

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Robustness of Group Delay Representations for Noisy Speech Signals

Sree Hari Krishnan P. ¹

IDIAP Research Institute, 1920 Martigny, Switzerland

R. Padmanabhan *

*Department of Computer Science and Engineering, Indian Institute of Technology
Madras, India*

Hema A. Murthy

*Department of Computer Science and Engineering, Indian Institute of Technology
Madras, India*

Abstract

This paper demonstrates the robustness of group delay based features to additive noise. First, we analytically show the robustness of group delay based representations. The analysis makes use of the fact that, for minimum-phase signals, the group delay function can be represented in terms of the cepstral coefficients of the log-magnitude spectrum. Such a representation results in the speech spectrum dominating over the noise spectrum, both at low and high SNRs. Further, we experimentally demonstrate the robustness of the representation on a voice activity detection (VAD) task, comparing a group delay based VAD algorithm with standard VAD methods as well as a magnitude-spectrum based method.

Key words: Voice activity detection, group delay functions, G.729 B, AMR VAD

* Corresponding author.

Email addresses: Hari.Parthasarathi@idiap.ch (Sree Hari Krishnan P.), padman@cse.iitm.ac.in (R. Padmanabhan), hema@cse.iitm.ac.in (Hema A. Murthy).

¹ Work was done at: Department of Computer Science and Engineering, Indian Institute of Technology Madras, India. Joint first authors appear in random order.

1 Introduction

All practical speech processing systems must have robustness as an important characteristic. The system must be able to perform satisfactorily under conditions where the input speech signal is distorted by effects that may be acoustic, articulatory or phonetic. Distortions caused by acoustic variations include those by additive or convolutive noise, channel and environmental effects. Development of robust feature extraction methods form a crucial part of building reliable speech processing systems.

Most short-term feature representations for speech are derived from the magnitude spectrum [eg. [Davis and Mermelstein \(1980\)](#); [Hermansky \(1990\)](#)]. There has been a growing series of studies which show that the short-term phase spectrum contains important information [[Alsteris and Paliwal \(2006\)](#); [Paliwal et al. \(2011\)](#)]. However, extracting information from the phase spectrum is not as straightforward as extracting information from the magnitude spectrum. Unlike the magnitude spectrum, the phase spectrum does not explicitly exhibit the system resonances. Further, signal processing difficulties (such as wrapping of the phase spectrum) are encountered while working directly with the phase spectrum [[Alsteris and Paliwal \(2006\)](#)]. The information in the phase spectrum has been utilised mostly by processing its derivative, the group delay function. The group delay function has been used in applications like signal reconstruction [[Yegnanarayana et al. \(1984\)](#)], formant extraction [[Murthy and Yegnanarayana \(1991\)](#)] and spectrum estimation [[Yegnanarayana and Murthy \(1992\)](#)]. Researchers have also come up with conventional short time feature representations derived from the group delay function [[Zhu and Paliwal \(2004\)](#); [Hegde et al. \(2007b\)](#)]. Moreover, it is shown that for applications like speech and speaker recognition, magnitude-based and phase-based features achieve comparable performance, and combining them results in improvement [[Zhu and Paliwal \(2004\)](#); [Hegde et al. \(2007a\)](#)].

This paper demonstrates that a group delay based representation is robust against additive noise. Although this property has been observed before (see [[Parthasarathi et al. \(2007\)](#)]), no attempt has yet been made to provide a mathematical analysis of the robustness property. In this article, we first show analytically how the group delay spectrum tends to follow the speech spectrum (as opposed to the noise spectrum) when speech is corrupted by additive noise. For minimum phase systems, log-magnitude and phase spectra are related through cepstral coefficients (see Section 2.) Using this property, we can represent the group delay spectrum in terms of speech and noise components at different signal to noise ratios.

The result of the analysis is validated by an experimental voice activity detection (VAD) task. A simple group delay based VAD algorithm is compared

1 to different VAD methods in various environments with additive noise. The
 2 experiments further corroborate the noise robustness properties of group delay
 3 functions.

4 The rest of this paper is organised as follows. In section 2, we briefly describe
 5 group delay functions. We then show analytically in section 3 that group delay
 6 functions are robust to additive noise. In section 4 we briefly describe the
 7 various representations of group delay functions and how features are derived
 8 from them. In section 5, we describe the experimental evidence supporting the
 9 analysis. Finally, we conclude in Section 6.
 10
 11
 12
 13
 14

15 2 A brief overview of group delay functions

16 The Fourier transform of a signal $x(n)$ can be represented as
 17
 18

$$19 X(\omega) = |X(\omega)|e^{j\theta(\omega)} \quad (1)$$

20
 21
 22
 23 The group delay function $\tau(\omega)$ of a signal is defined as the negative derivative
 24 of the continuous phase spectrum $\theta(\omega)$:
 25
 26

$$27 \tau(\omega) = -\frac{d(\theta(\omega))}{d\omega} \quad (2)$$

28
 29 Starting from Equation 1, the group delay function can be computed directly
 30 from the signal as:
 31
 32

$$33 \tau(\omega) = \frac{X_R(\omega)\hat{X}_R(\omega) + X_I(\omega)\hat{X}_I(\omega)}{|X(\omega)|^2} \quad (3)$$

