

SUVING: AUTOMATIC SILENCE /UNVOICED/VOICED CLASSIFICATION OF SPEECH

Mark Greenwood, Andrew Kinghorn

Department of Computer Science, University of Sheffield
Regent Court, 211 Portobello St, Sheffield S14DP, UK
{u7mag, u7awrk}@dcs.shef.ac.uk

ABSTRACT

This paper is concerned with labelling sections of speech samples based on whether they are silence, voiced or unvoiced speech. The labelling is done using calculations over the speech samples; zero crossing and short-term energy functions. These functions complement each other and as such can be used more accurately together to label the parts of speech. The results of applying these functions to ten speech samples are compared to the result of the same samples having been manually labelled, to produce a percentage accuracy for each speech file. This study found that the average percentage accuracy of the algorithms implemented, over all ten-speech samples was about 65%, and concludes that the accuracy could be slightly improved through the use of a more accurate windowing function.

1. INTRODUCTION

A classification of speech into voiced or unvoiced sounds provides a useful basis for subsequent processing, for example fundamental frequency estimation, formant extraction or syllable marking. A three-way classification into silence/unvoiced/voiced (hence the title, SUVing) extends the possible range of further processing to tasks such as stop consonant identification and endpoint detection for isolated utterances. Strictly, speech sounds such as voiced fricatives (e.g. "z") can have characteristics of both voiced and unvoiced sources simultaneously which makes classification more difficult, so for the purposes of this study we will assume that a 3-way classification is sufficient for the needs of any further processing.

We approach this problem of SUVing from two different directions: zero crossings and short-term energy. These two methods compliment each other well, and prevent us from having to rely heavily on one single method to label the different parts of speech.

We do not provide an interface for viewing the results of this study instead the results of the SUVing are stored in a format that allows them to be loaded into the *slt* tool which is included as part of [2].

2. SUVING USING ZERO CROSSINGS

The notion of zero-crossings is defined to be:

"the number of times in a sound sample that the amplitude of the sound wave changes sign"

For a 10ms sample of *clean* speech, the zero-crossing rate is approximately 12 for voiced speech and 50 for unvoiced speech [1]. For clean speech the zero-crossing rate should also be useful for detecting regions of silence, as the zero-crossing rate should be zero.

Unfortunately, very few sound samples are recordings of perfectly clean speech. This means that often there is some level of background noise, that interferes with the speech, meaning that silent regions actually have quite a high zero-crossings rate as the signal changes from just one side of zero amplitude to the other and back again. For this reason a tolerance threshold is included in the function that calculates zero-crossings to try and alleviate this problem. The thresholds work by removing any zero-crossings, which do not both start and end a certain amount from the zero value. In this study we have used a threshold of 0.001. This means that any zero-crossings that start and end in the range of x , where $-0.001 < x < 0.001$, are not included in the total number of zero-crossings for that window. This enables us to filter out most of the zero-crossings that occur during silent regions of the sample due only to background noise.

In this study to calculate zero-crossings we used a 10ms non-overlapping rectangular window. This does not produce such good zero-crossing results as an overlapping hamming window would, but since we are not interested in the fine details, this method works well when used to SUV a speech sample.

3. SUVING USING SHORT-TERM ENERGY

Short-term energy allows us to calculate the amount of energy in a sound at a specific instance in time, and is defined in Equation 3-1.

$$E_n = \sum_{m=n-N+1}^n (x(m)w(n-m))^2$$

Equation 3-1: Short-Term Energy (w is the window, n is the sample that the window is centered on, and N is the window size [1]).

Unfortunately, unlike zero-crossings there are no standard values of short-term energy for specific window sizes. Short-term energy is purely dependent upon the energy in the signal, which changes depending on how the sound was recorded. For example, if a person is recorded saying the same phrase twice, one while whispering and once while shouting, then the short-term energy values will be vastly different, although the zero-crossing values should be roughly the same. This means that you have to inspect the recorded speech files to determine at what level to make the distinction between voiced and unvoiced speech.

There is one thing that is standard though, and that is that short-term energy is higher for voiced than un-voiced speech, and should also be zero for silent regions in clean recording of clean speech.

In a similar way to zero-crossings we calculate the short-term energy using a 10ms non-overlapping rectangular window. This, again, is not as accurate as using an overlapping hamming window but it is adequate for the SUV labelling of speech.

4. SUVING USING BOTH METHODS

From the descriptions of the methods that are used to SUV label a speech signal, in this study, it should be clear that the two methods compliment each other well.

For voiced speech short-term energy is high and zero-crossings are low, and for un-voiced speech the opposite is true, short-term energy is low and zero-crossings are high. This can be seen clearly in Figure 4-1

In a perfect world, all speech samples would be clean and then Table 4-1 could be used to classify the speech as silence, un-voiced or voiced.

