

Variable Span Filtering for Speech Enhancement

ICASSP, 20–25 March 2016

Jesper Rindom Jensen^{*,1}, Jacob Benesty^{1,2},
and Mads Græsbøll Christensen¹

*jrj@create.aau.dk

¹Audio Analysis Lab,
AD:MT
Aalborg University
Denmark

²INRS-EMT
University of Quebec
Canada

Partly funded by the Danish Council for Independent Research, grant ID: DFF – 1337-00084, and the Villum Foundation.



AALBORG UNIVERSITY
DENMARK

Agenda



Introduction

Signal Model and Problem Formulation

Joint Diagonalization

Filtering

Performance

Optimal Filter Designs

Experimental Results

Conclusions



Introduction

- ▶ Noise reduction/enhancement is essential in many multichannel applications, such as hearing-aids.
- ▶ Have been tackled using many methods, e.g., linear filtering, spectral subtractive, and subspace methods.
- ▶ We propose a new class of methods that combines the linear filtering and subspace approaches.
- ▶ This is achieved by designing filters using a joint diagonalization.
- ▶ These variable span filters give explicit control over noise reduction versus signal distortion.



Signal Model

Model:

M sensors captures a convolved source signal and noise:

$$y_m(t) = g_m(t) * s(t) + v_m(t) = x_m(t) + v_m(t), \quad (1)$$

$m = 1, \dots, M$, where

g_m : m 'th acoustic impulse response,

s : desired source,

v_m : noise captured by sensor m ,



Problem Formulation

Goal

Extract x_1 from y_m , $m = 1, \dots, M$ with little/no distortion and residual noise using optimal filters.

A few assumptions to facilitate the task:

1. x_m and v_m uncorrelated, zero mean, stationary, real and broadband,
2. sensor signals are aligned wrt. the source direction.



Frequency Domain Model

Using the STFT, we get:

$$Y_m(k, n) = X_m(k, n) + V_m(k, n), \quad (2)$$

$m = 1, \dots, M$, where

k & n : frequency and time indices,

Y, X, V : STFTs of y, x and v at k 'th frequency.

In vector format:

$$\mathbf{y}(k, n) = [Y_1(k, n) \quad Y_2(k, n) \quad \dots \quad Y_M(k, n)]^T \quad (3)$$

$$= \mathbf{x}(k, n) + \mathbf{v}(k, n). \quad (4)$$



Observation

Common assumption:

$$X_m(k, n) = G_m(k)S(k, n), m = 1, 2, \dots, M, \quad (5)$$

where G_m and S are STFTs of g_m and s .

Observation:

- ▶ only valid with inf. long windows $N_{\text{win}} \rightarrow \infty$ or periodic sources,
- ▶ window length limited in practice,
- ▶ leads to $\text{rank}(\Phi_{\mathbf{x}}(k, n)) = P > 1$,
- ▶ not accounted for in conventional methods.



Interframe Correlation

Successive time frames (N) taken into account:

$$\underline{\mathbf{y}}(k, n) = [\mathbf{y}^T(k, n) \quad \mathbf{y}^T(k, n-1) \quad \cdots \quad \mathbf{y}^T(k, n-N+1)]^T \quad (6)$$

$$= \underline{\mathbf{x}}(k, n) + \underline{\mathbf{v}}(k, n). \quad (7)$$

Correlation matrix of $\underline{\mathbf{y}}(k, n)$:

$$\Phi_{\underline{\mathbf{y}}}(k, n) = \text{E} [\underline{\mathbf{y}}(k, n)\underline{\mathbf{y}}^H(k, n)] \quad (8)$$

$$= \Phi_{\underline{\mathbf{x}}}(k, n) + \Phi_{\underline{\mathbf{v}}}(k, n), \quad (9)$$

where

$\Phi_{\underline{\mathbf{x}}}(k, n)$: correlation matrix of $\underline{\mathbf{x}}$ of rank $P < MN$,

$\Phi_{\underline{\mathbf{v}}}(k, n)$: correlation matrix of $\underline{\mathbf{v}}$ of rank MN .



