

# MEL-GENERALIZED CEPSTRAL REGULARIZATION FOR DISCRIMINATIVE NON-NEGATIVE MATRIX FACTORIZATION

*Li Li<sup>1</sup>, Hirokazu Kameoka<sup>2</sup> and Shoji Makino<sup>1</sup>*

<sup>1</sup>University of Tsukuba, Japan

<sup>2</sup>NTT Communication Science Laboratories, NTT Corporation, Japan

## ABSTRACT

Non-negative matrix factorization (NMF) is effective in terms of signal recovery accuracy for single-channel speech enhancement task, while it does not directly lead to an enhancement in feature domain. To overcome this problem, we have previously proposed an extension of NMF which combines a NMF-model fitting criterion and a divergence measuring the NMF model of speech and the mel-generalized cepstral (MGC) representation of a pretrained prototype spectrum. The model has been shown effective both in increasing the signal recovery accuracy and feature domain enhancement. However one drawback is that the regularization term is formulated based on the NMF model of speech which may cause a degradation of the effect of regularization term since the enhanced speech is obtained using Wiener filtering. This paper proposes a novel formulation for MGC regularization and combines it with Discriminative NMF (DNMF) in order to achieve better speech enhancement performance. The experimental results revealed that the proposed method outperformed the perviously proposed model in terms of both the signal recovery accuracy and feature enhancement.

**Index Terms**— Discriminative non-negative matrix factorization, mel-generalized cepstral representation, speech enhancement, single-channel

## 1. INTRODUCTION

Speech enhancement is a technique for recovering the speech signal from an observed noisy speech signal. Since the presence of noise can significantly degrade the quality of speech transmission systems and the performance of applications such as speech recognition and speech conversion, many efforts have been devoted to increasing the performance of speech enhancement over recent decades.

For monaural speech enhancement task, non-negative matrix factorization (NMF) [1, 2] is a powerful approach attracted a lot of attention since it has been proposed. Although in the recent years, deep neural networks based machine learning approaches have showed the incredible

capability for various supervised audio signal process tasks including monaural speech enhancement [3, 4], NMF still remains attractive under unsupervised setting or only a limited training data set available.

Given an observed magnitude or power spectrogram of a mixture signal, NMF aims to approximate it as the sum of speech and noise spectrogram models, which are represented as a non-negative linear sum of the pretrained basis spectra scaled by time-varying amplitudes. The underlying speech components can be separated out using the Wiener filter constructed by the estimated power spectrograms of speech and noise. Although NMF is shown to be effective in terms of signal recovery accuracy, one drawback is that NMF does not directly lead to an enhancement in the feature domain (e.g., cepstral domain) or in terms of perceived quality. To overcome this drawback, we have previously proposed a NMF framework using mel-generalized cepstral regularization (MGCRNMF) [5], which combines an NMF-based model fitting criterion with a divergence measure between the estimated NMF model of speech and the mel-generalized cepstral (MGC) representation [6] of a prototype spectrum in a pretrained codebook. In [5], we have shown the effectiveness of the MGC regularization in terms of increasing both signal recovery accuracy and cepstral domain enhancement. However, MGCRNMF considers the regularization based on NMF model while the enhanced speech is finally separated out by Wiener filtering, which may lead to a reduction in the effect of the regularization term. To address this problem, this paper proposes a novel formulation of mel-generalized cepstral regularization which measures a divergence between a prototype spectrum represented by MGC representation and the estimated enhanced speech spectra obtained by Wiener filter directly.

On the other hand, Weninger recently proposed a basis training approach called discriminative NMF (DNMF) [7], which trains the basis spectra in such a way that the output of the Wiener filter becomes as close to the spectrogram of each of the training examples as possible so that the separated signals become optimal at test time. Since the basis spectra trained by DNMF also takes the Wiener filter into account, it makes us to believe that the basis spectra trained by DNMF

This work was supported by xxx.

are more appropriate for the proposed formulation than those trained by conventional NMF way. Furthermore, DNMF has been proved it can provide greater separation capability than conventional NMF [7, 8], which further motivates us to connect DNMF with mel-generalized cepstral regularization in order to achieve better performance of speech enhancement.

The remaining part of the paper proceeds as follows: we review the model of NMF with MGC regularization for speech enhancement task in sec. 2. In sec. 3, we introduce the new formulation for MGC regularization and derive the update rules of parameters based on majorization-minimization (MM) principle. In the experimental section (sec. 4), we define the data set, investigate the hyperparameters of the model and compare the proposed method with established methods. We conclude this work in sec. 5.