34 where subscripts R and I respectively denote the real part and imaginary
 35 part, $x(n) \leftrightarrow X(\omega)$ and $\hat{x}(n) \leftrightarrow \hat{X}(\omega)$ are Fourier transform pairs, and
 36 $\hat{x}(n) = nx(n)$.
 37
 38

39 For a minimum-phase signal, it can be shown that the log magnitude and con-
 40 tinuous phase spectra are related as follows [Oppenheim and Schaffer (2000)]:
 41
 42

$$43 \ln |X(\omega)| = \frac{1}{2}c(0) + \sum_{n=1}^{\infty} c(n) \cos(n\omega) \quad (4)$$

$$44 \theta(\omega) = -\sum_{n=1}^{\infty} c(n) \sin(n\omega) \quad (5)$$

where $c(n)$ are the cepstral coefficients. Taking the negative derivative of Equation 5, we get the group delay function for minimum phase signals:

$$\tau(\omega) = \sum_{n=1}^{\infty} nc(n) \cos(n\omega) \quad (6)$$

Equations 4 and 5 show that for minimum-phase signals, the log magnitude and phase are related through cepstral coefficients. Also, from Equation 6, we find that the group delay function is the Fourier transform of the weighted cepstrum.

3 Robustness of group delay functions to additive noise

In this section, we analytically show that group delay functions for minimum-phase signals are robust to additive, uncorrelated noise.

Let $x[n]$ denote a clean, minimum-phase speech signal, which is degraded by uncorrelated, zero-mean, additive noise $v[n]$ with variance $\sigma^2(\omega)$. Then, the noisy speech, $y[n]$, can be expressed as,

$$y[n] = x[n] + v[n] \quad (7)$$

Taking the Fourier transform, we have

$$Y(\omega) = X(\omega) + V(\omega) \quad (8)$$

Multiplying by corresponding complex conjugates and taking the expectation, we have the power spectrum

$$P_Y(\omega) = P_X(\omega) + \sigma^2(\omega) \quad (9)$$

where $P_Y(\omega) = |Y(\omega)|^2$, $P_X(\omega) = |X(\omega)|^2$ and we have used the fact that the expectation of noise is zero. The power spectra of the resulting noisy speech signal can be related to noise power and (clean) speech power in one of three mutually exclusive frequency regions:

- (1) High noise power regions where $P_X(\omega) \ll \sigma^2(\omega)$
- (2) High signal power regions where $P_X(\omega) \gg \sigma^2(\omega)$ and
- (3) Equal power case where $P_X(\omega) \approx \sigma^2(\omega)$.

The power spectrum of the noisy speech signal in each case is represented by $P_Y^n(\omega)$, $P_Y^s(\omega)$ and $P_Y^e(\omega)$ respectively. We analyse the group delay representation of noisy speech in the three cases mentioned above.

3.1 High noise power spectral regions ($P_Y^n(\omega)$)

In this subsection, we consider frequencies ω such that $P_X(\omega) \ll \sigma^2(\omega)$, i.e., regions where the noise power is higher than signal power. From Equation 9 we have

$$\begin{aligned} P_Y^n(\omega) &= P_Y(\omega) \quad \forall \omega \quad \text{s.t.} \quad P_X(\omega) \ll \sigma^2(\omega) \\ &= P_X(\omega) + \sigma^2(\omega) \\ &= \sigma^2(\omega) \left(1 + \frac{P_X(\omega)}{\sigma^2(\omega)} \right) \end{aligned}$$

Taking logarithms on both sides, using the Taylor series expansion of $\ln(1 + \frac{P_X(\omega)}{\sigma^2(\omega)})$, and ignoring the higher order terms,

$$\begin{aligned} \ln(P_Y^n(\omega)) &= \ln \left[\sigma^2(\omega) \left(1 + \frac{P_X(\omega)}{\sigma^2(\omega)} \right) \right] \\ &\approx \ln(\sigma^2(\omega)) + \frac{P_X(\omega)}{\sigma^2(\omega)} \end{aligned} \quad (10)$$

Expanding $P_X(\omega)$ as a Fourier series ($P_X(\omega)$ is a periodic, continuous, function of ω with a period $\omega_0 = 2\pi$),

$$\ln(P_Y^n(\omega)) \approx \ln(\sigma^2(\omega)) + \frac{1}{\sigma^2(\omega)} \left[\frac{d_0}{2} + \sum_{k=1}^{\infty} d_k \cos\left(\frac{2\pi}{\omega_0} \omega k\right) \right] \quad (11)$$

where, d_k are the Fourier series coefficients in the expansion of $P_X(\omega)$. Since $P_X(\omega)$ is an even function, coefficients of the sine terms are zero.

