

## Faculté Polytechnique



### SIGNAL PROCESSING'S PROJECT

How Useful Is Machine Learning ?  
A Speaker Classification Project

LEKEMO ATONFACK Rosabelle  
JIMENEZ MILLAN Sebastian



Under the direction of Mr DUTOIT and of Mr EL HADDAD

Academic Year: 2018-2019

## 0.1 Introduction

Machine learning is a kind of artificial intelligence based on the construction of computer programs that can understand, deduce and improve themselves through the experience of the data provided to it.

This project thus consist of taking over the concept of machine learning which is building systems based on data. In order to best acquire this concept, we will focus on the implementation rule-based systems using signal processing techniques and then compare them to machine learning-based techniques.

Within the limit of our study, we will therefore be focused on speaker classification, which means trying to recognize some speaker's voices.

## 0.2 SIGNAL PRE-PROCESSING

The libraries used for the functions :

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import scipy.signal as sig
4 from scipy.io.wavfile import read
5 from numpy import random as rnd
6 from scikit_talkbox_lpc import lpc_ref
```

Function that returns normalized data input

```
1 def normalize(data):
2     x=np.max(np.abs(data))
3     ampl=data/x
4     return ampl
```

Function that splits given that into frames

— signal = data to split  
— ms = size of each frame in ms  
— step = shifting step in ms

```
1 def frames(signal,ms,step):
2     wsize=ms*16
3     split=[]
4     etape=step*16
5     for i in range(0,len(signal),etape):
6         window=[]
7         for j in range(i,wsize+i,1):
8             if j<len(signal):
9                 window.append(signal[j])
10            if len(window)==wsize:
11                split.append(window)
12            else:
13                for k in range(wsize-len(window)):
14                    window.append(0)
15                split.append(window)
16    return split
```

# FEATURES EXTRACTION ALGORITHMS

## Signal Energy

Function which outputs the energy of a given signal, using the formula  $E = \sum |x(i)|^2$  :

```
1 def energy(signal):
2     Energy=sum(np.abs(np.array(signal))**2)
3     return Energy
```

## Pitch

### Autocorrelation-Based Pitch Estimation System

Function that estimates pitch for frames of a signal with autocorrelation, calculating the energy of each frame and classifying it as voiced or unvoiced depending on the threshold given. Returns energy, voiced as 1 unvoiced as 0, and frequency of each frame.

```
1 def acorrpitch(sound, threshold, wsize, step):
2     ns=normalize(sound)           #start normalizing signal
3     sp=frames(ns, wsize, step)    #divide it in frames
4     Energy=[]; voice=[]; f0=[]    #energy, voiced/unvoiced and frequency
5     for e in range(len(sp)):
6         en=energy(sp[e])          #calculate energy of the frame
7         Energy.append(en)         #append energy to energy vector
8     Energy=normalize(Energy)       #normalize energy
9     for i in range(len(sp)):
10        if Energy[i]>threshold:
11            voice.append(1)         #voiced frame==1
12            c=plt.xcorr(sp[i], sp[i], maxlags=320) #autocorrelation
13            pks=sig.find_peaks(c[1]) #find peaks on autocorrelation vector
14            p=[]                   #vector p which will have the real values, not the index
15            for j in range(len(pks[0])):
16                if pks[0][j]<287 or pks[0][j]>351: #frequency range
17                    p.append(c[1][pks[0][j]])      #append to vector p
18            s=np.sort(p)              #sort the peaks from lowest to highest
19            p2=s[-1]                  #second highest peak
20            for k in range(320):
21                if p2==c[1][k]:        #find position of second highest peak
22                    d=320-k            #highest peak is always in position 320
23                f=(16000)/d             #frequency=fe/N, N is distance between peaks
24                f0.append(f)           #append frequency to frequency vector
25            else:
26                voice.append(0)        #unvoiced frame==0
27                f0.append(0)           #frequency is 0 for unvoiced frames
28    return np.array(voice), np.array(f0)
```

## Cepstrum-Based Pitch Estimation System

Function that estimates pitch for frames of a window with cepstrum, calculating the energy of each frame and classifying it as voiced or unvoiced depending on the threshold given. Calculates frames on a rectangular window and on a hamming window. Returns energy, voiced as 1/unvoiced as 0, and frequency of each frame on a rectangular window and a hamming window

