

App development with audio applications from m-file to app

Daniel JANKOVIC

August 12, 2015

1 Part 1

Voice Activity Detector (VAD) is a technique used in signal processing to detect the presence of human voice in a signal. It can be an energy detector that indicates speech when the energy of the filtered signal exceeds a predefined threshold. Considered an important technology in speech based communication, today there are various types of applications that use it. Therefore a wide variety of VAD algorithms have been developed to provide the needed features.

There are different kind of stand-alone commercial baby monitors on the market today. From the most basic, that use one-way radio communication, to advance two-way communication monitors that use signal processing to transmit audio when a predefined threshold has been reached. It is also possible to find baby video monitors that broadcast both audio and video when the sensors notice movement. Since most of the monitor applications rely on radio signals to communicate between the units there is a probability that the signal will weaken or possibly not even reach the receiver because it needs to pass through multiple walls of varying thickness. The signal could also be effected by other applications. As the stand-alone monitor focuses on reliability (among other important sale strategies), little is known about the security features. It is possible to assume that the communication is unencrypted, at least in some products, and therefore introduces a potential risk for intrusion of peoples privacy.

To resolve the issues brought up above, an application such as the *Baby Activity Detector* (BAD) can be made more portable, versatile and secure with the help of todays smart-phone technology and VAD. There are many VAD algorithms to choose from and they all have their strengths and weaknesses. Complex algorithms such as Linear Predictive coding (LPC), mel-frequency cepstrum (MFC) are very powerful but quite difficult to grasp and also to implement, they can be considered out of scope for this course. The following VAD algorithms are easy to implement and can be, when combined, quite robust for the task of a basic BAD. The simple short-time energy algorithm calculates the energy levels for each frame to detect voice, unvoiced or silenced regions. Voiced regions will have higher energy levels, however, the algorithm does not

take unwanted noise into account which means that we can have false indication of voice detection. In order to remove the noise from the signal, spectral subtraction can be preformed. In the case of BAD, the threshold needs to be adjusted so that unforeseeable sound is not interpret as the infants cry. Zero-crossing rate (ZCR), is the rate at which a signal changes from plus to minus and back. The higher the rate the higher the frequency which indicates possible voice activity. According to [1] the cry sound that an infant makes has a fundamental frequency of 250-600 Hz (pitch). To be able to use ZCR together with the information above, it is necessary to extract the pitch from the signal in order to match the frequency interval.

The main task of BAD is to detect infant activity, an alternate algorithm is proposed in [1]. It describes an cry detection algorithm that is build up by three main stages. i) *VAD*, a statistical model-based detector [5] is used for detecting sections with sufficient audio acitivity. It also helps to reduce the power consumption. ii) *Classification*, uses k-nearest neighbours (k-NN) algorithm [6] to label each frame as either 'cry' (1), close enough, or 'no cry' (0). iii) Post-processing, validation stage in order to reduce false-negative errors. The idea of having devoted algorithm to detect infant cry is a winning concept for a BAD application, according to the authors it even had promising results in low SNR. Despite simplicity of the algorithm many of the features required to implement were mentioned earlier to be out of scope for this project.

An algorithm that might be of interest [3] suggests a new approach to speech enhancement, without the help of VAD technology. The signal is divided into multiple sub-bands and an noise floor level estimate is calculated simultaneously as the short-time average. The goal is to boost the sub-bands with high Signal-to-Noise Ration (SNR) instead of to suppress the lower. This algorithm has great potential to reduce the noise levels when analyzing incoming signals to the BAD application.

