

# A Priori SNR Estimation Using Weibull Mixture Model

12. ITG Fachtagung Sprachkommunikation

Aleksej Chinaev, Jens Heitkaemper, Reinhold Haeb-Umbach

Department of Communications Engineering  
Paderborn University

7. Oktober 2016

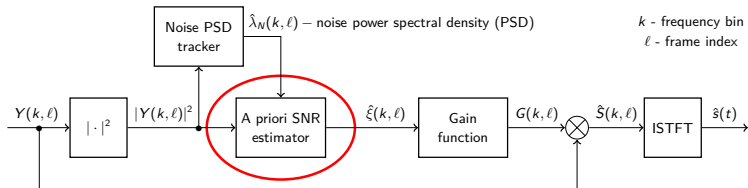
# Table of contents

- 1 Problem formulation and motivation
- 2 A priori SNR estimation based on Weibull mixture model
- 3 Experimental evaluation
- 4 Conclusions and outlook

# Problem formulation and motivation

- Single-channel clean speech  $s(t)$  contaminated by an additive noise  $n(t)$ :

$$y(t) = s(t) + n(t) \xrightarrow{STFT} Y(k, \ell) = S(k, \ell) + N(k, \ell)$$



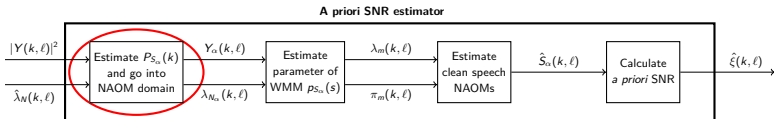
- A priori SNR  $\xi(k, \ell) = \frac{\lambda_S(k, \ell)}{\lambda_N(k, \ell)}$  - a key component in enhancement system  
 $\lambda_S(k, \ell) = E [|S(k, \ell)|^2]$  - clean speech PSD,  $\lambda_N(k, \ell) = E [|N(k, \ell)|^2]$  - noise PSD
- Motivated by a generalized spectral subtraction (GSS) denoising  $|Y(k, \ell)|^\alpha$  for  $\alpha \in \mathbb{R}_{>0}$  not restricted to  $(\alpha = 1)$  or  $(\alpha = 2)$  with assumption

$$|Y(k, \ell)|^\alpha = |S(k, \ell)|^\alpha + |N(k, \ell)|^\alpha$$

# Table of contents

- 1 Problem formulation and motivation
- 2 A priori SNR estimation based on Weibull mixture model**
- 3 Experimental evaluation
- 4 Conclusions and outlook

# Normalized $\alpha$ -order magnitude (NAOM) domain



- Normalize  $|Y(k, \ell)|^\alpha$  to a root of an averaged power  $P_{S_\alpha}(k)$  of  $|S(k, \ell)|^\alpha$

$$Y_\alpha(k, \ell) = \frac{|Y(k, \ell)|^\alpha}{\sqrt{P_{S_\alpha}(k)}} = S_\alpha(k, \ell) + N_\alpha(k, \ell) \quad \text{with} \quad P_{S_\alpha}(k) = \frac{1}{L} \sum_{\ell=1}^L |S(k, \ell)|^{2\alpha}$$

- Statistical models independent of speaker loudness
- Normalized energy of clean speech NAOMs  $E[S_\alpha^2(k)] = 1$
- $S_\alpha(k, \ell)$  &  $N_\alpha(k, \ell)$  – realizations of random variables  $S_\alpha(k)$  &  $N_\alpha(k)$

Estimate  $S_\alpha(k, \ell)$  from  $Y_\alpha(k, \ell)$  given models for  $S_\alpha(k)$  &  $N_\alpha(k)$

# Modeling of noise NAOM coefficients $N_\alpha(k, \ell)$

- $N(k, \ell) \sim \mathcal{N}_c(n; 0, \lambda_N(k, \ell))$

- $N_\alpha(k, \ell)$  – Weibull distributed

$$p_{N_\alpha(k, \ell)}(n) = \text{Weib}(n; \lambda_{N_\alpha}(k, \ell), \alpha)$$