```
1 def cepstrumpitch(sound, threshold, wsize, step):
2     ns=normalize(sound) #start normalizing signal
3     sp=frames(ns, wsize, step) #divide it in frames
4     Energy=[]; voice=[]; f0=[]; hf0=[] #energy, voiced/unvoiced and frequency
5     hamm=np.hamming(wsize*16) #hamming window with size of a frame
6     for e in range(len(sp)):
7         en=energy(sp[e]) #calculate energy of the frame
8         Energy.append(en) #append energy to energy vector
9     Energy=normalize(Energy) #normalize energy
10    for i in range(len(sp)):
11        hf=hamm*sp[i] #hamming window on frame
12        if Energy[i]>threshold:
13            voice.append(1) #voiced frame==1
14            w=np.fft.fft(sp[i]) #FFT on rect window frame
15            wlog=np.log10(np.abs(w)) #log spectrum on rect window frame
16            hw=np.fft.fft(hf) #FFT on hamming window frame
17            hwlog=np.log10(np.abs(hw)) #log spectrum on hamming window frame
18            cepstrum=np.fft.ifft(wlog) #cepstrum of rect window frame
19            hcep=np.fft.ifft(hwlog) #cepstrum of hamming window frame
20            cut=cepstrum[50:500] #to cut from 50 to 500Hz
21            hcut=hcep[50:500] #cut in hamming window frame
22            for j in range(len(cut)):
23                if max(cut)==cut[j]: #find peak in cut vector
24                    f=j+50 #frequency is position of the max value
25            f0.append(f) #append frequency to frequency vector
26            for k in range(len(hcut)):
27                if max(hcut)==hcut[k]: #find peak in cut vector
28                    hf=k+50 #position of the max value
29            hf0.append(hf) #append frequency to frequency vector
30        else:
31            voice.append(0) #unvoiced frame==0
32            f0.append(0) #frequency is 0 for unvoiced frames
33            hf0.append(0) #hamming frequency vector
34    return Energy, np.array(voice), np.array(f0), np.array(hf0)
```

## Other functions

Function that selects a random sentence from a given gender either 'male' or 'female'

```
1 def selectsound(gender):
2     ab=rnd.choice(['a', 'b'])
3     if ab=='a':
4         rm0=rnd.randint(1,594)
5         selection=ab,rm0
6     if rm0<10:
7         fs, Mono=read('/Users/Sebas/Documents/UMONS/'+gender+'/wav/
arctic_a000'+str(rm0)+'.wav')
```

```

8         if rm0>=10 and rm0<100:
9             fs ,Mono=read ( '/Users/Sebas/Documents/UMONS/' +gender+' /wav/
arctic_a00 '+str(rm0)+' .wav' )
10            if rm0>=100:
11                fs ,Mono=read ( '/Users/Sebas/Documents/UMONS/' +gender+' /wav/ arctic_a0
'+str(rm0)+' .wav' )
12            else :
13                rm1=rnd.randint(1,540)
14                selection=ab,rm1
15                if rm1<10:
16                    fs ,Mono=read ( '/Users/Sebas/Documents/UMONS/' +gender+' /wav/
arctic_b000 '+str(rm1)+' .wav' )
17                if rm1>=10 and rm1<100:
18                    fs ,Mono=read ( '/Users/Sebas/Documents/UMONS/' +gender+' /wav/
arctic_b00 '+str(rm1)+' .wav' )
19                if rm1>=100:
20                    fs ,Mono=read ( '/Users/Sebas/Documents/UMONS/' +gender+' /wav/ arctic_b0
'+str(rm1)+' .wav' )
21            return fs ,Mono, selection

```

Function that select a random sound from a given gender 'male' or 'female', normalizes signal, splits it in frames and calculates energy.

Plots the signal in temporal domain and normalized energy by frames.

Returns normalized signal, energy and the number of the wav file