A quick search on the net gives significant amount of hits for smart-phone based BAD applications. The techniques vary, from bluetooth to Wi-Fi and 3-/4G solutions. The award winning application *Baby Monitor 3G* [7] is a feature rich cross platform application that solves most the of issues brought up in this report. It supports both Wi-Fi and 3G/LTE networks, ability to transfer high quality live video, adjustable microphone sensitivity, talk-back functionality and guarantees both reliability and privacy. It can be assumed that the Android based BAD application *Dormi* [8] offers similar features as Baby Monitor 3G, even if Baby Monitor 3G's feature list provides barely any deeper information. It is noteworthy that *Domri* have both *Smart noise detection* and *Adaptive audio enhancement* as sales pitch, which from a engineering point of view is very attractive.

Following is a purposed algorithm for an BAD application

```

while(true){
    % Record audio can place it into the register %
    Get frame from the register
    Divide the signal into different sub-bands with FFT
    Calculate the total short-time energy average
    Calculate noise level for each sub-bands
    if energy average above threshold
        Calculate the gain for each sub-band
        Boost the sub-band with high SNR
        % Extract fundamental frequency (pitch)
        Count the ZCR under 1 seconds
        if the ZCR is within 250-600 Hz
            Possible infant activity detected!
            (Occurs only first time, and needs to be reseted)
    % Send frames %
        Start broadcasting the sound to the receiver
    end
    else
        Reset ZCR
        Dismiss and get next frame from register
    end
}

```

2 Part 2

Three various baby and noise sounds were provided. The sample frequency for all of the recorded files are 8 kHz. In Figure 1 and Figure 3 the frequency spectrum is plotted for baby respective noise sounds. Because the power level is unproportionally low for a specific audio recording an enhanced plot is displayed in Figure 2. From the plots it can be seen that the baby sound spectrum is between ~ 300 -2500 Hz and the noise sound spectrum is below 200 Hz and between ~ 1300 -2100 Hz. To get the wanted effect from the advanced algorithm a bandpass filter with a band between 300-1300 Hz is the optimal solution.

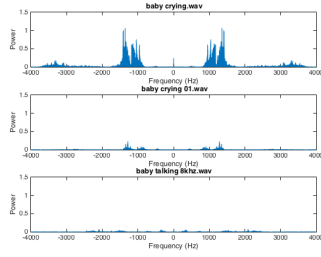


Figure 1: Frequency spectrum for baby sound

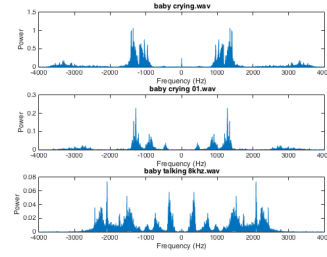


Figure 2: Frequency spectrum for baby sound, scaled y-axis

A quick and simple way to implement a BAD algorithm is to measure the power from the input signal. Mathematically described, the power equation is given by the following formula:

$$P(n) = \frac{1}{N} \sum_{k=0}^{N-1} x^2(n-k)$$

where $P(n)$ is a vector. Since the BAD application is calculating the power of the input signal in real-time, this equation becomes impossible to implement. Hence, the usage of the *recursive averaging* algorithm. It is, in perspective to memory usage, a light and hardware friendly algorithm that calculates the average power level of the signal. Given the equation below,

$$P(n) = \alpha P(n-1) + (1-\alpha)x^2(n)$$

instead of performing the calculation for one sample at the time, a block of samples (simulating an interval) can be calculated and summed. The old result is saved to be used as $P(n-1)$. The α is a constant between 0 and 1 and is related to the following formula:

$$\alpha = \frac{1}{T_s F_s}$$

Despite that the α can be calculated, the get the best results further tweaking and testing is required. By having two different α values, when the power

