

Speech Recognition System using the DTW Method

Anirudha Bhat Nekkare (211EC268)

*Department of Electronics and Communication Engineering
National Institute of Technology Karnataka
Surathkal, India
anirudh.211ec213@nitk.edu.in*

Savio Kuttikal Santhosh (211EC247)

*Department of Electronics and Communication Engineering
National Institute of Technology Karnataka
Surathkal, India
saviosanthosh.211ec247@nitk.edu.in*

Abstract—Speech recognition is a fundamental component of modern human-computer interaction systems. Dynamic Time Warping (DTW) has emerged as a powerful technique for speech recognition, especially in scenarios where traditional methods like Hidden Markov Models (HMMs) face challenges such as variability in speech patterns and non-linear temporal alignments. This paper presents a speech recognition system based on DTW, focusing on its effectiveness in handling diverse speech inputs. The system incorporates pre-processing techniques for feature extraction, such as Mel Frequency Cepstral Coefficients (MFCCs), to represent speech signals effectively. DTW is then applied on the test dataset to compare these feature representations with reference templates, enabling robust recognition even in noisy environments or with limited training data. The scalability and adaptability of DTW make it a promising approach for real-world applications requiring accurate and robust speech recognition capabilities.

Index Terms—Feature Extraction, Feature Matching, Mel Frequency Cepstral Coefficients (MFCC), Dynamic Time Warping (DTW), Discrete Cosine Transform (DCT)

I. INTRODUCTION

Speech recognition plays a pivotal role in various applications, from virtual assistants to hands-free control systems. Traditional approaches to speech recognition, such as Hidden Markov Models (HMMs), have been widely used but often encounter challenges with variability in speech patterns and non-linear temporal alignments. Dynamic Time Warping (DTW) has emerged as a robust alternative, particularly effective in scenarios where precise alignment of speech sequences is crucial. By allowing flexible matching of temporal sequences, DTW can accommodate variations in speech speed and timing, making it suitable for real-world applications with diverse speech inputs.

In this paper, we present a speech recognition system leveraging DTW as the core technique. The system is designed to handle a wide range of speech inputs while maintaining high accuracy and reliability. Key components of the system include pre-processing techniques like Mel Frequency Cepstral Coefficients (MFCCs) for feature extraction, which provide a compact yet informative representation of speech signals. DTW is then applied to compare these feature vectors with reference templates, enabling robust recognition even in noisy environments or with limited training data.

The rest of this paper is organized as follows: Section II details the implementation of our DTW-based speech recognition system, including data preprocessing, feature extraction, and

the DTW algorithm itself. Section III presents experimental results and performance evaluations on the dataset, demonstrating the efficacy of our approach.

II. IMPLEMENTATION

A. Dataset

The dataset comprises forty recordings, four recordings of each digit from 0-9. For training and testing purposes, we partitioned the dataset into standard training and testing datasets, with the training dataset consisting of thirty recordings and the remaining ten going to the testing dataset. The goal is to develop an automatic speech recognition system capable of accurately identifying digits irrespective of speech flow, loudness, intonation and intensity of overtones.

B. Data pre-processing

Before training the model, other tasks such as data loading and preprocessing tasks, including reading data from WAV files and dividing the data into training test sets were handled.

All speech signals were normalized to have zero mean and unit variance. This step ensured consistent signal amplitude across recordings, reducing the impact of amplitude variations on feature extraction and recognition.

The normalized speech signals were segmented into frames using a fixed-length windowing technique. Each frame typically spanned 20-30 milliseconds, with a 10-millisecond overlap between consecutive frames. This frame segmentation process facilitated the extraction of time-varying features from the speech signals.

Mel Frequency Cepstral Coefficients (MFCCs) were extracted from each frame to capture spectral characteristics and temporal dynamics of the speech signals. The MFCCs were computed using a standard feature extraction pipeline, including the application of a Mel filterbank, logarithmic scaling, Discrete Cosine Transform, and liftering to emphasize relevant frequency components.

The extracted MFCC features were normalized to have zero mean and unit variance across frames within each utterance. This normalization step enhanced the robustness of the feature representation by mitigating variations in speech intensity and loudness.

Categorical classification accuracy was employed as the performance metric for evaluating the model.

C. Methodology

Reference templates for each digit are constructed using a training dataset comprising examples of each digit spoken by multiple speakers. These templates capture the expected patterns for each digit, accounting for variations in speech due to different speakers, speaking rates, and accents.

DTW is a dynamic programming algorithm that measures the similarity between two sequences with varying lengths and alignments. It finds the optimal alignment between a query sequence (input speech signal) and a reference sequence (digit template) by minimizing the total distance or cost between corresponding elements. The algorithm works by creating a matrix where each cell represents the cumulative distance between corresponding elements of the two sequences at that point. The goal is to find the path through this matrix that minimizes the total distance. DTW allows for flexible matching of sequences by considering local alignments and allowing for warping or stretching of the sequences in the time domain. This flexibility is crucial for handling variations in speech timing and duration.

We recognized a spoken digit by preprocessing the input speech signal to extract MFCC feature vectors, computing the DTW distance between the input feature sequence and each reference template (one for each digit), and selecting the digit with the lowest DTW distance as the recognized digit, indicating the closest match between the input speech and the reference templates.

III. RESULTS AND DISCUSSIONS

The performance of the proposed Dynamic Time Warping (DTW)-based digit recognition system was evaluated using a small dataset comprising spoken digits (0-9) across varying conditions.

The system achieved a remarkable accuracy of 100%, indicating perfect recognition of spoken digits across all test samples. Precision, recall, and F1 score values of 1.0 were also obtained for each digit class, demonstrating the system's ability to correctly classify all instances of each digit without any false positives or false negatives.

REFERENCES

- [1] Senin, Pavel. (2009). Dynamic Time Warping Algorithm Review.
- [2] Majeed, Sayf & HUSAIN, HAFIZAH & Samad, Salina & Idbeaa, Tarik. (2015). Mel frequency cepstral coefficients (Mfcc) feature extraction enhancement in the application of speech recognition: A comparison study. *Journal of Theoretical and Applied Information Technology*. 79. 38-56.
- [3] Z. K. Abdul and A. K. Al-Talabani, "Mel Frequency Cepstral Coefficient and its Applications: A Review," in *IEEE Access*, vol. 10, pp. 122136-122158, 2022, doi: 10.1109/ACCESS.2022.3223444.
- [4] Bhadrageiri Jagan Mohan and Ramesh Babu N., "Speech recognition using MFCC and DTW," 2014 International Conference on Advances in Electrical Engineering (ICAEE), Vellore, India, 2014, pp. 1-4, doi: 10.1109/ICAEE.2014.6838564.

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

▼ Function to read .wav file

```
1 import wave
2 import numpy as np
3 import pylab as pl
4
5 def wavread(filename):
6     # Read wav file
7     print('Reading WAV file ',filename,'---\n')
8     wavefile = wave.open(filename, 'r')
9
10    #Function to read four types of information from wav files. numframes indicates how many frames were read in total
11    nchannels = wavefile.getnchannels()          # Returns number of audio channels (1 for mono, 2 for stereo).
12    sample_width = wavefile.getsampwidth()        # Returns sample width in bytes.
13    framerate = wavefile.getframerate()          # Returns sampling frequency.
14    numframes = wavefile.getnframes()            # Returns number of audio frames.
15
16    #print("channel",nchannels)
17    #print("sample_width",sample_width)
18    #print("framerate",framerate)
19    #print("numframes",numframes)
20
21    str_data = wavefile.readframes(numframes)     # Reads and returns at most n frames of audio, as a bytes object.
22    wavefile.close()
23
24    #Convert waveform data to array
25    wave_data = np.fromstring(str_data, dtype=np.short)
26    #print(len(wave_data))
27    time = np.arange(0, numframes) * (1.0 / framerate)
28    #print(len(time))
29
30    return wave_data
```

▼ Pre-emphasis

```
1 def pre_emphasis(signal,coefficient=0.95):
2     '''
3     Pre-emphasize the signal
4     Parameter meaning:-
5     signal: original signal
6     coefficient: emphasis coefficient, default is 0.95
7     '''
8     print('Pre emphasis \n')
9     return numpy.append(signal[0],signal[1:]-coefficient*signal[:-1])
```

▼ Signal framing and windowing

```

1 def enframe(signal, nw, inc, winfunc):
2     '''
3     Convert audio signals into frames.
4     Parameter meaning:
5     signal:original audio model
6     nw: The length of each frame (here refers to the length of the sampling point, that is, the sampling frequency multiplied by the
7     inc: interval between adjacent frames (same as defined above)
8     '''
9     print('Signal framing and windowing \n')
10    signal_length=len(signal) #total signal length
11    if signal_length<nw: #If the signal length is less than the length of one frame, the number of frames is defined as 1
12        nf=1
13    else: #Otherwise, calculate the total length of the frame
14        nf=int(np.ceil((1.0*signal_length-nw+inc)/inc))
15    pad_length=int((nf-1)*inc+nw) #The total flattened length of all frames added up
16    zeros=np.zeros((pad_length-signal_length,)) #Insufficient lengths are padded with 0s, similar to the extended array operation in
17    pad_signal=np.concatenate((signal,zeros)) #The padded signal is recorded as pad_signal
18    indices=np.tile(np.arange(0,nw),(nf,1))+np.tile(np.arange(0,nf*inc,inc),(nw,1)).T #It is equivalent to extracting the time point
19    indices=np.array(indices,dtype=np.int32) #Convert indices into matrices
20    frames=pad_signal[indices] #Get frame signal
21    win=np.tile(winfunc,(nf,1)) #window window function, the default value here is 1
22    #print(np.shape(frames*win))
23    return frames*win #Return frame signal matrix

1 def sgn(n):
2     # Returns the sign of frames, used in ZCR Calculation
3     if n>=0:
4         return 1
5     if n<0:
6         return -1
7
8 def energy(frames):
9     #print(frames)
10    frames=frames/np.amax(np.absolute(frames)) # classified as [-1,1]
11    nframe=np.shape(frames)[0] # Number of frames
12    lframe=np.shape(frames)[1] # Energy of frames
13    energy=np.power(frames,2) # Sum energy along each frame
14    e=energy.sum(axis=1)
15    time = np.arange(0, nframe)
16    pl.plot(time,e)
17    pl.xlabel("Time")
18    pl.ylabel("Energy of signal")
19    pl.show()
20    return e
21
22 def zcr(frames):
23     nframe=np.shape(frames)[0] # Number of frames
24     lframe=np.shape(frames)[1] # Length of each frame
25     zcr=np.zeros(nframe) # Initialize zero crossing rate array
26     for i in range(nframe):
27         zframe=0
28         for j in range(1,lframe):
29             zframe=zframe+0.5*(abs(sgn(frames[i,j]))-sgn(frames[i,j-1]))
30         zcr[i]=zframe/(lframe-1) # Compute zero crossing rate for each frame
31     time = np.arange(0, nframe)
32     pl.plot(time,zcr)
33     pl.xlabel("Time")
34     pl.ylabel("Zero crossing Rate")
35     pl.show()
36     return zcr
37
38 def vioceextrac(frames):
39     print('Extracting active voice \n')
40     nframe=np.shape(frames)[0] # Number of frames
41     lframe=np.shape(frames)[1]
42     mh=10 # High energy threshold
43     ml=2 # Low energy threshold
44     zs=0.18 # Zero crossing rate threshold
45     a1=0
46     a2=0
47     status=0
48     count=0 # This variable counts the number of consecutive frames with energy levels below mh
49     e=energy(frames) # Compute energy of frames
50     z=zcr(frames) # Compute zero crossing rate of frames
51     mh=min(mh,np.amax(e)/4) # Adjust high energy threshold
52     ml=min(ml,np.amax(e)/8) # Adjust low energy threshold
53     #print(mh,ml)
54     for i in range(nframe):
55         # a1, a2, a2t: These variables are used to track the indices of frames where active
56         if status==0 and e[i]>mh:
57             a1=i
58             status=1
59         if status==1 and e[i]<ml:
60             a2t=i

```

```

59         status=2
60         if status==2 and e[i]<mh:
61             count=count+1
62         if status==2 and e[i]>mh:
63             count=0
64             status=1
65         if status==2 and count>30:
66             a2=a2t
67     #print(a1,a2)
68     b1=a1-1
69     b2=a2+1
70     while e[b1]>m1:
71         b1=b1-1
72     while e[b2]>m1:
73         b2=b2+1
74     #print(b1,b2)
75     c1=b1-1
76     c2=b2+1
77     while z[c1]>=(3*zs):
78         c1=c1-1
79     while z[c2]>=(3*zs):
80         c2=c2+1
81     #print(c1,c2)
82     frames_a=frames[c1:c2,:]
83     return(frames_a)

```

b1, b2: These variables are used to adjust the boundaries of the active voice segm

c1, c2: These variables are used to further refine the boundaries of the active voi

status: This variable tracks the current state of processing the frames. It transit

Return frames containing active voice segments

▼ Feature Extraction (using MFCC)

```

1 import numpy
2 import scipy.io.wavfile
3 from scipy.fftpack import dct
4
5 def mfcc(frames,NFFT=512,sample_rate=16000):
6     print('Computing MFCC features \n')
7     #frame energy spectrum
8     mag_frames = numpy.absolute(numpy.fft.rfft(frames, NFFT)) # Magnitude of the FFT
9     pow_frames = ((1.0 / NFFT) * ((mag_frames) ** 2)) # Power Spectrum
10    energy=numpy.sum(pow_frames,1) #Sum the energy spectrum of each frame
11    energy=numpy.where(energy==0,numpy.finfo(float).eps,energy) #Adjust the places where the energy is 0 to eps, which facilitates :
12    #print(numpy.shape(mag_frames),numpy.shape(pow_frames))
13    #filter bank
14    nfilt=40
15
16    low_freq_mel = 0
17    high_freq_mel = (2595 * numpy.log10(1 + (sample_rate / 2) / 700)) # Convert Hz to Mel
18    mel_points = numpy.linspace(low_freq_mel, high_freq_mel, nfilt + 2) # Equally spaced in Mel scale
19    hz_points = (700 * (10**(mel_points / 2595) - 1)) # Convert Mel to Hz
20    bin = numpy.floor((NFFT + 1) * hz_points / sample_rate)
21
22    fbank = numpy.zeros((nfilt, int(numpy.floor(NFFT / 2 + 1))))
23    for m in range(1, nfilt + 1):
24        f_m_minus = int(bin[m - 1]) # left
25        f_m = int(bin[m]) # center
26        f_m_plus = int(bin[m + 1]) # right
27
28        for k in range(f_m_minus, f_m):
29            fbank[m - 1, k] = (k - bin[m - 1]) / (bin[m] - bin[m - 1])
30        for k in range(f_m, f_m_plus):
31            fbank[m - 1, k] = (bin[m + 1] - k) / (bin[m + 1] - bin[m])
32    filter_banks = numpy.dot(pow_frames, fbank.T)
33    filter_banks = numpy.where(filter_banks == 0, numpy.finfo(float).eps, filter_banks) # Numerical Stability
34    filter_banks = 20 * numpy.log10(filter_banks) # dB
35
36    #mfcc
37    num_ceps = 13
38    appendEnergy=1
39
40    mfcc = dct(filter_banks, type=2, axis=1, norm='ortho')[:, 1 : (num_ceps + 1)] # Keep 2-13
41    if appendEnergy:
42        mfcc[:,0]=numpy.log(energy) #Take only 2-13 coefficients and replace the first one with the logarithm of the energy
43    #mean-normalized
44    mfcc -= (numpy.mean(mfcc, axis=0) + 1e-8)
45
46
47    return mfcc_delta_delta(mfcc)
48
49 def derivate(feats,big_theta=2,cep_num=13):
50     '''General transformation formula for calculating first-order coefficients or acceleration coefficients
51     Parameter Description:
52     feat: MFCC array or first-order coefficient array
53     big_theta: the big theta in the formula, the default is 2. This parameter controls the range of the calculation
54     '''
55     result=numpy.zeros(feats.shape) #result
56     denominator=0 #Denominator
57     for theta in numpy.linspace(1,big_theta,big_theta):
58         denominator=denominator+theta**2
59     denominator=denominator*2 #Calculate the value of the denominator
60     for row in numpy.linspace(0,feats.shape[0]-1,feats.shape[0]):
61         tmp=numpy.zeros((cep_num,))
62         numerator=numpy.zeros((cep_num,)) #numerator
63         for t in numpy.linspace(1,cep_num,cep_num):
64             a=0
65             b=0
66             s=0
67             for theta in numpy.linspace(1,big_theta,big_theta):
68                 if (t+theta)>cep_num:
69                     a=0
70                 else:
71                     a=feat[int(row)][int(t+theta-1)]
72                 if (t-theta)<1:
73                     b=0
74                 else:
75                     b=feat[int(row)][int(t-theta-1)]
76                 s+=theta*(a-b)
77             numerator[int(t-1)]=s
78         tmp=numerator*1.0/denominator
79         result[int(row)]=tmp
80     return result
81
82