For a minimum-phase signal, the group delay function can be computed in terms of the cepstral coefficients of the log-magnitude spectrum, as given in [Yegnanarayana et al. (1984)],

$$\begin{aligned} \log |X(\omega)| &= \frac{c_0}{2} + \sum_{k=1}^{\infty} c_k \cos(\omega k) \\ \tau(\omega) &= \sum_{k=1}^{\infty} k c_k \cos(\omega k) \end{aligned} \quad (12)$$

where, $\tau(\omega)$ is the minimum-phase group delay function and a_k are the cepstral coefficients. From Equation 12, it can be observed that the group delay function can be obtained from the log-magnitude response by ignoring the dc term, and by multiplying each coefficient with k . Applying this observation to

Equation 11, we get the group delay function as:

$$\tau_{Y^n}(\omega) \approx \frac{1}{\sigma^2(\omega)} \sum_{k=1}^{\infty} k d_k \cos(\omega k) \quad (13)$$

This expression shows that the group delay function is inversely proportional to the noise power ($\sigma^2(\omega)$) in regions where noise power is greater than the signal power.

3.2 High signal power spectral regions ($P_Y^s(\omega)$)

Now consider frequencies ω such that $P_X(\omega) \gg \sigma^2(\omega)$. Starting with Equation 9, and following the steps similar to those in the previous subsection:

$$\ln(P_Y^s(\omega)) \approx \ln(P_X(\omega)) + \frac{\sigma^2(\omega)}{P_X(\omega)} \quad (14)$$

Since $P_X(\omega)$ is non-zero, continuous, and periodic in ω , $\frac{1}{P_X(\omega)}$ is also periodic and continuous. Consequently, $\ln(P_X(\omega))$ and $\frac{1}{P_X(\omega)}$ can be expanded using Fourier series, giving

$$\ln(P_Y^s(\omega)) \approx \frac{d_0 + \sigma^2(\omega) e_0}{2} + \sum_{k=1}^{\infty} (d_k + \sigma^2(\omega) e_k) \cos(\omega k)$$

Using Equation 12, and following the steps in the previous case, we obtain the expression for the group delay function as,

$$\tau_{Y^s}(\omega) \approx \sum_{k=1}^{\infty} k (d_k + \sigma^2(\omega) e_k) \cos(\omega k) \quad (15)$$

where d_k and e_k are the Fourier series coefficients of $\ln(P_X(\omega))$ and $\frac{1}{P_X(\omega)}$ respectively. It is satisfying to observe that if $\sigma^2(\omega)$ is negligible, the group delay function can be expressed solely in terms of the log-magnitude spectrum.

3.3 Signal power \approx noise power regions ($P_Y^e(\omega)$)

For frequencies ω such that $P_X(\omega) \approx \sigma^2(\omega)$, we again start with Equation 9, and follow steps similar to those in the previous subsections, except in this case we do not need the Taylor series expansion:

$$\begin{aligned}
P_Y^e(\omega) &\approx 2P_X(\omega) \\
\ln(P_Y^e(\omega)) &\approx \ln 2 + \ln(P_X(\omega))
\end{aligned}
\tag{16}$$

Expanding $\ln(P_X(\omega))$ as a Fourier series, since it is a periodic, continuous, function of ω with a period 2π , the group delay function can be computed as,

$$\tau_{Y^e}(\omega) \approx \sum_{k=1}^{\infty} k d_k \cos(\omega k)
\tag{17}$$

where d_k are the Fourier series coefficients of $\ln(P_X(\omega))$.

3.4 Behaviour of minimum-phase group delay functions in noise

From Equations 13, 15, and 17, the estimated group delay functions are summarised respectively for the three cases:

$$\tau(\omega) \approx \begin{cases} \frac{1}{\sigma^2(\omega)} \sum_{k=1}^{\infty} k d_k \cos(\omega k) & \text{for } \omega : P_X(\omega) \ll \sigma^2(\omega), \\ \sum_{k=1}^{\infty} k (d_k + \sigma^2(\omega) e_k) \cos(\omega k) & \text{for } \omega : P_X(\omega) \gg \sigma^2(\omega), \\ \sum_{k=1}^{\infty} k d_k \cos(\omega k) & \text{for } \omega : P_X(\omega) \approx \sigma^2(\omega) \end{cases}
\tag{18}$$

From Equation 18, we note that the group delay function of a minimum-phase signal (which is corrupted with additive noise), is *inversely* proportional to the noise power at frequencies corresponding to high noise regions. On the other hand, for low noise regions, the group delay function becomes *directly* proportional to the signal power. In other words, the group delay function tends to follow the magnitude spectrum of the signal, rather than that of the noise. This indicates that the group delay function of a minimum-phase signal preserves the peaks and valleys of the magnitude spectrum well in the presence of additive noise. The group delay representation is thus robust and information like formant peaks are well retained.