```

1 def randsound (gender , wsize , step ) :
2     fs ,Mono,sound=selectsound (gender )
3     t=np.linspace (0 ,(len (Mono)-1)/fs ,len (Mono))
4     Mono=normalize (Mono)
5     sp=frames (Mono, wsize , step )
6     E=[]
7     print ( 'Number of frames: ' +str (len (sp)))
8     for i in range (len (sp)) :
9         en=energy (sp [i])
10        E.append (en)
11    Energy=normalize (E)
12    plt.figure ()
13    plt.suptitle (gender+' '+sound[0]+str (sound [1]) , fontsize=16)
14    plt.subplot (211)
15    plt.plot (t, Mono)
16    plt.xlabel ( 'time ' )
17    plt.ylabel ( 'pressure ' )
18    plt.title ( 'Signal ' )
19    plt.subplot (313)
20    plt.plot (Energy)
21    plt.xlabel ( 'frame ' )
22    plt.ylabel ( 'energy ' )
23    plt.title ( 'Energy ' )
24    return Mono, Energy , sound

```

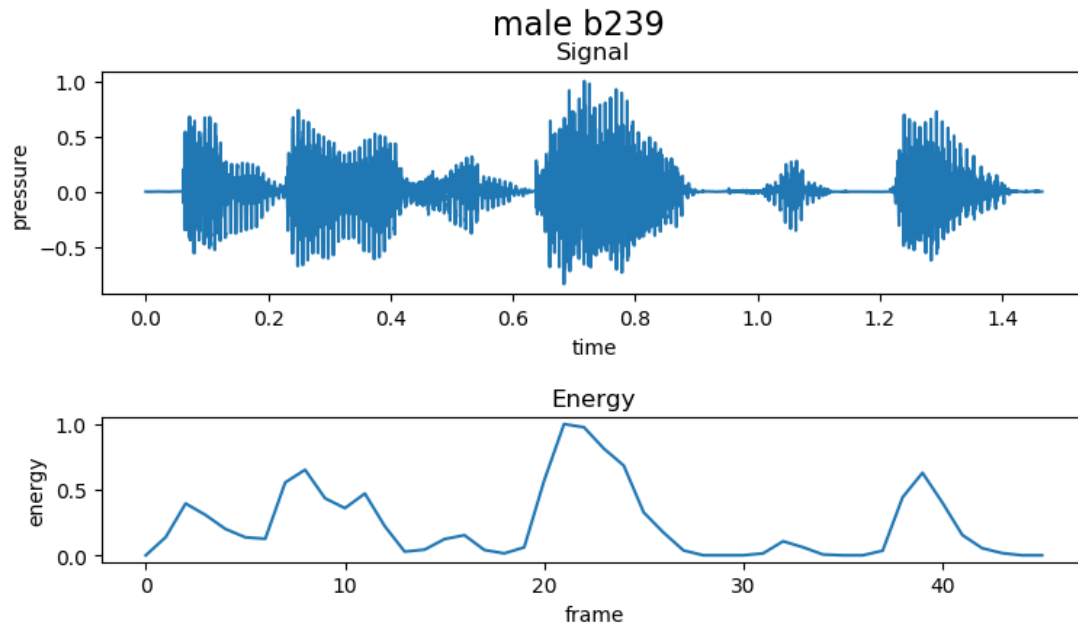


FIGURE 1 – Example of a random sound chosen from BLT with plots of temporal domain and energy, using the function randsound

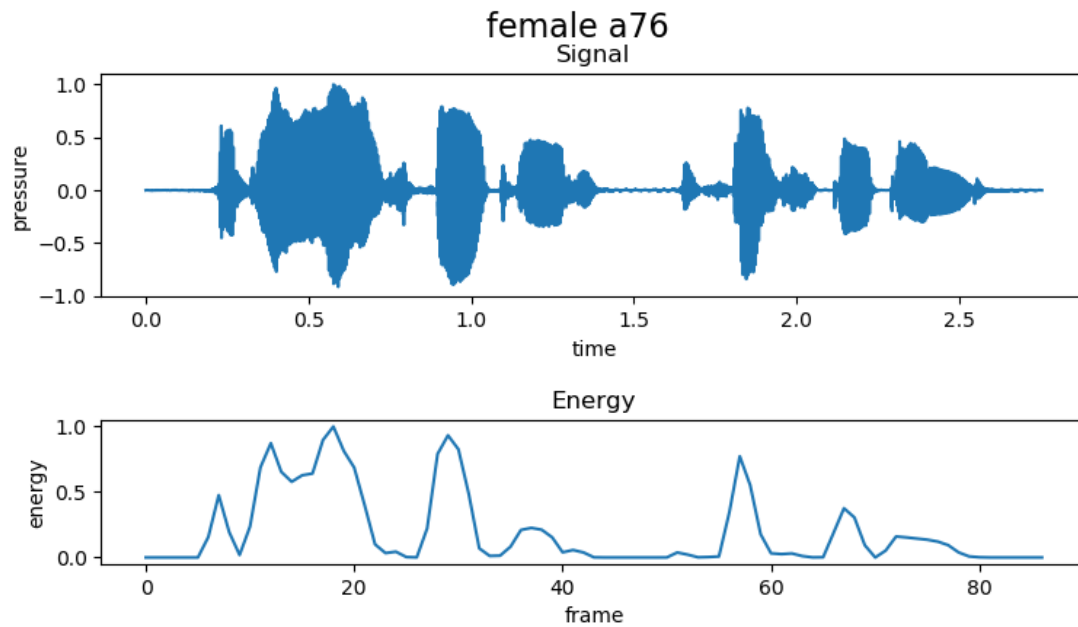


FIGURE 2 – Example of random sound chosen from SLT with plots of temporal domain and energy, using the function randsound

## Formants

Function that calculates the formants of the frames of a given signal, with *wsz* being the size of each frame in *ms* and the step the shifting step.

Returns a list with the formants on each frame.