Zero-Crossings	Short-Term Energy	Label
approx. 12	High	Voiced
approx. 50	Low	Un-Voiced
0	0	Silence

Table 4-1: Perfect world labelling scheme.

Unfortunately sampled speech is never perfectly clean, usually containing some level of background noise, and so the labelling scheme in Table 4-2 is used in this study to label the speech samples.

Another problem, apart from background noise is that it is often difficult to detect silent regions of speech samples due to the fact that the short-term energy for a breath can quite easily be confused with the short-term energy of a fricative sound [3].

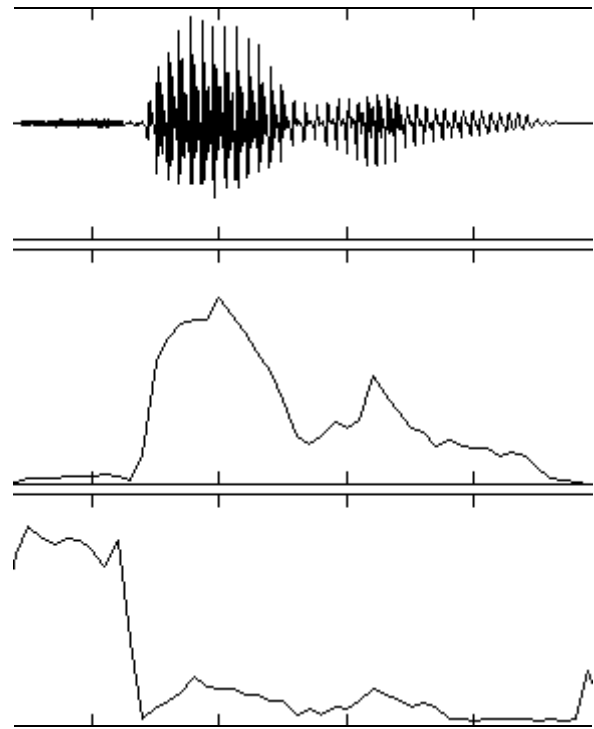


Figure 4-1: Showing the waveform, short-term energy, and zero crossings for the word “seven”. Produced using the *timedom* tool which is included as part of [2].

Zero-Crossings	Short-Term Energy	Label
approx 0	approx 0	Silence
High	Low	Un-voiced
Low	High	Voiced
approx. 0	High	Voiced
High	High	Voiced
Low	Low	Voiced
approx. 0	Low	Un-voiced
Low	approx. 0	Silence
High	approx. 0	?

Table 4-2: Real world labelling scheme.

The one obvious anomaly in Table 4-2 is the last labelling scheme that labels a window as ‘?’. This is because it is impossible to get this arrangement of zero-crossings and short-term energy, in a speech signal. We decided that it made more sense to label this anomaly as a ‘?’ than to try and fudge the zero-crossings and short-term energy values to make them fit the criteria of a different label. We also realised that it would be a useful debugging aid to know that a window could not be successfully labelled by the labelling function.

5. RESULTS

Two people independent of each other, manually labelled the ten sound files with silence, unvoiced or voiced. The results of the manual labelling were then

7. CONCLUSIONS

The parts of speech labelling produced using the algorithms outlined in this study are reasonably accurate for well recorded, fairly clean speech but are not nearly as accurate for quiet recordings of speech.

The accuracy, of the algorithms outlined in this study, could be improved in two ways. Firstly more time could be spent on tweaking the cut-off values used by the algorithms to label the different parts of speech. The problem with this, however, is that if the values are fine tuned for one speech sample it is unlikely that they will be as accurate on other speech samples.

The other possible way of increasing the accuracy of the algorithms would be to use an over-lapping hamming window, when calculating the zero-crossings and short-term energy. This would, however, mean that many more calculations were necessary for each speech file, which would drastically increase the time taken to label an entire speech file. If, however, the speed of SUVing is not an issue, then this method of improving the algorithms is preferred to fine tuning the cut-off values.

The Matlab code of the algorithms outlined in this paper, and the manual and automatic transcriptions of the speech samples, for use with the *slt* tool (provided as part of [2]), can be found on the WWW at:
<http://www.dcs.shef.ac.uk/~u7mag/com325/>

REFERENCES

- [1] M. Cooke. *COM325: Computer Speech & Hearing* (Lecture Notes). Presented at the University of Sheffield, 1999
- [2] M. Cooke, G. Brown, S. Wrigley, and D. Ellis. *Matlab Auditory Demos, version 2.0*. Available Dec 1999 from:
www.dcs.shef.ac.uk/~martin/MAD/docs/mad.htm
- [3] B. Gold, and N. Morgan. *Speech and Audio Signal Processing*. John Wiley & Sons, 2000