ADAPTIVE SUPPRESSION OF NON-STATIONARY NOISE BY USING THE VARIATIONAL BAYESIAN METHOD

Takuya Yoshioka and Masato Miyoshi

NTT Communication Science Laboratories, NTT Corporation 2-4, Hikari-dai, Seika-cho, Soraku-gun, Kyoto 619-0237, Japan

ABSTRACT

This paper proposes an adaptive noise suppression method for nonstationary noise based on the Bayesian estimation method. The following conditions are assumed: (1) Speech and noise samples are statistically independent, and they follow auto-regressive (AR) processes. (2) The prior distribution of the parameters of the noise AR model of a current frame is identical to the posterior distribution of those parameters calculated in the previous frame. Under these conditions, the proposed method approximates the joint posterior distribution of the AR model parameters and the speech samples by using the variational Bayesian method. Furthermore, we describe an efficient implementation by assuming that all involved covariance matrices have the Toeplitz structure. The proposed method was tested on real speech and noise signals and compared with other noise suppression methods.

Index Terms— Noise suppression, Bayesian estimation, variational Bayesian method, auto-regressive process

1. INTRODUCTION

The task of suppressing noise in degraded speech signals observed at microphones has been of great interest for decades. Noise suppressors working in realistic environments must be able to cope with non-stationary noise. Traditional noise suppression methods that are based on voice activity detectors (VADs) are unable to suppress nonstationary noise very much because they update noise estimates only during periods when speech is absent. As a result, adaptive noise suppression methods that can update noise estimates even during speech activity have been studied in the last decade.

One approach to adaptive noise suppression is based on statistical models of speech and noise [1, 2, 3]. The statistical model based approach estimates the model parameters or the speech signals from the observed noisy speech signals by using the maximum likelihood (ML) estimation method, the Bayesian estimation method, or the like. Auto-regressive (AR) models are usually employed for both the speech and noise models. The statistical model based approach seems to produce less distorted speech signals than the minimum statistics based approach [4]. This may be partly because the minimum statistics based approach estimates noise power spectral components independently for each frequency bin, and partly because this approach involves spectral smoothing, which makes speech signals reverberant.

The statistical model based methods are further classified into two categories. One requires the speech and noise models to be trained in advance [1, 2]. The methods in this category fail to recover the speech signals contaminated by noise signals out of the training data. The other estimates both the model parameters and the speech signals online without the need for prior training of the models [3]. However, the methods in this category are, as far as we know, based on the ML estimation method, which is in principle inferior to the Bayesian estimation method.

In this paper, we propose an adaptive noise suppression method that performs the Bayesian estimation of the model parameters and the speech signals. That is, the proposed method calculates the joint posterior distribution of the AR parameters and the speech signals on a frame-by-frame basis. The proposed method assumes that the speech and noise signal samples follow auto-regressive (AR) processes. Moreover, the prior distribution of the parameters of the noise AR model of a current frame is assumed to be identical to the posterior distribution of those parameters calculated in the previous frame. Since it is difficult to derive the exact posterior distribution, in reality, the proposed method approximates the true posterior distribution by using the variational Bayesian method [5]. Furthermore, we describe an efficient implementation of the proposed method, resulting in only a slightly higher computational cost than the ML estimation method. Note that the proposed method can be distinguished from the method described in [6], which is also based on the variational Bayesian method, in that the proposed method does not assume the noise to be stationary.

2. BAYESIAN APPROACH TO ADAPTIVE NOISE SUPPRESSION

2.1. Task formulation of adaptive noise suppression

Let $s_t(n)$, $v_t(n)$, and $x_t(n)$ denote a speech sample, a noise sample, and an observed noisy speech sample, respectively, in the *t*-th short time frame of length *N*. The vector of the noisy speech samples contained in the *t*-th frame, $\boldsymbol{x}_t = [x_t(N), \cdots, x_t(1)]^T$, is represented by

$$\boldsymbol{x}_t = \boldsymbol{s}_t + \boldsymbol{v}_t, \tag{1}$$

where $\boldsymbol{s}_t = [s_t(N), \dots, s_t(1)]^T$, $\boldsymbol{v}_t = [v_t(N), \dots, v_t(1)]^T$, and superscript T is the transpose operator. The noise suppression task addressed in this paper is to estimate the speech samples of the t-th frame, \boldsymbol{s}_t , from the observed noisy speech samples up to the t-th frame, $\boldsymbol{x}_1^t = \{\boldsymbol{x}_u\}_{1 \le u \le t}$.