## 2. NMF WITH MEL-GENERALIZED CEPSTRAL REGULARIZATION

### 2.1. NMF for speech enhancement

Given an observed power spectrogram of a noisy speech signal  $\mathbf{Y} = (Y_{\omega,t})_{\Omega \times T} \in \mathbb{R}^{\geq 0, \Omega \times T}$ , where  $\omega$  and  $t$  are frequency and time indices, we consider approximating it by the sum of speech and noise components,  $X_{\omega,t} = X_{\omega,t}^s + X_{\omega,t}^n$ , where  $X_{\omega,t}^s$  and  $X_{\omega,t}^n$  are represented by the non-negative linear combination of  $K_s$  speech basis spectra  $W_{1,\omega}^s, \dots, W_{K_s,\omega}^s$  and  $K_n$  noise basis spectra  $W_{1,\omega}^n, \dots, W_{K_n,\omega}^n$ :

$$X_{\omega,t}^s = \sum_{k=1}^{K_s} W_{k,\omega}^s H_{k,t}^s, \quad X_{\omega,t}^n = \sum_{k=1}^{K_n} W_{k,\omega}^n H_{k,t}^n. \quad (1)$$

In a supervised setting, a concentrated basis matrix  $\mathbf{W} = [\mathbf{W}^s \ \mathbf{W}^n]$  are pretrained using speech and noise training samples respectively.  $\mathbf{H} = [\mathbf{H}^s; \mathbf{H}^n]$  are the variables to be estimated at test time. NMF leads to different optimization problems according to the definition of the measure of the difference between  $\mathbf{Y}$  and  $\mathbf{X} = (X_{\omega,t})_{\Omega \times T}$ . Here we use the generalized Kullback Leibler (KL) divergence

$$\mathcal{D}_{KL}(\mathbf{Y}|\mathbf{X}) = \sum_{\omega,t} \left( Y_{\omega,t} \log \frac{Y_{\omega,t}}{X_{\omega,t}} - Y_{\omega,t} + X_{\omega,t} \right) \quad (2)$$

as a goodness-of-fit criterion. Once  $X_{\omega,t}^s$  and  $X_{\omega,t}^n$  are estimated, the enhanced speech can be separated out using the Wiener filter constructed with the estimated power spectrogram of speech and noise

$$\hat{\mathbf{X}}^s = \frac{\mathbf{W}^s \mathbf{H}^s}{\mathbf{W} \mathbf{H}} \otimes \mathbf{Y}, \quad (3)$$

where  $\dot{\cdot}$  and  $\otimes$  here are element-wise operations.

### 2.2. Mel-generalized cepstral regularization

When estimating the speech spectrogram  $X_{\omega,t}^s$  in a mixture spectrogram  $Y_{\omega,t}$ , we would want to ensure that the features

of  $X_{\omega,t}^s$  in cepstral domain are also enhanced. With this motivation, we have proposed a penalty term in [5] defined as the Itakura-Saito (IS) divergence [9] between  $X_{\omega,t}^s$  and  $S_{\omega,t}(\theta)$

$$\mathcal{J}(\mathbf{H}^s, \boldsymbol{\theta}) = \sum_{\omega,t} \left( \frac{X_{\omega,t}^s}{S_{\omega,t}(\theta)} - \log \frac{X_{\omega,t}^s}{S_{\omega,t}(\theta)} - 1 \right). \quad (4)$$

Here  $S_{\omega,t}(\theta) = \beta_{t,r_t} \mu_{\omega,r_t}$  is a scaled pretrained prototype spectrum chosen from  $I$  prototypes, which are trained using clean speech samples by  $k$ -means. At each iteration, we find the prototype spectrum  $\mu_i$  closest to  $\mathbf{X}_t^s$  in terms of the IS divergence.  $r_t \in \{1, \dots, I\}$  denotes a cluster indicator variable, describing to which of the  $I$  clusters the  $t$ -th speech spectrum is assigned. To eliminate the scaling indeterminacy, we invoke a scaling parameter  $\beta_{t,r_t}$ , which can be obtained as

$$\hat{\beta}_{t,i} = \frac{1}{\Omega} \sum_{\omega} \frac{X_{\omega,t}^s}{\mu_{\omega,i}}. \quad (5)$$

when  $X_{\omega,t}^s$  and  $\mu_{\omega,i}$  are given. Thus, the parameter  $\boldsymbol{\theta} = \{r_t, \beta_{t,r_t}\}$  consists of a set of cluster indicator variables and corresponding scaling parameters. Note the prototype spectra  $\mu_i = [\mu_{1,i}, \dots, \mu_{\Omega,i}]^T$  here are represented by MGC representation, which is a parametric model for spectral envelopes of speech described by  $M+1$  coefficients and two hyperparameters  $\gamma$  and  $\alpha$ :

$$\begin{aligned} \mu_{\omega} &= l_{\gamma}^{-1} \left( \sum_{m=0}^M c(m) \Psi_{\alpha}^m(e^{j\omega}) \right) \\ &= \begin{cases} \left( 1 + \gamma \sum_{m=0}^M c(m) \Psi_{\alpha}^m(e^{j\omega}) \right)^{1/\gamma} & (0 < |\gamma| \leq 1) \\ \exp \sum_{m=0}^M c(m) \Psi_{\alpha}^m(e^{j\omega}) & (\gamma = 0) \end{cases}. \end{aligned} \quad (6)$$

