

A Statistically Principled and Computationally Efficient Approach to Speech Enhancement using Variational Autoencoders: Supporting Document

Manuel Pariente, Antoine Deleforge, Emmanuel Vincent

► To cite this version:

Manuel Pariente, Antoine Deleforge, Emmanuel Vincent. A Statistically Principled and Computationally Efficient Approach to Speech Enhancement using Variational Autoencoders: Supporting Document. [Research Report] RR-9268, INRIA. 2019, pp.1-8. hal-02089062

HAL Id: hal-02089062

<https://hal.inria.fr/hal-02089062>

Submitted on 8 Apr 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



A Statistically Principled and Computationally Efficient Approach to Speech Enhancement using Variational Autoencoders : Supporting document

Manuel Pariente, Antoine Deleforge, Emmanuel Vincent

**RESEARCH
REPORT**

N° 9268

April 2019

Project-Team MULTISPEECH



A Statistically Principled and Computationally Efficient Approach to Speech Enhancement using Variational Autoencoders : Supporting document

Manuel Pariente, Antoine Deleforge, Emmanuel Vincent

Project-Team MULTISPEECH

Research Report n° 9268 — April 2019 — 8 pages

Abstract: Recent studies have explored the use of deep generative models of speech spectra based of variational autoencoders (VAEs), combined with unsupervised noise models, to perform speech enhancement. These studies developed iterative algorithms involving either Gibbs sampling or gradient descent at each step, making them computationally expensive. This paper proposes a variational inference method to iteratively estimate the power spectrogram of the clean speech. Our main contribution is the analytical derivation of the variational steps in which the encoder of the pre-learned VAE can be used to estimate the variational approximation of the true posterior distribution, using the very same assumption made to train VAEs. Experiments show that the proposed method produces results on par with the aforementioned iterative methods using sampling, while decreasing the computational cost by a factor 36 to reach a given performance.

Key-words: Speech enhancement, variational autoencoders, variational Bayes, non-negative matrix factorization

**RESEARCH CENTRE
NANCY – GRAND EST**

615 rue du Jardin Botanique
CS20101
54603 Villers-lès-Nancy Cedex

1 Model

This document provides the details of the analytical derivation of the algorithm presented in [1]. We first remind the statistical model used in [1].

We use f to denote the frequency index and t to denote the time frame index. Independently for $(f, t) \in \{0, \dots, F-1\} \times \{0, \dots, N-1\}$, the single channel observation of the mixture is modeled by

$$x_{ft} = s_{ft} + n_{ft} + \epsilon_{ft}, \quad (1)$$

where x_{ft} , s_{ft} and n_{ft} are the short term Fourier transform (STFT) coefficients of the mixture, the speech source and the noise source respectively, and $\epsilon_{ft} \sim \mathcal{N}_c(0, \sigma_\epsilon^2)$ is introduced to prevent $p(x_{ft}|s_{ft}, b_{ft})$ from being a Dirac, we have

$$x_{ft}|s_{ft}, b_{ft} \sim \mathcal{N}_c(s_{ft} + b_{ft}, \sigma_\epsilon^2), \quad (2)$$

where \mathcal{N}_c denotes the univariate proper complex Gaussian defined in Section 3. σ_ϵ^2 will be set to 0, once the variational updates are obtained.

The noise STFT coefficients are modelled using Non-Negative Matrix Factorization (NMF) [2] :

$$n_{ft} \sim \mathcal{N}_c(0, (\mathbf{W}\mathbf{H})_{ft}). \quad (3)$$

with $\mathbf{W} \in \mathbb{R}_+^{F \times K}$, $\mathbf{H} \in \mathbb{R}_+^{K \times N}$, K being the rank of the NMF model.

We define $\mathbf{s}_t = [s_{1t}, \dots, s_{Ft}]^T$, $\mathbf{s} = [\mathbf{s}_1, \dots, \mathbf{s}_N]$, $\mathbf{n} = [n_{1t}, \dots, n_{Ft}]^T$, $\mathbf{n}_t = [\mathbf{n}_1, \dots, \mathbf{n}_N]$, $\mathbf{x}_t = [x_{1t}, \dots, x_{Ft}]^T$, $\mathbf{x} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ and $\mathbf{H}_t = [H_{1t}, \dots, H_{Kt}]^T$, where $(\cdot)^T$ denotes the transpose operator.

The speech STFT coefficients are modeled using a variational autoencoder (VAE) [3], we have

$$\mathbf{z}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{1}), \quad (4)$$

$$s_{ft}|\mathbf{z}_t; \theta \sim \mathcal{N}_c(0, \sigma_f^2(\mathbf{z}_t)), \quad (5)$$

$$z_{l,t}|\mathbf{s}_t; \phi \sim \mathcal{N}(\tilde{\mu}_l(|\mathbf{s}_t|^2), \tilde{\sigma}_l^2(|\mathbf{s}_t|^2)), \quad (6)$$

where $\mathbf{z}_t = [z_{1t}, \dots, z_{Lt}]^T$ is the latent variable and $L < F$.