If we are to solve the noise suppression task, we must introduce some constraints on the speech and noise. In this paper, we assume the following conditions.

(a1) The speech sample $s_t(n)$ follows an auto-regressive (AR) process of order P. Thus, the probability density function (PDF) of $s_t(n)$ conditioned on the past P samples $\boldsymbol{s}_t^P(n-1) = [s_t(n-1), \cdots, s_t(n-P)]^T$, the regression coefficient vector $\boldsymbol{a}_t = [a_{t,1}, \cdots, a_{t,P}]^T$, and the innovation variance

 σ_t^2 is given by

$$p(s_t(n)|\boldsymbol{s}_t^P(n-1), \boldsymbol{a}_t, \sigma_t^2) = \mathcal{N}\{s_t(n); \boldsymbol{s}_t^P(n-1)^T \boldsymbol{a}_t, \sigma_t^2\}, \qquad (2)$$

where $\mathcal{N}\{\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}\}$ denotes the PDF of a random variable \boldsymbol{x} following the normal distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. We hereafter represent the parameters of this speech model by $\Phi_t = \{\boldsymbol{a}_t, \sigma_t^2\}$.

(a2) The noise sample $v_t(n)$ follows an AR process of order Q. Thus, the PDF of $v_t(n)$ conditioned on the past Q samples $\boldsymbol{v}_t^Q(n-1) = [v_t(n-1), \cdots, v_t(n-Q)]^T$, the regression coefficient vector $\boldsymbol{b}_t = [b_{t,1}, \cdots, b_{t,Q}]^T$, and the innovation variance γ_t^2 is given by

$$p(v_t(n)|\boldsymbol{v}_t^Q(n-1), \boldsymbol{b}_t, \gamma_t^2) = \mathcal{N}\{v_t(n); \boldsymbol{v}_t^Q(n-1)^T \boldsymbol{b}_t, \gamma_t^2\}.$$
 (3)

We hereafter represent the parameters of this noise model by $\Psi_t = \{ \boldsymbol{b}_t, \gamma_t^2 \}.$

(a3) $s_t(m)$ and $v_u(n)$ are mutually independent random variables for any t, u, m, and n. Φ_t and Ψ_u are also random variables that are independent of each other.

Under the above assumptions, the conditional distribution of x_t and the conditional distribution of s_t are characterized respectively by

$$p(\boldsymbol{x}_t | \boldsymbol{s}_t, \Psi_t) = \prod_{n=1}^N \mathcal{N}\{\boldsymbol{x}_t(n); \boldsymbol{s}_t(n) + (\boldsymbol{x}_t^Q(n-1) - \boldsymbol{s}_t^Q(n-1))^T \boldsymbol{b}_t, \gamma_t^2\}$$
(4)

$$p(\boldsymbol{s}_t | \Phi_t) = \prod_{n=1}^N \mathcal{N}\{s_t(n); \boldsymbol{s}_t^P (n-1)^T \boldsymbol{a}_t, \sigma_t^2\}.$$
 (5)

In addition to the above, we further assume the following conditions as regards the dynamics of the parameters of the speech and noise models.

(a4) The parameters of the noise model vary so slowly that the prior distribution of Ψ_t given \boldsymbol{x}_1^{t-1} is identical to the posterior distribution of Ψ_{t-1} given \boldsymbol{x}_1^{t-1} , *i.e.*

$$p(\Psi_t | \boldsymbol{x}_1^{t-1}) = p(\Psi_{t-1} | \boldsymbol{x}_1^{t-1}).$$
(6)

Furthermore, this distribution is assumed to be the conjugate distribution as follows:

$$p(\Psi_t | \boldsymbol{x}_1^{t-1}) = \mathcal{N} \Big\{ \boldsymbol{b}_t; \boldsymbol{\nu}_{v,t-1}, \left(\frac{\zeta_{v,t-1} \Xi_{v,t-1}}{\gamma_t^2} \right)^{-1} \Big\} \\ \times \chi^{-2} \{ \gamma_t^2; \rho_{v,t-1}, \lambda_{v,t-1} \},$$
(7)

where $\chi^{-2}{\xi^2; \rho, \lambda}$ denotes the PDF of a random variable ξ^2 following the scaled inverse chi square distribution with degree of freedom ρ and scale parameter λ . The choice of the conjugate prior is for mathematical tractability.