```

```

83 def mfcc_delta(feats):
84     '''Calculate 13 MFCC+13 first-order differential coefficients
85     '''
86     result=derivate(feats) #Call the derive function
87     result=numpy.concatenate((feats,result),axis=1)
88     return result
89
90
91 def mfcc_delta_delta(feats):
92     '''Calculate 13 MFCCs + 13 first-order differential coefficients + 13 acceleration coefficients, a total of 39 coefficients
93     '''
94     result1=derivate(feats)
95     result2=derivate(result1)
96     result3=numpy.concatenate((feats,result1),axis=1)
97     result=numpy.concatenate((result3,result2),axis=1)
98     return result

```

▼ DTW Principle

```

1 import numpy as np
2
3 def dist(feats1,feats2):
4     n=np.shape(feats1)[0]
5     m=np.shape(feats2)[0]
6     d=np.zeros((n,m))
7     for i in range(n):
8         for j in range(m):
9             d[i,j]=np.sqrt(np.sum(np.square(feats1[i,:]-feats2[j,:])))
10    return d
11
12 def dtw(dist):
13
14     realmax=1.79E+308
15     n=np.shape(dist)[0]
16     m=np.shape(dist)[1]
17     D=np.ones((n+1,m+1))*realmax
18     D[0,0]=0
19     #print(D)
20     for i in range(1,n+1):
21         for j in range(1,m+1):
22             D[i,j]=dist[i-1,j-1]+min(D[i-1,j],D[i,j-1],D[i-1,j-1])
23     return D[n,m]
24
25 def score(feats1,feats2):
26     return dtw(dist(feats1,feats2))

```

Calculates the Euclidean distance between feats1 and feats2

Maximum values of numerical operations

dimensions of the distance matrix

initializes a matrix D with dimensions (n+1) x (m+1) filled with realmax

sets the starting point of the DTW matrix to 0

calculates the DTW score (similarity) between two sets of features

▼ Training

```

1 import scipy.signal as signal
2 import numpy
3 import os
4
5 modelob = open("/content/drive/MyDrive/model.txt",'w')
6
7 for line in os.listdir("/content/drive/MyDrive/train"):
8     if line.endswith(".wav"):
9         #Get training file name
10        label=line[0]
11        filename="/content/drive/MyDrive/train/"+line
12        #print(filename)
13        #file reading
14        wavearr=wavread(filename)
15        #Signal pre-emphasis
16        wavearr_pre=pre_emphasis(wavearr,0.98)
17        #Select window function
18        winfunc = signal.hamming(240)
19        #Signal framing
20        n=enframe(wavearr_pre,240,80,winfunc)
21        #vad
22        na=voiceextract(n)
23        #Extract features
24        feat=mfcc(na,512)
25        #Save features and criteria to the output model file
26        print('Output results of',filename,' \n')
27        print(label, file = modelob)
28        print(feat, file = modelob)
29        print(" ")
30        #print(numpy.shape(feat))

```

```
30 print('omp: %s' % omp.get_status())  
31 print('\n Training complete \n')  
32 modelob.close()
```

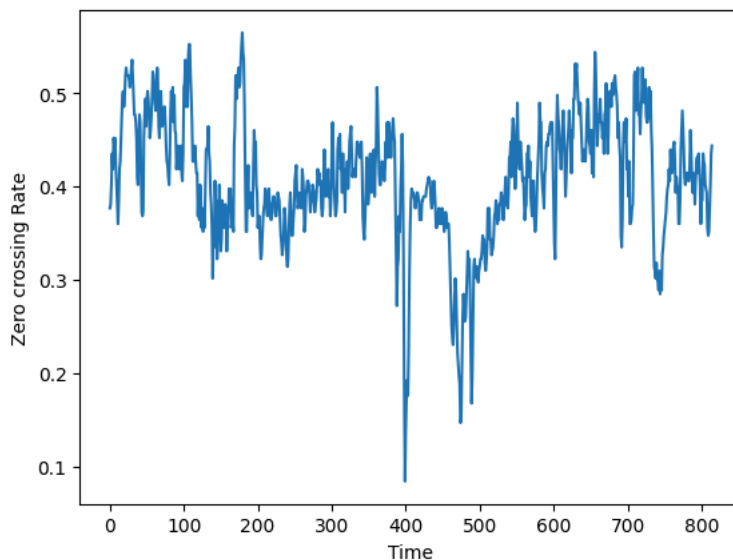
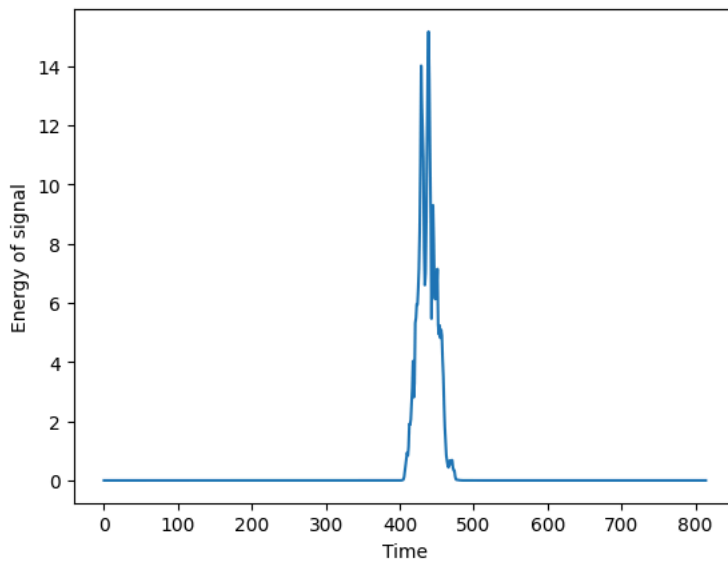

Reading WAV file /content/drive/MyDrive/train/1a.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice

```
<ipython-input-2-38ef7c5dd080>:25: DeprecationWarning: The binary mode of fromstring is deprecated, as it behaves surprisingly on  
wave_data = np.fromstring(str_data, dtype=np.short)  
<ipython-input-11-ffa5ce279020>:18: DeprecationWarning: Importing hamming from 'scipy.signal' is deprecated and will raise an error in the future  
winfunc = signal.hamming(240)
```



Computing MFCC features

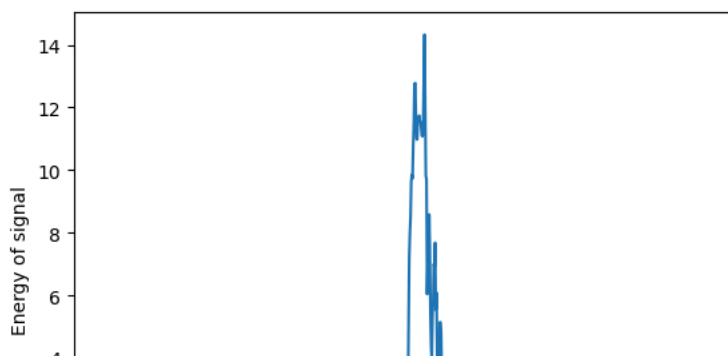
Output results of /content/drive/MyDrive/train/1a.wav

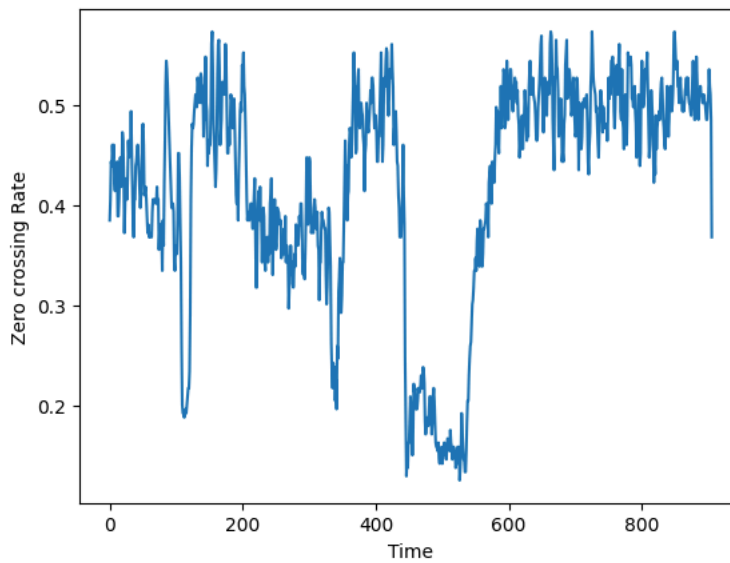
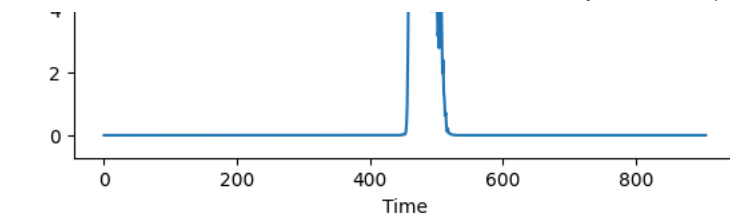
Reading WAV file /content/drive/MyDrive/train/2a.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice





Computing MFCC features

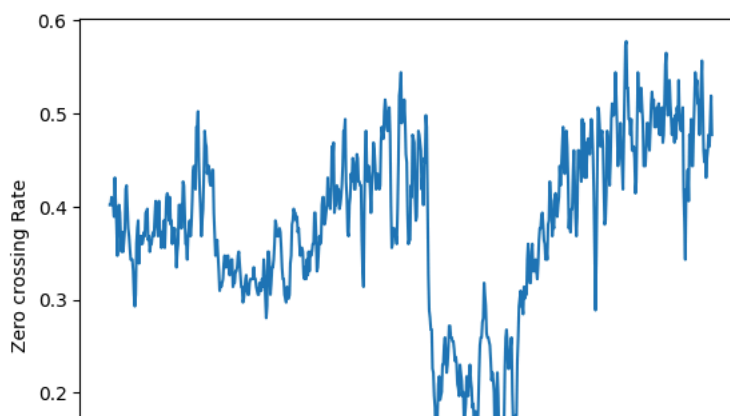
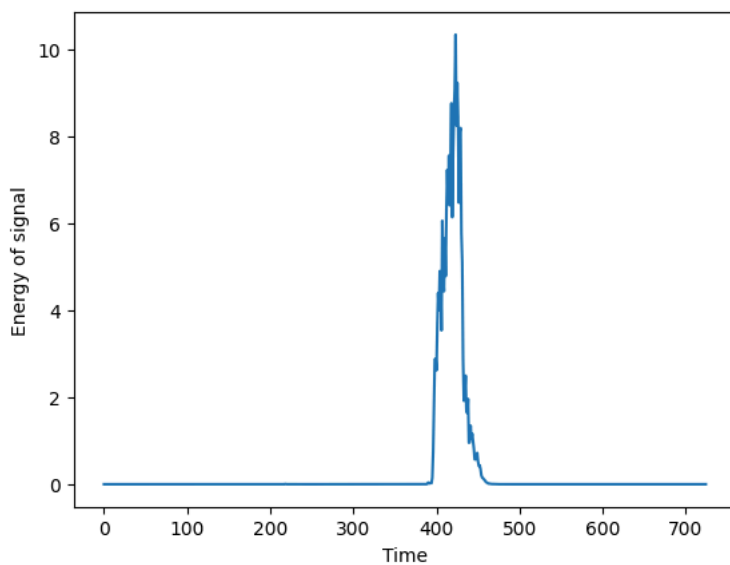
Output results of /content/drive/MyDrive/train/2a.wav

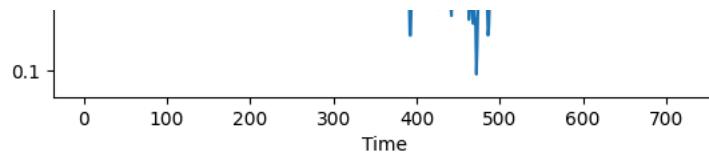
Reading WAV file /content/drive/MyDrive/train/0c.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice





Computing MFCC features

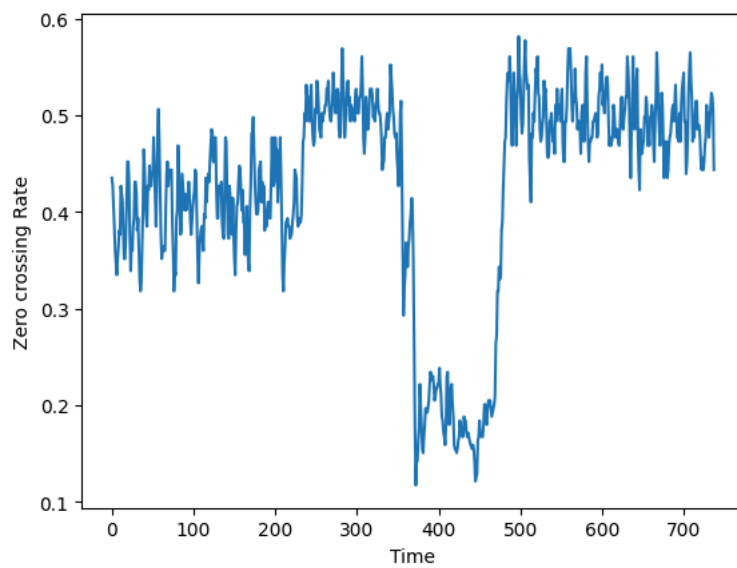
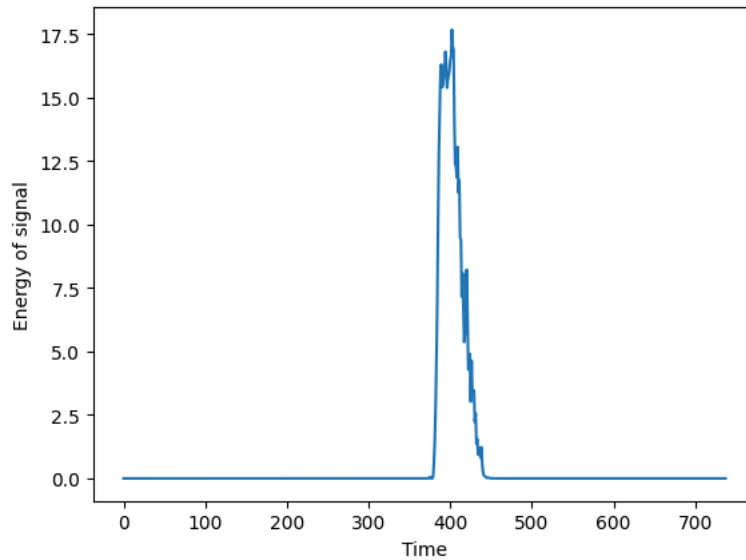
Output results of /content/drive/MyDrive/train/0c.wav