4 Representations of group delay functions

Although the previous section showed the robustness of group delay functions in additive noise, it was done under the assumption of the signal being minimum-phase. Furthermore, to gainfully employ group delay functions for a speech processing task, features need to be extracted. This section endeavours to summarise literature relevant in handling both issues.

In general, the group delay function can be effectively utilised only if the signal under consideration is a minimum-phase signal [Hegde et al. (2007b)].

1 A signal is defined as a minimum phase signal only if both itself and its inverse
2 are causal and stable. The group delay function becomes unstable if the roots
3 of the transfer function are on the unit circle. Since speech is produced by
4 a stable system, the poles of the system are well within the unit circle. The
5 possible presence of zeros near the unit circle lead to an unstable group delay
6 function. For speech signals, the zeros are introduced by pitch, nasals and the
7 short-time analysis window function. Further, zeros are also introduced by
8 the glottal return phase (the time interval between the most negative value of
9 the glottal flow derivative and glottal closure [Plumpe et al. (1999)] when the
10 glottal flow is modelled as a linear filter [Doval et al. (2003)]. The presence
11 of zeros makes the denominator term in Equation 3 vanish, leading to an ill-
12 behaved group delay function. Thus, the analysis presented in the previous
13 section applies to minimum-phase signals.

14
15 To compute the group delay function for non-minimum phase signals, two
16 approaches can be used:

17 4.1 Minimum phase group delay function (MINGD)

18
19 In this method, a minimum-phase equivalent is derived from the non minimum-
20 phase signal using the root cepstrum [Murthy and Yegnanarayana (1991)]. In
21 [Nagarajan et al. (2003)], it was shown that for an all-pole, non minimum-
22 phase signal, the causal portion of the inverse Fourier transform of the magni-
23 tude spectrum is a minimum-phase signal. The group delay function of this sig-
24 nal is then utilised. The MINGD of a signal is thus derived from its minimum-
25 phase component.

26 4.2 Modified group delay function (MODGD)

27
28 In this method, the zeros near the unit circle in Equation 3 are suppressed
29 by smoothing the magnitude spectrum $X(\omega)$ [Murthy and Gadde (2003);
30 Hegde et al. (2007b)]. The resulting group delay function is well behaved.
31 The MODGD $\tau_m(\omega)$ is defined as

$$32 \tau_m(\omega) = \left(\frac{\tau_s(\omega)}{|\tau_s(\omega)|} \right) (|\tau_s(\omega)|)^\alpha \quad (19)$$

33 where

$$34 \tau_s(\omega) = \frac{X_R(\omega)\hat{X}_R(\omega) + X_I(\omega)\hat{X}_I(\omega)}{|S(\omega)|^{2\gamma}} \quad (20)$$

35 and $|S(\omega)|$ is a smoothed version of $|X(\omega)|$, while the parameters α and γ are

introduced to control the dynamic range. The length of the cepstral smoothing window is controlled by the parameter lifter_ω .

The MINGD of a signal is extracted from the magnitude spectrum of the signal, whereas, the MODGD can be extracted directly from the signal. Studies such as [Hegde et al. (2007b)] have shown the equivalence between MINGD and MODGD in a least squares sense, when the signal is not minimum phase. Furthermore, it is also demonstrated in [Hegde et al. (2007b)] the application of the MODGD for tasks such as syllable recognition, language identification and speaker recognition. Given these reasons, we validate the robustness property of group delay functions on the MODGD by experimental evaluation.

4.3 Converting the MODGD into features

To convert the MODGD into meaningful features, a decorrelation operation is performed using homomorphic processing. This is achieved by applying a discrete cosine transform (DCT) on the MODGD. Thus, the MODGD is converted into cepstral features and the resulting feature representation is called the modified group delay feature (MODGDF.) The decorrelation implied in the MODGDF features allow the use of diagonal covariances for modelling the feature distribution.

The extraction of MODGDF features from the signal are now described as follows [Hegde et al. (2007b)]:

- (1) Compute the DFT of the speech signal $x[n]$ as $X[k]$.
- (2) Next, the DFT of the signal $n x[n]$ is computed as $\hat{X}[k]$.
- (3) Compute the cepstrally smoothed spectra of $X[k]$ and denote it as $S[k]$. The parameter lifter_ω is used to control the length of the window in the cepstral domain.
- (4) Compute the MODGD as:

$$\tau_m[k] = \left(\frac{\tau[k]}{|\tau[k]|}\right)(|\tau[k]|)^\alpha \quad (21)$$

where $\tau[k] = \frac{X_R[k]\hat{X}_R[k]+X_I[k]\hat{X}_I[k]}{|S[k]|^{2\gamma}}$ where the parameters α and γ are used to control the dynamic range of the MODGD.

- (5) Compute the MODGDF features by taking the DCT:

$$c[n] = \sum_{k=0}^{k=N_f} \tau_m[k] \cos(n(2k+1)\pi/N_f) \quad 0 \leq n < N_c \quad (22)$$

where N_f is the DFT size and N_c are the number of cepstral coefficients.