The coefficients  $\mathbf{c} = [c(0), \dots, c(M)]^T$  are called the MGC coefficients (MGCC). The function  $l_{\gamma}^{-1}(\cdot)$  is the inverse of the generalized logarithmic function

$$l_{\gamma}(\omega) = \begin{cases} (\omega^{\gamma} - 1)/\gamma & (0 < |\gamma| \leq 1) \\ \log \omega & (\gamma = 0) \end{cases}, \quad (7)$$

parameterized by  $\gamma$ .  $\Psi_{\alpha}(z)$  is an all-pass function given by

$$\Psi_{\alpha}(z) = \frac{z^{-1} - \alpha}{1 - \alpha z^{-1}}, \quad (8)$$

which can be seen as a frequency warping function parameterized by  $\alpha$ . Here,  $\alpha$  must satisfy  $|\alpha| < 1$ . Note that MGC representation takes the all-pole spectral model and the cepstral representation as special cases when  $(\gamma, \alpha) = (-1, 0)$  and  $(\gamma, \alpha) = (0, 0)$  respectively. When the sampling frequency is 16 kHz, the phase characteristic of the all-pass function becomes a good approximation to the mel scale with  $\alpha = 0.42$  and to the bark scale with  $\alpha = 0.55$  [10].

### 2.3. Objective function of MGCRNMF

MGCRNMF considers an optimization problem of minimizing a combined objective function of (2) and (4)

$$\mathcal{F}(\mathbf{H}, \boldsymbol{\theta}) = \mathcal{D}_{KL}(\mathbf{Y}|\mathbf{X}) + \lambda \mathcal{J}(\mathbf{H}^s, \boldsymbol{\theta}), \quad (9)$$

where  $\lambda \geq 0$  weighs the importance of the mel-generalized cepstral regularization term relative to the NMF cost.

However, since the enhanced speech finally separated out using a Wiener filter, the formulation of MGC regularization based on NMF model of speech  $\mathbf{W}^s \mathbf{H}^s$  may cause a reduction of the effect of the regularization term during the filtering process.

## 3. PROPOSED METHOD

### 3.1. A new formulation of MGC regularization

To address the problem mentioned above, we introduce a new formulation of MGC regularization which measures IS divergence between the pretrained prototypes and the enhanced speech spectra obtained using a Wiener filter directly.

$$\tilde{\mathcal{J}}(\mathbf{H}^s, \boldsymbol{\theta}) = \sum_{\omega, t} \left( \frac{\hat{X}_{\omega, t}^s}{S_{\omega, t}(\boldsymbol{\theta})} - \log \frac{\hat{X}_{\omega, t}^s}{S_{\omega, t}(\boldsymbol{\theta})} - 1 \right), \quad (10)$$

where,  $\hat{X}_{\omega, t}^s$  are the enhanced speech spectra obtained by Wiener filtering (3). With the new regularization term, we can easily obtain the objective function of the proposed method

$$\tilde{\mathcal{F}}(\mathbf{H}, \boldsymbol{\theta}) = \mathcal{D}_{KL}(\mathbf{Y}|\mathbf{X}) + \lambda \tilde{\mathcal{J}}(\mathbf{H}^s, \boldsymbol{\theta}). \quad (11)$$

Similarly,  $\lambda \geq 0$  here is a weight parameter to measure the importance of the regularization term relative to the NMF cost.

### 3.2. Relation to DNMF

Instead of applying NMF to speech and noise training samples respectively to train the basis spectra, Weninger [7] proposed directly using the reconstruction error of the separated signals as an objective function for the basis training

$$\begin{aligned} & \text{minimize } f(\mathbf{W}, \mathbf{H}) = \mathcal{D}_{KL} \left( \mathbf{T} \left| \frac{\mathbf{W}^s \mathbf{H}^s}{\mathbf{W} \mathbf{H}} \otimes \mathbf{Y} \right. \right) \quad (12) \\ & \text{subject to } \forall k, \sum_{\omega} W_{\omega, k} = 1, \end{aligned}$$

where  $\mathbf{T}$  denotes the spectrograms of clean speech training samples. This framework is called discriminative NMF (DNMF) by analogy with the discriminative models for classification or regression. DNMF trains basis spectra which are optimal to construct a Wiener filter instead of basis spectra which represents spectrograms of clean train samples well, which is more appropriate for the proposed formulation of MGC regularization.