Joint Diagonalization

Joint diagonalization of correlation matrices:

$$\mathbf{B}^H(k, n)\mathbf{\Phi}_{\underline{\mathbf{x}}}(k, n)\mathbf{B}(k, n) = \mathbf{\Lambda}(k, n), \quad (10)$$

$$\mathbf{B}^H(k, n)\mathbf{\Phi}_{\underline{\mathbf{v}}}(k, n)\mathbf{B}(k, n) = \mathbf{I}_{MN}, \quad (11)$$

where

B: full rank, $MN \times MN$ matrix,

Λ : diagonal matrix with P real, positive entries (sorted),

\mathbf{I}_{MN} : identity matrix with dimensions $MN \times MN$.

Λ and **B** are eigenvalue and -vector matrices of $\mathbf{\Phi}_{\underline{\mathbf{v}}}^{-1}\mathbf{\Phi}_{\underline{\mathbf{x}}}$, i.e.,

$$\mathbf{\Phi}_{\underline{\mathbf{v}}}^{-1}(k, n)\mathbf{\Phi}_{\underline{\mathbf{x}}}(k, n)\mathbf{B}(k, n) = \mathbf{B}(k, n)\mathbf{\Lambda}(k, n). \quad (12)$$



Filtering

Desired signal X_1 estimated from $\underline{\mathbf{y}}$ through filtering:

$$\underline{Z}(k, n) = \underline{\mathbf{h}}^H(k, n)\underline{\mathbf{y}}(k, n), \quad (13)$$

where

$$\underline{\mathbf{h}}(k, n) = [\mathbf{h}^T(k, n) \quad \mathbf{h}^T(k, n-1) \quad \dots \quad \mathbf{h}^T(k, n-N+1)]^T. \quad (14)$$

With \mathbf{B} as basis, the filter is

$$\underline{\mathbf{h}}(k, n) = \mathbf{B}(k, n)\underline{\mathbf{a}}(k, n), \quad (15)$$

where $\underline{\mathbf{a}}$ is a filter in \mathbf{B} , and

$$\underline{\mathbf{a}}(k, n) = [A_1(k, n) \quad \dots \quad A_{MN}(k, n)]^T. \quad (16)$$



Filtering

Using $\underline{\mathbf{a}}$, the signal estimate is

$$Z(k, n) = \underline{\mathbf{a}}^H(k, n)\mathbf{B}^H(k, n)\underline{\mathbf{x}}(k, n) + \underline{\mathbf{a}}^H(k, n)\mathbf{B}^H\underline{\mathbf{v}}(k, n) \quad (17)$$

$$= X_{\text{fd}}(k, n) + V_{\text{rn}}(k, n). \quad (18)$$

Variance of Z :

$$\phi_Z(k, n) = \underline{\mathbf{a}}^H(k, n)\mathbf{\Lambda}(k, n)\underline{\mathbf{a}}(k, n) + \underline{\mathbf{a}}^H(k, n)\underline{\mathbf{a}}(k, n) \quad (19)$$

$$= \phi_{X_{1,\text{fd}}}(k, n) + \phi_{V_{1,\text{rn}}}(k, n) \quad (20)$$



Performance

Output SNR:

$$\text{oSNR}[\underline{\mathbf{h}}(k, n)] = \frac{\phi_{X_{1,\text{fd}}}(k, n)}{\phi_{V_{1,\text{m}}}(k, n)} = \frac{\underline{\mathbf{a}}^H(k, n)\underline{\boldsymbol{\Lambda}}(k, n)\underline{\mathbf{a}}(k, n)}{\underline{\mathbf{a}}^H(k, n)\underline{\mathbf{a}}(k, n)} \quad (21)$$

Signal reduction factor:

$$\xi_{\text{sr}}[\underline{\mathbf{h}}(k, n)] = \frac{\phi_{X_1}(k, n)}{\phi_{X_{1,\text{fd}}}(k, n)} = \frac{\phi_{X_1}(k, n)}{\underline{\mathbf{a}}^H(k, n)\underline{\boldsymbol{\Lambda}}(k, n)\underline{\mathbf{a}}(k, n)} \quad (22)$$