5 Experimental evaluation

In earlier sections, we analysed the effect of additive noise on the MINGD and how the MODGD is an approximation of the MINGD. The previous section described the process of converting the MODGD into cepstral features. In this section, we corroborate the analysis presented in section 3 with experimental results using the MODGDF.

5.1 Voice activity detection

A voice activity detection (VAD) algorithm distinguishes speech and non-speech regions in a signal. VAD is used as a preprocessing stage in many speech applications, including automatic speech recognition (ASR) and speech coding. An effective VAD subsystem improves performance significantly, and achieves better resource utilisation. When viewed as a pattern recognition problem, VAD has to discriminate two classes, (noisy) speech and non-speech. Non-speech may be silence, or acoustic background noise. In high SNR conditions, most VAD algorithms work well. But in low SNR conditions, VAD performance deteriorates considerably, with many detection errors.

Many VAD algorithms use short term energy [Evangelopoulos and Maragos (2005); Parthasarathi et al. (2007)]. Energy-based algorithms usually do not work well in low SNR conditions. Other approaches to VAD use statistical methods [Sohn et al. (1999); Nemer et al. (2001); Li et al. (2005)], which make use of the assumption that the statistics of noise and speech are different. Recently, to perform VAD non-stationary noise, long-term analysis windows have been utilised [Ghosh et al. (2011)].

Standard VAD algorithms include the G.729-Annex B algorithm [Benyassine et al. (1997)] (recommended by the International Telecommunication Union) and the adaptive multi rate (AMR) VAD 1 and 2 [GSM 06.94 (1999)] (recommended by the European Telecommunications Standards Institute).

5.2 Utilising MODGDF for VAD

In this section, we develop a VAD algorithm which utilises MODGDF features for representing the speech signal. The proposed MODGDF-based VAD algorithm is compared to a baseline algorithm as well as standard VAD algorithms.

The use of Gaussian mixture models (GMMs) is popular for modelling generic

1 probability distributions. In its application to VAD, the underlying principle
2 is that the distributions of speech and non-speech are sufficiently separate in
3 the MODGDF space. In the training phase, MODGDF features extracted from
4 speech and non-speech are used to build GMMs for each class, using framewise
5 labelled data. In the testing phase, each frame of the test utterance is classified
6 as either speech or non-speech. This method is called MODGD-VAD.
7

8 The MODGDF features were extracted as given in section 4.3. The speech
9 signal was pre-emphasised, after using a frame size of 25 ms and a frame shift of
10 10 ms and applying a Hamming window. As specified in [Hegde et al. (2007b)]
11 the MODGDF parameters α and γ were set to 0.4 and 0.9 respectively, and
12 the parameter lifter_{ω} was set to 8. 13-dimensional MODGDF features were
13 extracted. Cepstral mean subtraction is performed to remove channel effects.
14 Two 64-mixture GMMs were built for speech and non-speech respectively,
15 each initialised using k-means clustering, followed by five iterations of the
16 expectation maximisation (EM) algorithm. Diagonal covariance matrices were
17 utilised.
18

19 The performance of MODGD-VAD is compared to a baseline GMM-based
20 algorithm using the Mel frequency cepstra (MFCC). This method is called
21 MFCC-VAD. MFCC features were computed as described in [Murthy et al.
22 (1999)]. A bank of 40 filters, which were spaced between 0 Hz to 4000 Hz are
23 used to compute filterbank energies by integration. The shape of the filters was
24 set by setting the shape constant T_r to 0.0 and the filter warping constant F_{ω}
25 was set to 0.4. 13-dimensional MFCC coefficients were extracted, and cepstral
26 mean subtraction is performed. 64-mixture GMMs were built for speech and
27 non-speech, initialised by k-means clustering, followed by five EM iterations.
28 As in the case for MODGDF, diagonal covariances were used.
29
30
31
32
33
34
35
36

37 5.3 Reference VAD systems

38
39
40

41 In addition to the baseline MFCC-VAD, the performance of the MODGD-
42 VAD algorithm is compared to two standard VAD algorithms. We provide a
43 brief description of them below.
44

45 **G.729B VAD:** ITU's G.729B VAD [Benyassine et al. (1997)] makes use of
46 spectral distortion, full band and low band energy and zero crossings to make
47 a VAD decision every 10 ms. These parameters are estimated as the difference
48 from the respective running averages. Piecewise linear hyperplanes are used to
49 separate speech and non-speech in the 4-dimensional parameter space, followed
50 by a four stage hangover scheme.
51
52
53

54 **AMR VAD options 1 and 2:** ETSI's AMR option 1 and 2 [GSM 06.94
55 (1999)] are sub-band based methods. AMR 1 uses nine non-uniform sub-bands
56
57
58

1 and AMR 2 uses sixteen. Both methods estimate sub-band SNRs using a first-
2 order auto-regressive model. VAD decisions are made every 20 ms using an
3 adaptive threshold, followed by a hangover scheme.
4
5
6
7