Reading WAV file /content/drive/MyDrive/train/2b.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice



Computing MFCC features

Output results of /content/drive/MyDrive/train/2b.wav

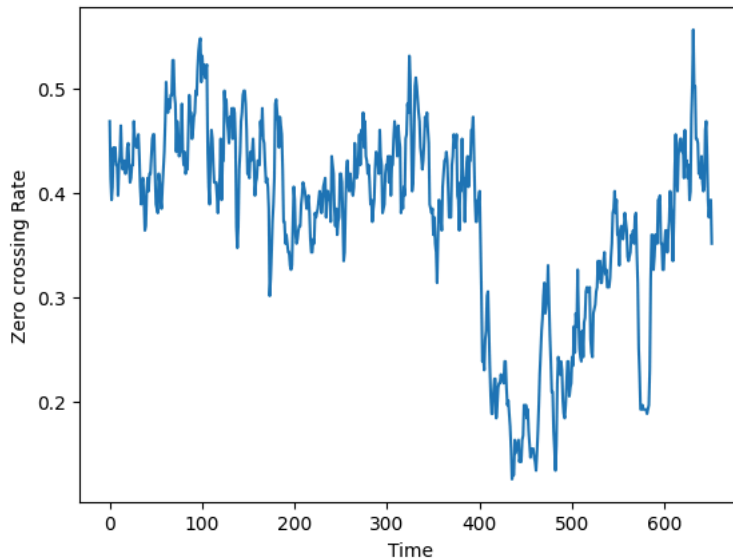
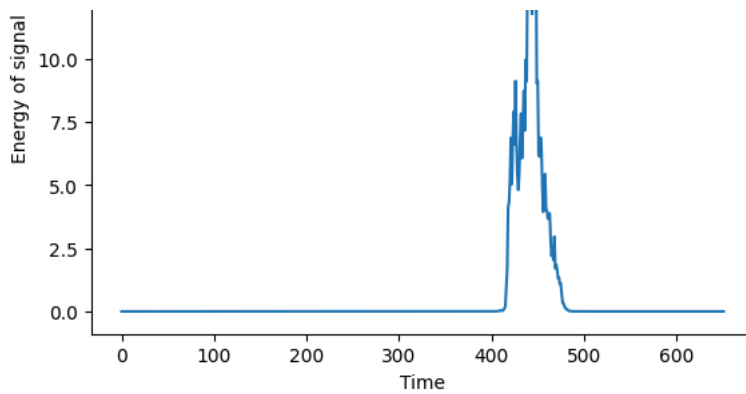
Reading WAV file /content/drive/MyDrive/train/0b.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice





Computing MFCC features

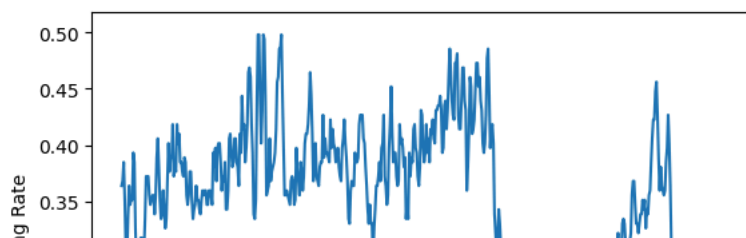
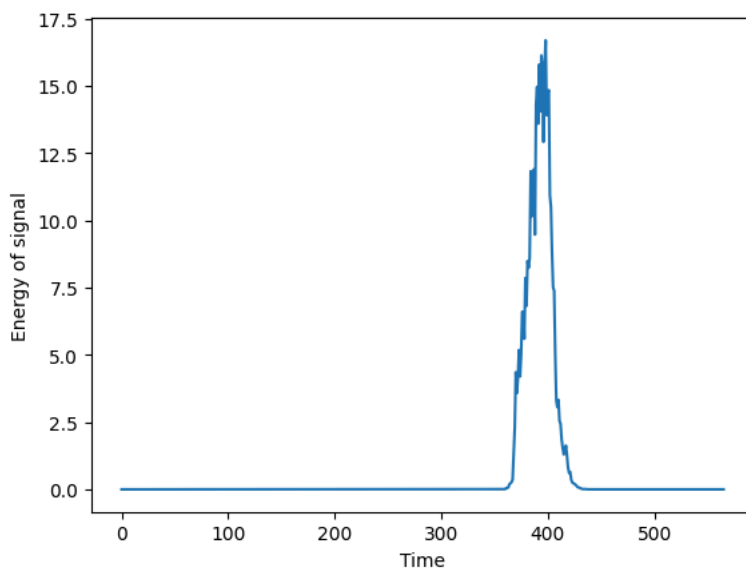
Output results of /content/drive/MyDrive/train/0b.wav

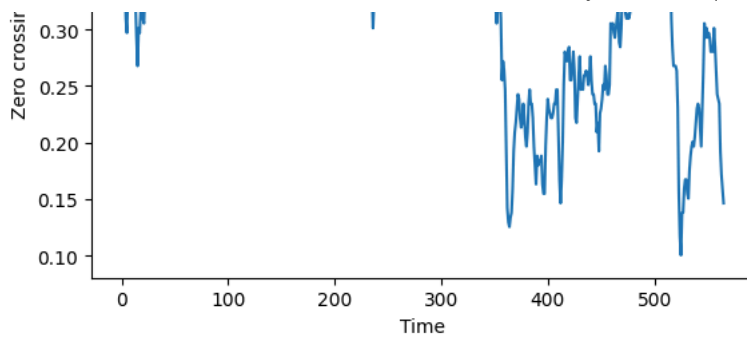
Reading WAV file /content/drive/MyDrive/train/0a.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice





Computing MFCC features

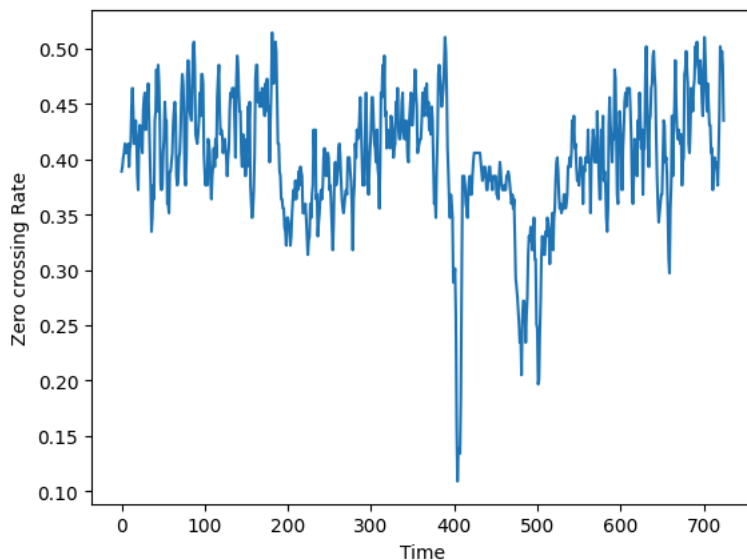
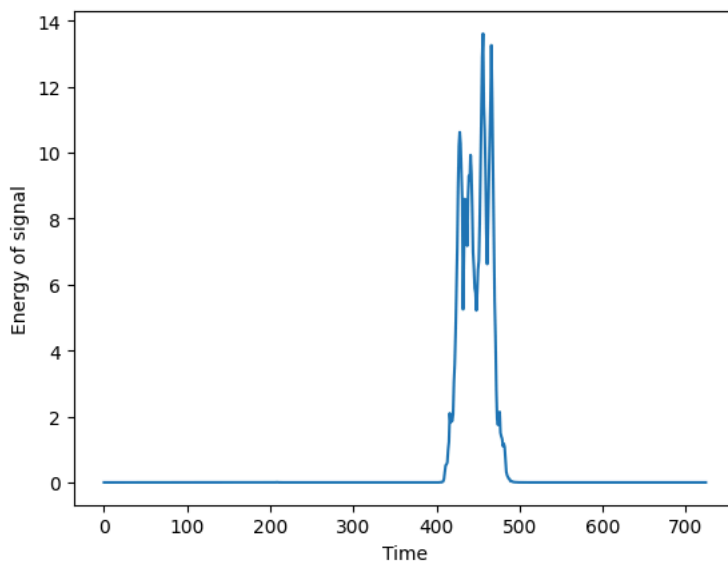
Output results of /content/drive/MyDrive/train/0a.wav

Reading WAV file /content/drive/MyDrive/train/1b.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice



Computing MFCC features

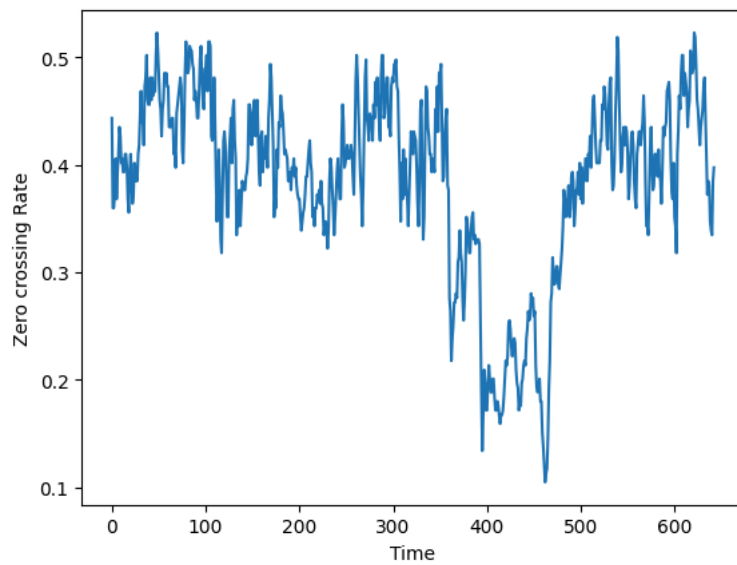
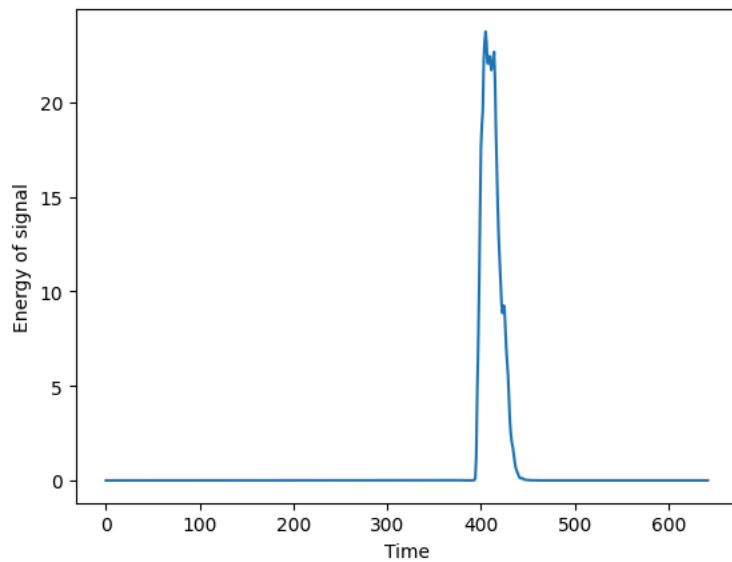
Output results of /content/drive/MyDrive/train/1b.wav

Reading WAV file /content/drive/MyDrive/train/3b.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice



Computing MFCC features

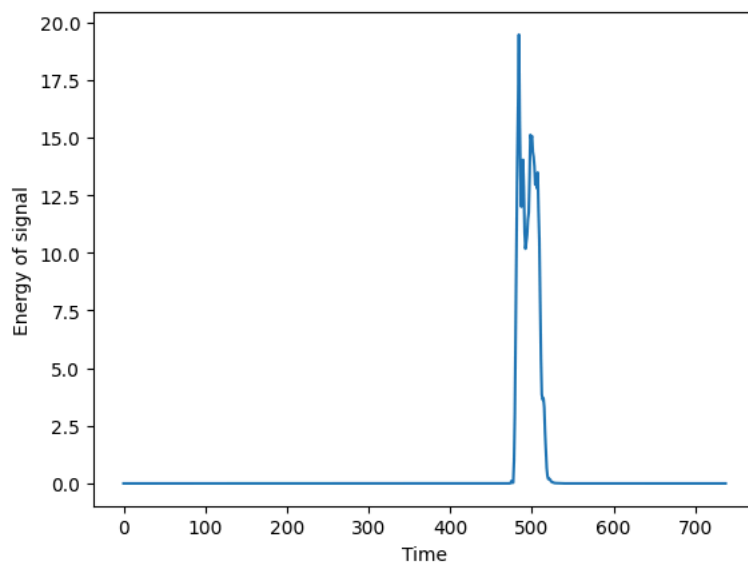
Output results of /content/drive/MyDrive/train/3b.wav

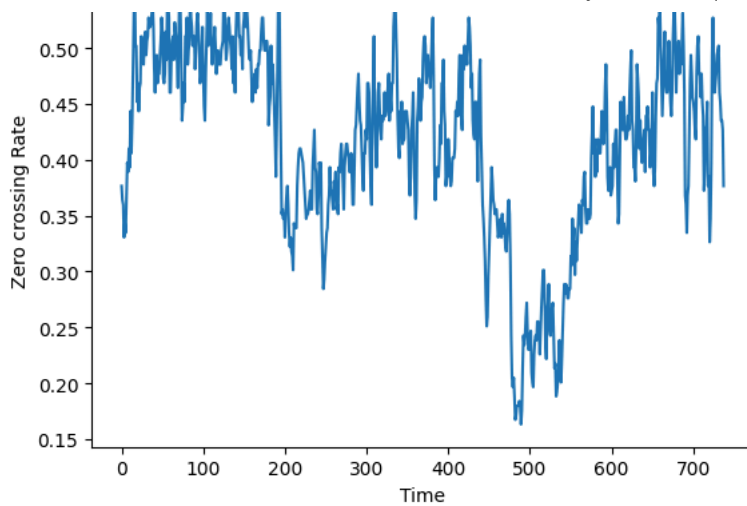
Reading WAV file /content/drive/MyDrive/train/3a.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice





Computing MFCC features

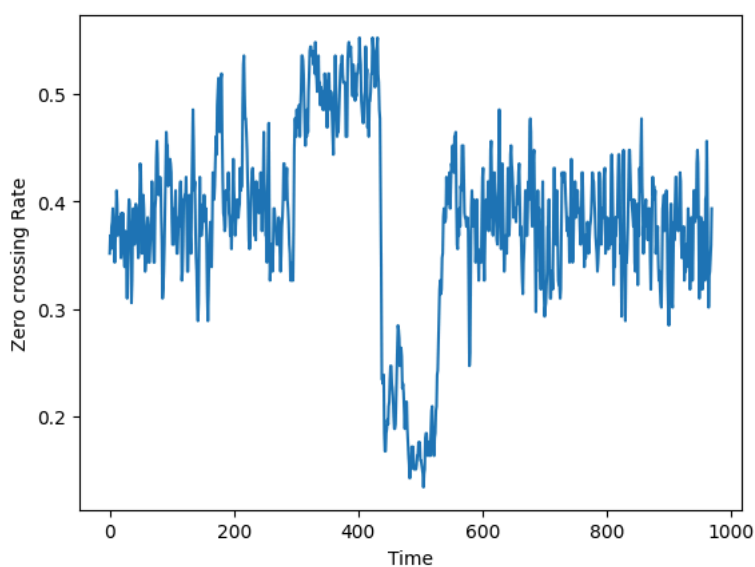
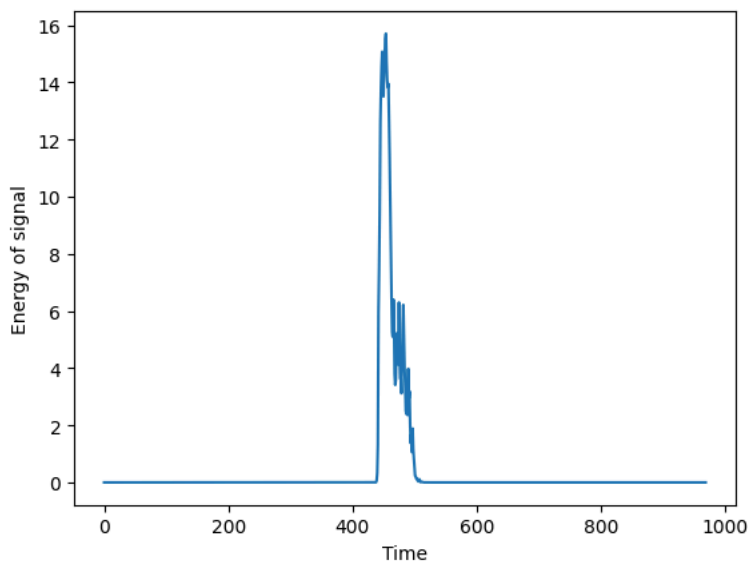
Output results of /content/drive/MyDrive/train/3a.wav

Reading WAV file /content/drive/MyDrive/train/2c.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice



Computing MFCC features

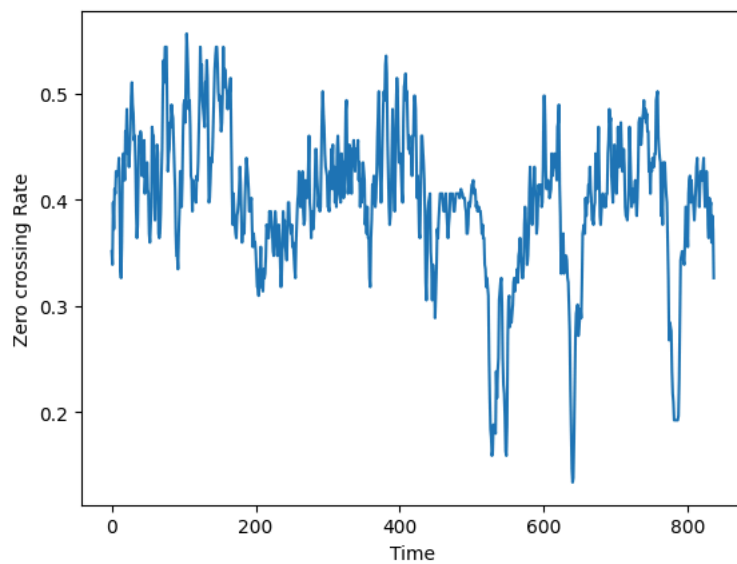
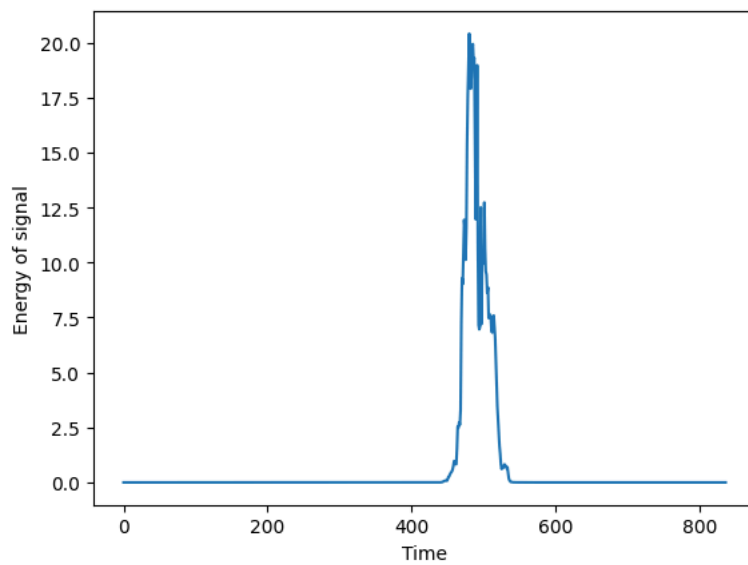
Output results of /content/drive/MyDrive/train/2c.wav

Reading WAV file /content/drive/MyDrive/train/1c.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice



Computing MFCC features

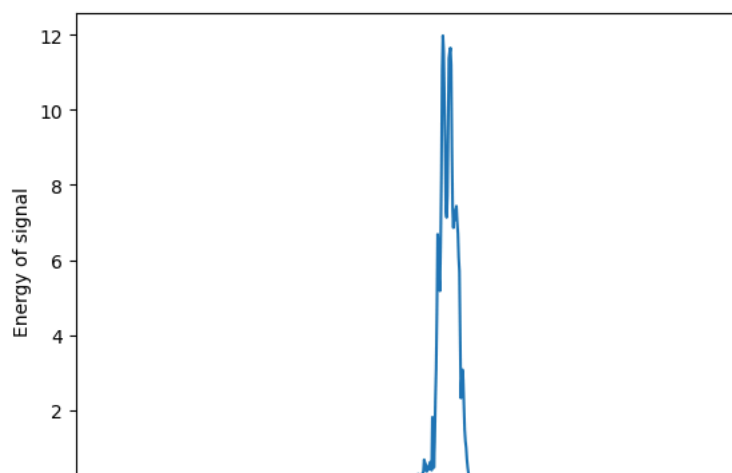
Output results of /content/drive/MyDrive/train/1c.wav

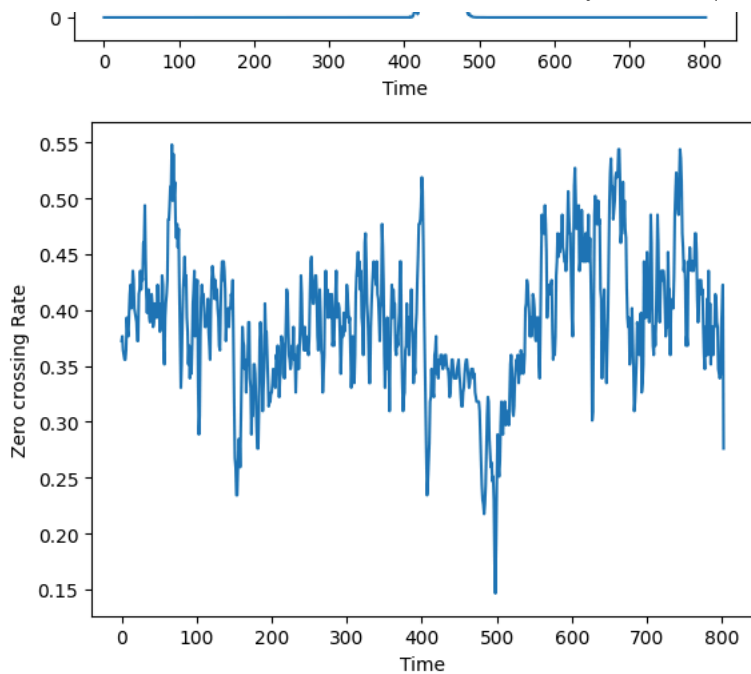
Reading WAV file /content/drive/MyDrive/train/7b.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice





Computing MFCC features

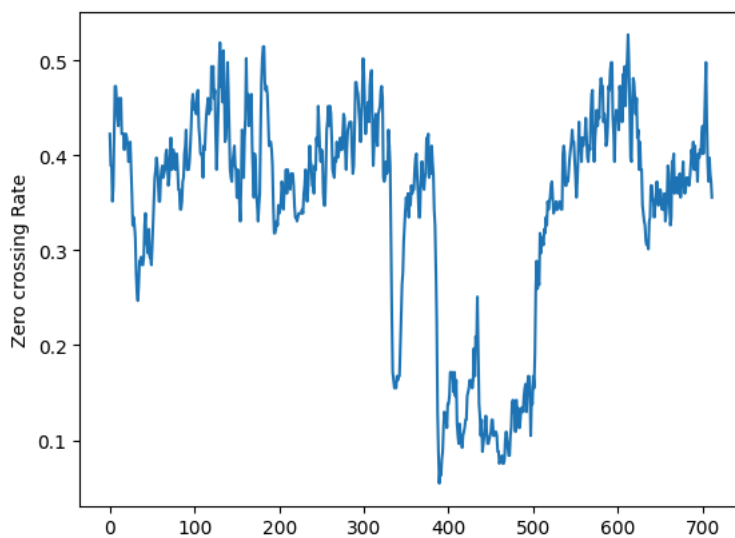
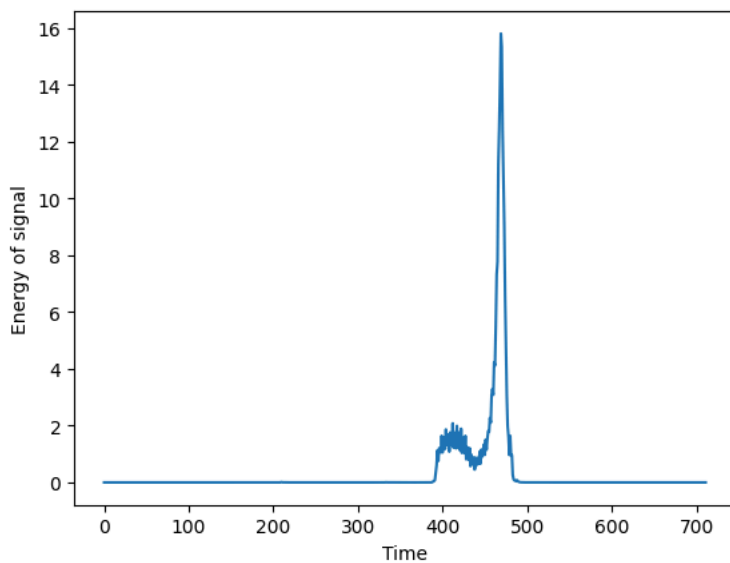
Output results of /content/drive/MyDrive/train/7b.wav

Reading WAV file /content/drive/MyDrive/train/5c.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice



Time

Computing MFCC features

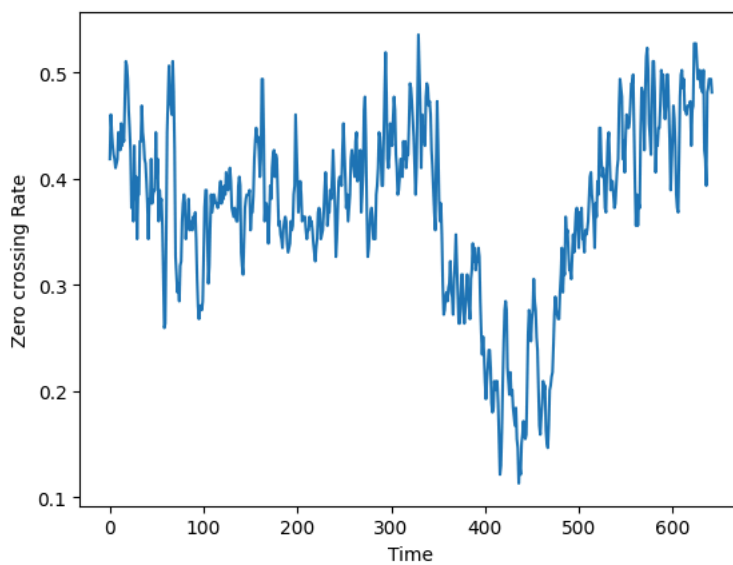
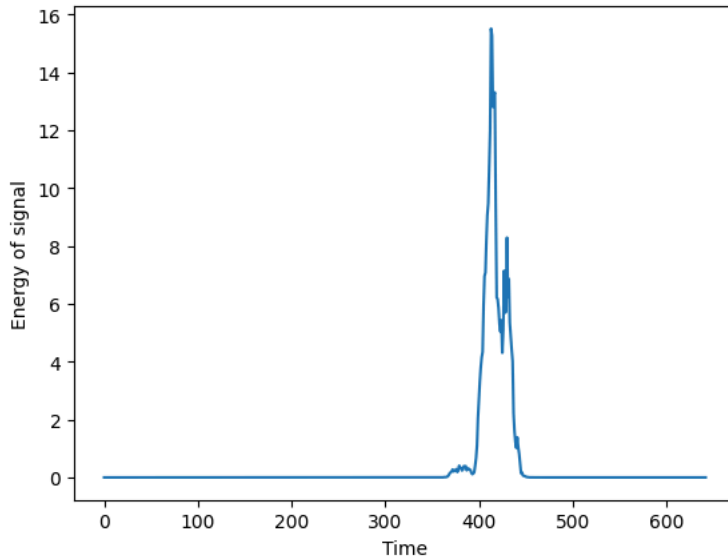
Output results of /content/drive/MyDrive/train/5c.wav

Reading WAV file /content/drive/MyDrive/train/4c.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice



Computing MFCC features

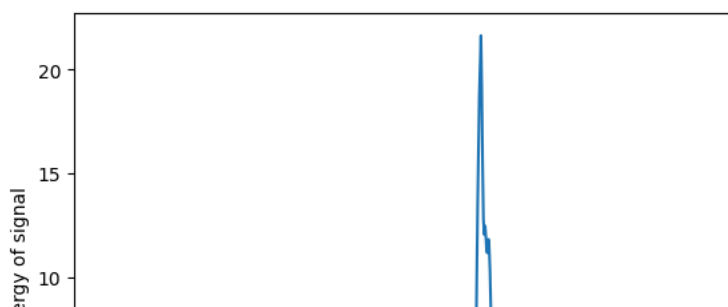
Output results of /content/drive/MyDrive/train/4c.wav

