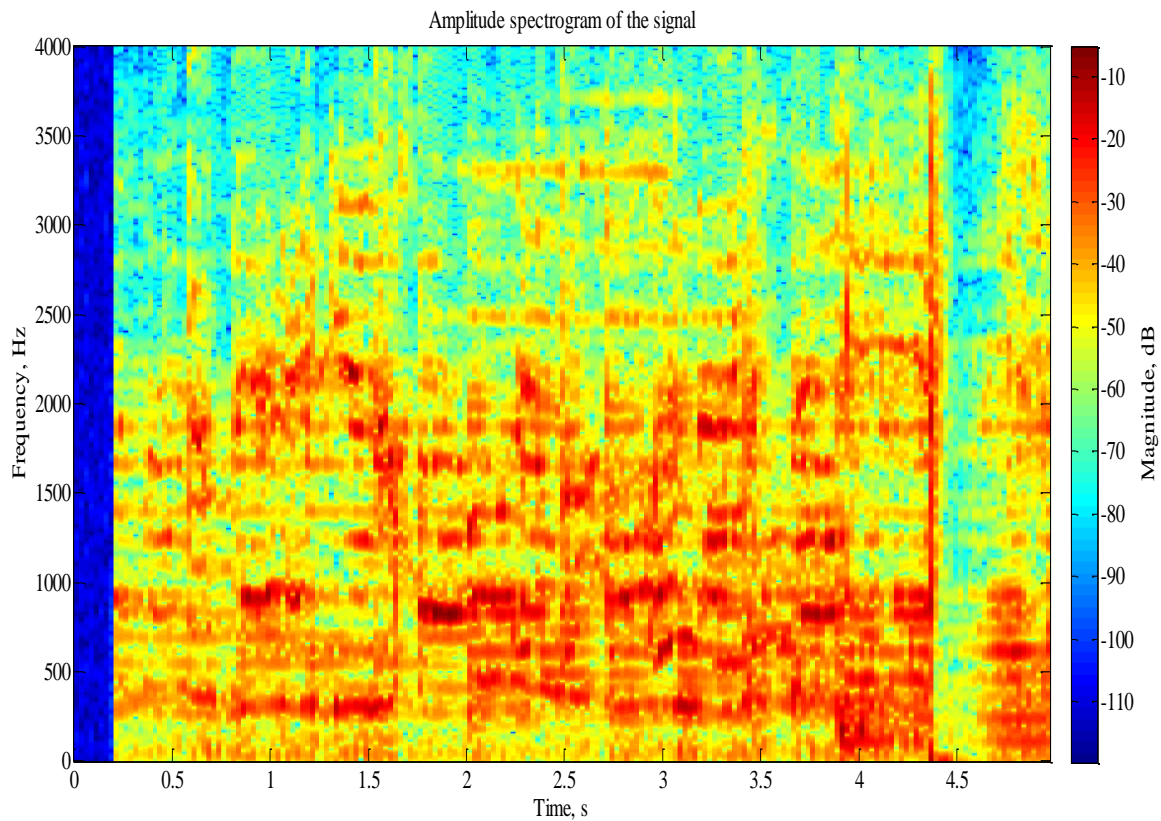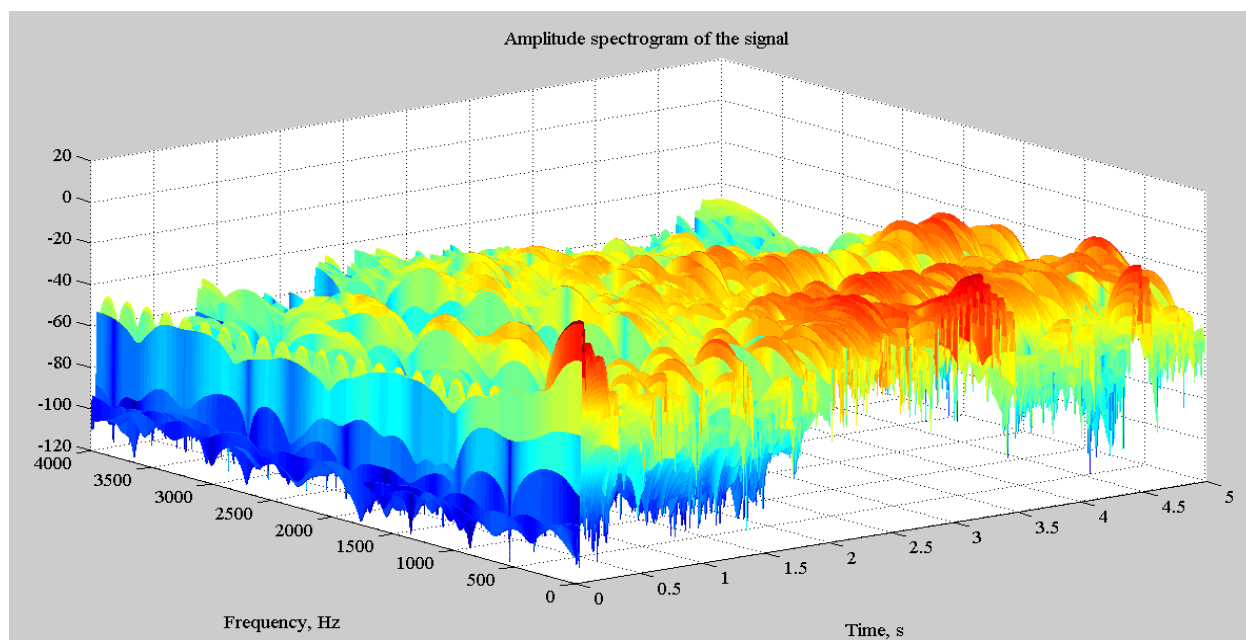