### 3.3. Update rules and algorithm

Although minimizing the objective function including the regularization term (10) directly is analytically difficult, we can derive a computationally efficient algorithm to find a locally optimal solution based on majorization-minimization (MM) principle [11, 12].

Suppose  $F(\Theta)$  is an objective function that we wish to minimize with respect to  $\Theta$ . Majorization-minimization principle considers to construct a ‘‘majorizer’’  $F^+(\Theta, \alpha)$  defined as a function satisfying  $F(\Theta) = \min_{\alpha} F^+(\Theta, \alpha)$ , where  $\alpha$  is called an auxiliary parameter. An algorithm that consists of iteratively minimizing  $F^+(\Theta, \alpha)$  with respect to  $\Theta$  and  $\alpha$  is guaranteed to converge to a stationary point of the objective function. It should be noted that this concept is adopted in many existing algorithms [1, 13].

Here, we derive a majorizer for the objective function (11) with respect to  $\mathbf{H}^s$  and  $\mathbf{H}^n$  when target MGC representation  $S_{\omega, t}(\hat{\boldsymbol{\theta}})$  with  $\boldsymbol{\theta}$  fixed to  $\hat{\boldsymbol{\theta}}$ . First,  $\mathcal{D}_{KL}(\mathbf{Y}|\mathbf{X})$  involves a ‘‘log-of-sum’’ form of  $W_{k, \omega} H_{k, t}$ . Since the negative logarithm function is a convex function, we can invoke Jensen’s inequality to construct an upper bound of  $\mathcal{D}_{KL}(\mathbf{Y}|\mathbf{X})$  having a ‘‘sum-of-logs’’ form in the same way as [1]

$$\begin{aligned} \mathcal{D}_{KL}(\mathbf{Y}|\mathbf{X}) &\leq \mathcal{D}_{KL}^+(\mathbf{Y}|\mathbf{X}) \quad (13) \\ \mathcal{D}_{KL}^+(\mathbf{Y}|\mathbf{X}) &\stackrel{c}{=} \sum_{\omega, t} \left( -Y_{\omega, t} \sum_k \zeta_{k, \omega, t} \log \frac{W_{k, \omega} H_{k, t}}{\zeta_{k, \omega, t}} + X_{\omega, t} \right), \end{aligned}$$

where  $\stackrel{c}{=}$  denotes equality up to a constant term and  $\zeta_{k, \omega, t}$  is a positive weight that sums to unity,  $\sum_k \zeta_{k, \omega, t} = 1$ . It can be shown that equality of (13) holds if and only if

$$\zeta_{k, \omega, t} = \frac{W_{k, \omega} H_{k, t}}{\sum_{k'=1}^K W_{k', \omega} H_{k', t}}. \quad (14)$$

Then, we focus on the regularization term

$$\tilde{\mathcal{J}}(\mathbf{H}^s; \hat{\boldsymbol{\theta}}) \stackrel{c}{=} \sum_{\omega, t} \left( \frac{Y_{\omega, t} G_{\omega, t}^s}{S_{\omega, t}(\hat{\boldsymbol{\theta}}) G_{\omega, t}} - \log G_{\omega, t}^s + \log G_{\omega, t} \right), \quad (15)$$

where  $G_{\omega, t}^s = \sum_{k=1}^{K_s} W_{k, \omega}^s H_{k, t}^s$  and  $G_{\omega, t} = \sum_{k=1}^K W_{k, \omega} H_{k, t}$ . To construct an upper bound for the first term of (15), we can invoke the Lemma 1 introduced in [8]

$$\frac{G_{\omega, t}^s}{G_{\omega, t}} \leq \frac{\tau_{\omega, t} G_{\omega, t}^{s, 2}}{2} + \frac{1}{2\tau_{\omega, t} G_{\omega, t}^2}. \quad (16)$$

The equality of (16) holds if and only if

$$\tau_{\omega, t} = \frac{1}{G_{\omega, t}^s G_{\omega, t}}. \quad (17)$$

Since a quadratic function is convex, we can apply Jensen’s inequality to  $G_{\omega, t}^{s, 2}$ , which yields

$$G_{\omega, t}^{s, 2} \leq \sum_{k=1}^{K_s} \frac{W_{\omega, k}^s H_{k, t}^{s, 2}}{\alpha_{k, \omega, t}}, \quad (18)$$

where  $\alpha_{k,\omega,t} > 0$  is also a positive number that sums to unity, i.e.,  $\sum_k \alpha_{k,\omega,t} = 1$ . The equality of (18) holds if and only if

$$\alpha_{k,\omega,t} = \frac{W_{k,\omega}^s H_{k,t}^s}{\sum_{k'=1}^{K_s} W_{k',\omega}^s H_{k',t}^s}. \quad (19)$$