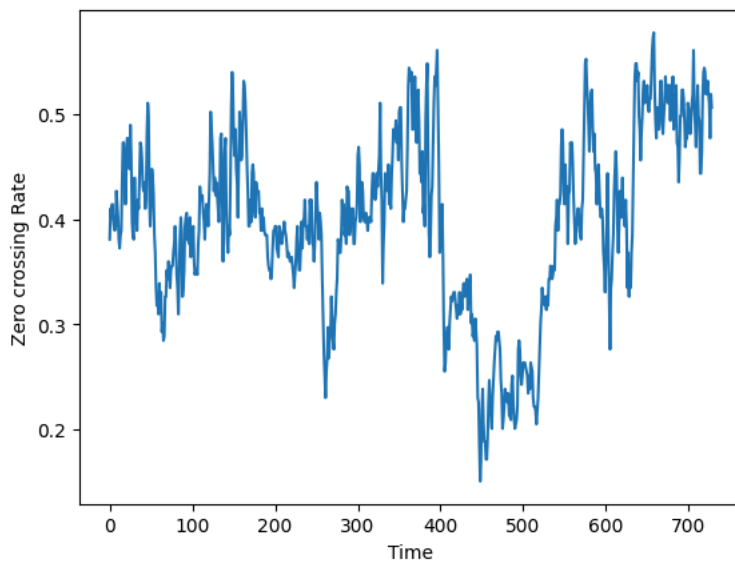
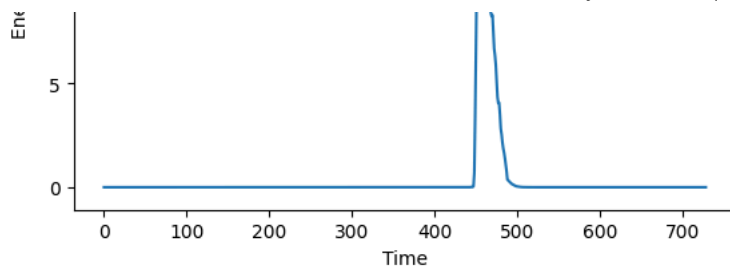
Reading WAV file /content/drive/MyDrive/train/3c.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice





Computing MFCC features

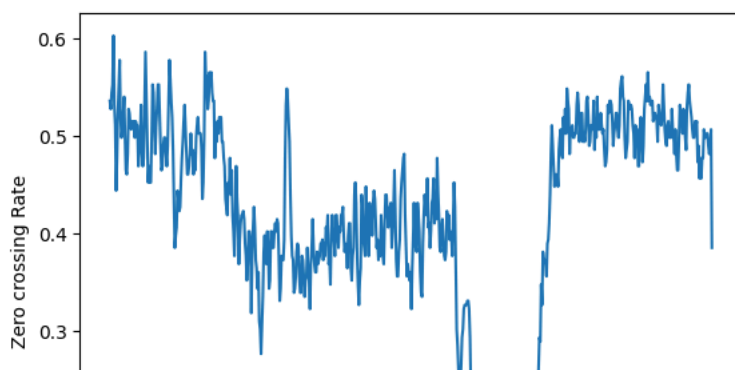
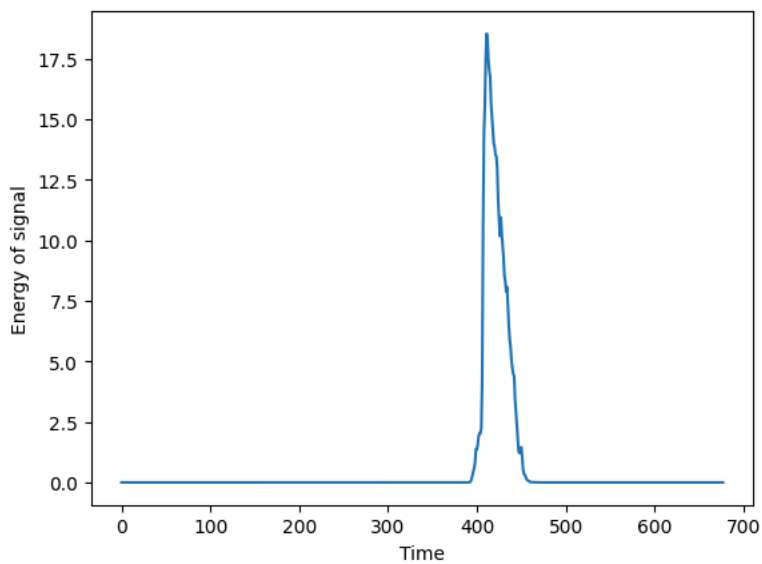
Output results of /content/drive/MyDrive/train/3c.wav

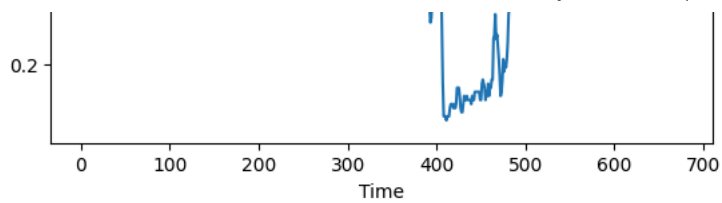
Reading WAV file /content/drive/MyDrive/train/8c.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice





Computing MFCC features

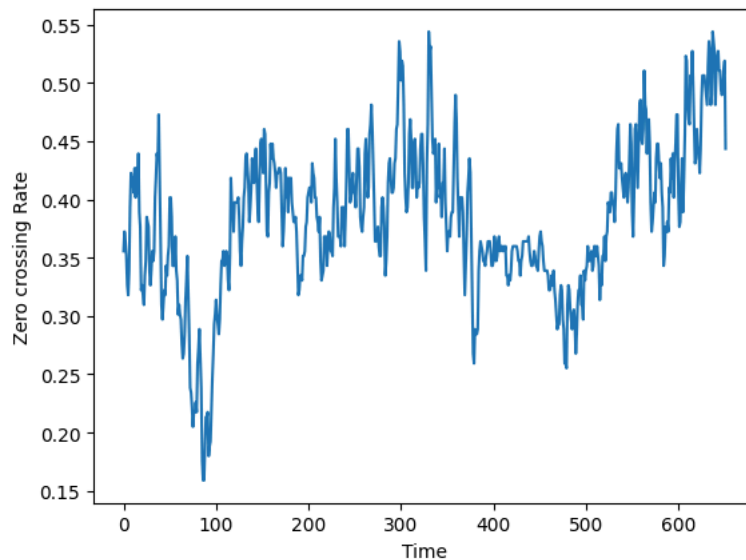
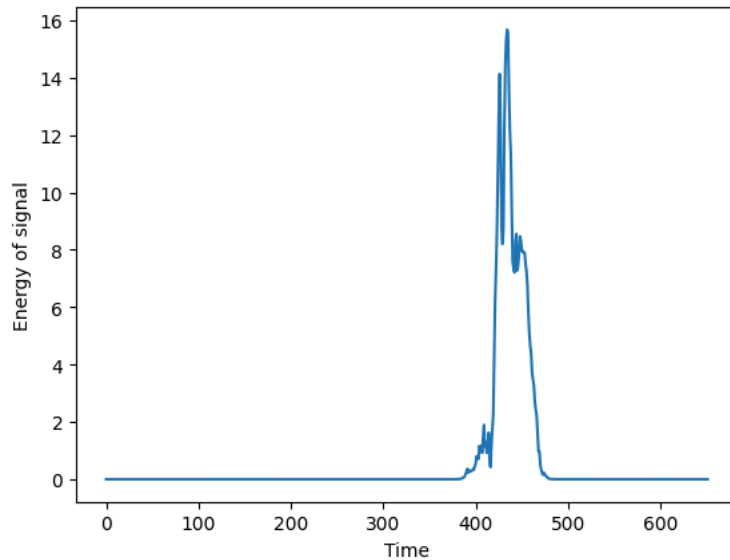
Output results of /content/drive/MyDrive/train/8c.wav

Reading WAV file /content/drive/MyDrive/train/7a.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice



Computing MFCC features

Output results of /content/drive/MyDrive/train/7a.wav

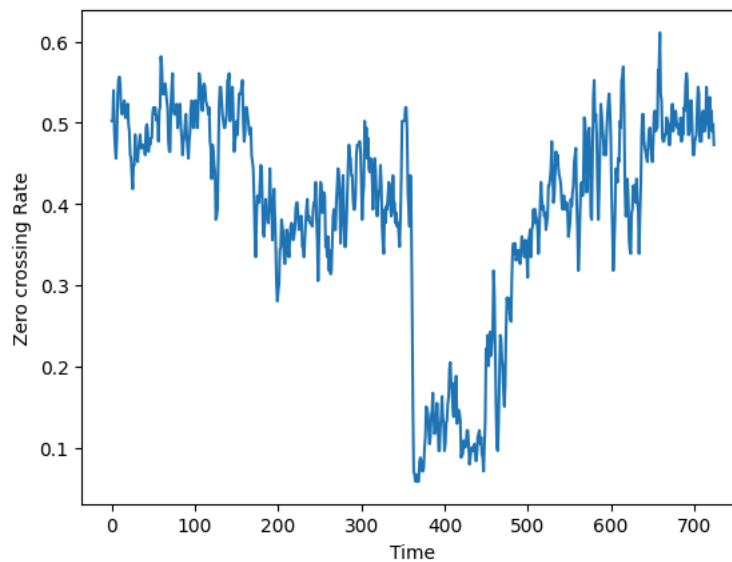
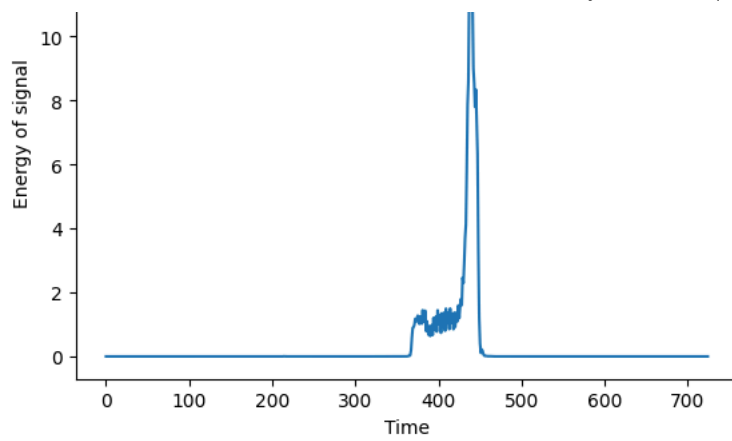
Reading WAV file /content/drive/MyDrive/train/5a.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice





Computing MFCC features

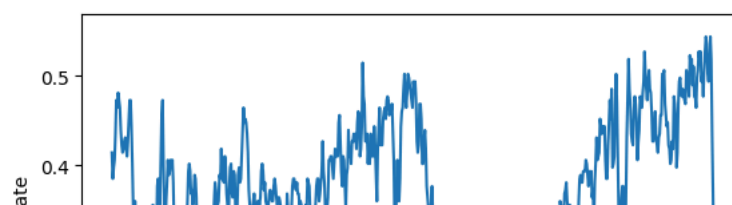
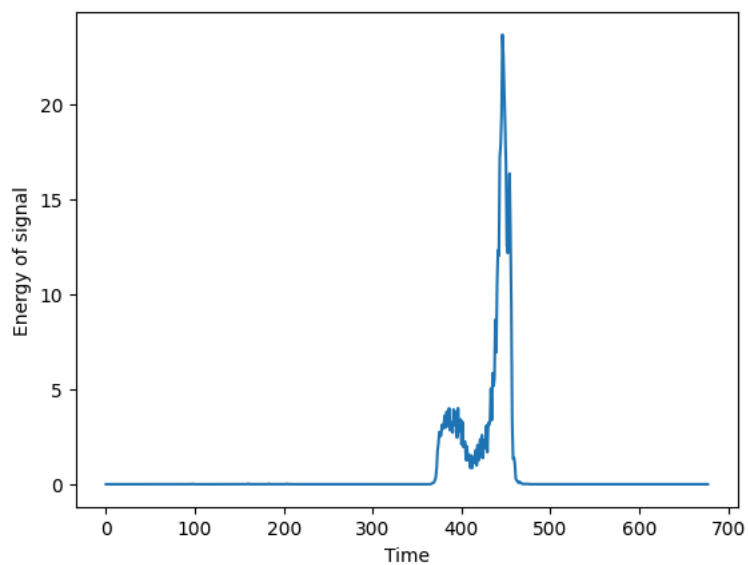
Output results of /content/drive/MyDrive/train/5a.wav

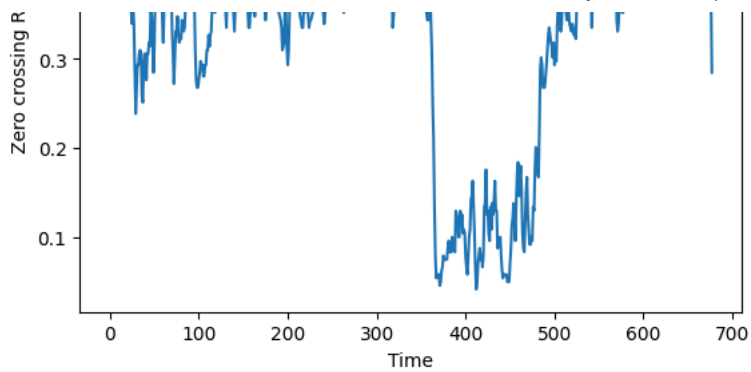
Reading WAV file /content/drive/MyDrive/train/5b.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice





Computing MFCC features

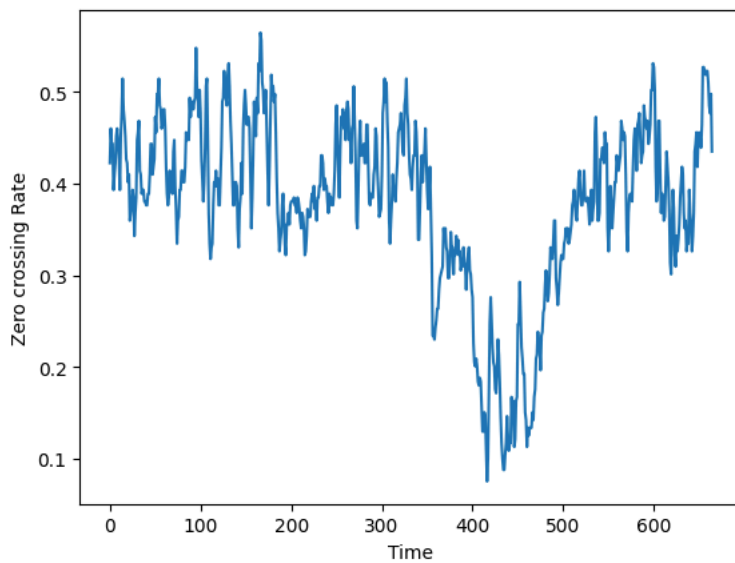
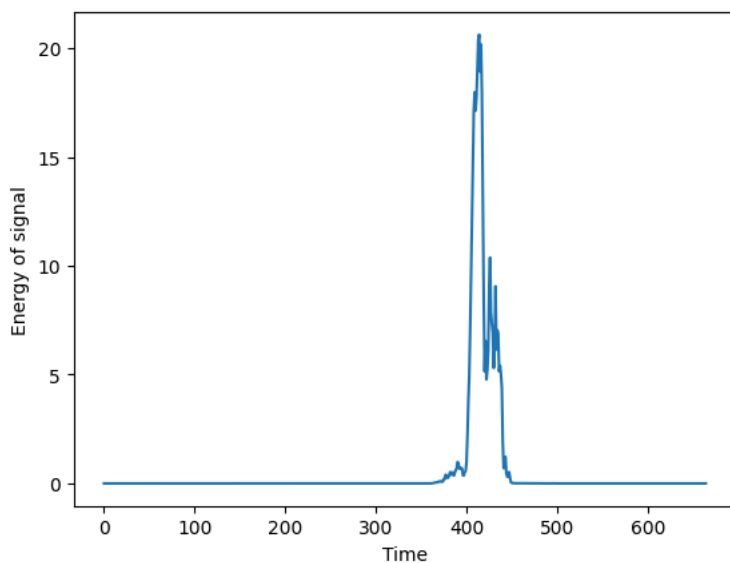
Output results of /content/drive/MyDrive/train/5b.wav

