

Improved noise power spectral density tracking by a MAP-based postprocessor

Aleksej Chinaev, Alexander Krueger, Dang Hai Tran Vu, Reinhold Haeb-Umbach

University of Paderborn, Germany

March 28th, 2012

Computer Science, Electrical Engineering and Mathematics



Communications Engineering Prof. Dr.-Ing. Reinhold Häb-Umbach



Table of Contents

- 1 Introduction and problem formulation
- 2 MAP-based noise PSD estimation
- 3 Experimental framework
- 4 Performance evaluation
- 5 Summary and outlook





A. Chinaev, A. Krueger, D.H. Tran Vu, R. Haeb-Umbach



Introduction



Motivation

- Noise PSD estimation is a key component
 - to speech enhancement
 - to robust automatic speech recognition

Basic assumptions

- Many sophisticated algorithms rely on two assumptions:
 - Noise 'more stationary' than speech
 - Noise-only time-frequency bins at regular intervals

Here

- MAP-based (MAP-B) postprocessor
 - Estimate of noise power even if speech is dominant in time-frequency bin
 - Initial estimate of the current speech power required







MAP-B as a postprocessor



Improved noise power spectral density tracking by a MAP-based postprocessor



A. Chinaev, A. Krueger, D.H. Tran Vu, R. Haeb-Umbach

Problem formulation

Observed

$$\mathbf{Y}_l = \mathbf{N}_l + \mathbf{X}_l$$

- We consider \mathbf{N}_l as target process 'corrupted' by clean speech \mathbf{X}_l of known variance $\sigma^2_{\mathbf{X},l}$

Our goal

- Estimation of noise variance $\sigma_{{\bf N},l+1}^2$ at frame l+1, from the noisy observation ${\bf y}_{l+1}$ given
 - the current speech power $\sigma^2_{\mathbf{X},l+1}$
 - \blacktriangleright a priori PDF $p_{\sigma^2_{\mathbf{N}}}(\sigma^2)$ of variance $\sigma^2_{\mathbf{N},l}$

Approach

• Maximum A Posteriori (MAP) estimation





MAP-based noise PSD estimation

Bayesian variance estimation: $\mathbf{Y}_l = \mathbf{N}_l$

If $\sigma^2_{\mathbf{X},l+1} = 0$ (uncorrupted observ.) and $\sigma^2_{\mathbf{N},l+1} = \sigma^2_{\mathbf{N}}$ (stationary) then the textbook problem:

• Scaled inverse chi-square distribution

$$p_{\sigma_{\mathbf{N}}^2}(\sigma^2) \propto (\sigma^2)^{-\frac{\nu_l+2}{2}} \cdot e^{-\frac{\nu_l \lambda_l^2}{2\sigma^2}}$$

with the degrees of freedom ν_l and the scale factor λ_l^2 is conjugate prior to normal observation PDF

$$p_{\mathbf{Y}_{l+1}|\sigma_{\mathbf{N}}^{2}}(\mathbf{y}_{l+1}|\sigma^{2}) = \frac{1}{\pi\sigma^{2}} \cdot e^{-\frac{|\mathbf{y}_{l+1}|^{2}}{\sigma^{2}}}$$

- Parameter update $\nu_{l+1} = \nu_l + 2$ and $\lambda_{l+1}^2 = \frac{2}{\nu_l+2} |\mathbf{y}_{l+1}|^2 + \frac{\nu_l}{\nu_l+2} \lambda_l^2$
- MAP-estimate of variance

$$\hat{\sigma}_{\mathbf{N},l+1}^2 = \operatorname*{argmax}_{\sigma^2} \left[p_{\sigma_{\mathbf{N}}^2 | \mathbf{Y}_{l+1}}(\sigma^2 | \mathbf{y}_{l+1}) \right] = \frac{\nu_{l+1}}{\nu_{l+1}+2} \cdot \lambda_{l+1}^2$$





Extension to non-stationary noise

Still 'uncorrupted' noise: $\mathbf{Y}_l = \mathbf{N}_l$

If $\sigma^2_{\mathbf{N},l+1}$ is time-variant then:

- The parameter ν_l is kept at a constant value $\nu_{l+1}=\nu_l=\nu_0$
- This results in recursive smoothing of variance estimate

$$\hat{\sigma}_{\mathbf{N},l+1}^2 = (1-\alpha) \cdot \hat{\sigma}_{\mathbf{N},l}^2 + \alpha \cdot |\mathbf{y}_{l+1}|^2, \text{ where } \alpha = \frac{2}{\nu_0 + 4}$$

Choice of ν_0

• Trade-off between tracking ability and estimation error in stationary noise







Observation 'corrupted' by speech: $\mathbf{Y}_l = \mathbf{N}_l + \mathbf{X}_l$

• If
$$\sigma_{\mathbf{X},l+1}^2 \neq 0$$
 and is known then the posterior PDF
 $p_{\sigma_{\mathbf{X}}^2|\mathbf{Y}_{l+1}}(\sigma^2|\mathbf{y}_{l+1}) \propto (\sigma_{\mathbf{X},l+1}^2 + \sigma^2)^{-1} \cdot (\sigma^2)^{-\frac{\nu_l+2}{2}} \cdot e^{-\left(\frac{|\mathbf{y}_{l+1}|^2}{\sigma_{\mathbf{X},l+1}^2 + \sigma^2} + \frac{\nu_l \lambda_l^2}{2\sigma^2}\right)},$
is no longer a conjugate prior for the observation PDF.