8 5.4 Dataset used for the evaluation 9

10
11
12 The dataset used for the study consisted of 426 speech utterances (213 female,
13 213 male), selected from the TIMIT training set (phonetically balanced, read
14 speech [Garofolo (1993)]). The utterances were downsampled to 8kHz, because
15 the software implementations of the G.729B VAD and AMR VAD we used
16 required so. Framewise VAD decisions manually labelled on these utterances
17 formed the reference labels. These utterances were then concatenated in sets
18 of three and silence was inserted between them so that the ratio of speech
19 to non-speech frames was 40% to 60%. This models the average activity in a
20 typical telephone conversation [Beritelli et al. (2001)]. Three different types
21 of noise (babble, pink and white) from the NOISEX-92 database [Varga et al.
22 (1992)] were added to these files, resulting in four datasets, each having an
23 SNR of 0, 5, 10 and 15 dB. It is to be noted that SNR computation is over
24 the whole utterance, including silences.
25
26
27
28
29

30 Six such utterances at 10 dB SNR (comprising of about two minutes each of
31 speech and non-speech) were used as training data. Using framewise speech
32 and non-speech reference labels, the training data was used to estimate param-
33 eters the respective GMMs. On the other hand, the total test dataset consists
34 of about thirty hours of noisy data.
35
36
37
38
39
40

41 5.5 Evaluation 42 43 44

45 The evaluation data for each noise type consists of four SNR levels. Framewise
46 speech/non-speech outputs of the five VAD methods considered (G.729 B,
47 AMR 1, AMR 2, MFCC-VAD and MODGD-VAD) are compared with the
48 reference labels. The performance evaluation is based on the metrics defined
49 in [Li et al. (2005)], namely:
50
51
52

- 53 • P_{cs} : Percentage of correct speech identification,
- 54 • P_{cn} : Percentage of correct non-speech identification and
- 55 • P_f : Percentage of error (misses and false alarms.)
56
57
58
59
60
61
62
63
64
65

Table 1. Performance of VAD methods in different noise environments.

Method	G.729B			AMR 1			AMR 2			MFCC-VAD			MODGD-VAD		
	P_{cn}	P_{cs}	P_f	P_{cn}	P_{cs}	P_f	P_{cn}	P_{cs}	P_f	P_{cn}	P_{cs}	P_f	P_{cn}	P_{cs}	P_f
SNR															
Babble noise environment															
0	65.3	57.3	37.8	47.8	95.7	33.0	44.5	95.6	34.9	88.7	78.1	15.6	95.7	79.3	10.9
5	65.4	68.2	33.4	56.9	97.5	26.8	42.7	99.1	34.6	95.0	83.7	9.5	95.9	88.4	7.1
10	65.3	77.4	29.8	72.2	97.2	17.7	45.4	99.7	32.8	97.3	88.0	6.4	95.0	93.8	5.4
15	65.8	86.4	25.9	84.5	97.1	10.4	57.0	99.5	25.9	96.7	91.4	5.4	93.4	96.6	5.3
Avg.	65.5	72.3	31.7	65.3	96.9	22.0	47.4	98.5	32.1	94.5	85.3	9.2	95.0	89.6	7.1
White noise environment															
0	89.5	54.2	24.5	87.5	67.5	20.4	94.7	89.0	7.5	86.3	76.7	17.5	90.6	85.8	11.4
5	89.4	66.8	19.5	90.2	84.6	11.9	93.6	96.4	5.2	94.5	82.5	10.3	91.9	90.2	10.1
10	89.3	77.4	15.4	92.4	93.4	7.1	93.2	97.9	4.8	97.1	87.9	6.6	92.6	94.1	6.8
15	93.1	87.4	9.1	95.6	95.7	4.3	93.0	98.7	4.6	96.3	92.7	5.16	90.6	97.9	6.4
Avg.	90.3	71.5	17.1	91.4	85.3	10.9	93.6	95.5	5.5	93.5	85.0	9.9	91.0	92.0	8.3
Pink noise environment															
0	89.5	60.5	22.0	89.6	69.2	18.5	95.0	91.4	6.4	77.3	82.4	20.6	91.3	80.4	13.1
5	88.8	71.4	18.1	92.6	86.2	9.9	93.6	97.2	4.9	92.8	86.1	9.9	94.3	85.2	9.4
10	89.3	81.4	13.8	94.7	93.0	5.8	93.4	98.3	4.6	97.0	90.5	5.6	95.1	91.0	6.5
15	88.4	88.6	11.5	95.6	96.6	3.9	92.2	99.2	4.9	93.0	94.6	6.3	94.4	95.0	5.3
Avg.	89.0	75.5	16.3	93.1	86.3	9.5	93.6	96.5	5.2	90.0	88.4	10.6	93.8	87.9	8.6
All	81.6	73.1	21.7	83.7	89.5	14.1	78.2	96.8	14.3	92.7	86.2	9.9	93.3	89.8	8.2

5.6 Results and discussion

The experimental performance of the proposed MODGD-VAD and other methods are given in Table 1.