- Shape parameter  $\alpha \in \mathbb{R}_{>0}$
- Scale parameter

$$\lambda_{N_\alpha}(k, \ell) = \frac{\lambda_N(k, \ell)}{\alpha \sqrt[\alpha]{P_{S_\alpha}(k)}} \in \mathbb{R}_{>0}$$

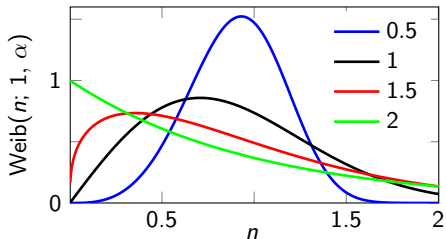
- Model  $N_\alpha(k)$  with Weibull PDF

$$p_{N_\alpha(k)}(n) = \text{Weib}(n; \lambda_{N_\alpha}(k), \alpha)$$

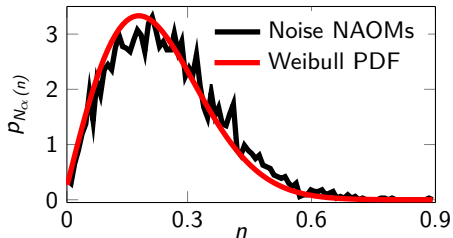
$$\text{with } \lambda_{N_\alpha}(k) = \frac{1}{L} \sum_{\ell=1}^L \lambda_{N_\alpha}(k, \ell)$$

- NAOM coefficients of white noise signal and estimated  $p_{N_\alpha(k)}(n)$

Weibull PDF for  $\lambda = 1$  and different  $\alpha$



Histogram and Weibull PDF for  $\alpha = 0.7$



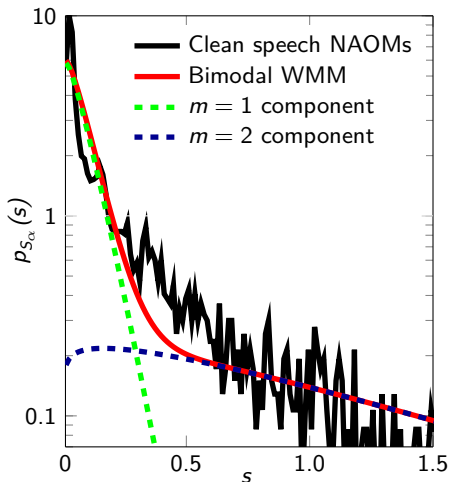
# Modeling of NAOM coefficients of clean speech $S_\alpha(k, \ell)$

- $S(k, \ell) \sim \mathcal{N}_c(n; 0, \lambda_S(k, \ell))$
- Bimodal Weibull mixture model (WMM) to model  $S_\alpha(k)$

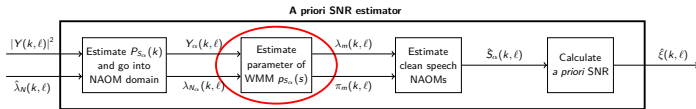
$$p_{S_\alpha(k)}(s) = \sum_{m=1}^2 \pi_m(k) \cdot \text{Weib}(s; \lambda_m(k), \beta)$$

- $m = 1$  : silence
- $m = 2$  : activity
- $\pi_m(k) \in [0, 1]$ : weights
- $\lambda_m(k)$ : scale parameters
- $\beta$ : shape parameter
- $\beta \neq \alpha$  : additional degree of freedom in the model
- Clean speech NAOMs & estimated WMM ( $\alpha = 0.7$ ;  $\beta = 2.5$ )

