A Polynomial Subspace Projection Approach for the Detection of Weak Voice Activity

Imperial College London



Vincent W. Neo, Stephan Weiss, Patrick A. Naylor SSPD 2022

A Polynomial Subspace Projection Approach for the Detection of Weak Voice Activity - 1/32

Outline

1. Introduction

Voice Activity Detection Motivations

2. Background

Multichannel Signal Model Polynomial Matrices and Polynomial EVD

- 3. PEVD Preprocessor for VAD
- 4. Experiment and Results
- 5. Conclusion

Introduction

A Polynomial Subspace Projection Approach for the Detection of Weak Voice Activity - 3/32

What is Voice Activity Detection (VAD)



What is Voice Activity Detection (VAD)



What is Voice Activity Detection (VAD)



Detection of voice activity is important for many applications:

- Speech enhancement in hearing aids, telecommunications
- Automatic speech recognition (ASR) systems
- Robot audition

Detection of voice activity is important for many applications:

- Speech enhancement in hearing aids, telecommunications
- Automatic speech recognition (ASR) systems
- Robot audition
- Main challenges:
 - Background noise
 - Interfering sources
 - Reverberation

• Exploit differences in noise and speech distributions

- Exploit differences in noise and speech distributions
- \Rightarrow Challenging to measure signal statistics in very noisy environments

- Exploit differences in noise and speech distributions
- \Rightarrow Challenging to measure signal statistics in very noisy environments

Machine learning-based methods [Google 2021; Zhang2016; lvry2019]

• Speech feature extraction for classification

- Exploit differences in noise and speech distributions
- \Rightarrow Challenging to measure signal statistics in very noisy environments

Machine learning-based methods [Google 2021; Zhang2016; Ivry2019]

- Speech feature extraction for classification
- \Rightarrow Feature extraction becomes difficult in adverse acoustic environments

- Exploit differences in noise and speech distributions
- \Rightarrow Challenging to measure signal statistics in very noisy environments

Machine learning-based methods [Google 2021; Zhang2016; Ivry2019]

- Speech feature extraction for classification
- \Rightarrow Feature extraction becomes difficult in adverse acoustic environments
 - Weak Transient Signal Detection Using PEVD [Weiss2021]
 - Exploits multichannel signal processing to amplify weak transient signals

- Exploit differences in noise and speech distributions
- \Rightarrow Challenging to measure signal statistics in very noisy environments

Machine learning-based methods [Google 2021; Zhang2016; Ivry2019]

- Speech feature extraction for classification
- \Rightarrow Feature extraction becomes difficult in adverse acoustic environments
 - Weak Transient Signal Detection Using PEVD [Weiss2021]
 - Exploits multichannel signal processing to amplify weak transient signals

This Talk: PEVD-based Multichannel Preprocessing for VAD

Background

A Polynomial Subspace Projection Approach for the Detection of Weak Voice Activity - 9/32

Multichannel Signal Model

The received signal at the q-th sensor with time index n is

$$x_q(n) = \sum_{p=1}^p \mathbf{h}_{p,q}^T(n) \mathbf{s}_p(n)$$

where

- $\mathbf{h}_{p,q}(n)$ is the room impulse response from pth source to qth microphone modelled as a Jth order FIR filter,
- $\mathbf{s}_p(n)$ is the *p*th localized source signal.

The data vector collected from Q microphones:

$$\mathbf{x}(n) = [x_1(n), x_2(n), \dots, x_Q(n)]^T \in \mathbb{R}^Q$$
.

Imperial College

London

Space-time Covariance Polynomial Matrix

Imperial College London

Assuming stationarity, the space-time covariance matrix is

$$\mathbf{R}(\tau) = \mathbb{E}[\mathbf{x}(n)\mathbf{x}^T(n-\tau)] \in \mathbb{R}^{Q \times Q} ,$$

where (i, j)th element is the correlation function $r_{ij}(\tau) = \mathbb{E}[x_i(n)x_j(n-\tau)]$ and τ is the time-shift.

Z-transform of $\mathbf{R}(\tau)$ is a para-Hermitian polynomial matrix

$$\mathcal{R}(z) = \sum_{\tau = -W}^{W} \mathbf{R}(\tau) z^{-\tau},$$

where $\mathbf{R}(\tau) \approx 0$ for $|\tau| > W$, calligraphic \mathcal{R} for polynomial matrices and regular \mathbf{R} for matrices.