We can use the fact that  $1/x^2$  is convex in the first quadrant and use Jensen's inequality to obtain a majorizer:

$$\frac{1}{G_{\omega,t}^2} \leq \sum_{k=1}^K \frac{\xi_{k,\omega,t}^3}{W_{k,\omega}^2 H_{k,t}^2}, \quad (20)$$

where  $\xi_{k,\omega,t} > 0$  and  $\sum_k \xi_{k,\omega,t} = 1$ . It can be proved that the equality of this inequality holds if and only if

$$\xi_{k,\omega,t} = \frac{W_{k,\omega} H_{k,t}}{\sum_{k'=1}^K W_{k',\omega} H_{k',t}}. \quad (21)$$

By substituting (18) and (20) into (16), the majorizer for the first term can be written as

$$\frac{G_{\omega,t}^s}{G_{\omega,t}} \leq \sum_{k=1}^{K_s} \frac{\tau_{\omega,t} W_{k,\omega}^s H_{k,t}^s}{2\alpha_{k,\omega,t}} + \sum_{k=1}^K \frac{\xi_{k,\omega,t}^3}{2\tau_{\omega,t} W_{k,\omega}^2 H_{k,t}^2}. \quad (22)$$

As regards the second term, Jensen's inequality can be invoked again since  $-\log G_{\omega,t}^s$  is convex in  $G_{\omega,t}^s$ ,

$$-\log G_{\omega,t}^s \leq -\sum_{k=1}^{K_s} \gamma_{k,\omega,t} \log \frac{W_{k,\omega}^s H_{k,t}^s}{\gamma_{k,\omega,t}}, \quad (23)$$

where  $\gamma_{k,\omega,t}$  is a positive weight that sums to unity. The equality of (23) holds if and only if

$$\gamma_{k,\omega,t} = \frac{W_{k,\omega}^s H_{k,t}^s}{\sum_{k'=1}^{K_s} W_{k',\omega}^s H_{k',t}^s}. \quad (24)$$

The third term  $\log G_{\omega,t}$  is concave in  $G_{\omega,t}$ . Hence, we can use the fact that a tangent line to the graph of a differentiable concave function lies entirely above the graph:

$$\log G_{\omega,t} \leq \sum_{k=1}^K \frac{W_{k,\omega} H_{k,t}}{\eta_{\omega,t}} + \log \eta_{\omega,t} - 1, \quad (25)$$

where  $\eta_{\omega,t}$  is an arbitrary positive number. The equality of this inequality holds if and only if

$$\eta_{\omega,t} = G_{\omega,t}. \quad (26)$$

From (22), (23) and (25), we can construct a majorizer for the regularization term as

$$\begin{aligned} \tilde{\mathcal{J}}(\mathbf{H}^s; \hat{\boldsymbol{\theta}}) &\leq \tilde{\mathcal{J}}^+(\mathbf{H}^s, \boldsymbol{\Gamma}; \hat{\boldsymbol{\theta}}) \\ &= \sum_{k,\omega,t} \frac{\tau_{\omega,t} Y_{\omega,t} W_{k,\omega}^s H_{k,t}^s}{2\alpha_{k,\omega,t} S_{\omega,t}(\hat{\boldsymbol{\theta}})} + \sum_{k,\omega,t} \frac{\xi_{k,\omega,t}^3 Y_{\omega,t}}{2\tau_{\omega,t} S_{\omega,t}(\hat{\boldsymbol{\theta}}) W_{k,\omega}^2 H_{k,t}^2} \end{aligned}$$

$$- \sum_{k,\omega,t} \gamma_{k,\omega,t} \log \frac{W_{k,\omega}^s H_{k,t}^s}{\gamma_{k,\omega,t}} + \sum_{k,\omega,t} \frac{W_{k,\omega} H_{k,t}}{\eta_{\omega,t}} + d,$$

where  $\boldsymbol{\Gamma} = \{\zeta_{k,\omega,t}, \tau_{\omega,t}, \gamma_{k,\omega,t}, \eta_{\omega,t}, \alpha_{k,\omega,t}, \xi_{k,\omega,t}\}$  denotes a set of all the auxiliary variables and  $d$  denotes a constant term. The upper bound for the objective function can be easily obtained by combining the majorizers for each term as

$$\tilde{\mathcal{F}}^+(\mathbf{H}, \boldsymbol{\Gamma}; \hat{\boldsymbol{\theta}}) = \mathcal{D}_{KL}^+(\mathbf{Y}|\mathbf{X}) + \lambda \tilde{\mathcal{J}}^+(\mathbf{H}^s, \boldsymbol{\Gamma}; \hat{\boldsymbol{\theta}}). \quad (27)$$