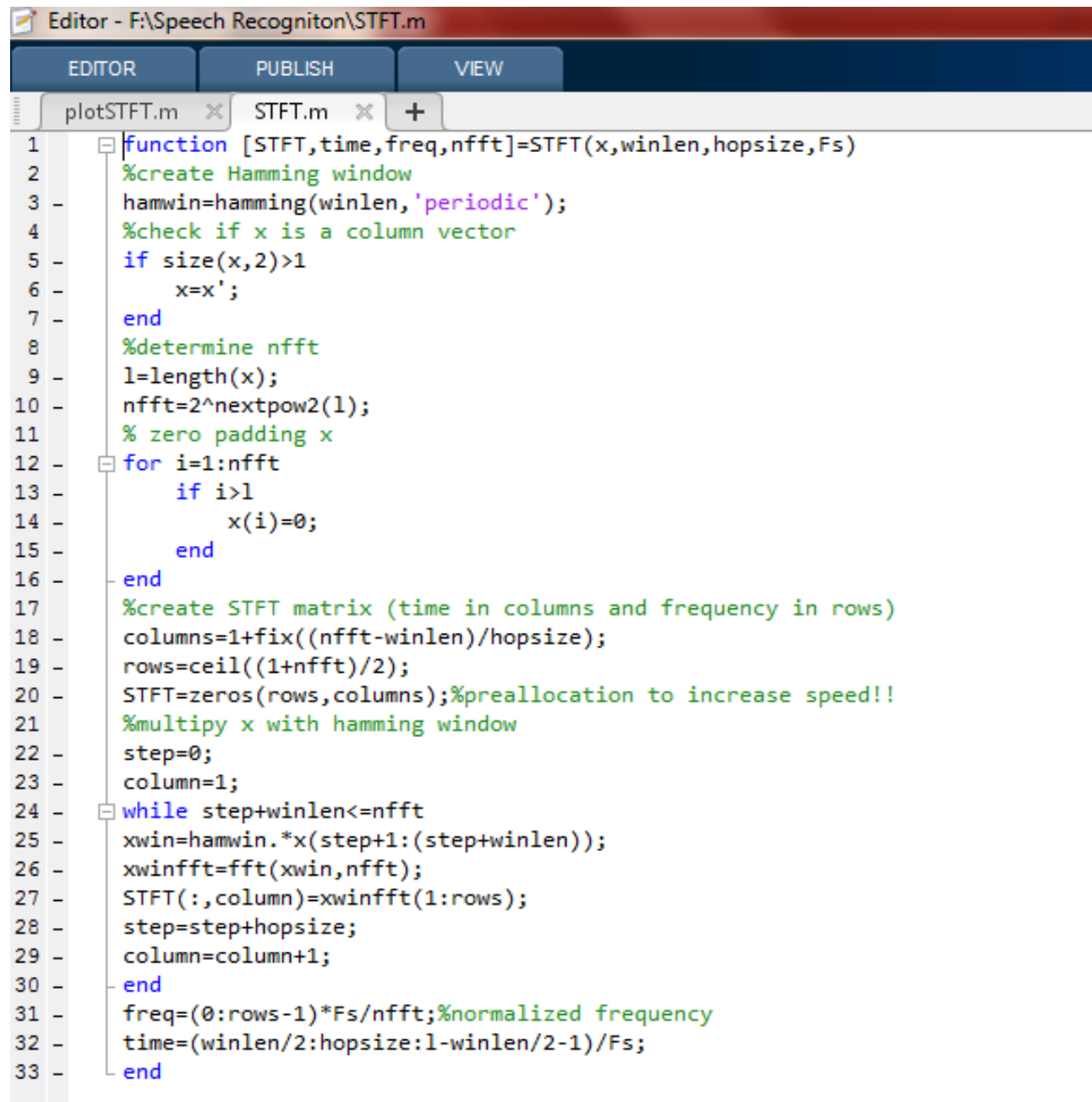
Figure 1.4