(a5) The prior distribution of the speech model parameters Φ_t given \boldsymbol{x}_1^{t-1} is free from \boldsymbol{x}_1^{t-1} and is independent of the frame index *t*. Specifically, the prior distribution is given by

$$p(\Phi_t) = \mathcal{N}\left\{\boldsymbol{a}_t; \boldsymbol{\nu}_s', \left(\frac{\zeta_s'\Xi_s'}{\sigma_s^2}\right)^{-1}\right\} \chi^{-2}\left\{\sigma_t^2; \rho_s', \lambda_s'\right\}, \quad (8)$$

where we let ζ' and ρ' be a very small value to make the prior distribution non-informative.

Under the above conditions, we shall estimate s_t from x_1^t .

2.2. Bayesian approach

In this paper, we approach the above task based on the Bayesian estimation. The goal of the Bayesian estimation is to estimate the joint posterior distribution of the speech samples \boldsymbol{s}_t and the parameters $\Theta_t = \{\Phi_t, \Psi_t\}$ of the *t*-th frame, $p(\boldsymbol{s}_t, \Theta_t | \boldsymbol{x}_1^t)$. The posterior distribution can be factorized as

$$p(\boldsymbol{s}_t, \Theta_t | \boldsymbol{x}_1^t) \propto p(\boldsymbol{x}_t | \boldsymbol{s}_t, \Psi_t) p(\boldsymbol{s}_t | \Phi_t) p(\Phi_t) p(\Psi_t | \boldsymbol{x}_1^{t-1}).$$
(9)

The terms of the right hand side of (9) have been already defined in (4), (5), (8), and (7). Moreover, by integrating out \boldsymbol{s}_t and Φ_t from (9), we obtain the posterior distribution for the *t*-th frame, $p(\Psi_t | \boldsymbol{x}_t^1)$, which is used as the prior distribution for the (t + 1)th frame, $p(\Psi_{t+1} | \boldsymbol{x}_t^1)$. Importantly, with the variational approximation described in the next section, $p(\Psi_t | \boldsymbol{x}_t^1)$ has the same form as $p(\Psi_t | \boldsymbol{x}_t^{t-1})$ by virtue of $p(\Psi_t | \boldsymbol{x}_t^{t-1})$ being conjugate. Thus, in principle, every time noisy speech samples of a new frame are observed, we can obtain the posterior distribution of the speech samples and the parameters of that frame adaptively.

3. ALGORITHM BASED ON VARIATIONAL BAYESIAN METHOD

3.1. Variational Bayesian method

We employ the variational Bayesian method, which approximates the posterior distribution $p(s_t, \Theta_t | \boldsymbol{x}_1^t)$ because it is difficult to derive the analytic form of $p(\boldsymbol{s}_t, \Theta_t | \boldsymbol{x}_1^t)$. The difficulty arises from the fact that \boldsymbol{s}_t is a hidden variable in the model that generates observed variable \boldsymbol{x}_t . The variational Bayesian method approximates the true posterior distribution $p(\boldsymbol{s}_t, \Theta_t | \boldsymbol{x}_1^t)$ by using the hypothetical posterior distribution of the following factorized form:

$$q(\boldsymbol{s}_t, \Theta_t) = q(\boldsymbol{s}_t)q(\Theta_t). \tag{10}$$

 $q(\mathbf{s}_t)$ and $q(\Theta_t)$ are calculated so that the Kullback-Leibler (KL) divergence between $p(\mathbf{s}_t, \Theta_t | \mathbf{x}_1^t)$ and $q(\mathbf{s}_t, \Theta_t)$ is minimized.