The proposed MODGD-VAD achieves the least average VAD error (8.2%) across all noise types. Also, in babble noise, MODGD-VAD also achieves the least VAD error (7.1%) compared to the other methods. Among the three noise types, it achieves the best non-speech detection (92.0%) in white noise, whereas the best speech detection is achieved in babble noise (95.0%).

The performance of MFCC-VAD is similar, although the algorithm attains a higher average VAD error. At high noise levels (0 dB SNR), MFCC-VAD has a 6% performance reduction in compared MODGD-VAD. Similar to MODGD-VAD, MFCC-VAD achieves the least VAD error and the best speech detection in babble noise (9.2% and 94.5% respectively), whereas the best non-speech detection is achieved in pink noise (88.4.%) In particular, both these supervised methods perform better than the standard methods in babble noise.

G.729B has the highest average VAD error across all environments. The average speech detection rate is similar across various noise types, although the non-speech detection is poor in babble noise. Among the three noise types considered, G.729 obtains the least VAD error and best non-speech detection in pink noise (16.3% and 75.5%), whereas the best speech detection is achieved in white noise (90.3%)

The AMR VAD 2 algorithm achieves the best speech detection performance in pink and white noise (93.6%). In babble noise, it has good speech detection performance, although at the cost of noise detection. The AMR VAD 1 algorithm on the other hand, fares better in babble noise (average VAD error of 22.0% when compared to 32.1% obtained by AMR VAD 2.) Both the AMR algorithms achieve the least VAD error in pink noise (9.5% for AMR 1 and 5.2% for AMR 2), and the best non-speech detection in babble noise (96.9% for AMR 1 and 98.5% for AMR 2.)

As expected, all VAD methods considered show improvement in average error as the SNR increases. In babble noise, MODGD-VAD outperforms all other VAD methods considered. In pink noise and white noise, it performs better than AMR option 1 and is comparable to AMR option 2. Although G.729B achieves the highest average VAD error, it fares better than AMR VAD option 2 in babble noise. In babble noise, MODGD-VAD gives the least VAD error (7.1%), whereas in both white and pink noise AMR VAD option 2 gives the least error (5.5% and 5.2% respectively.)

The supervised VAD methods (MFCC-VAD) and (MODGD-VAD) fare bet-

1 ter than the unsupervised methods in all conditions. Being non-stationary in
2 nature, babble noise degrades performance of the standard VAD methods con-
3 siderably. The GMM is able to capture the non-stationarity and separate the
4 speech and non-speech in a better manner.

5
6 The standard VAD algorithms implement a hang-over scheme in which the
7 VAD decision of the current frame is based on that of the earlier frames. To
8 reduce the misdetection of speech, the hang-over scheme usually delays the
9 transition from speech to non-speech, at the cost of increased false alarm rate
10 [Sohn et al. (1999)]. It is to be noted that a hang-over scheme was not used
11 in the GMM-based methods.
12
13
14

15 6 Conclusions

16
17
18
19
20 In this paper, we demonstrated the robustness of group delay based repre-
21 sentations to additive noise. Analysis of group delay representations in noise
22 showed that the group delay spectrum tends to follow the signal spectrum,
23 rather than the noise spectrum. To support the analysis, a GMM-based voice
24 activity detector was developed using modified group delay based features. The
25 modified group delay features achieved performance comparable to or better
26 than standard VAD methods in various environments. The performance was
27 also compared to an MFCC feature based GMM voice activity detector, which
28 achieved an average error rate slightly more than the MODGD-based system.
29
30
31

32 We provided a mathematical explanation to the noise robustness of group de-
33 lay functions. The experiments also demonstrated that features derived from
34 group delay functions separate noisy-speech and non-speech in a reliable man-
35 ner. Thus, the use of the phase-based features is as appealing as magnitude-
36 based features for this important task.
37
38
39

40 41 References

- 42
43
44 Alsteris, L. D., Paliwal, K. K., 2006. Short-time phase spectrum in speech
45 processing: a review and some experimental results. *Digital Signal Process.*
46 17, 578–616.
47
48 Benyassine, A., Shlomot, E., Su, H.-Y., 1997. ITU Recommendation G.729
49 Annex B: A Silence Compression Scheme for use with G.729 optimized for
50 V.70 digital simultaneous voice and data applications. *IEEE Comm. Mag.*
51 35 (9), 64–73.
52
53 Beritelli, F., Casale, S., Ruggeri, G., 2001. Performance evaluation and com-
54 parison of ITU-T/ETSI voice activity detectors. In: *Proc. ICASSP*. Vol. 3.
55 pp. 1425–1428.
56
57
58
59
60
61
62
63
64
65