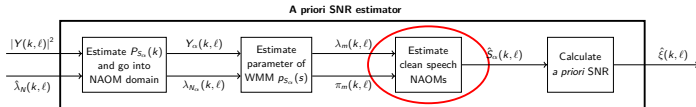
Histogram and estimated WMM



# Estimation of WMM parameters and clean speech NAOMs



- Set  $\lambda_1(k)$  acc. to  $\xi_{\min}$  usually used in a *a priori* SNR estimation [Cappe 94]
- Expectation Maximization algorithm to estimate  $\lambda_2(k)$ ,  $\pi_m(k)$ 
  - After EM, weights  $\pi_m(k)$  are corrected with the constraint  $E[S_\alpha^2(k)] = 1$



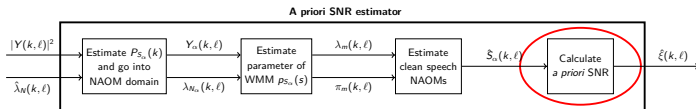
- Maximum a posteriori (MAP) estimation:

$$\hat{S}_\alpha^{\text{MAP}}(k, \ell) = \underset{s}{\operatorname{argmax}} p_{S_\alpha(k)} | Y_\alpha(k, \ell)(s|y)$$

- $Y_\alpha(k, \ell)$  is a realisation of random variable  $Y_\alpha(k) = S_\alpha(k) + N_\alpha(k)$
- Approximative computationally efficient solution for  $\beta = \alpha = 1$



# Calculation of *a priori* SNR and causal implementation



- Go back into domain of power spectral density by calculating

$$\hat{\xi}(k, \ell) = \max \left( \frac{\left[ \hat{S}_\alpha(k, \ell) \cdot \sqrt{P_{S_\alpha}(k)} \right]^{\frac{2}{\alpha}}}{\lambda_N(k, \ell)}, \xi_{\min} \right)$$

## Causal implementation of WMM-based *a priori* SNR estimators

- Calculate  $P_{S_\alpha}(k)$  and  $\lambda_{N_\alpha}(k)$  in a causal way
- Causal EM for  $\lambda_2(k)$  and  $\pi_2(k)$  with one EM-iteration per time frame
- Note, parameters  $\alpha$  and  $\beta$  have to be set appropriately  $\rightarrow$  optimization

# Table of contents

- 1 Problem formulation and motivation
- 2 A priori SNR estimation based on Weibull mixture model
- 3 Experimental evaluation**
- 4 Conclusions and outlook

# Experimental evaluation

## Data and setup

- Clean speech: *Wall Street Journal* database 16 kHz (male and female)
- 7 different noise types of *Noisex92* database: *white*, *pink*, *f16*, *hfchannel*, *factory-1*, *factory-2*, *babble*
- Input global SNR from  $-5$  dB up to 25 dB in 5 dB steps

## Spectral speech enhancement framework

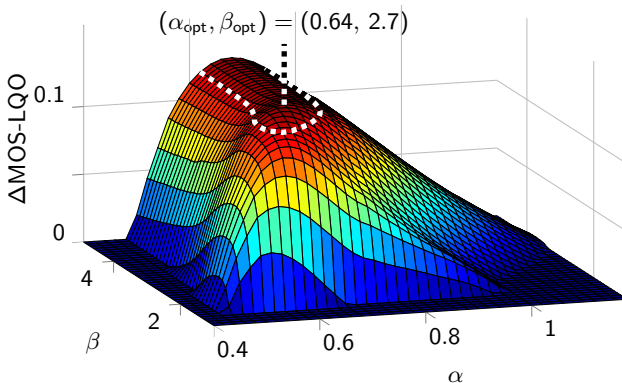
- Noise PSD tracking using *Minimum statistics* approach [Martin 01]
- A priori SNR estimation with  $\xi_{\min} = -18$  dB [Cappe 94]
  - Proposed WMM-based approach with Wiener filter
  - Reference approach: *Decision Directed* [Ephraim 84]