Such $q(s_t)$ and $q(\Theta_t)$, which are called optimal variational posterior distributions, can be proven to satisfy

$$q(\boldsymbol{s}_t) \propto \exp\{\langle \log p(\boldsymbol{x}_t, \boldsymbol{s}_t | \Theta_t) \rangle_{q(\Theta_t)}\}$$
(11)

$$q(\Theta_t) \propto p(\Theta_t | \boldsymbol{x}_1^{t-1}) \exp\{\langle \log p(\boldsymbol{x}_t, \boldsymbol{s}_t | \Theta_t) \rangle_{q(\boldsymbol{s}_t)}\}, \quad (12)$$

where the complete data likelihood $p(\boldsymbol{x}_t, \boldsymbol{s}_t | \Theta_t)$ is obtained by multiplying (4) and (5). It is obvious that these two conditions are dependent on each other. Therefore, the variational Bayesian method first calculates $q(\boldsymbol{s}_t)$ according to (11) for a fixed $q(\Theta_t)$, and then calculates $q(\Theta_t)$ for a fixed $q(\boldsymbol{s}_t)$ according to (12). $q(\boldsymbol{s}_t)$ and $q(\Theta_t)$ are obtained by repeating these two steps until convergence. Below we derive specific forms of the optimal variational posterior distributions $q(\Theta_t)$ and $q(\boldsymbol{s}_t)$.

3.2. Optimal variational posterior distribution of parameters

It is shown that substituting (4) and (5) into (12) leads to

$$q(\Theta_t) = \mathcal{N} \Big\{ \boldsymbol{a}_t; \boldsymbol{\nu}_{s,t}, \left(\zeta_{s,t} \frac{\Xi_{s,t}}{\sigma_t^2} \right)^{-1} \Big\} \chi^{-2} \{ \sigma_t^2; \rho_{s,t}, \lambda_{s,t} \} \\ \times \mathcal{N} \Big\{ \boldsymbol{b}_t; \boldsymbol{\nu}_{v,t}, \left(\zeta_{v,t} \frac{\Xi_{v,t}}{\gamma_t^2} \right)^{-1} \Big\} \chi^{-2} \{ \gamma_t^2; \rho_{v,t}, \lambda_{v,t} \}.$$
(13)

In (13), the hyperparameters of the speech model (the first line of (13)) are defined as follows.

$$\zeta_{s,t} = N + \zeta'_s \tag{14}$$

$$\Xi_{s,t} = \frac{NR_{s,t} + \zeta_s'\Xi_s'}{N + \zeta_s'} \tag{15}$$

$$\boldsymbol{\nu}_{s,t} = \Xi_{s,t}^{-1} \boldsymbol{\xi}_{s,t} \tag{16}$$

$$\rho_{s,t} = N + \rho'_s \tag{17}$$

$$\lambda_{s,t} = \lambda'_s + Nr_{s,t} + {\zeta'_s \boldsymbol{\nu}'_s}^T \Xi'_s \boldsymbol{\nu}'_s - (N + {\zeta'_s}) \boldsymbol{\xi}_{s,t}^T \Xi_{s,t}^{-1} \boldsymbol{\xi}_{s,t}, \quad (18)$$

where

$$\boldsymbol{\xi}_{s,t} = \frac{N\boldsymbol{r}_{s,t} + \zeta'_s \Xi'_s \boldsymbol{\nu}'_s}{N + \zeta'_s}.$$
(19)

 $R_{s,t}$, $r_{s,t}$, and $r_{s,t}$ that appear in (14) to (19) are respectively the P-th order autocovariance matrix, the one-sample delayed autocovariance vector of order P, and the variance of the speech samples $\{s_t(n)\}_{1 \le n \le N}$ expected on $q(s_t)$:

$$R_{s,t} = \frac{1}{N} \sum_{n=1}^{N} \langle \boldsymbol{s}_{t}^{P}(n-1) \boldsymbol{s}_{t}^{P}(n-1)^{T} \rangle_{q(\boldsymbol{s}_{t})}$$
(20)

$$\boldsymbol{r}_{s,t} = \frac{1}{N} \sum_{n=1}^{N} \langle s_t(n) \boldsymbol{s}_t^P(n-1) \rangle_{q(\boldsymbol{s}_t)}$$
(21)

$$r_{s,t} = \frac{1}{N} \sum_{n=1}^{N} \langle s_t(n)^2 \rangle_{q(s_t)}.$$
 (22)