Figure 1.5

STFT can also be shown using mesh in MATLAB which gives 3-D plot of STFT as shown in Figure 1.5.

MATLAB code for SFTF is shown in Figure 1.6 and Figure 1.7 shows MATLAB code for plotting STFT of a signal.

```matlab
Editor - F:\Speech Recogniton\STFT.m
EDITOR        PUBLISH        VIEW
plotSTFT.m  X   STFT.m  X   +
1    function [STFT,time,freq,nfft]=STFT(x,winlen,hopsize,Fs)
2    %create Hamming window
3    hamwin=hamming(winlen,'periodic');
4    %check if x is a column vector
5    if size(x,2)>1
6        x=x';
7    end
8    %determine nfft
9    l=length(x);
10   nfft=2^nextpow2(l);
11   % zero padding x
12   for i=1:nfft
13       if i>l
14           x(i)=0;
15       end
16   end
17   %create STFT matrix (time in columns and frequency in rows)
18   columns=1+fix((nfft-winlen)/hopsize);
19   rows=ceil((1+nfft)/2);
20   STFT=zeros(rows,columns);%preallocation to increase speed!!
21   %multipy x with hamming window
22   step=0;
23   column=1;
24   while step+winlen<=nfft
25   xwin=hamwin.*x(step+1:(step+winlen));
26   xwinfft=fft(xwin,nfft);
27   STFT(:,column)=xwinfft(1:rows);
28   step=step+hopsize;
29   column=column+1;
30   end
31   freq=(0:rows-1)*Fs/nfft;%normalized frequency
32   time=(winlen/2:hopsize:l-winlen/2-1)/Fs;
33   end
```

Figure 1.6

EDITOR        PUBLISH        VIEW

plotSTFT.m   ✕   STFT.m   ✕   +

```matlab
1    function plotSTFT(x,winlen,hopsize,Fs)
2    xmax=max(abs(x));
3    l=length(x);
4    scaledx=x/xmax;
5    K = sum(hamming(winlen, 'periodic'))/winlen;
6    [stft,t,f,nfft]=STFT(scaledx,winlen,hopsize,Fs);
7    stft = abs(stft)/winlen/K;
8    stft=stft(:,1:(ceil((l-winlen)/hopsize)));%HEll YEAH!!!!
9    if rem(nfft, 2)                        % odd nfft excludes Nyquist point
10       stft(2:end, :) = stft(2:end, :).*2;
11   else                                   % even nfft includes Nyquist point
12       stft(2:end-1, :) = stft(2:end-1, :).*2;
13   end
14   stft = 20*log10(stft + 1e-6);
15   imagesc(t, f, stft);
16   set(gca,'YDir','normal');
17   set(gca, 'FontName', 'Times New Roman', 'FontSize', 14);
18   xlabel('Time, s');
19   ylabel('Frequency, Hz');
20   title('Amplitude spectrogram of the signal');
21   handl = colorbar;
22   set(handl, 'FontName', 'Times New Roman', 'FontSize', 14);
23   ylabel(handl, 'Magnitude, dB');
24   % maxstft=max(abs(stft));
25   % valstft=abs(stft)/maxstft;
26   % [m,n]=size(valstft);
27   % colormap=zeros(m,n,3);
28   % for i=1:m
29   %     for j=1:n
30   %         for k=1:2
31   %     colormap(i,j,k)=.5;
32   %         end
33   %     end
34   % end
35   % mesh(t,f,stft);
36   % set(gca,'YDir','normal');
37   % set(gca, 'FontName', 'Times New Roman', 'FontSize', 14);
38   % xlabel('Time, s');
39   % ylabel('Frequency, Hz');
40   % title('Amplitude spectrogram of the signal');
```