```
1 def formants(signal, wsz, step):
2     sp=frames(signal, wsz, step)           #divide signal in frames
3     hamm=np.hamming(wsz*16)               #hamming window with size of a frame
4     coeff=[]; roots=[]; raiz=[]; angle=[]; freq=[];
5     for i in range(len(sp)):
6         filtered=sig.lfilter([1, -0.67], [1], sp[i]) #high pass filter on frame
7         hf=hamm*filtered                   #hamming window on frame
8         coeff.append(lpc_ref(hf, 10))      #LPC coefficients of each frame
9         roots.append(np.roots(coeff[i]))   #roots of the LPC
10    for j in range(len(roots)):
11        l=len(roots[j])                   #number of roots
12        erase=[k for k in range(1, l, 2)]  #delete half of roots
13        r=np.delete(roots[j], erase)
14        raiz.append(r)                   #save roots
15        angle.append([])                 #this empty lists are to append here the formants
16        freq.append([])                  #and then append the list of formants of each frame
17        for z in range(int(l/2)):         #angle for each root
18            angle[j].append(np.arctan(np.imag(raiz[j][z])/np.real(raiz[j][z])))
19            fn=abs((angle[j][z]*16000)/(2*np.pi)) #calculate frequency
20            freq[j].append(fn)            #save frequency on formant list
21        freq[j].sort()                   #sort frequencies from lowest ot highest
22    return np.array(freq)
```

### 0.2.1 MFCC

Function that outputs the power spectrum of the signal, the energy and finally the MFCCs which are the coefficients representing the vocal tract activities.

```
1 import scipy as sp
2 from scipy.fftpack import dct
3
4 def MFCC(signal, ms, step, f0):
5
6     mfccs=[]
7     P=[]
8     y=[]
9     const_a=0.97
10    a=[]
11    a[0]=1
12    a[1]=0
13    b=[]
14    b[0]=1
15    b[1]=-const_a
16    NTFD=512
17    fbanks=[]
18
19    sigFiltered=scipy.signal.lfilter(a, b, signal)
20    split=frames(sigFiltered, ms, step)
21
```