As shown later, $q(s_t)$ is a normal distribution with mean μ_t and covariance matrix Υ_t . Hence, the expected product of $s_t(n)$ and $s_t(m)$ that constitute $R_{s,t}$, $r_{s,t}$ and $r_{s,t}$ can be calculated as

$$\langle s_t(m)s_t(n)\rangle_{q(s_t)} = \mu_t(m)\mu_t(n) + v_{m,n},$$
 (23)

where $\mu_t(n)$ is the (N - n + 1)-th component of μ_t and $v_{m,n}$ is the (N - m + 1, N - n + 1)-th component of Υ_t . Recall and note that the elements of \boldsymbol{s}_t are arranged in time-descending order.

The hyperparameters of the noise model are defined in the same way as those of the speech model given by (14) to (19). To take the nonstationarity of noise into account, however, we propose multiplying the hyperparameters of the prior distribution of the noise model parameters by forgetting factor α . Thus, the update equations for the hyperparameters of the noise model are as follows.

$$\zeta_{v,t} = N + \alpha \zeta_{v,t-1} \tag{24}$$

$$\Xi_{v,t} = \frac{NR_{v,t} + \alpha\zeta_{v,t-1} \Xi_{v,t-1}}{N + \alpha\zeta_{v,t-1}}$$
(25)

$$\nu_{v,t} = \Xi_{v,t}^{-1} \boldsymbol{\xi}_{v,t}$$
(26)

$$\rho_{v,t} = N + \alpha \rho_{v,t-1}$$
(27)

$$\lambda_{v,t} = \alpha \lambda_{v,t-1} + Nr_{v,t} + \alpha \zeta_{v,t-1} \boldsymbol{\nu}_{v,t-1}^T \Xi_{v,t-1} \boldsymbol{\nu}_{v,t-1}$$

$$-(N+\alpha\zeta_{v,t-1})\boldsymbol{\xi}_{v,t}^{T}\Xi_{v,t}\boldsymbol{\xi}_{v,t}, \qquad (28)$$

where

$$\boldsymbol{\xi}_{v,t} = \frac{N\boldsymbol{r}_{v,t} + \alpha \zeta_{v,t-1} \Xi_{v,t-1} \boldsymbol{\nu}_{v,t-1}}{N + \alpha \zeta_{v,t-1}}$$
(29)

 $R_{v,t}$, $r_{v,t}$, and $r_{v,t}$ in (24) to (29) are respectively the Q-th order autocovariance matrix, the one-sample delayed autocovariance vector of order Q, and the variance of the noise samples $\{x_t(n) -$ $s_t(n)\}_{1 \le n \le N}$ expected on $q(s_t)$. These are given in the same way as (20) to (22).

To summarize this subsection, once the posterior distribution of the speech samples, $q(s_t)$, is given, the posterior distribution of the parameters is updated by updating its hyperparameters according to (14) to (18) and (24) to (28).

3.3. Optimal variational posterior distribution of speech samples

Next, let us derive the update procedure for $q(\mathbf{s}_t)$ given $q(\Theta_t)$. By substituting (4) and (5) into (11), we finally have

$$q(\boldsymbol{s}_t) = \mathcal{N}\{\boldsymbol{s}_t; \boldsymbol{\mu}_t, \boldsymbol{\Upsilon}_t\}.$$
(30)

The covariance matrix Υ_t and the mean μ_t are defined respectively as follows:

$$\Upsilon_t = \left(\frac{\rho_{s,t}}{\lambda_{s,t}} A_t^T A_t + \frac{1}{\zeta_{s,t}} \Omega_{s,t} + \frac{\rho_{v,t}}{\lambda_{v,t}} B_t^T B_t + \frac{1}{\zeta_{v,t}} \Omega_{v,t}\right)^{-1} (31)$$

$$\boldsymbol{\mu}_{t} = \left(\frac{\rho_{v,t}}{\lambda_{v,t}} B_{t}^{T} B_{t} + \frac{1}{\zeta_{v,t}} \Omega_{v,t}\right) \Upsilon_{t} \boldsymbol{x}_{t}, \qquad (32)$$

where A_t and B_t are N-th order upper triangular Toeplitz matrices whose first rows are respectively

$$\underbrace{[1, -\boldsymbol{\nu}_{s,t}^{T}, \underbrace{0, \cdots, 0}_{N,t}]}_{P+1}, \text{ and } \underbrace{[1, -\boldsymbol{\nu}_{v,t}^{T}, \underbrace{0, \cdots, 0}_{N,t}]}_{N-Q-1}.$$
(33)