Figure 1.7

**Filterbanks**:

In signal processing, a filter bank is an array of band-pass filters that separates the input signal into multiple components, each one carrying a single frequency sub-band of the original signal. Mel frequency filter bank is a triangular filter banks which are equally spaced in Mel frequency domain. The relation between Mel frequency and conventional frequency measured in Hz is as Equ1.1. The filterbank used in this project is shown in Figure 1.8.

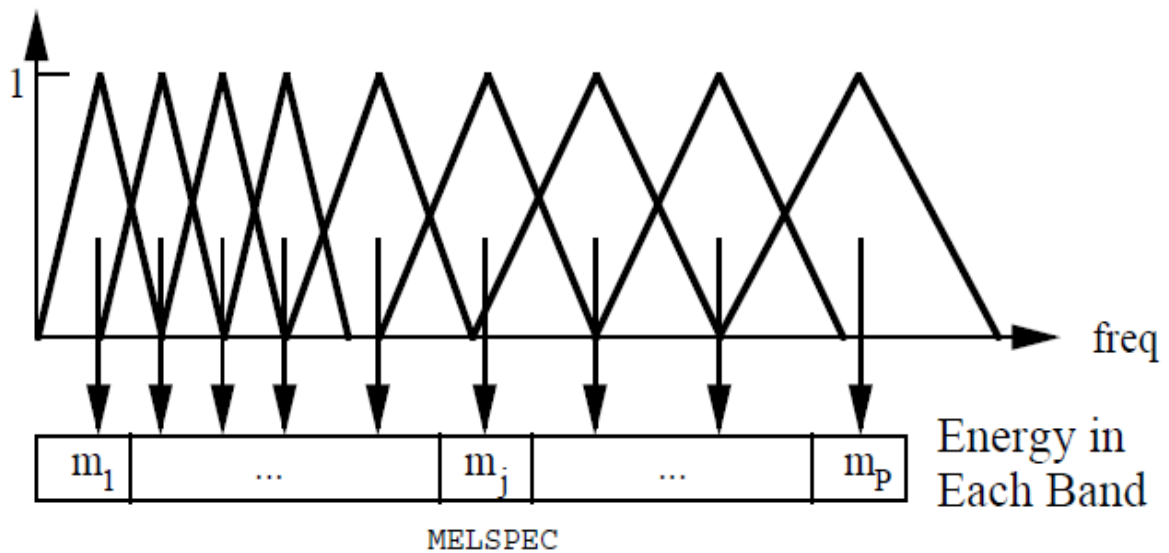$$\text{Equ 1.1} : Mel(f) = 2595 \log_{10}(1 + \frac{f}{700})$$



Figure 1.8

MATLAB code to create Mel frequency filter bank is shown in Figure 1.9. Note that **hz2mel** and **mel2hz** are function handles that are defined according to Equ1.1.

```
Editor - F:\Speech Recogniton\Trifilterbank.m

    EDITOR          PUBLISH          VIEW

  Trifilterbank.m  ✕   +
 1     ⊟ function [Fbank,freq,cfreq]=Trifilterbank(fnumber,flength,range,fs,hz2mel,mel2hz)
 2 -       f_min=0;
 3 -       f_low=range(1);
 4 -       f_high=range(2);
 5 -       f_max=0.5*fs;
 6 -       freq=linspace(f_min,f_max,flength);
 7 -       cfreq=mel2hz(hz2mel(f_low)+[0:fnumber+1]*((hz2mel(f_high)-hz2mel(f_low))/(fnumber+1)));
 8 -       Fbank=zeros(fnumber,flength);
 9 -    ⊟ for i=1:fnumber
10 -           j=freq>=cfreq(i)&freq<=cfreq(i+1); % up-slope
11 -           Fbank(i,j)=(freq(j)-cfreq(i))/(cfreq(i+1)-cfreq(i));
12 -           j = freq>=cfreq(i+1)&freq<=cfreq(i+2); % down-slope
13 -           Fbank(i,j)=(cfreq(i+2)-freq(j))/(cfreq(i+2)-cfreq(i+1));
14 -     ⊢ end
15       % for k=1:fnumber
16       %      plot(freq,Fbank(k,:));
17       %      hold on
18       % end
19 -     ⊢ end
```

Figure 1.9

**Discrete cosine transform (DCT)**:

Discrete cosine transform expresses a signal in terms of cosine function. Its application in signal processing is for signal compaction. In this project we use DCT2 formula shown in Equ1.2.