```

22     for i in range(len(split)):
23         fenetre=scipy.signal.hamming(len(split[i]))    #each window taken has
the same lenght as the frame
24         yframe=split[i]*fenetre                        #apply each window to each
frame
25         #y.append(yframe)
26         P.append((numpy.linalg.norm(numpy.fft(yframe))^2)/NTFD)
27
28
29     pow_frames=frames(P,ms,step)
30
31     fbanks=filter_banks(pow_frames, f0, nfilt = 40, NFFT=512)
32     y=scipy.fftpack.dct(fbanks, type=2, axis=1, norm='ortho')
33
34
35     en=y[0]      #the first output is the energy
36     for i in range(1,12,1):
37         mfccs[i-1]=y[i]    #the 12 oders are the mfccs
38
39
40     return P,en,mfccs

```

## Rule Based System

To implement the classification system, we create a set of rules for classifying the sentences of the speakers. For the parameters, we use size of a frame =  $50ms$ , shifting step =  $16ms$  and threshold = 0.10. The threshold chosen at 0.10, is after observing the plotted energy along with the temporal plot of the sound, finding that the sounds we want to find frequencies from are mostly above this level.

We select randomly 15 sentences from each speaker, we calculate the average frequency on each sentence selecting only the voiced frames. We then calculate the average frequency for each speaker. Finally the average frequency of both speakers is calculated.

```

1  'Rule Based System'
2
3
4  import speaker as spk
5  import numpy as np
6
7
8  #Parameters
9  wsize=50          #size of each frame in ms, real size is wsize*16
10 step=16           #shifting step in ms, real step is step*16
11 threshold=0.10    #threshold for clasifying frame as voiced or unvoiced
12
13 #Lists for the selected sounds and extracted features
14 speaker_male=[]
15 speaker_female=[]
16 num_male=[]
17 num_female=[]
18 voice_male=[]
19 voice_female=[]
20 freq_male=[]

```

```

21 freq_female=[]
22
23 #Select 15 random sentences from each speaker
24 for i in range(15):                                     #Sentences
25     MonoM,energiaF ,numeroM=spk.randsound( 'Male',wsize ,step)    #Male
26     MonoF,energiaF ,numeroF=spk.randsound( 'Female',wsize ,step)  #Female
27     speaker_male.append(MonoM)    #Signal of selected male sentence
28     speaker_female.append(MonoF)  #Signal of selected female sentence
29     num_male.append(numeroM)      #File number for male sentences
30     num_female.append(numeroF)    #File number for female sentences
31     voiceM,frequencyM=spk.acorrpitch(MonoM,0.10,wsize ,step) #autocorrelation
32     voiceF,frequencyF=spk.acorrpitch(MonoF,0.10,wsize ,step) #autocorrelation
33     voice_male.append(voiceM)     #voiced or unvoiced classification
34     voice_female.append(voiceF)    #for each frame
35     freq_male.append(frequencyM)   #List of frequencies for male sentences
36     freq_female.append(frequencyF) #list of frequencies for male sentences
37
38 #Will calculate average frequency for each sentence
39 average_freq_male=[]
40 average_freq_female=[]
41
42 for i in range(15):
43     f=0; v=0                                             #frequency , number of voiced frames
44     l=len(voice_male[i])                                #number of frames of the signal
45     for j in range(l):
46         if voice_male[i][j]==1:                        #only select the voiced frames
47             f=f+freq_male[i][j]                        #sum of all frequencies
48             v=v+1                                       #number of frequencies summed
49     average_freq_male.append(f/v)                       #average frequency
50
51 #same for female sentences
52 for i in range(15):
53     f=0; v=0
54     l=len(voice_female[i])
55     for j in range(l):
56         if voice_female[i][j]==1:
57             f=f+freq_female[i][j]
58             v=v+1
59     average_freq_female.append(f/v)
60
61 #Calculate total average frequency for male and female sentences
62 tot_average_male=np.mean(average_freq_male)
63 tot_average_female=np.mean(average_freq_female)
64
65 print('Average frequency of male sentences: '+str(tot_average_male))
66 print('Average frequency of female sentences: '+str(tot_average_female))
67
68 #Calculate average frequency between both speakers
69 average_ensemble=(tot_average_male+tot_average_female)/2
70
71 print('Average frequency of all sentences: '+str(average_ensemble))

```

## Classification System 1

For the first rule based system, we work on the average frequency obtained and we use this as the threshold to classify new randomly chosen sentences. If the average frequency of a sentence is lower than this threshold, we classify it as male, and if it is higher we classify it as female. We then calculate the accuracy of the system.