The update rules for  $H_{k,t}$  can be obtained by setting at zeros the partial derivatives of the derived majorizer with respect to  $H_{k,t}^s$  and  $H_{k,t}^n$ . Thus, the update rules can be obtained as the positive solution of the following quartic and cubic equations:

$$\begin{aligned} \sum_{\omega} \frac{\tau_{\omega,t} Y_{\omega,t}}{2\alpha_{k,\omega,t} S_{\omega,t}(\hat{\boldsymbol{\theta}})} W_{k,\omega}^s H_{k,t}^s + \sum_{\omega} \frac{W_{k,\omega}^s H_{k,t}^s}{\eta_{\omega,t}} \\ - \sum_{\omega} \gamma_{k,\omega,t} H_{k,t}^s - \sum_{\omega} \frac{Y_{\omega,t} \xi_{k,\omega,t}^3}{2\tau_{\omega,t} S_{\omega,t}(\hat{\boldsymbol{\theta}}) W_{k,\omega}^2} = 0, \quad (28) \end{aligned}$$

$$\sum_{\omega} \frac{W_{k,\omega}^n H_{k,t}^n}{\eta_{\omega,t}} - \sum_{\omega} \frac{Y_{\omega,t} \xi_{k,\omega,t}^3}{2\tau_{\omega,t} S_{\omega,t}(\hat{\boldsymbol{\theta}}) W_{k,\omega}^2} = 0. \quad (29)$$

It is noteworthy that all the parameters can be updated in parallel using these update rules, which means this algorithm is well suited to parallel implementations. Furthermore, since each of the update rules consists of a negative 0th-order term and a negative 2nd-order term, it turns out that there is only one positive solution, implying that there is no need to solve a solution selection problem.

Algorithm. 1 shows the whole procedure.

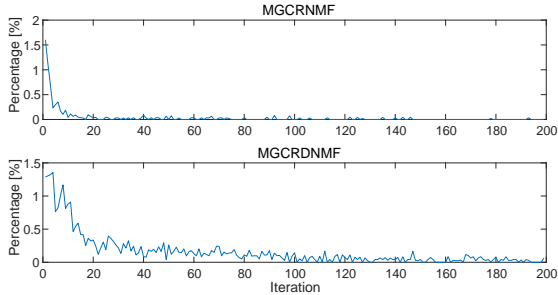
---

**Algorithm 1** Algorithm presented in subsec. 3.3

---

**Require:** pretrained speech basis  $\mathbf{W}$  and  $I$  MGC prototypes

- $\boldsymbol{\mu}$ , parameters  $\lambda$  and  $MaxIter$
  - 1: random initialize  $\mathbf{H}^s$  and  $\mathbf{H}^n$
  - 2: **for**  $iter = 1$  to  $MaxIter$  **do**
  - 3:   **if**  $iter \leq 50$  **then**
  - 4:     update  $\mathbf{H}^s$  and  $\mathbf{H}^n$  using SSNMF
  - 5:   **else**
  - 6:     calculate the enhanced speech  $\hat{\mathbf{X}}^s$  using (3)
  - 7:     **for** Frame  $t = 1$  to  $T$  **do**
  - 8:       compute  $\hat{\beta}_{t,i}$  using (5)
  - 9:        $\hat{r}_t = \arg \min_{r_t} \tilde{\mathcal{J}}(\mathbf{H}^s, \boldsymbol{\theta})$
  - 10:        $S_{\omega,t}(\hat{\boldsymbol{\theta}}) = \hat{\beta}_{t,\hat{r}_t} \boldsymbol{\mu}_{\omega,\hat{r}_t}$
  - 11:     **end for**
  - 12:     update auxiliary variables  $\boldsymbol{\Gamma}$  using (14), (17), (19), (21), (24) and (26)
  - 13:     update  $\mathbf{H}^s, \mathbf{H}^n$  by solving the equations (28) and (29)
  - 14:     **end if**
  - 15:   **end for**
  - 16: **end for**
-



**Fig. 1.** Percentage of the number of frames which shift the prototype during the updating of MGCRNMF (upper) and MGCRDNMF (bottom). The average number of 10 samples randomly selected under 5 noise types are shown in the figure.

## 4. EXPERIMENTS

### 4.1. Experimental conditions

To evaluate the effect of the proposed method for speech enhancement task, we tested supervised NMF (SNMF) [2], Discriminative NMF (DNMF) [8], NMF with mel-generalized cepstral regularization (MGCRNMF) [5] and the proposed method (MGCRDNMF) using the speech data excerpted from the ATR503 database [14] and 5 types of measured noise, respectively BusTerminal-5dB, Square-5dB, BowlingAlley-5dB, SubwayStation0dB and DepartmentStore0dB, excerpted from the ATR ambient noise sound database.