 $\Omega_{s,t}$ is defined as follows. Let $\Omega_{s,t}(n)$ denote the *N*-th order matrix whose *P*-th order submatrix beginning from the (N - n + 1, N - n + 1)-th component is equal to $\Xi_{s,t}^{-1}$ and whose components outside the submatrix are all zero. Now, $\Omega_{s,t}$ is defined as

$$\Omega_{s,t} = \sum_{n=1}^{N} \Omega_{s,t}(n-1).$$
(34)

The definition of $\Omega_{v,t}$ is similar to that of $\Omega_{s,t}$. Therefore, the posterior distribution of the speech samples, $q(s_t)$, is updated by updating its mean μ_t and covariance matrix Υ_t according to (32) and (31), respectively, when the posterior distribution of the parameters $q(\Theta_t)$ is given.

It is noteworthy that if we ignore $\Omega_{s,t}$ and $\Omega_{v,t}$, the mean μ_t is the minimum mean square error (MMSE) estimate of s_t under the condition that the parameters Θ_t are deterministic as $a_t = \nu_{s,t}$, $1/\sigma_t^2 = \rho_{s,t}/\lambda_{s,t}$, $b_t = \nu_{v,t}$, and $1/\gamma_t^2 = \rho_{v,t}/\lambda_{v,t}$. Furthermore, the covariance matrix Υ_t corresponds to the associated error covariance matrix. The terms $\Omega_{s,t}/\zeta_{s,t}$ and $\Omega_{v,t}/\zeta_{v,t}$ reflect the degree of uncertainty of the parameters Θ_t .

3.4. Efficient Implementation

Based on the above derivation, the proposed method is summarized as follows.

- (s1) Initialize the hyperparameters of the posterior distribution of the speech and noise models, $\zeta_{s,t}$, $\Xi_{s,t}$, $\boldsymbol{\nu}_{s,t}$, $\rho_{s,t}$, $\lambda_{s,t}$, $\zeta_{v,t}$, $\Xi_{v,t}$, $\boldsymbol{\nu}_{v,t}$, $\rho_{v,t}$, and $\lambda_{v,t}$.
- (s2) Update the mean μ_t and covariance matrix Υ_t of the posterior distribution of the speech samples according to (32) and (31), respectively.
- (s3) Update the hyperparameters according to (14) to (18) and (24) to (28).

(s4) Unless convergence is reached, return to (s2).

The above algorithm can be implemented efficiently in the following two respects.

• We suppose that both

$$\begin{bmatrix} * & \boldsymbol{\xi}_{s,t}^T \\ \boldsymbol{\xi}_{s,t} & \boldsymbol{\Xi}_{s,t} \end{bmatrix} \text{ and } \begin{bmatrix} * & \boldsymbol{\xi}_{v,t}^T \\ \boldsymbol{\xi}_{v,t} & \boldsymbol{\Xi}_{v,t} \end{bmatrix},$$
(35)

where * denotes a certain value, are symmetric Toeplitz matrices. Let us represent the first column of the former matrix by $[\bar{r}_{s,t}(0), \cdots, \bar{r}_{s,t}(P)]^T$. Then, $\boldsymbol{\nu}_{s,t}$ and $\lambda_{s,t}$ of (16) and (18) are calculated by applying the Levinson-Durbin algorithm [7] of order P to $\{\bar{r}_{s,t}(0), \cdots, \bar{r}_{s,t}(P)\}$. $\boldsymbol{\nu}_{v,t}$ and $\lambda_{v,t}$ of (26) and (28) can also be calculated by using the Q-th order Levinson-Durbin algorithm.