In order to maintain an efficient MAP estimation procedure we

- approximate the posterior PDF by a scaled inverse chi-squared distribution,
- and match its maximum $(\hat{\sigma}_{\mathbf{N},l+1}^2)$ with the maximum of the posterior PDF,
- which we calculate efficiently using a bisection and Newton approach.







Experimental framework

MAP-B as postprocessor

- First noise PSD estimator: Improved Minima Controlled Recursive Averaging (IMCRA) algorithm [Cohen, 2003]
- Gain function: Optimally-Modified Log-Spectral Amplitude (OM-LSA) estimator
 [Cohen, Berdugo, 2001]





Performance evaluation

Setup

- Clean speech: TIMIT database, sentences concatenated to 3 minutes length
- Artificially added noise from Noisex92 database:
 - ▶ Noise types: 'Stationary WGN', 'Triangular WGN', 'Babble' and 'Factory-1'
 - ▶ SNR values: -5,0,5,10,15 dB
- MAP-B estimator: we set $u_0 = 40$ corresponding to a time constant of \simeq 0.164 s

Reference noise PSD

• Recursive temporal smoothing

$$\sigma_{\mathbf{N},k,l}^2 = 0.95 \cdot \sigma_{\mathbf{N},k,l-1}^2 + 0.05 \cdot |\mathbf{N}_{k,l}|^2, \text{ with known noise periodogram } |\mathbf{N}_{k,l}|^2$$





Sample trajectories of noise variance estimates

• 'Babbble' noise at frequency bin k = 97 (3 kHz)



 MAP-B: continious update of noise variance estimate • 'Triangular WGN' (averaged over all frequency bins)



noise power



Quantitative evaluation

- Performance measures adopted from [Taghia et al., 2011]
- Fig.(a): minimum averaged log distance *LE_m* between the reference and estimated PSD
- MAP-B obtains lower error LE_m for all noise types and SNRs less than or equal to 5 dB or 10 dB
- Fig.(b): variance of the logarithmic difference *LE_v*
- MAP-B yields lower variance LE_v for all noise types and SNRs than the IMCRA









Speech enhancement

• Gain in perceptual speech quality (PESQ) scores

$$PESQ_{Gain} = PESQ_{MAP-B} - PESQ_{IMCRA}$$



• MAP-B has a favourable effect on speech quality for non-stationary noise types



Summary and outlook

Summary

- Proposed MAP-B estimator is able to track the noise statistics even if the speech is dominant
- Low computational complexity
- Single parameter ν_0
- Experimental evaluation: MAP-B obtains
 - Iower estimation error under Iow SNR conditions
 - Iower fluctuation of the estimated values under all tested environments
 - slightly improved speech quality for non-stationary noise types

Outlook

- Investigations about dependence of MAP-B algorithm performance:
 - on the first noise PSD estimator
 - \blacktriangleright on the degrees of freedom ν_0

Improved noise power spectral density tracking by a MAP-based postprocessor

A. Chinaev, A. Krueger, D.H. Tran Vu, R. Haeb-Umbach





Thank you for your attention!



Questions?

Computer Science, Electrical Engineering and Mathematics



Communications Engineering Prof. Dr.-Ing. Reinhold Häb-Umbach

Periodograms of cleen, noisy and enhanced signals

- Periodogramms of the spoken sentence 'Biblical scholars argue history' for 'factory-1' noise at an SNR of 5 dB
 - (a) cleen speech signal
 - (b) noisy speech signal
 - (c) enhanced speech signal based on IMCRA estimates
 - (d) enhanced speech signal based on MAP-B estimates
- MAP-B has a positive influence on periodogramm of enhanced signal particularly clearly seen in the highlighted window







Minimum Statistics (MS) instead of IMCRA

• Performance measures LE_m and LE_v for female speaker signals



- Fig.(a): MAP-B obtains lower estimation error LE_m than the MS for all noise types and SNRs
- Fig.(b): MAP-B yields lower variance LE_v than the MS for all tested setups







Minimum Statistics (MS) instead of IMCRA

• $PESQ_{Gain} = PESQ_{MAP-B} - PESQ_{NoiseEst}$

with 'NoiseEst' \in ['IMCRA'; 'MS'] for female speaker signals



• A favourable effect of MAP-B estimator on speech quality for non-stationary noise types is for MS smaller than for IMCRA