Equ 1.2: $C(i) = \sqrt{\dfrac{2}{N}} \sum_{j=1}^{N} \cos(\dfrac{\pi i}{N}(j - 0.5))$

In Equ1.2 N is the number of filter banks used filter bank design.

**MFCC extraction**:

MATLAB code to extract MFCC from voice is shown in Figure1.10.

```matlab
1     function [MFCC,FBE]=extractMFCC(signal,fs,Frametime,Frameshift,alpha,range,fnumber,cepcoef,lift)
2     signal=signal(:,1);%channel one!!!
3     l=length(signal);
4     %Handle functions
5     hz2mel = @( hz )( 1127*log(1+hz/700) );      % Hertz to mel
6     mel2hz = @( mel )( 700*exp(mel/1127)-700 ); % mel to Hertz
7     dctm = @( N, M )( sqrt(2.0/M) * cos( repmat([0:N-1].',1,M) ...
8                                        .* repmat(pi*([1:M]-0.5)/M,N,1) ) );
9     ceplifter = @(N,L)( 1+0.5*L*sin(pi*[0:N-1]/L) );
10    %First filtering signal with an FIR filter
11    fsignal=filter([1,-alpha],1,signal);
12    %Applying signal to short-time fourier transform
13    %first determining winlen and hopsize
14    winlen=fix((l/(Frametime*fs/1000)));
15    hopsize=Frameshift*fs/1000;
16    [stft,~,~,nfft]=STFT(fsignal,winlen,hopsize,fs);
17    K = sum(hamming(winlen, 'periodic'))/winlen;
18    stft=abs(stft)/winlen/K;
19    stft=stft(:,1:(ceil((l-winlen)/hopsize)));
20    if rem(nfft, 2)                          % odd nfft excludes Nyquist point
21        stft(2:end, :) = stft(2:end, :).*2;
22    else                                     % even nfft includes Nyquist point
23        stft(2:end-1, :) = stft(2:end-1, :).*2;
24    end
25    % Trianngular filter bank
26    flength=nfft/2+1;
27    [Fbank,~,~]=Trifilterbank(fnumber,flength,range,fs,hz2mel,mel2hz);
28    % filter bank energy
29    FBE=Fbank*stft;
30    %DCT computation
31    DCT=dctm(cepcoef,fnumber);
32    % Conversion of logFBEs to cepstral coefficients through DCT
33    CC=DCT*log(FBE);
34    % Cepstral lifter computation
35    lifter=ceplifter(cepcoef,lift);
36    % Cepstral liftering gives liftered cepstral coefficients
37    MFCC=diag(lifter)*CC; % ~ HTK's MFCCs
38    end
```
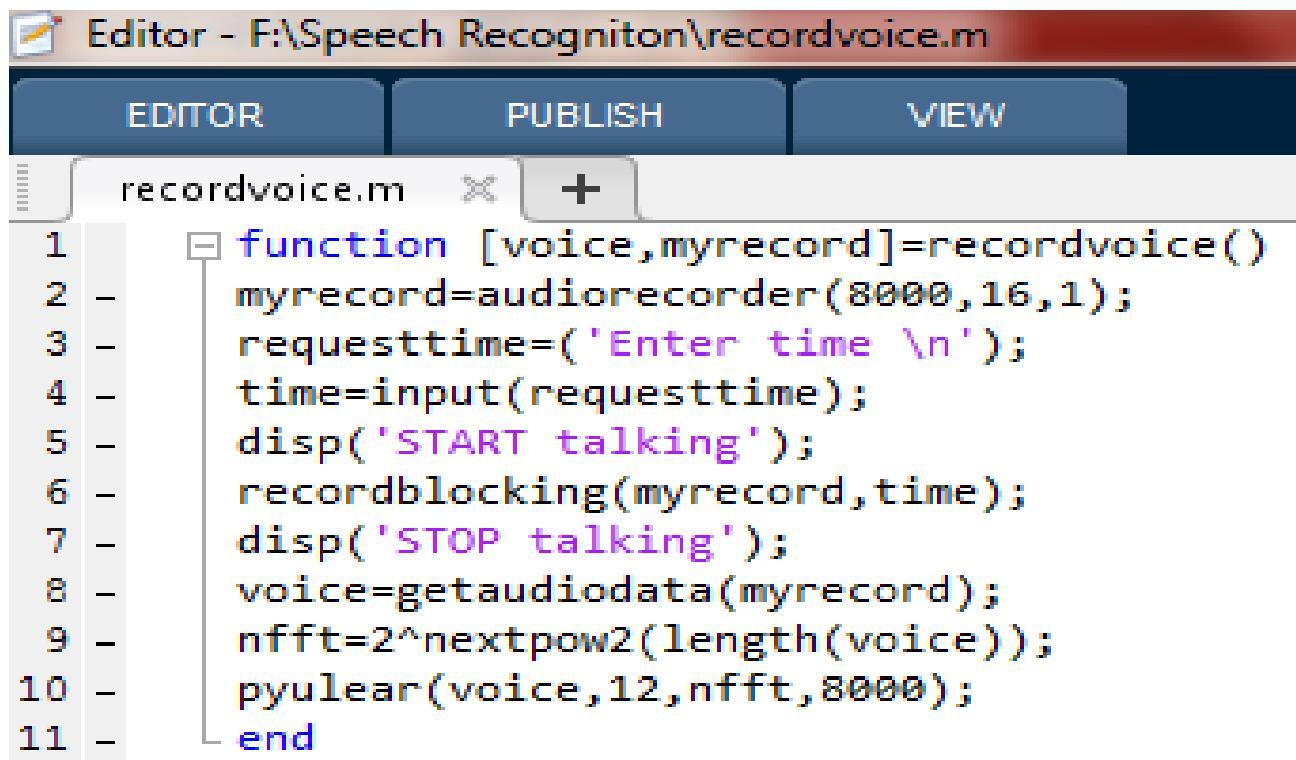
Figure 1.10