Mean Squared Error

Error given by $\mathcal{E}(k, n) = Z(k, n) - X_1(k, n)$, leads to MSE

$$J[\mathbf{a}'(k, n)] = E \left[|\mathcal{E}(k, n)|^2 \right] = \underbrace{J_{\text{ds}}[\mathbf{a}'(k, n)]}_{\text{distortion MSE}} + \underbrace{J_{\text{rs}}[\mathbf{a}'(k, n)]}_{\text{residual noise MSE}},$$

where

$$J_{\text{ds}}[\mathbf{a}'(k, n)] = E \left[|X_1(k, n) - \mathbf{a}'^H(k, n)\mathbf{B}'^H(k, n)\underline{\mathbf{x}}(k, n)|^2 \right] \quad (23)$$

$$J_{\text{rs}}[\mathbf{a}'(k, n)] = E \left[|\mathbf{a}'^H(k, n)\mathbf{B}'^H(k, n)\underline{\mathbf{v}}(k, n)|^2 \right]. \quad (24)$$



Tradeoff Filter

A general filter is obtained by solving:

$$\min_{\mathbf{a}'} J_{\text{ds}}[\mathbf{a}'(k, n)] \quad \text{s.t.} \quad J_{\text{rs}}[\mathbf{a}'(k, n)] = \beta \phi_{V_1}(k, n), \quad (25)$$

where $0 \leq \beta \leq 1$ controls the level of noise reduction.

Tradeoff filter design:

$$\underline{\mathbf{h}}_{\text{T}, \mu}(k, n) = \sum_{p=1}^P \frac{\underline{\mathbf{b}}_p(k, n) \underline{\mathbf{b}}_p^H(k, n)}{\mu + \lambda_p(k, n)} \Phi_{\underline{\mathbf{x}}}(k, n) \underline{\mathbf{i}}, \quad (26)$$

with μ a Lagrange multiplier adjusting noise level.



General Tradeoff Filter

Even more general filter obtained by using an arbitrary number of eigenvalues instead

$$\underline{\mathbf{h}}_{\mu, Q}(k, n) = \sum_{q=1}^Q \frac{\underline{\mathbf{b}}_q(k, n) \underline{\mathbf{b}}_q^H(k, n)}{\mu + \lambda_q(k, n)} \Phi_{\underline{\mathbf{x}}}(k, n) \mathbf{i}. \quad (27)$$

Observations:

- ▶ $\underline{\mathbf{h}}_{0,1}(k, n) = \underline{\mathbf{h}}_{\max}(k, n)$, (max SNR filter)
- ▶ $\underline{\mathbf{h}}_{1,P}(k, n) = \underline{\mathbf{h}}_W(k, n)$, (general Wiener filter)
- ▶ $\underline{\mathbf{h}}_{0,P}(k, n) = \underline{\mathbf{h}}_{\text{MVDR}}(k, n)$, (distortionless filter)
- ▶ $\underline{\mathbf{h}}_{0,Q}(k, n) = \underline{\mathbf{h}}_{\text{MD}}(k, n)$, (minimum distortion filter)
- ▶ $\underline{\mathbf{h}}_{\mu,P}(k, n) = \underline{\mathbf{h}}_{T,\mu}(k, n)$.



Experimental Results

Estimation of statistics

The statistics were estimated directly from the speech and noise signals, respectively.

It was conducted recursively using the following general equation for approximating the correlation matrix of a vector $\underline{\mathbf{a}}(k, n)$:

$$\widehat{\Phi}_{\underline{\mathbf{a}}}(k, n) = (1 - \xi)\widehat{\Phi}_{\underline{\mathbf{a}}}(k, n - 1) + \xi\underline{\mathbf{a}}(k, n)\underline{\mathbf{a}}(k, n)^H, \quad (28)$$

where ξ is a forgetting factor.

Forgetting factors for the signal and noise statistics estimation were chosen as $\xi_s = 0.05$ and $\xi_n = 0.05$, respectively.