# Optimization of $\alpha$ and $\beta$

- Speech quality maximization in terms of wide-band mean opinion score listening quality objective (MOS-LQO) with

$$\Delta \text{MOS-LQO} = \max(\text{MOS-LQO}_{\text{WMM}} - \text{MOS-LQO}_{\text{DD}}, 0)$$

- Averaging over genders, noise types and input global SNR values



# Final experimental results

- Clean speech: WSJ database signals other than used for optimization
- Estimation error – Itakura-Saito distance (ISD) and estimator's variance – logarithmic error variance (LEV): the smaller the better

Resulting ISD, LEV and MOS-LQO values averaged over noise types

| SNR, dB |     | -5          | 0           | 5           | 10          | 15          | 20          | 25          | AVG         |
|---------|-----|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| ISD     | DD  | 48.8        | 44.0        | 39.6        | 34.9        | 30.2        | 24.5        | 19.1        | 34.4        |
|         | WMM | <b>42.6</b> | <b>38.1</b> | <b>34.1</b> | <b>30.4</b> | <b>27.3</b> | <b>23.0</b> | <b>18.9</b> | <b>30.6</b> |
| LEV     | DD  | 53.1        | 49.0        | 46.4        | 45.1        | 45.5        | 47.4        | 50.5        | 48.1        |
|         | WMM | <b>45.6</b> | <b>43.9</b> | <b>42.6</b> | <b>41.1</b> | <b>39.0</b> | <b>37.0</b> | <b>35.9</b> | <b>40.7</b> |
| MOS-LQO | DD  | 1.11        | 1.30        | 1.63        | 2.09        | 2.57        | 3.00        | 3.39        | 2.16        |
|         | WMM | <b>1.18</b> | <b>1.46</b> | <b>1.77</b> | <b>2.13</b> | <b>2.62</b> | <b>3.16</b> | <b>3.61</b> | <b>2.28</b> |

# Conclusions and outlook

## Conclusions

- Novel causal *a priori* SNR estimator based on a bimodal Weibull mixture model for the normalized  $\alpha$ -order spectral magnitudes (NAOMs)
- Optimization of the proposed approach by maximization of speech quality
  - Power exponent  $\alpha_{\text{opt}} = 0.64$  smaller than 1 (spectral magnitudes)
  - Shape factor  $\beta_{\text{opt}} = 2.7$  – a heavier tailed Weibull distribution
- Compared to the wide-spread *Decision Directed* approach:
  - Reduced error and variance of the WMM-based *a priori* SNR estimator
  - Improvement of speech quality of the enhanced signals
  - Higher computational effort

## Outlook

- Reduction of computational effort – fixed speaker-independent models
- Development of model-based spectral enhancement using generalized (arbitrary) power exponent in the spirit of generalized spectral subtraction



**Thank you for your attention!**

**Questions?**

Paderborn University

Department of  
Communications Engineering

Web: [nt.upb.de](http://nt.upb.de)

# Resulting WMM parameter and audio samples

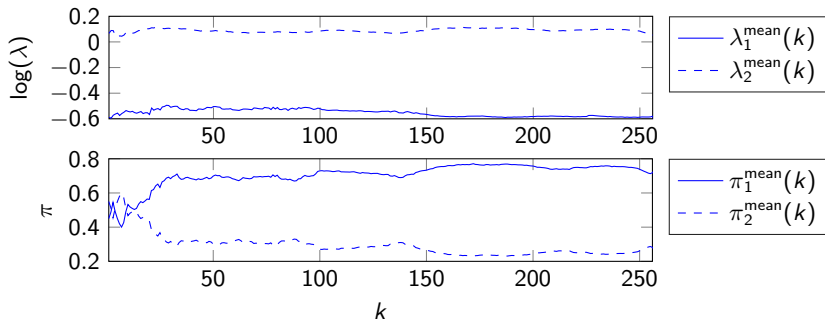


Figure : Resulting WMM parameter over frequency bins

- Exemplarily speech samples: **Noisy** **DD** **WMM**