**PSD(Power spectra density):**

Another way to extract information from speech signal is using PSD or power spectrum density. PSD shows the average power of signal in a frequency band. A useful and convenient way to acquire PSD of signal is using Yule-Walker method to estimate PSD of speech signal. MATLAB has the function **pyulear** which can give us PSD estimation. In this project, in addition to MFCC feature extraction, we also give PSD of each sample voice to ANN for clustering and recognition.

Note that order of **pyulear** in MATLAB is 12 for applications of speech processing. MATLAB code to acquire Yule-Walker PSD estimation of a signal is shown in Figure 1.11.



```matlab
function [voice,myrecord]=recordvoice()
myrecord=audiorecorder(8000,16,1);
requesttime=('Enter time \n');
time=input(requesttime);
disp('START talking');
recordblocking(myrecord,time);
disp('STOP talking');
voice=getaudiodata(myrecord);
nfft=2^nextpow2(length(voice));
pyulear(voice,12,nfft,8000);
end
```

Figure 1.11

**Voice Recognition**:

After extracting information from speech signal using MFCC and Yule-Walker PSD estimation, we give these data to neural network. We assess both methods and compare them.

MATLAB code for ANN voice recognition based on MFCC is shown in Figure1.12.

```matlab
MFCCneuralclassification.m  ×  +
1     function [result,myNet,tr]=MFCCneuralclassification(MFCCdata)
2         myNet = feedforwardnet(10 , 'trainlm');
3         m=500;
4         % myNet.trainparam.epochs=5000;
5         dataClass = zeros(5,m);
6         for i=1:m
7             if (MFCCdata{i,1} == 1)
8                 dataClass(1:5,i) = [1;0;0;0;0];
9             elseif(MFCCdata{i,1} == 2)