Simulation Setup

Parameters:

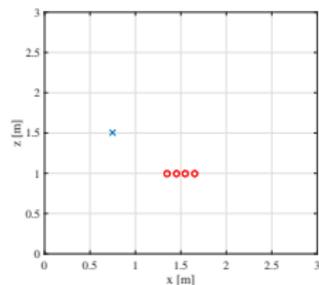
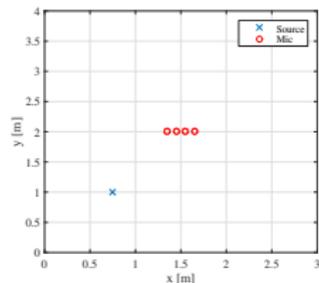
- ▶ sensor distance: 5 cm
- ▶ sound speed: 343 m/s
- ▶ reverb time: 0.2 s
- ▶ RIR length: 2,048
- ▶ mic type: omnidirectional

Signals:

Desired: speech (2 female and 2 male),

Noise: diffusive (babble) + sensor noise (white).

Room layout:

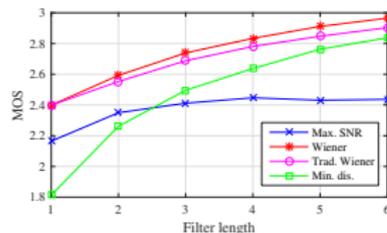
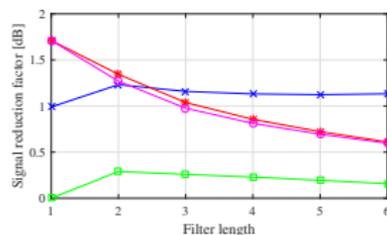
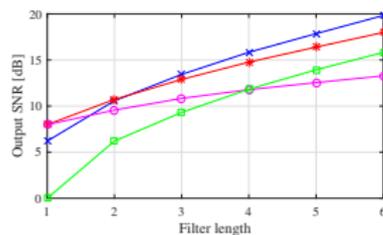


Results

Evaluation vs. filter length

Parameters:

- ▶ # of sensors: 3
- ▶ SDNR: 0 dB
- ▶ SSNR: 30 dB
- ▶ window length: 40
- ▶ FFT length: 64
- ▶ Assumed rank: 3

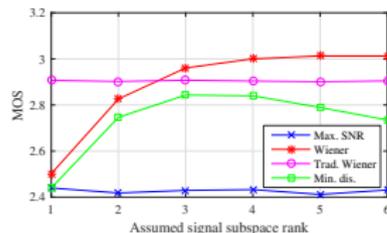
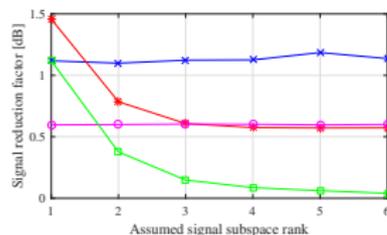
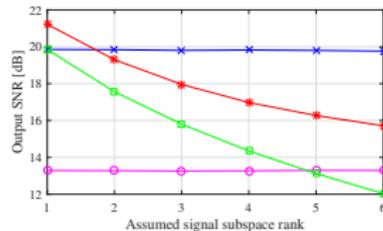


Results

Evaluation vs. filter rank

Parameters:

- ▶ # of sensors: 3
- ▶ SDNR: 0 dB
- ▶ SSNR: 30 dB
- ▶ window length: 40
- ▶ FFT length: 64
- ▶ Filter length: 6



Examples

Parameters:

of sensors: 3, SDNR: 0 dB, SSNR: 30 dB, window length: 40, FFT length: 64, filter length: 4.

Clean



Noisy



Max. SNR



Min. Dis.



Wiener



Trad. Wiener



Conclusions

- ▶ We considered the topic of multichannel speech enhancement.
- ▶ Proposed a new class of variable span filters (STFT domain).
- ▶ Unifies the filtering and subspace approaches.
- ▶ Provides explicit control over noise reduction versus signal distortion.
- ▶ Encompasses many well-known filter designs.
- ▶ Can outperform the traditional filtering counterparts according to experimental results (oSNR, PESQ, distortion).