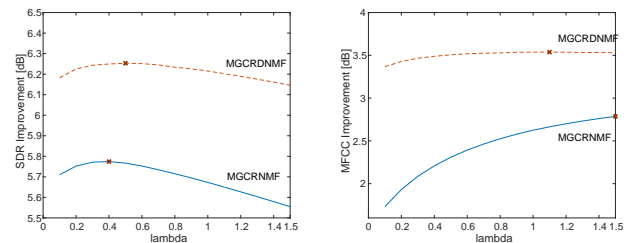
The test data were created by adding noise signals to clean speech signals with the signal-to-noise ratios (SNRs) of -5, 0 dB. All the audio signals were monaural and sampled at 16KHz. The STFT was computed using the Hanning window with 32ms long and 16ms overlap.

In the training phase, 200 utterances spoken by 2 male and 2 female speakers were used to train 40 speech basis spectra. For noise we used the same number of basis spectra. We run 200 iterations for SNMF basis training and 25 iterations for DNMF training with running 100 iterations NMF for initialization. The same training set was also used for  $k$ -means training. The cluster number was set at 1000. We used 20 order MGCCs with hyperparameter  $(\gamma, \alpha) = (-1, 0.42)$  since it has been shown in [5] that this hyperparameters setting can achieve relatively high performance stably under all conditions.

In the test phase, 50 utterances selected randomly from ATR503 database were used as test data set. All methods were run 200 iterations while MGCRNMF and the proposed method run 50 iterations NMF as an initialization. During the update, each frame should converge to one cluster, which means we do not need to figure out the closest prototype spectrum at every iteration. Fig. 1 shows the percentage of how many frames shift the cluster to another one among all the frames at every iteration. The average frame number of the selected test data was about 324 so that 1% of the whole

**Table 1.** A comparison of runtime [sec] between updating indicator variables at each iteration and updating them at the first iteration of MGCRNMF and the proposed method. The length of the test data was 5 seconds.

|          | w/ update | w/o update |
|----------|-----------|------------|
| MGCRNMF  | 135.5534  | 1.9471     |
| MGCRDNMF | 259.9060  | 126.9996   |



**Fig. 2.** Average SDR improvement [dB] (left) and MFCC improvement [dB] (right) achieved by MGCRNMF and the proposed method using 50 test samples with  $\lambda$  from 0.1 to 1.5 at 0.1 intervals. The points draw the maximum of the curves.

frames was about 3 frames, which means that there was only about 3 frames have not converged to the clusters after updating once. With this result, in the following experiments, we set the update iteration number of the indicator variables  $\{r_t\}$  at 1. Tab. 1 shows a comparison of the runtime between updating indicator variable every iteration and only updating it once using MGCRNMF and the proposed method. The programs were run in the MATLAB 2015b with Inter(R) Xeon E3-1505M V5 CPU @2.80GHz 64bit and 16.0 GB memory. The results show a significant improvement in runtime between w/o update and w/ update with a very small degradation of the performance about 0.05 dB. It is worth noting that without updating every iteration, MGCRNMF realized a real time computation.

We investigated the weight parameter  $\lambda$  during 0.1 to 1.5 at 0.1 intervals and the results are shown in Fig. 2. According to the Fig. 2, we set  $\lambda = 0.4$  for MGCRNMF and  $\lambda = 0.5$  for the proposed method this time in order to achieve the highest signal-to-distortion ratio improvement.

### 4.2. Objective evaluation

We used Signal-to-distortion ratios (SDRs), signal-to-interference ratios (SIRs) [15] and MFCC distance for the evaluation. Given two  $D$ -dimension MFCC sequences  $x[d]$  and  $y[d]$  calculated from  $N$  frequency bins, the MFCC distance is defined as follow:

$$Dist = \frac{20D}{N \ln 10} \sqrt{2 \sum_d^D (x[d] - y[d])^2}. \quad (30)$$

Tab. 2 shows the results of average SDR, SIR and MFCC distance improvement [dB] obtained using SNMF, DNMF,

**Table 2.** From top to bottom, there are respectively average SDR, SIR, MFCC Improvement [dB] evaluated under 5 noise conditions. The highest score of each term is shown in bold font type.

| Noise Type      | SNMF  | MGCRNMF | DNMF  | Proposed     |
|-----------------|-------|---------|-------|--------------|
| BusTerminal     | 10.71 | 11.22   | 11.57 | <b>11.94</b> |
| Square          | 6.19  | 6.45    | 6.76  | <b>6.88</b>  |
| BowlingAlley    | 3.37  | 3.40    | 4.01  | <b>4.30</b>  |
| SubwayStation   | 3.90  | 3.78    | 4.22  | <b>4.46</b>  |
| DepartmentStore | 4.73  | 4.95    | 4.76  | <b>4.97</b>  |