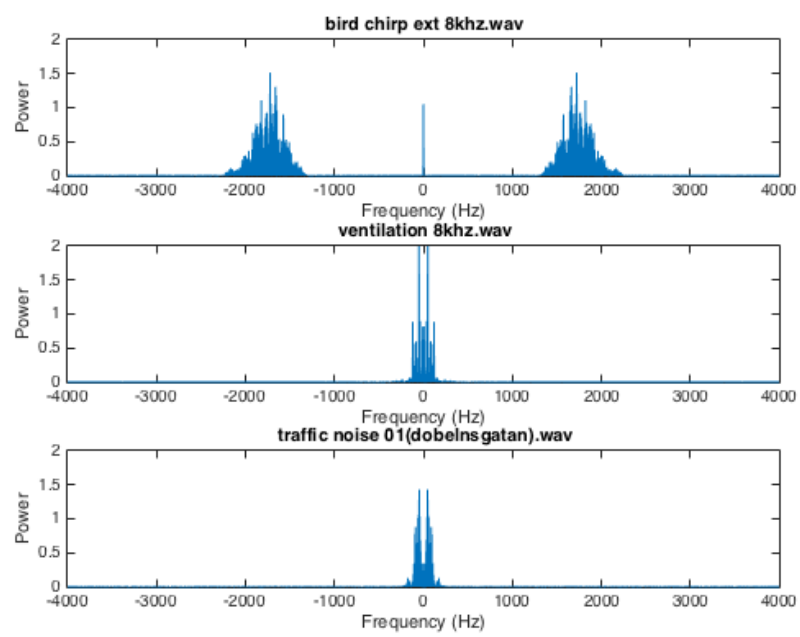


Figure 3: Frequency spectrum for noise sound files

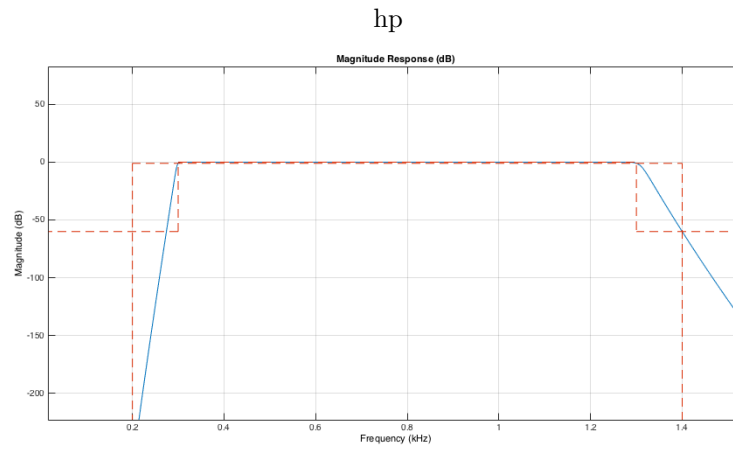


Figure 4: Butterworth 300-1300 Hz bandpass filter

is increasing respective decreasing, the algorithm can response faster and the decrease slower. That way, hard spikes are avoided and the average power is slowly fading, causing quicker reaction if high power signals enter the system.

The advanced algorithm is based upon the on the same, *recursive averaging*, algorithm but before calculating the average power the signal is first filtered through a Butterworth bandpass filter to remove unwanted frequencies. With the noise frequencies suppressed, the average power will focus on baby sound frequencies only.