Reading WAV file /content/drive/MyDrive/train/4a.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice



Computing MFCC features

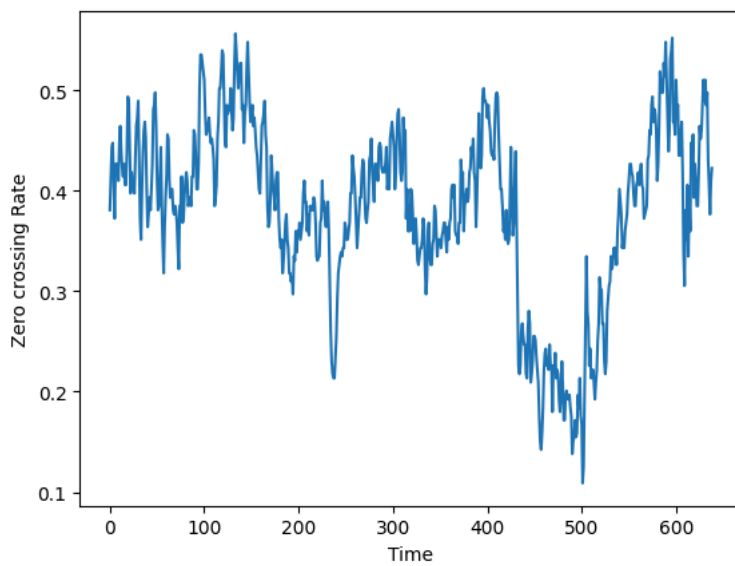
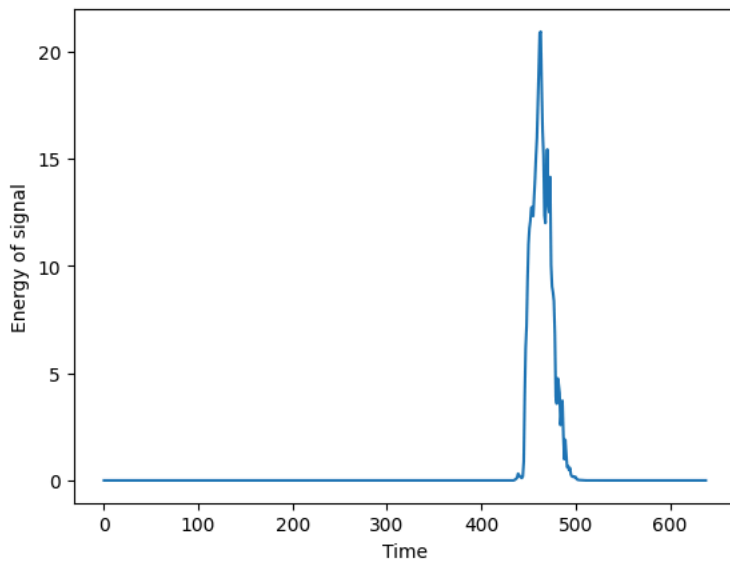
Output results of /content/drive/MyDrive/train/4a.wav

Reading WAV file /content/drive/MyDrive/train/6a.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice



Computing MFCC features

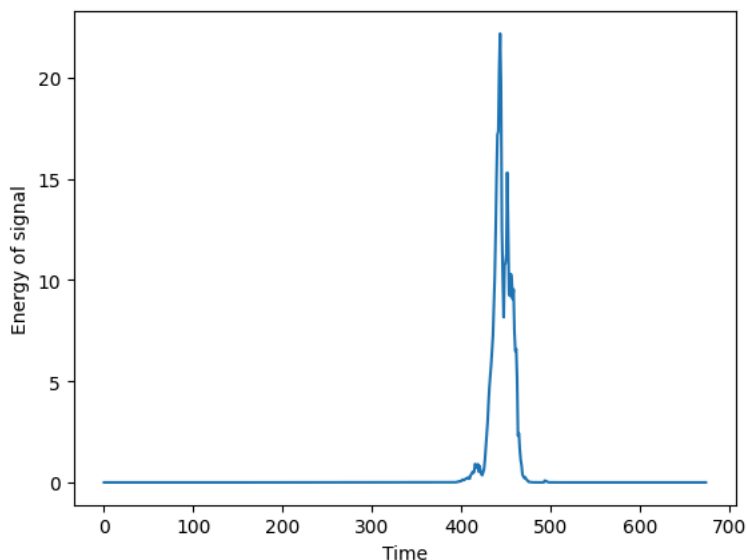
Output results of /content/drive/MyDrive/train/6a.wav

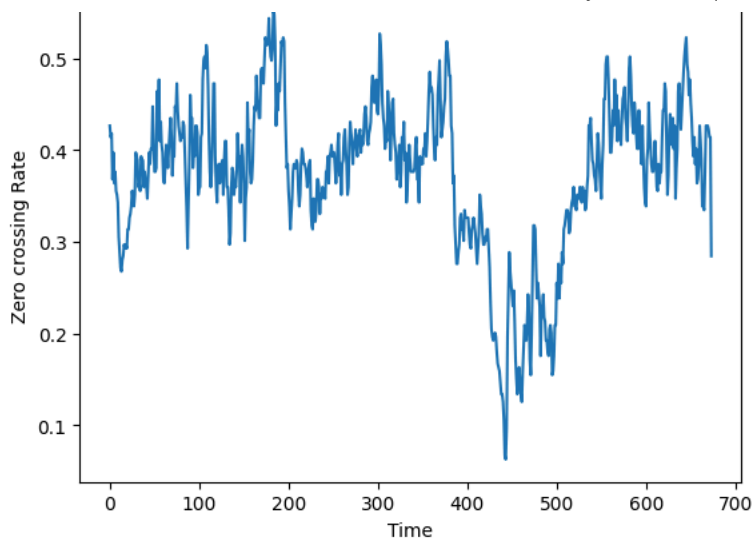
Reading WAV file /content/drive/MyDrive/train/4b.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice





Computing MFCC features

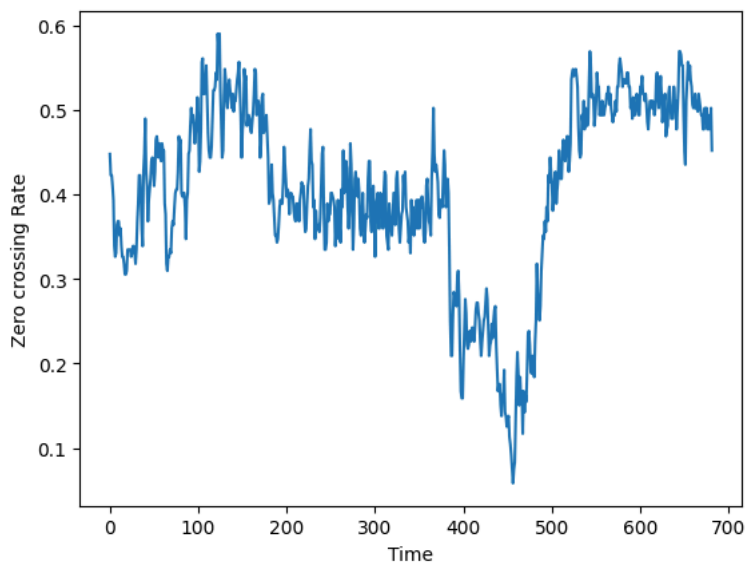
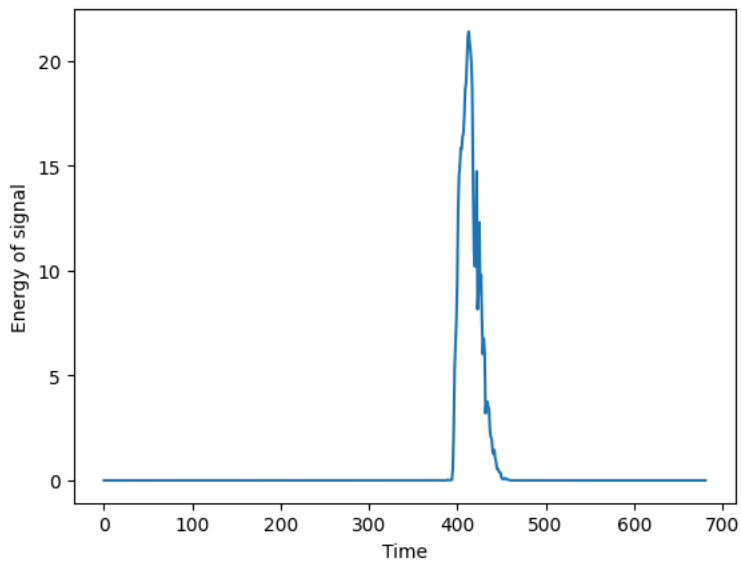
Output results of /content/drive/MyDrive/train/4b.wav

Reading WAV file /content/drive/MyDrive/train/6c.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice



Computing MFCC features

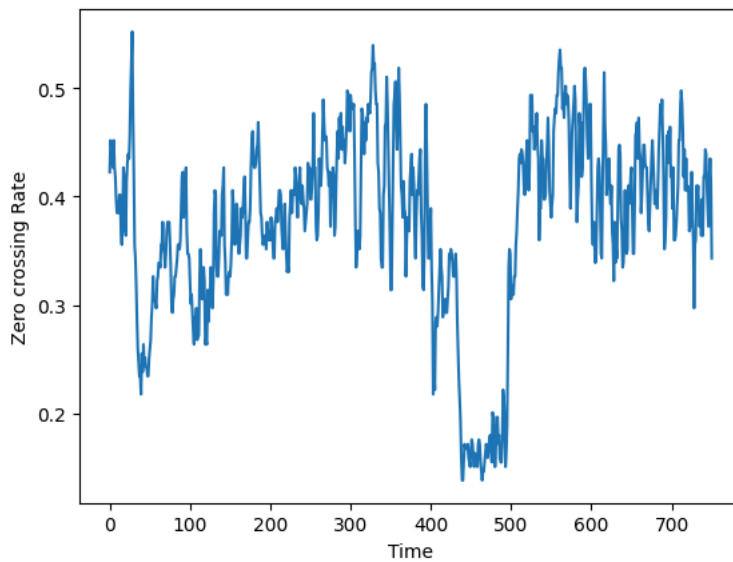
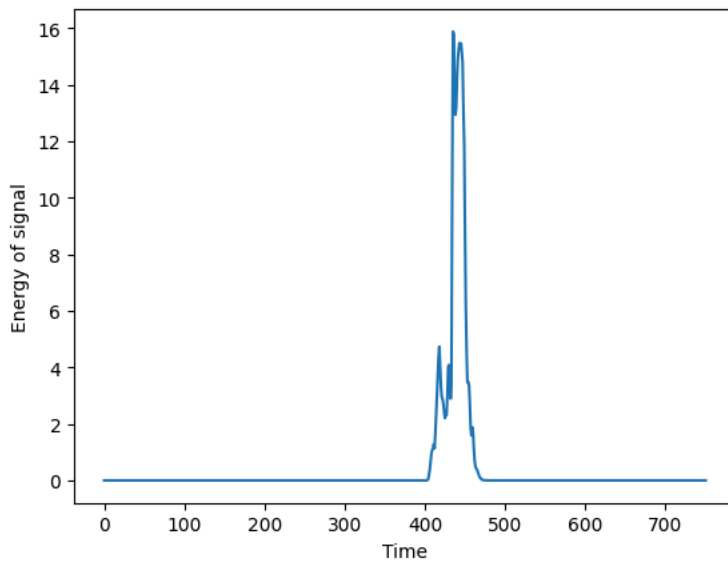
Output results of /content/drive/MyDrive/train/6c.wav

```
Reading WAV file /content/drive/MyDrive/train/8b.wav ---
```

```
Pre emphasis
```

```
Signal framing and windowing
```

```
Extracting active voice
```



```
Computing MFCC features
```

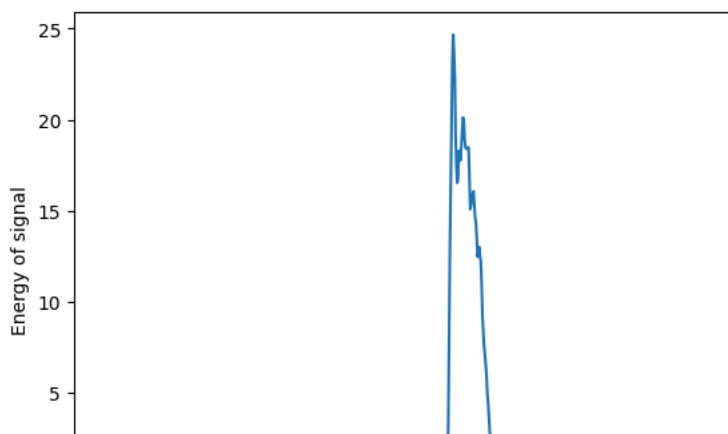
```
Output results of /content/drive/MyDrive/train/8b.wav
```

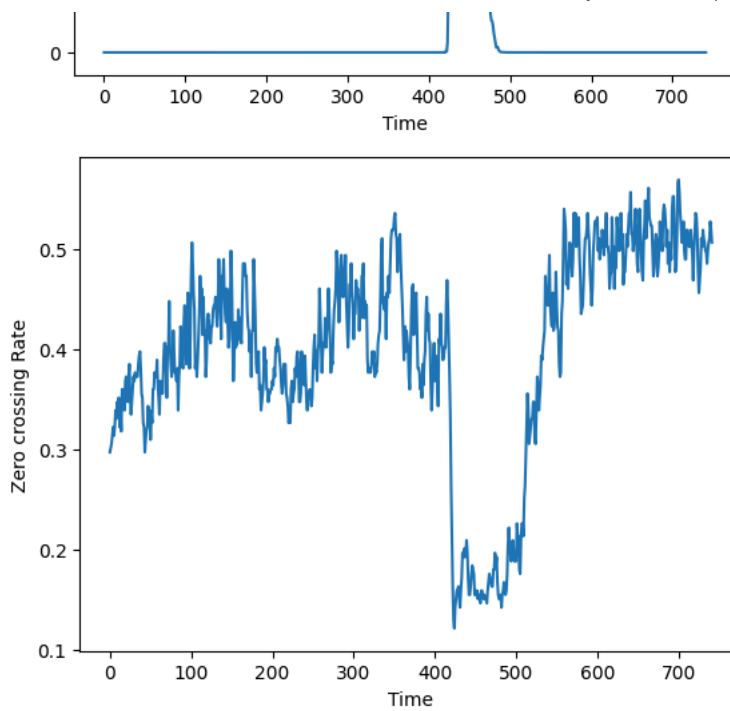
```
Reading WAV file /content/drive/MyDrive/train/8a.wav ---
```

```
Pre emphasis
```

```
Signal framing and windowing
```

```
Extracting active voice
```