---

| Noise Type      | SNMF  | MGCRNMF | DNMF  | Proposed     |
|-----------------|-------|---------|-------|--------------|
| BusTerminal     | 13.85 | 15.07   | 17.43 | <b>18.32</b> |
| Square          | 8.61  | 9.29    | 10.32 | <b>11.14</b> |
| BowlingAlley    | 5.27  | 5.77    | 6.72  | <b>7.40</b>  |
| SubwayStation   | 6.74  | 7.71    | 8.27  | <b>8.87</b>  |
| DepartmentStore | 7.09  | 7.89    | 8.92  | <b>10.42</b> |

---

| Noise Type      | SNMF | MGCRNMF | DNMF | Proposed    |
|-----------------|------|---------|------|-------------|
| BusTerminal     | 1.79 | 2.27    | 3.24 | <b>3.66</b> |
| Square          | 1.87 | 2.10    | 1.84 | <b>3.05</b> |
| BowlingAlley    | 1.32 | 1.97    | 2.81 | <b>3.13</b> |
| SubwayStation   | 1.81 | 2.52    | 2.85 | <b>3.51</b> |
| DepartmentStore | 1.79 | 2.27    | 2.36 | <b>3.22</b> |

MGCRNMF and the proposed method under 5 noise conditions. The proposed method outperformed the other methods under all the conditions in all the evaluation criteria.

## 5. CONCLUSION

This paper proposed a novel formulation for mel-generalized cepstral regularization to enhance speech in spectral and cepstral domain, which takes the Wiener filtering process into account. We combined the proposed regularization with Discriminative NMF basis training approach and derived a computationally efficient algorithm based on majorization-minimization principle. The experimental results showed that the proposed algorithm outperformed SNMF, DNMF and previously proposed MGCRNMF in SDR, SIR and MFCC distance improvements, which showed the effectiveness of the proposed method.

## 6. REFERENCES

- [1] D. D. Lee and H. S. Seung, "Algorithms for nonnegative matrix factorization," in *Adv. NIPS*, pp. 556–562, 2000.
- [2] P. Smaragdis, B. Raj and M. Shashanka, "Supervised and semi-supervised separation of sounds from single-channel mixtures," in *Proc. ICA 2007*, pp. 414–421, 2007.
- [3] Y. Xu, J. Du, L. R. Dai and C. H. Lee, "A regression approach to speech enhancement based on deep neural networks," *IEEE/ACM Transactions on Audio, Speech and Language Processing*, Vol. 23, No. 1, pp. 7–19, 2015.
- [4] J. R. Hershey, C. Zhuo, J. L. Roux and S. Watanabe, "Deep clustering: Discriminative embeddings for segmentation and separation," in *Proc. ICASSP*, pp. 31–35, 2015.
- [5] L. Li, H. Kameoka, T. Toda and S. Makino, "Speech enhancement using non-negative spectrogram models with mel-generalized cepstral regularization," in *Proc. INTERSPEECH*, 2017.
- [6] K. Tokuda, T. Kobayashi, T. Masuko and S. Imai, "Mel-generalized cepstral analysis—a unified approach to speech spectral estimation," in *ICSLP*, Vol. 94, pp. 18–22, 1994.
- [7] F. Weninger, J. L. Roux, J. R. Hershey, and S. Watanabe, "Discriminative NMF and its application to single-channel source separation," in *Proc. INTERSPEECH*, pp. 865–869, 2014.
- [8] L. Li, H. Kameoka, and S. Makino, "Discriminative non-negative matrix factorization with majorization-minimization," in *Proc. HSCMA*, pp. 141–145, 2017.
- [9] F. Itakura and S. Saito, "Analysis synthesis telephony based on the maximum likelihood method," in *Proc. the 6th International Congress on Acoustics*, pp. 17–20, 1968.
- [10] T. Masuko, "HMM-based speech synthesis and its applications," Institute of Technology, 2002.
- [11] J.D. Leeuw, and W.J. Heiser, "Convergence of correction matrix algorithms for multidimensional scaling," in *Geometric representations of relational data*, pp. 735–752, 1977.
- [12] D.R. Hunter, and K. Lange, "A tutorial on MM algorithms," *The American Statistician*, 58 (1), pp. 30–37, 2004.
- [13] M. Nakano, H. Kameoka, J. Le Roux, Y. Kitano, N. Ono, and S. Sagayama, "Convergence-guaranteed multiplicative algorithms for non-negative matrix factorization with beta-divergence," in *Proc. MLSP*, pp. 283–288, 2010.
- [14] A. Kurematsu, K. Takeda, Y. Sagisaka, S. Katagiri, H. Kuwabara, and K. Shikano, "ATR Japanese speech database as a tool of speech recognition and synthesis," *Speech Communication*, vol. 9, pp. 357–363, 1990.
- [15] E. Vincent, R. Gribonval, and C. Fevotte, "Performance measurement in blind audio source separation," *IEEE Transactions on Audio, Speech, and Language Processing*, Vol. 14, No. 4, pp. 1462–1469, 2016.