- 1 Davis, S., Mermelstein, P., 1980. Comparison of parametric representations
2 for monosyllabic word recognition in continuously spoken sentences. *IEEE*
3 *Trans. Acoust., Speech Signal Process.* 28 (4), 357–366.
- 4 Doval, B., d’Alessandro, C., Henrich, N., 2003. The voice source as a
5 causal/anticausal linear filter. In: *Proc. ISCA Tutorial and Research Work-*
6 *shop on Voice Quality: Functions, Analysis and Synthesis.*
- 7 Evangelopoulos, G., Maragos, P., 2005. Speech event detection using multi-
8 band modulation energy. In: *Proc. Interspeech.* pp. 685–688.
- 9
10 Garofolo, J. S., 1993. Timit acoustic-phonetic continuous speech corpus. Lin-
11 guistic Data Consortium, Philadelphia.
- 12 Ghosh, P. K., Tsiartas, A., Narayanan, S., 2011. Robust voice activity detec-
13 tion using long-term signal variability. *IEEE Trans. Audio, Speech, Lang.*
14 *Process.* 19 (3), 600–613.
- 15 GSM 06.94, 1999. Digital Cellular Telecommunications System (Phase 2+);
16 Voice Activity Detector (VAD) for Adaptive Multi Rate (AMR) Speech
17 Traffic Channels.
- 18
19 Hegde, R. M., Murthy, H. A., Gadde, V. R., 2007a. Significance of joint fea-
20 tures derived from the modified group delay function in speech processing.
21 *EURASIP Jnl. Audio, Speech, Music Process.*
- 22
23 Hegde, R. M., Murthy, H. A., Gadde, V. R. R., 2007b. Significance of the
24 Modified Group Delay Feature in Speech Recognition. *IEEE Trans. Audio,*
25 *Speech Lang. Process.* 15 (1), 190–202.
- 26
27 Hermansky, H., 1990. Perceptual linear predictive (plp) analysis of speech.
28 *Jnl. Acoust. Soc. America* 87 (4), 1738–1752.
- 29
30 Li, K., Swamy, M. N. S., Ahmad, M. O., 2005. An improved voice activity
31 detection using higher order statistics. *IEEE Trans. Speech Audio Process.*
32 13 (5), 965–974.
- 33
34 Murthy, H., Beaufays, F., Heck, L., Weintraub, M., 1999. Robust text-
35 independent speaker identification over telephone channels. *IEEE Trans.*
36 *Speech Audio Process.* 7 (5), 554–568.
- 37
38 Murthy, H. A., Gadde, V. R. R., 2003. The modified group delay function and
39 its application to phoneme recognition. In: *Proc. ICASSP. Vol. 1.* pp. 68–71.
- 40
41 Murthy, H. A., Yegnanarayana, B., 1991. Formant extraction from minimum
42 phase group delay function. *Speech Comm.* 10 (3), 209–221.
- 43
44 Nagarajan, T., Prasad, V. K., Murthy, H. A., 2003. Minimum phase signal
45 derived from root cepstrum. *IEE Electronics Lett.* 39 (12), 941– 942.
- 46
47 Nemer, E., Goubran, R., Mahmoud, S., 2001. Robust voice activity detec-
48 tion using higher-order statistics in the LPC residual domain. *IEEE Trans.*
49 *Speech Audio Process.* 9 (3), 217–231.
- 50
51 Oppenheim, A. V., Schaffer, R. W., 2000. *Discrete-time Signal Processing.*
52 Prentice-Hall.
- 53
54 Paliwal, K. K., Wojcicki, K., Shannon, B., 2011. The importance of phase in
55 speech enhancement. *Speech Commun.* 53, 465–494.
- 56
57 Parthasarathi, S. H. K., Rajan, P., Murthy, H. A., 2007. Voice activity detec-
58 tion using group delay processing on buffered short term energy. In: *Proc.*

National Conference Comm.

- 1 Plumpe, M. D., Quatieri, T. F., Reynolds, D. A., 1999. Modeling of the glottal
2 flow derivative waveform with application to speaker identification. IEEE
3 Trans. Speech Audio Process. 7 (5), 569–586.
4
5 Sohn, J., Kim, N. S., Sung, W., 1999. A statistical model-based voice activity
6 detection. IEEE Signal Process. Lett. 6, 1–3.
7
8 Varga, A. P., Steeneken, H. J. M., Tomlinson, M., Jones, D., 1992. The
9 NOISEX-92 study on the effect of additive noise on automatic speech recog-
10 nition. Tech. Report DRA Speech Research Unit.
11
12 Yegnanarayana, B., Murthy, H. A., 1992. Significance of group delay functions
13 in spectrum estimation. IEEE Trans. Signal Process. 40, 2281–2289.
14
15 Yegnanarayana, B., Saikia, D. K., Krishnan, T. R., Jun. 1984. Significance of
16 group delay functions in signal reconstruction from spectral magnitude or
17 phase. IEEE Trans. Acoust., Speech, Signal Process., 610–622.
18
19 Zhu, D., Paliwal, K., 2004. Product of power spectrum and group delay func-
20 tion for speech recognition. In: Proc. ICASSP. Vol. 1.
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65