10                dataClass(1:5,i) = [0;1;0;0;0];
11            elseif(MFCCdata{i,1} == 3)
12                dataClass(1:5,i) = [0;0;1;0;0];
13            elseif(MFCCdata{i,1} == 4)
14                dataClass(1:5,i) = [0;0;0;1;0];
15            elseif(MFCCdata{i,1} == 5)
16                dataClass(1:5,i) = [0;0;0;0;1];
17            end
18        end
19        trainRand=randperm( m,floor(7*m/10) ) ;
20        traindata=MFCCdata(trainRand,:);
21        dataClass1 = dataClass(:,trainRand);
22        testdata=MFCCdata;
23        testdata(trainRand,:)=[];
24        testClass = dataClass ;
25        testClass(:,trainRand) = [] ;
26        [m1,~] = size(traindata);
27        trainvoice = zeros(1287,m1);
28        trainoutput=zeros(5,m1);
29        for i=1:m1
30            trainvoice(1:1287,i) = reshape(traindata{i,2},1287,1);
31            trainoutput(1:5,i) = dataClass1(1:5,i);
32        end
33
34        tic;
35        [myNet,tr] = train(myNet,trainvoice,trainoutput);
36        toc;
37        view(myNet);
38        %%Evaluation of ANN
39        [m2,~] = size(testdata);
40        testvoice = zeros(1287,m2);
```
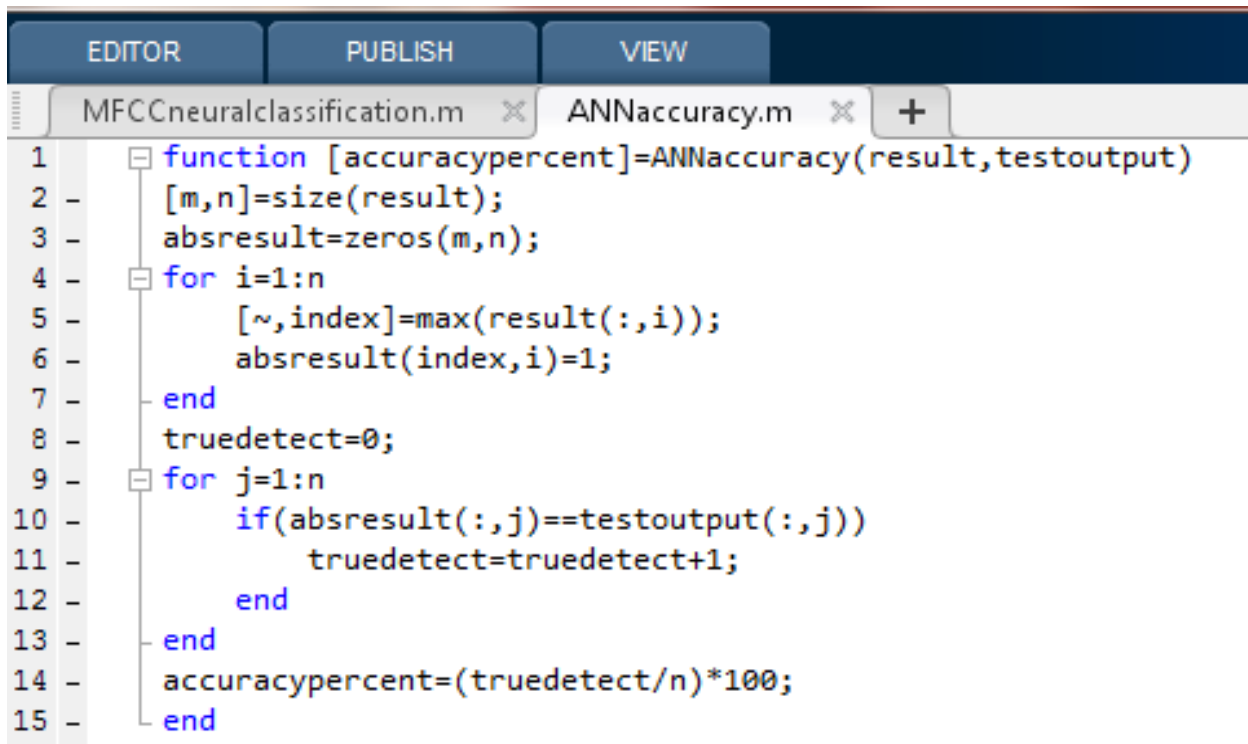
```
41 -    testoutput=zeros(5,m2);
42 - ┌ for i=1:m2
43 -    testvoice(1:1287,i) = reshape(testdata{i,2},1287,1);
44 -    testoutput(1:5,i)=testClass(1:5,i);
45 - ├ end
46 -    result=myNet(testvoice);
47 -    save('var.mat','testoutput','result');
48 - └ end
```

Figure 1.12

In this ANN, number of epochs is considered 5000 epochs but since the training method is trainlm, which means in every epoch all neuron coefficients are updated, training normally takes 10 to 14 epochs. Number of layers is 1 and there are 10 neurons in the code above.