Computing MFCC features

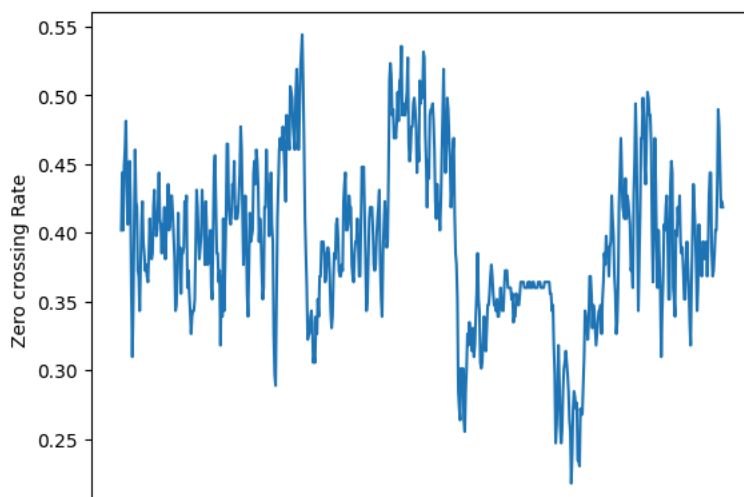
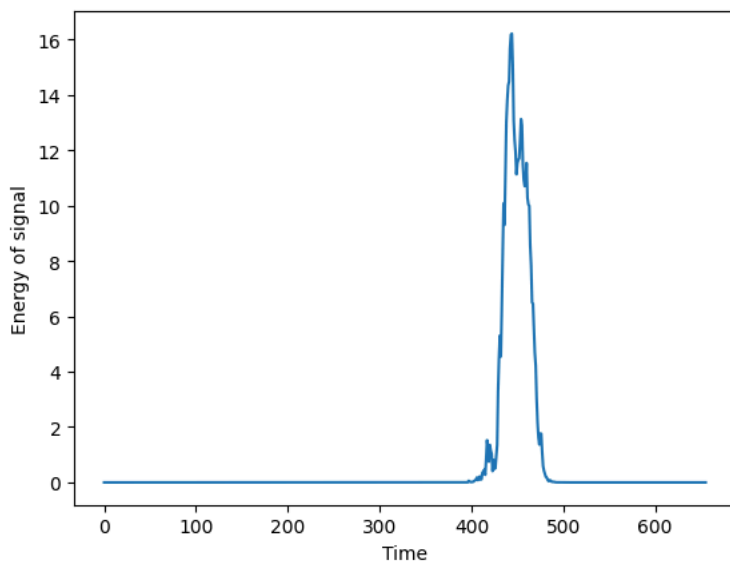
Output results of /content/drive/MyDrive/train/8a.wav

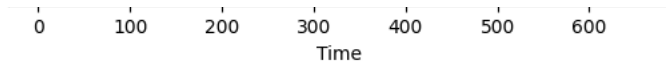
Reading WAV file /content/drive/MyDrive/train/7c.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice





Computing MFCC features

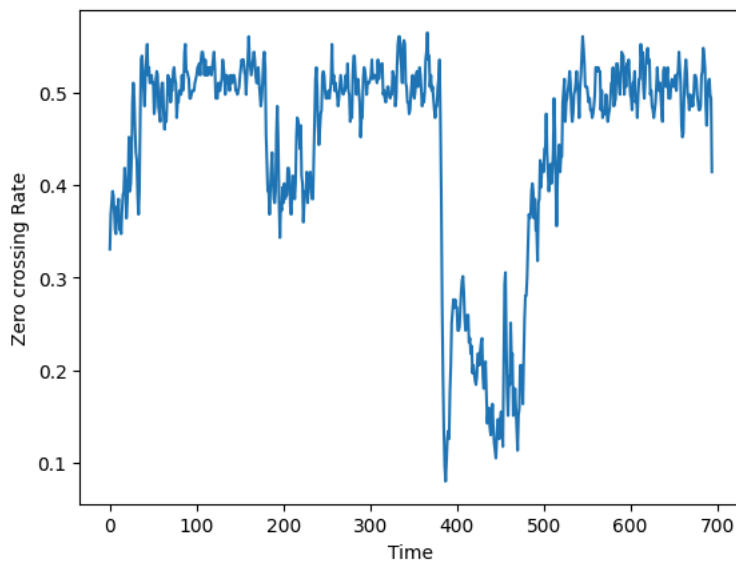
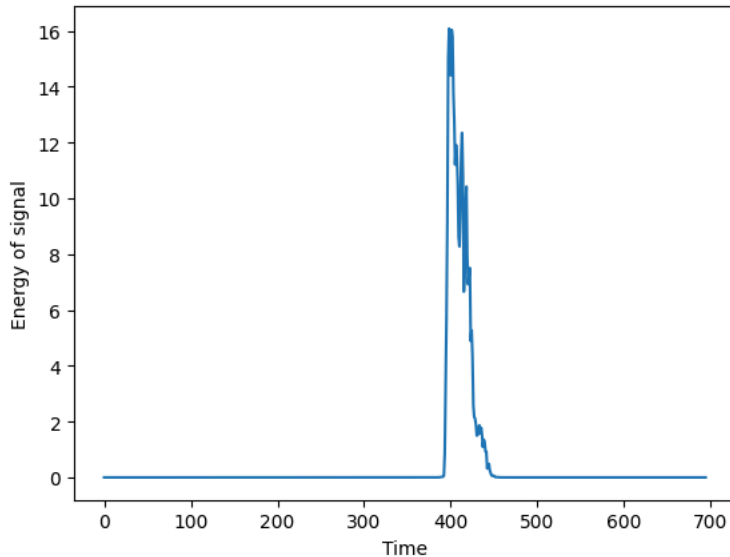
Output results of /content/drive/MyDrive/train/7c.wav

Reading WAV file /content/drive/MyDrive/train/6b.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice



Computing MFCC features

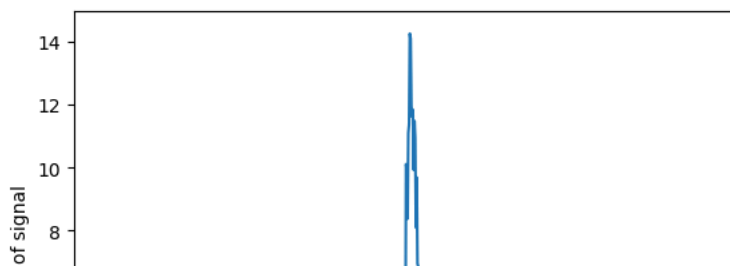
Output results of /content/drive/MyDrive/train/6b.wav

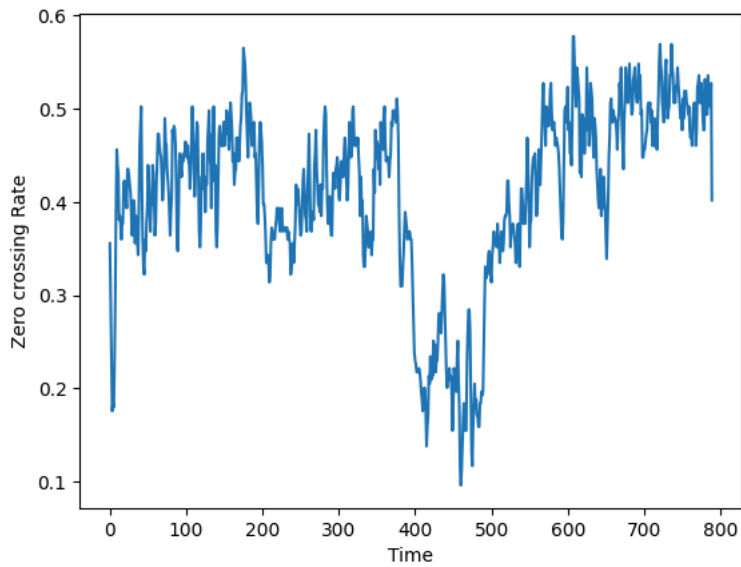
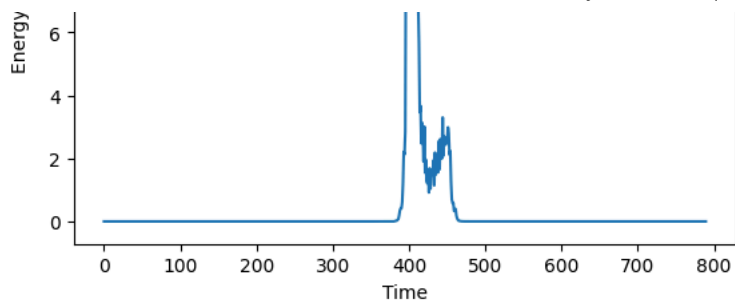
Reading WAV file /content/drive/MyDrive/train/9a.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice





Computing MFCC features

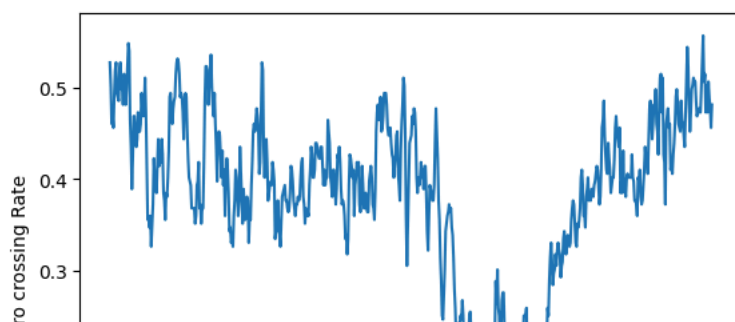
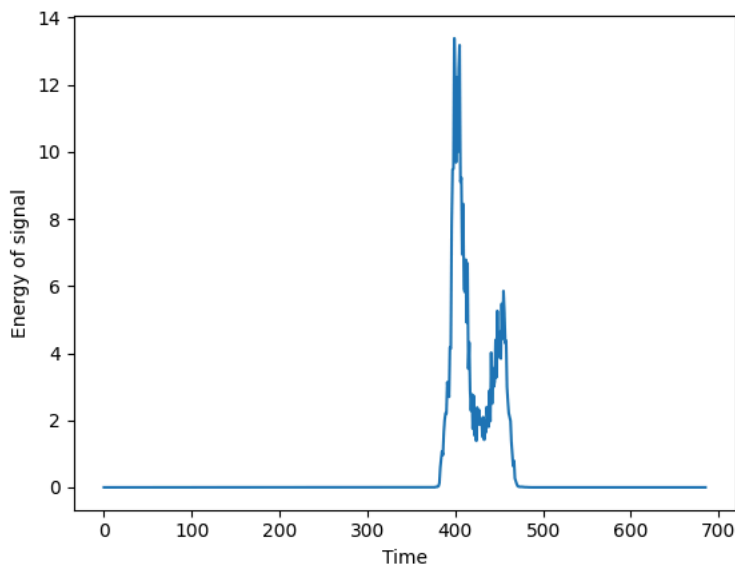
Output results of /content/drive/MyDrive/train/9a.wav

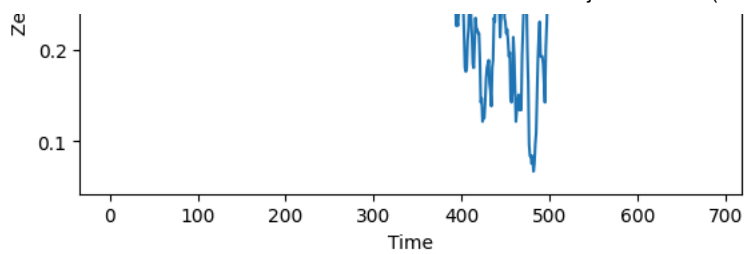
Reading WAV file /content/drive/MyDrive/train/9b.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice





Computing MFCC features

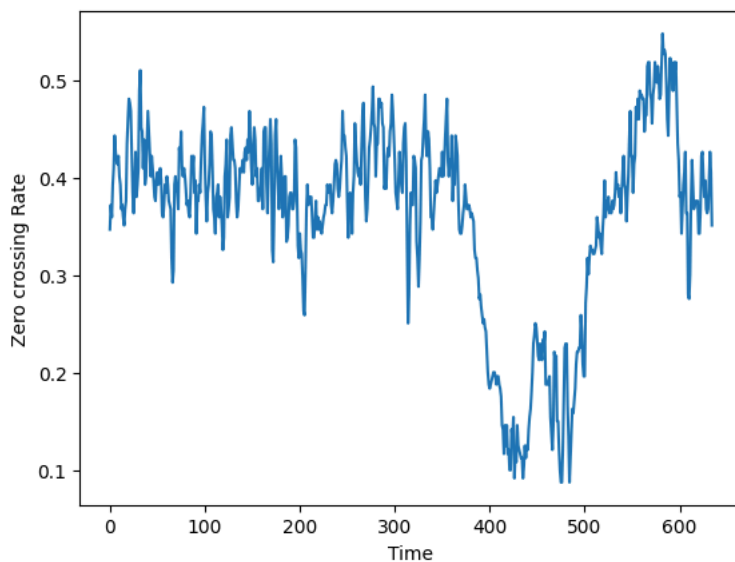
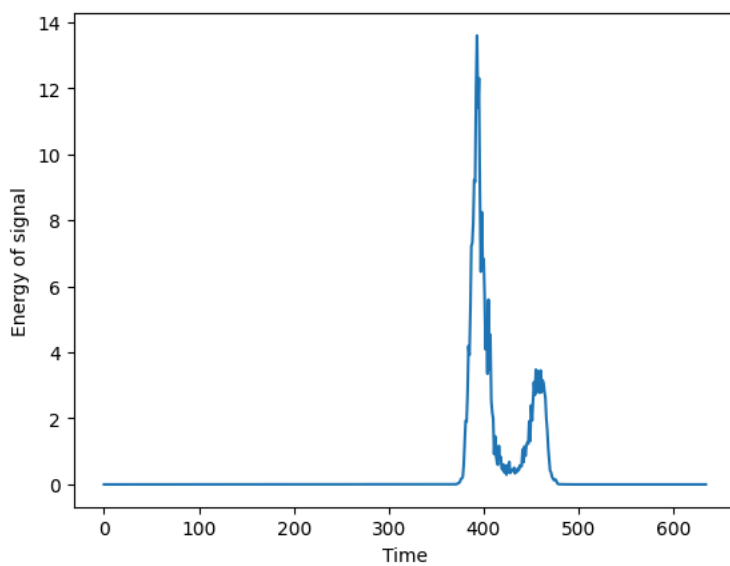
Output results of /content/drive/MyDrive/train/9b.wav

Reading WAV file /content/drive/MyDrive/train/9c.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice



Computing MFCC features

Output results of /content/drive/MyDrive/train/9c.wav

Training complete

Testing

```
1 import re
2 import scipy.signal as signal
3 import numpy as np
```

```

3 import numpy
4 import os
5
6 total = 0
7 acc = 0
8 for line in os.listdir("/content/drive/MyDrive/test"):
9     if line.endswith(".wav"):
10         total += 1
11         label=line[0]
12         print("Actual Label : ", label)
13         filename="/content/drive/MyDrive/test/"+line
14         #file reading
15         wavearr=wavread(filename)
16         #Signal pre-emphasis
17         wavearr_pre=pre_emphasis(wavearr,0.98)
18         #Select window function
19         winfunc = signal.hamming(240)
20         #Signal framing
21         n=enframe(wavearr_pre,240,80,winfunc)
22         #vad
23         na=vioceextrac(n)
24         #Extract features
25         feattest=mfcc(na,512)
26         #print(numpy.shape(feattest))
27
28         modelob=open("/content/drive/MyDrive/model1.txt")
29         print(' DTW \n')
30         i=0
31
32         dist1=[]
33         label=[]
34         mframes=modelob.readlines()
35         #print(len(frames))
36         for line1 in mframes:
37             i+=1
38             mline=line1.strip()
39             if i==1 and len(mline)==1:
40                 label.append(mline)
41
42                 featarr=[]
43                 nf=0
44
45                 if i>1 and len(mline)==1:
46                     #Regenerate numpy matrix and unify specifications
47                     feat=numpy.array(featarr,dtype=numpy.float64)
48                     feattrain=feat.reshape(nf,m)
49                     #print(feattrain)
50                     featarr=[]
51                     #Calculate the dtw distance from the test signal to each signal in the model
52                     dist1.append(score(feattrain,feattest))
53                     #print(nf,m)
54                     nf=0
55                     label.append(mline)
56
57                     if i!=len(mframes) and len(mline)>1:
58                         #print(mline)
59                         mline=mline.strip('[]')
60                         mline=mline.strip()
61                         #print(mline)
62                         mlinearr=re.split('\s+',mline)
63                         #print(len(mlinearr))
64                         #print(mlinearr)
65                         featarr.append(mlinearr)
66                         m=len(mlinearr)
67                         nf+=1
68
69                     if i==len(mframes) and len(mline)>1:
70                         #print(mline)
71                         mline=mline.strip('[]')
72                         mline=mline.strip()
73                         #print(mline)
74                         mlinearr=re.split('\s+',mline)
75                         #print(len(mlinearr))
76                         #print(mlinearr)
77                         featarr.append(mlinearr)
78                         m=len(mlinearr)
79                         nf+=1
80
81                         #Regenerate numpy matrix and unify specifications
82                         feat=numpy.array(featarr,dtype=numpy.float64)
83                         feattrain=feat.reshape(nf,m)
84                         #print(feattrain)
85                         featarr=[]
86                         #Calculate the dtw distance from the test signal to each signal in the model
87                         dist1.append(score(feattrain,feattest))
88                         #print(nf,m)

```

No. of FVs extracted from each model

```
86     if dist1 != []:
87         print('Min distance matching \n')
88         #Search for the smallest among all distances
89         labelnum=dist1.index(min(dist1))
90         print('Output \n')
91         #Output the label corresponding to the minimum distance signal as the recognition result.
92         print(label[labelnum])
93         print("_____")
94         if line[0] == label[labelnum]:
95             acc += 1
96 print("Test Accuracy :", (acc/total)*100, "%")
97 modelob.close
98 print('\n Recognition complete \n')
```

Actual Label : 2

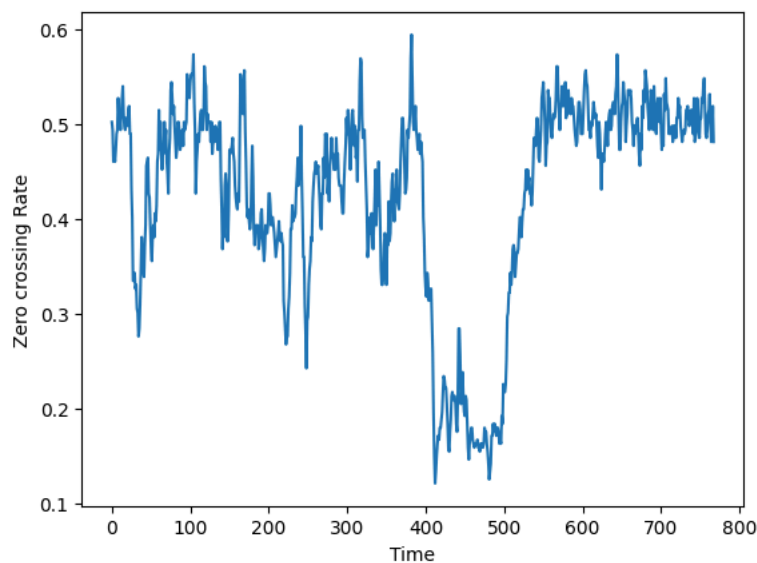
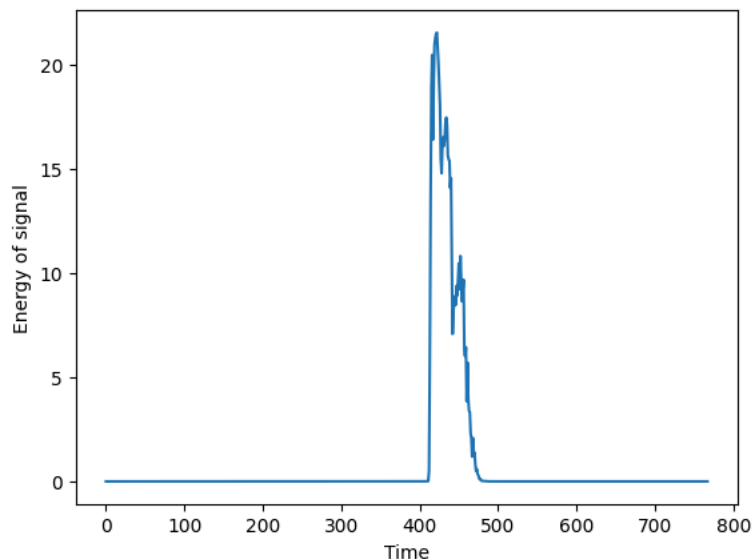
Reading WAV file /content/drive/MyDrive/test/2t.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice

```
<ipython-input-2-38ef7c5dd080>:25: DeprecationWarning: The binary mode of fromstring is deprecated, as it behaves surprisingly on
wave_data = np.fromstring(str_data, dtype=np.short)
<ipython-input-12-347b22e80d44>:19: DeprecationWarning: Importing hamming from 'scipy.signal' is deprecated and will raise an error in the future
winfunc = signal.hamming(240)
```



Computing MFCC features

DTW

Min distance matching

Output

2

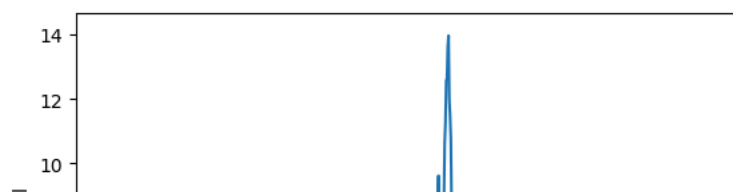
Actual Label : 1

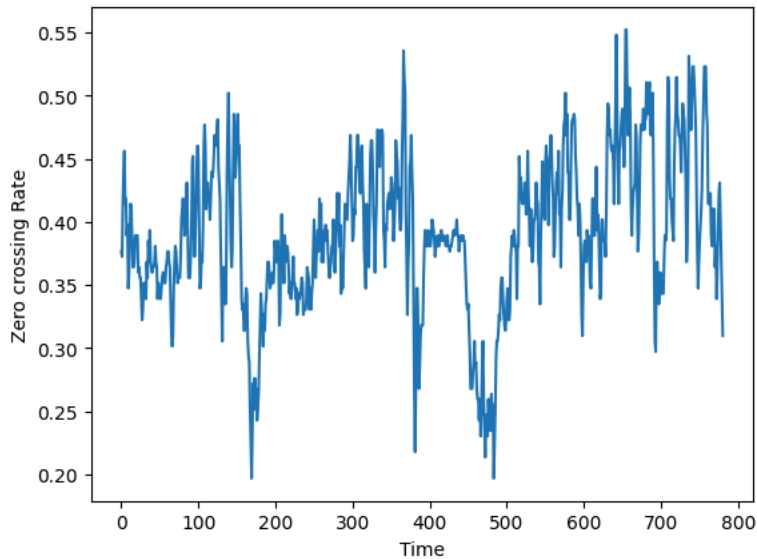
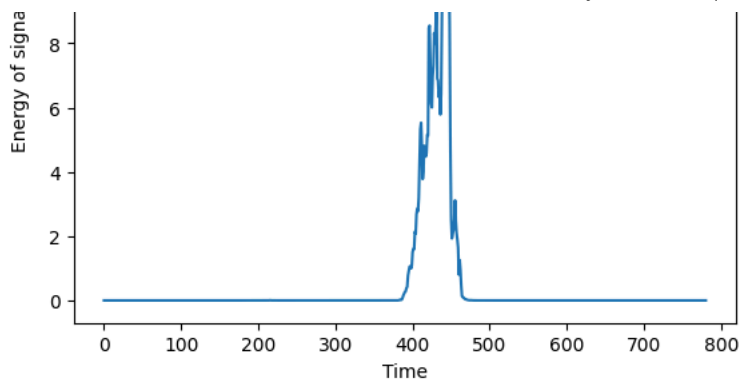
Reading WAV file /content/drive/MyDrive/test/1t.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice





Computing MFCC features

DTW

Min distance matching

Output

1

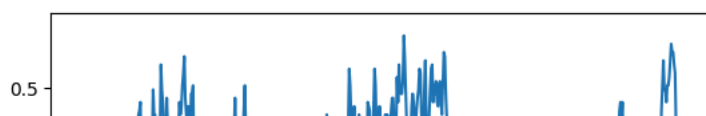
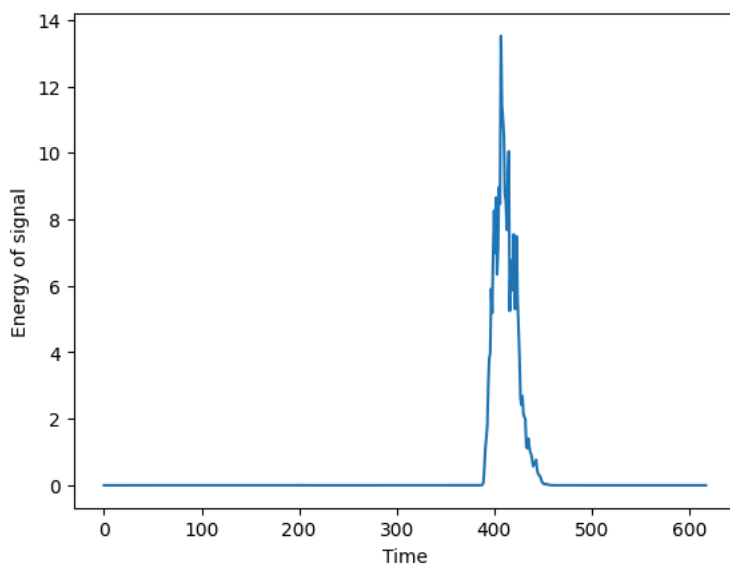
Actual Label : 0

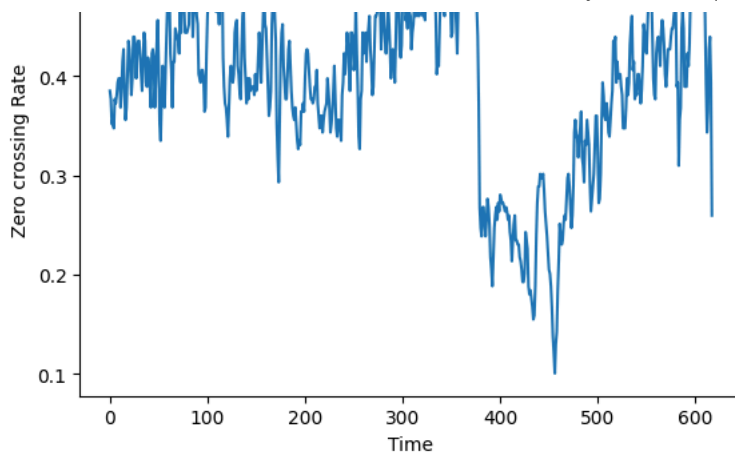
Reading WAV file /content/drive/MyDrive/test/0t.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice





Computing MFCC features

DTW

Min distance matching

Output

0

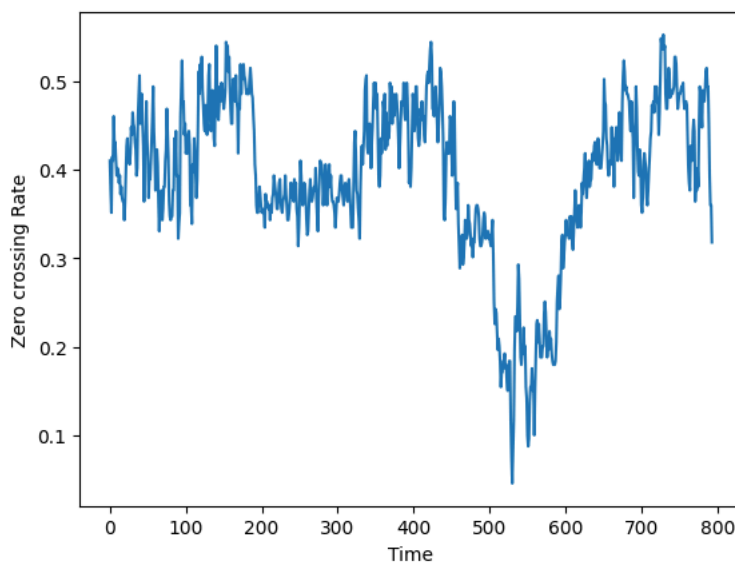
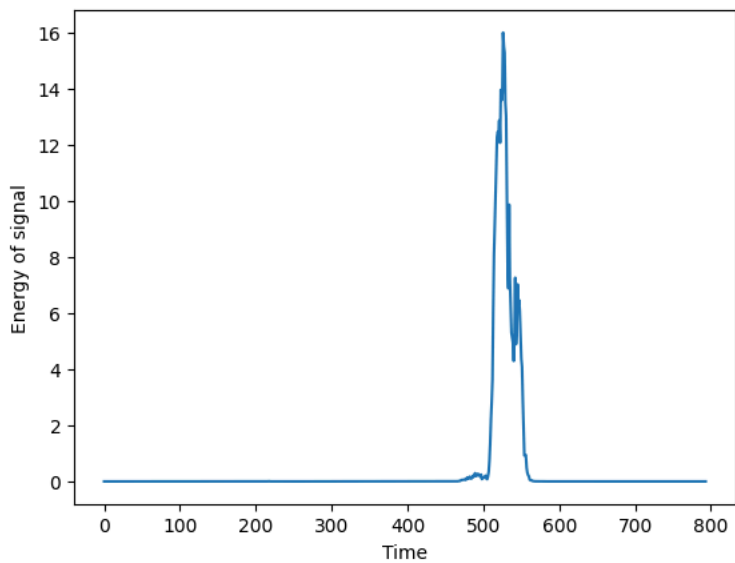
Actual Label : 4

Reading WAV file /content/drive/MyDrive/test/4t.wav ---

Pre emphasis

Signal framing and windowing

Extracting active voice



Computing MFCC features