We suppose that A^T_t A_t, Ω_{s,t}, B^T_t B_t, and Ω_{v,t} are all Toeplitz matrices. Then, μ_t and Υ_t of (32) and (31) can be calculated in the frequency domain. Moreover, based on the result described in [8], Ω_{s,t} is found to be calculated from the (P-1)-th order regression coefficients for {τ_{s,t}(0), · · · , τ_{s,t}(P - 1)}. Surprisingly and fortunately, the (P-1)-th order regression coefficients have already been obtained in the calculation of ν_{s,t}. The same holds for Ω_{v,t}.

4. EXPERIMENTAL RESULTS

We conducted experiments to evaluate the performance of the proposed method. Japanese utterances of five males and five females were taken from the JNAS database. The sampling frequency was 8 kHz. The speech signals of these utterances were contaminated by eight types of additive noise with a signal to noise ratio (SNR) of 10 dB, which were taken from the AURORA-2 database. The system parameters were set as follows: the frame size N was 256 samples, the frame shift was 128 samples, the window function was a Hanning window, the order of the speech AR model, P, was 12, the order of the noise AR model, Q, was 6, and the forgetting factor α was 0.9.

The proposed method was compared with the minimum statistics based method of [4] and the ML estimation based method. The ML estimation based method is derived by forcing $\Omega_{s,t}$ and $\Omega_{v,t}$ of (31) and (32) to be zero matrices. The noise suppression performance was evaluated by using the segmental SNR (SSNR) and the segmental Itakura-Saito distortion (SISD) measure. The SSNR and SISD were calculated by using median SNRs and ISDs for each frame as in [3] to mitigate the influence of outliers.

Figs. 1 and 2 show the experimental results. It can be seen that the proposed method achieved the best performance in terms of both measures. The reason for the high SISDs of the minimum statistics based method may be that the method is likely to introduce spectral distortion owing to spectral smoothing. Hence, it can be concluded that the proposed method is able to suppress non-stationary noise while reducing spectral distortion.

5. CONCLUSION

In this paper, we have described a Bayesian method for estimating the parameters of speech and noise AR models and speech signals. The variational Bayesian method was used to obtain the approximate posterior distribution of the parameters and the speech signals. An efficient implementation was also developed by assuming the



Fig. 1. Average SSNRs for each noise type. Greater SSNRs indicate better performance.



Fig. 2. Average SISDs for each noise type. Smaller SISDs indicate better performance.

Toeplitz structure for all involved covariance matrices. Experimental results showed that the proposed method performed better than alternative methods in terms of SSNR and SISD.

6. REFERENCES

- Y. Ephraim, "A Bayesian estimation approach for speech enhancement using hidden Markov models," *IEEE Trans. Signal Process.*, vol. 40, no. 4, pp. 725–735, 1992.
- [2] S. Srinivasan, J. Samuelsson, and B. Kleijn, "Codebook-based Bayesian speech enhancement for nonstationary environments," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 15, no. 2, pp. 441–452, 2007.
- [3] S. Gannot, D. Burshtein, and E. Weinstein, "Iterative and sequential Kalman filter-based speech enhancement algorithms," *IEEE Trans. Speech, Audio Process.*, vol. 6, no. 4, pp. 373–385, 1998.
- [4] R. Martin, "Noise power spectral density estimation based on optimal smoothing and minimum statistics," *IEEE Trans. Speech, Audio Process.*, vol. 9, no. 5, pp. 504–512, 2001.
- [5] S. R. Waterhouse, D. MacKay, and A. J. Robinson, "Bayesian methods for mixture of experts," in *Adv. Neural Inform. Process. Syst.* 8, 1995, pp. 351–357.
- [6] H. Attias, J. C. Platt, A. Acero, and L. Deng, "Speech denoising and dereverberation using probabilistic models," in *Adv. Neural Inform. Process. Syst.* 13, 2000, pp. 758–764.
- [7] J. Durbin, "The fitting of time-series models," *Rev. Inst. Int. Stat.*, vol. 28, no. 3, pp. 233–243, 1960.
- [8] W. F. Trench, "An algorithm for the inversion of finite Toeplitz matrices," J. Soc. Indust. Appl. Math., vol. 12, no. 3, pp. 515– 522, 1964.