Polynomial Matrix Eigenvalue Decomposition

The PEVD of $\Re(z)$ is [Weiss2018a; Weiss2018b]

$$\mathcal{R}(z) = \mathcal{U}(z)\Lambda(z)\mathcal{U}^{P}(z) , \qquad (1)$$

Imperial College

London

where $\Lambda(z), \mathcal{U}(z)$ contain the eigenvalues and eigenvectors and $\mathcal{R}^{P}(z) = \mathcal{R}^{H}(1/z^{*})$.

Subspace decomposition using PEVD:

$$\mathcal{R}(z) = \begin{bmatrix} \mathcal{U}_s(z) & \mathcal{U}_{\perp}(z) \end{bmatrix} \begin{bmatrix} \Lambda_s(z) & \mathbf{0} \\ \mathbf{0} & \Lambda_{\bar{s}}(z) \end{bmatrix} \begin{bmatrix} \mathcal{U}_s^P(z) \\ \mathcal{U}_{\perp}^P(z) \end{bmatrix}, \quad (2)$$

associated with signal, $\{\cdot\}_s$ and orthogonal complement, $\{\cdot\}_{\perp}$ subspaces.

Example: Polynomial Matrix from ST-Covariance

Imperial College London



Example: PEVD Algorithm

Algorithm converges when $|g| < 1.68 \times 10^{-2}$

Example: PEVD Algorithm Outputs





PEVD Algorithms

Iterative PEVD algorithms approximating (1) include:

- Second-order Sequential Best Rotation (SBR2) [McWhirter2007]
- Sequential Matrix Diagonalization (SMD) [Redif2015]
- Householder PEVD [Neo2019]
- Fixed-order approximate PEVD [Tkacenko 2010]
- Multiple-shift SBR2/SMD [Wang2015; Corr2014b]
- Causality-constrained Multiple-shift SMD [Corr2014a]

PEVD Preprocessor for VAD

A Polynomial Subspace Projection Approach for the Detection of Weak Voice Activity - 17/32

Ambient Acoustics Subspace Characterization

Imperial College London



For L estimated signal components, $\mathcal{U}_s(z) \in \mathbb{C}^{Q \times L}$ and $\mathcal{U}_{\perp}(z) \in \mathbb{C}^{Q \times (Q-L)}$,

$$\mathcal{U}_s(z)\mathcal{U}_s^P(z) + \mathcal{U}_\perp(z)\mathcal{U}_\perp^P(z) = \mathbf{I}$$
.

The component associated with $\mathcal{U}_{\perp}(z) \bullet \mathcal{O} \mathbf{U}(n)$ can be recovered using

$$\mathbf{y}(n) = \sum_{k} \sum_{m} \mathbf{U}_{\perp}(k) \mathbf{U}_{\perp}^{H}(k-m) \mathbf{x}(n-m) .$$

This is equivalent to $\mathbf{x}(n)$ with the $\mathcal{U}_s(z)$ component removed.





Experiment and Results

A Polynomial Subspace Projection Approach for the Detection of Weak Voice Activity - 22/32

Setup: Male Speaker in a Measured Room [Kayser2009]

Imperial College London



Comparative algorithms:

- 1. Sohn [Sohn1999]
- 2. WebRTC [Google 2021] : G0, G3 (Least to most aggressive)
- 3. Proposed (PEVD+Sohn): R1, R2, R5, R7 (different rank estimates)

Comparative algorithms:

- 1. Sohn [Sohn1999]
- 2. WebRTC [Google 2021] : G0, G3 (Least to most aggressive)
- 3. Proposed (PEVD+Sohn): R1, R2, R5, R7 (different rank estimates)

Evaluation measures [Tharwat 2018] :

- Label evaluation metrics
 - Correct labels: True Positive (TP), True Negative (TN)
 - Wrong labels: False Positive (FP), False Negative (FN)
- Overall scores: F1, Balanced Accuracy (BACC)

Comparative algorithms:

- 1. Sohn [Sohn1999]
- 2. WebRTC [Google 2021] : G0, G3 (Least to most aggressive)
- 3. Proposed (PEVD+Sohn): R1, R2, R5, R7 (different rank estimates)

Evaluation measures [Tharwat 2018] :

- Label evaluation metrics
 - Correct labels: True Positive (TP), True Negative (TN)
 - Wrong labels: False Positive (FP), False Negative (FN)
- Overall scores: F1, Balanced Accuracy (BACC)
- \implies Focus on first microphone in the results.