3 Appendix

```
close all
clear all

% Read the audio files and extract the sample rate
%fileName1=('baby-signals/baby-crying.wav');
fileName1=('noise-signals/bird.chirp_ext.8khz.wav');
[xOrig1, fsOrig1] = audioread(fileName1);
%fileName2=('baby-signals/baby-crying-01.wav');
fileName2=('noise-signals/ventilation.8khz.wav');
[xOrig2, fsOrig2] = audioread(fileName2);
%fileName2=('baby-signals/baby-talking.8khz.wav');
fileName2=('noise-signals/traffic-noise-01(dobelnskatan).wav');
[xOrig3, fsOrig3] = audioread(fileName2);

% disp('Baby-crying'), disp(fsOrig1);
% disp('Baby-crying-01'), disp(fsOrig2);
% disp('Baby-talking.8khz'), disp(fsOrig3);

% Get the length of the audio file
m1 = length(xOrig1);
m2 = length(xOrig2);
m3 = length(xOrig3);
% Find a proper transform length
n1 = pow2(nextpow2(m1));
n2 = pow2(nextpow2(m2));
n3 = pow2(nextpow2(m3));

% Preform FFT on the different audio files
% and rearrange so that they are zero-centered
% First audio file to process
y1 = fft(xOrig1,n1);
f1 = (0:n1-1)*(fsOrig1/n1);
power1 = y1.*conj(y1)/n1;

y1c = fftshift(y1);
fs1c = (-n1/2:n1/2-1)*(fsOrig1/n1);
power1c = y1c.*conj(y1c)/n1;

% Second audio file to process
y2 = fft(xOrig2,n2);
f2 = (0:n2-1)*(fsOrig2/n2);
power2 = y2.*conj(y2)/n2;

y2c = fftshift(y2);
fs2c = (-n2/2:n2/2-1)*(fsOrig2/n2);
power2c = y2c.*conj(y2c)/n2;

% Third audio file
y3 = fft(xOrig3,n3);
f3 = (0:n3-1)*(fsOrig3/n3);
power3 = y3.*conj(y3)/n3;

y3c = fftshift(y3);
```

```

fs3c = (-n3/2:n3/2-1)*(fsOrig3/n3);
power3c = y3c.*conj(y3c)/n3;

% Plot the all FFT into one figure
figure
plot1 = subplot(311)
plot(fs1c,power1c)
xlabel('Frequency (Hz)')
ylabel('Power')
title('\bf bird chirp ext 8khz.wav')
plot2 = subplot(312)
plot(fs2c,power2c)
xlabel('Frequency (Hz)')
ylabel('Power')
title('\bf ventilation 8khz.wav')
plot3 = subplot(313)
plot(fs3c,power3c)
xlabel('Frequency (Hz)')
ylabel('Power')
title('\bf traffic noise 01(dobelnskatan).wav')

%linkaxes([plot1,plot2,plot3],'y');

% fax.bins1 = [0:N1-1];
% fax.bins2 = [0:N2-1];
% fax.bins3 = [0:N3-1];
%
% fax.Hz1 = fax.bins1*fsOrig1/N1;
% fax.Hz2 = fax.bins2*fsOrig1/N2;
% fax.Hz3 = fax.bins3*fsOrig3/N3;
%
% figure(1)
% subplot(311)
% plot(fax.Hz1, abs(fft(xOrig1)));
% subplot(312)
% plot(fax.Hz2, abs(fft(xOrig2)));
% subplot(313)
% plot(fax.Hz3, abs(fft(xOrig3)));

% segmentlen=100;
% noverlap=90;
% NFFT=128;

%spectrogram(xOrig1,segmentlen,noverlap,NFFT,fsOrig1,'yaxis')
%cceps(xOrig1)

```


References

- [1] R. Cohen, Y. Lavner, *Infant Cry Analysis and Detection*, 2012
- [2] E. Verteletskaya, K. Sakhnov, *Voice Activity Detection for Speech Enhancement Applications*, 2010
- [3] N. Westerlund, M. Dahl, *Speech Enhancement using an Adaptive Gain Equalizer*, 2003
- [4] R. Narayanam, *An Efficient Peak Valley Detection based VAD Algorithm for Robust Detection of Speech Auditory Brainstem Responses*, 2013
- [5] J. Sohn, N.S. Kim, W. Sung, *A Statistical Model-Based Voice Activity Detection*, 1999
- [6] O. Sutton, *A Statistical Model-Based Voice Activity Detection*, 2012
- [7] TappyTaps, <https://www.babymonitor3g.com>,
- [8] Sleekbit, <http://dormi.sleekbit.com/index.html>,
- [9] D. Jankovic, M. Johansson, M. Lichota, *Adaptive Gain Control in Digital Signal Processors*, 2015