For determining the accuracy of ANN we need to compare result of ANN with the test data that we used for testing ANN. MATLAB code for accuracy assessment of ANN is shown in Figure1.13.

```
                EDITOR          PUBLISH          VIEW

   MFCCneuralclassification.m  ✕   ANNaccuracy.m  ✕   +
 1       ┌ function [accuracypercent]=ANNaccuracy(result,testoutput)
 2 -       [m,n]=size(result);
 3 -       absresult=zeros(m,n);
 4 -     ┌ for i=1:n
 5 -           [~,index]=max(result(:,i));
 6 -           absresult(index,i)=1;
 7 -     ├ end
 8 -       truedetect=0;
 9 -     ┌ for j=1:n
10 -           if(absresult(:,j)==testoutput(:,j))
11 -               truedetect=truedetect+1;
12 -           end
13 -     ├ end
14 -       accuracypercent=(truedetect/n)*100;
15 -     └ end
```

Figure 1.13

We also train the ANN using PSD estimation data. MALAB code for this ANN is shown in Figure1.14.

```matlab
function [result,myNet,tr]=PSDneuralclassification(voicedata)
for i=1:500
  voicedata{i,1}=filter([1,-.97],1,voicedata{i,1});
end
[m,n]=size(voicedata);
myNet = feedforwardnet(10 , 'trainb');
% myNet.trainparam.epochs=5000
dataClass=zeros(5,m);
for i=1:m
    if (voicedata{i,2} == 1)
        dataClass(1:5,i) = [1;0;0;0;0];
      elseif(voicedata{i,2} == 2)
        dataClass(1:5,i) = [0;1;0;0;0];
      elseif(voicedata{i,2} == 3)
        dataClass(1:5,i) = [0;0;1;0;0];
      elseif(voicedata{i,2} == 4)
        dataClass(1:5,i) = [0;0;0;1;0];
      elseif(voicedata{i,2} == 5)
        dataClass(1:5,i) = [0;0;0;0;1];
    end
end
PSDdata=cell(m,n);
nfft=8192;
for i=1:m
    PSDdata{i,2}=pyulear(voicedata{i,1},12,nfft,8000);
    PSDdata{i,1}=voicedata{i,2};
end
trainRand=randperm( m,floor(7*m/10) ) ;
traindata=PSDdata(trainRand,:);
dataClass1 = dataClass(:,trainRand);
testdata=PSDdata;
testdata(trainRand,:)=[];
testClass = dataClass ;
testClass(:,trainRand) = [] ;
[m1,~] = size(traindata);
trainvoice = zeros(4097,m1);
trainoutput = zeros(5,m1);
for i=1:m1
    trainvoice(1:4097,i) = reshape(traindata{i,2},4097,1);
    trainoutput(1:5,i) = dataClass1(1:5,i);
```

```
41 -      end
42 -      tic;
43 -      [myNet,tr] = train(myNet,trainvoice,trainoutput);
44 -      toc;
45 -      view(myNet);
46        %%% ANN assessment
47 -      [m2,~] = size(testdata);
48 -      testvoice = zeros(4097,m2);
49 -      testoutput=zeros(5,m2);
50 -  ⊟  for i=1:m2
51 -      testvoice(1:4097,i) = reshape(testdata{i,2},4097,1);
52 -      testoutput(1:5,i)=testClass(1:5,i);
53 -      end
54 -      result=myNet(testvoice);
55 -      save('var.mat','testoutput','result');
56 -      end
```

Figure 1.20

**Deduction:**

We tested both ANNs with voice data and the result showed that MFCC based trained ANNs work considerably more precise than PSD based trained ANNs. When using MFCC based training for ANN, we used both trainb and trainlm method for training ANN. The trainb method trains the ANN more rapidly than trainlm method does but in some cases it has less accuracy than trainlm method. In MFCC based training both trainb and trainlm method lead to an astonishing 97% accuracy when they're tested, so we can surely say that MFCC based voice recognition is significantly more robust and precise than PSD based Voice recognition.

**References:**

I.    Mathworks.com
II.   HTK book by Steve Young et al. Cambridge University Engineering Department
III. htk.eng.cam.ac.uk
IV. Discrete time signal processing, Alan V Oppenheim. 3$^{rd}$ edition.