VAD Performance for -20 dB SIR F16 Cockpit Noise

Imperial College London



A Polynomial Subspace Projection Approach for the Detection of Weak Voice Activity - 25/32

VAD Performance for -20 dB SIR F16 Cockpit Noise

Imperial College London

Method	TP	ΤN	FP	FN	F1	BACC
Sohn	130	241	38	185	0.538	0.638
R1	136	249	30	179	0.565	0.662
R2	158	244	35	157	0.622	0.688
R5	148	247	32	167	0.598	0.678
R7	136	224	55	179	0.538	0.617
G0	315	0	279	0	0.693	0.500
G3	315	0	279	0	0.693	0.500

A Polynomial Subspace Projection Approach for the Detection of Weak Voice Activity - 26/32

VAD Performance for -20 dB SIR F16 Cockpit Noise

Imperial College London

Method	TP	TN	FP	FN	F1	BACC
Sohn	130	241	38	185	0.538	0.638
R1	136	249	30	179	0.565	0.662
R2	158	244	35	157	0.622	0.688
R5	148	247	32	167	0.598	0.678
R7	136	224	55	179	0.538	0.617
G0	315	0	279	0	0.693	0.500
G3	315	0	279	0	0.693	0.500

Other results in the paper:

- Since G0, G3 always predict the presence of speech, F1 scores significantly decrease when the speech segment is short.
- Tested on destroyer noise at various SIR from -30 dB to 20 dB.

Conclusion

A Polynomial Subspace Projection Approach for the Detection of Weak Voice Activity - 27/32

Conclusion

- PEVD-based multi-microphone preprocessing for VAD
 - Characterize the ambient acoustics using PEVD to generate multichannel syndrome signals, which are microphone signals without ambient acoustics
 - Apply single channel VAD to each microphone
- Performance of proposed PEVD-based approach
 - Almost always improves F1 and BACC scores over the single channel method even in adverse environments, i.e. -30 dB SIR
 - Consistent performance unaffected by length of speech segments

References

diagonalisation for parahermitian matrices". In: Proc. Eur. Signal Process. Conf. (EUSIPCO), pp. 1277–1281.
diagonalisation for parahermitian matrices". In: Proc. IEEE/SP Workshop on Statistical Signal Process. Pp. 844–848.
Process. 11.5, pp. 498–505.
Sogie (2021). Weby IC Voice Activity Detector.
narrowband telephone networks and speech codecs. Recommendation P.862. Int. Telecommun. Union (ITU-T).
voice and data applications. Recommendation. Int. Telecommun. Union (ITU-T).
Sel. Topics Signal Process. 13.2, pp. 254–264.
 behind-the-ear head-related and binaural room impulse responses". In: EURASIP J. on Advances in Signal Process. 2009.1, p. 298605.

References

In: IEEE Trans. Signal Process. 55.5, pp. 2158–2169.
eigenvalue decomposition". In: Proc. IEEE Int. Conf. on Acoust., Speech and Signal Process. (ICASSP), pp. 8043–8047.
Realt, S., S. Weiss, and J. G. McWhiter (Jan. 2015). Sequential matrix diagonalisation algorithms for polynomial EVD of parametricitian matrices". In: IEEE Trans. Signal Process. 63.1, pp. 81–89. Schuller, B. M., E. & Berner, and J. & Delanardi (Ann. 2018). "Presence southing a pathog sequence of the polynomial EVD of parametricitian contents of polynomial evolution contents
processing algorithms". In: Proc. IEEE Int. Conf. on Acoust., Speech and Signal Process. (ICASSP), pp. 351–355.
 6.1, pp. 1–3. Therest A. (Aug. 2018). "Checification recomment methods". In Applied Computer and International Vice and Vice 2018.
In: Proc. IEEE Int. Conf. on Acoust., Speech and Signal Process. (ICASSP), pp. 4074–4077.
 study the effect of additive noise on speech recognition systems". In: Speech Commun. 3.3, pp. 247–251.

References

Very C. J. Vamarichi and S. King (Nev. 2013). "The voice bark comus design collection and data analysis of a large regional accent speech
database". In: Conf. Asian Spoken Language Research and Evaluation.
EVD". In: Proc. Eur. Signal Process. Conf. (EUSIPCO), pp. 844–848.
subspace approach". In: Sensor Signal Process. for Defence Conf. (SSPD).
matrix". In: IEEE Trans. Signal Process. 66.10, pp. 2659–2672.
Decomposition of a Parahermitian Matrix". In: IEEE Trans. Signal Process. 66.23, pp. 6325–6327.
IEEE/ACM Trans. Audio, Speech, Language Process. 24.2, pp. 252–264.



Thank you

Listening Examples: https://vwn09.github.io/research/pevd-vad

A Polynomial Subspace Projection Approach for the Detection of Weak Voice Activity - 32/32