```
1 'System 1'
2
3 #System 1's threshold will rely on the average frequency of all the sentences
4 #After calculating the average frequency of all sentences in both speakers ,
5 #it will classify a sentence as BLT(male) speaker if the average frequency
6 #is below the total average frequency , and as SLT(female) if it is higher
7
8 #New selection of sentences
9 sentence_male_1=[]
10 sentence_female_1=[]
11 voice_male_1=[]
12 voice_female_1=[]
13 freq_male_1=[]
14 freq_female_1=[]
15
16 for i in range(15):
17     MonoM,energiaF ,numeroM=spk.randsound( 'Male',wsize,step)    #male
18     MonoF,energiaF ,numeroF=spk.randsound( 'Female',wsize,step)  #female
19     sentence_male_1.append(MonoM)                                #Signal
20     sentence_female_1.append(MonoF)                             #Lists
21     voiceM,frequencyM=spk.acorrpitch(MonoM,0.10,wsize,step) #Pitch estimation
22     voiceF,frequencyF=spk.acorrpitch(MonoF,0.10,wsize,step) #by autocorrelation
23     voice_male_1.append(voiceM)                                #voiced
24     voice_female_1.append(voiceF)                             #or unvoiced
25     freq_male_1.append(frequencyM)                             #Frequencies
26     freq_female_1.append(frequencyF)                           #lists
27
28 #Will calculate average frequency for each sentence
29 average_freq_male_1=[]
30 average_freq_female_1=[]
31
32 for i in range(15):
33     f=0; v=0
34     l=len(voice_male_1[i])
35     for j in range(l):
36         if voice_male_1[i][j]==1:
37             f=f+freq_male_1[i][j]
38             v=v+1
39     average_freq_male_1.append(f/v)
40
41 #same for female sentences
42 for i in range(15):
43     f=0; v=0
44     l=len(voice_female_1[i])
45     for j in range(l):
46         if voice_female_1[i][j]==1:
47             f=f+freq_female_1[i][j]
48             v=v+1
49     average_freq_female_1.append(f/v)
```

```

50
51 #Classification lists
52 classification_male_1=[]
53 classification_female_1=[]
54
55
56 #go through all sentences
57 for i in range(15):
58     if average_freq_male_1[i]<=average_ensemble:      #if freq < average,
59         classification_male_1.append(1)                #1==true
60     else:                                              #otherwise 0==false
61         classification_male_1.append(0)
62     if average_freq_female_1[i]>average_ensemble:      #for female
63         classification_female_1.append(1)              #the condition
64     else:                                              #freq > average
65         classification_female_1.append(0)
66
67 #accuracy of male and female classification , and system total
68 accuracy_male_1=np.mean(classification_male_1)
69 accuracy_female_1=np.mean(classification_female_1)
70 accuracy1=np.mean(classification_male_1+classification_female_1)
71
72 print('System 1')
73 print('Accuracy for male speaker: '+str(accuracy_male_1))
74 print('Accuracy for female speaker: '+str(accuracy_female_1))
75 print('System Accuracy: '+str(accuracy1))

```

These are some results from the classification of the first system, testing on 15 sentences of each speaker :

```

Accuracy for male speaker: 0.8666666666666667
Accuracy for female speaker: 1.0
System Accuracy: 0.9333333333333333

```

FIGURE 3 – Results 1

```

Accuracy for male speaker: 0.9333333333333333
Accuracy for female speaker: 1.0
System Accuracy: 0.9666666666666667

```

FIGURE 4 – Results 2

## 0.3 Conclusion

In this project, our study was about the implementation of rule-based systems in order to figure out the machine learning concept. We therefore focused on speaker classification which meant recognizing speaker's voices.

All a long this project, we first of all had to pre-process the signals by normalizing them and splitting them into frames. Then, we had to extract some of their features which were their energy and an estimation of the pitch. After that, we had to characterize the shape of the signal at a certain instant by outputting the resonances frequencies of the vocal tract. And finally, we define rules to build a speaker recognition system which are, in this case, rules to classify males and females.

Arrived to the end of this project, we were certainly able to classify speakers. An important point to remind is that, our rules classification had been based on our data and will then obviously showed a good accuracy. This might not be the case on others data that have never been seen by our algorithm.