We remind that σ_f^2 ; $\tilde{\mu}_l$ and $\tilde{\sigma}_l^2$, are non-linear functions implemented by deep neural networks (DNNs) with parameters θ and ϕ respectively. They are learned by maximizing the marginal log-likelihood $\log p_\theta(\mathbf{s})$. All the STFT frames are considered to be independent, we have

$$\log p_\theta(\mathbf{s}) = \sum_t \log p_\theta(\mathbf{s}_t). \quad (7)$$

We can then maximize the marginal likelihoods of individual STFT frames, we write

$$\log p_\theta(\mathbf{s}_t) = \mathcal{D}_{KL}(q_\phi(\mathbf{z}_t|\mathbf{s}_t)||p_\theta(\mathbf{z}_t|\mathbf{s}_t)) + \mathcal{L}(\theta, \phi; \mathbf{s}_t), \quad (8)$$

$$\mathcal{L}(\theta, \phi; \mathbf{s}_t) = \mathbb{E}_{q_\phi(\mathbf{z}_t|\mathbf{s}_t)}[\log p_\theta(\mathbf{s}_t|\mathbf{z}_t)] - \mathcal{D}_{KL}(q_\phi(\mathbf{z}_t|\mathbf{s}_t)||p(\mathbf{z}_t)), \quad (9)$$

where $\mathcal{D}_{KL}(q||p) = -\mathbb{E}_q[\log(p/q)]$ denotes the Kullback-Leibler (KL) divergence. Based on (4), (5) and (6), we rewrite $\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{s}_t)$:

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}) \stackrel{c}{=} \sum_f \mathbb{E}_{q_\phi(\mathbf{z}_t | \mathbf{s}_{ft})} [d_{IS}(|\mathbf{s}_{ft}|^2, \sigma_f^2(\mathbf{z}_t))] + \frac{1}{2} \sum_l \log(\tilde{\sigma}_l^2(|\mathbf{s}_t|^2) - \tilde{\mu}_l^2(|\mathbf{s}_t|^2) - \tilde{\sigma}_l^2(|\mathbf{s}_t|^2)),$$

where $d_{IS}(x; y) = x/y - \log(x/y) - 1$ denotes the Itakura-Saito divergence, and $\stackrel{c}{=}$ denotes equality up to a constant.

2 Inference

Given an observation $\mathbf{X} = \{x_{ft}\}_{(ft) \in (B)}$, our goal is now to maximize the log-likelihood of \mathbf{X} given the mixture model (1), the generative model of speech ((4) and (5)) and the noise's NMF model (3).

We consider $\mathbf{y}_t = \{\mathbf{s}_t, \mathbf{n}_t, \mathbf{z}_t\}$ to be the set of latent variables, $\boldsymbol{\Theta}_t = \{\mathbf{W}, \mathbf{H}_t\}$ the parameters of the model. We introduce a variational distribution $r(\mathbf{y}_t)$ and write the following decomposition:

$$\log p(\mathbf{x}_t; \boldsymbol{\Theta}_t) - \mathcal{D}_{KL}(r(\mathbf{y}_t) || p(\mathbf{y}_t | \mathbf{x}_t; \boldsymbol{\Theta}_t)) = \mathcal{L}(r, \boldsymbol{\Theta}_t), \quad (10)$$

where $\mathcal{L}(r, \boldsymbol{\Theta}_t)$ follows

$$\mathcal{L}(r, \boldsymbol{\Theta}_t) = \mathbb{E}_{r(\mathbf{y}_t; \boldsymbol{\Theta}_t)} \left[\log \frac{p(\mathbf{x}_t, \mathbf{y}_t; \boldsymbol{\Theta}_t)}{r(\mathbf{y}_t; \boldsymbol{\Theta}_t)} \right]. \quad (11)$$

We suppose that the variation distribution $r(\mathbf{y}_t)$ factorizes as

$$r(\mathbf{s}_t, \mathbf{n}_t, \mathbf{z}_t) = r(\mathbf{s}_t, \mathbf{n}_t) r(\mathbf{z}_t) = \prod_f r(s_{ft}, n_{ft}) \prod_l r(z_{lt}). \quad (12)$$

Given the independence of s_t and n_t , and n_t and z_t , we can write the complete data likelihood as

$$p(\mathbf{x}_t, \mathbf{y}_t; \boldsymbol{\Theta}_t) = p(\mathbf{x}_t | \mathbf{s}_t, \mathbf{n}_t) p(\mathbf{s}_t | \mathbf{z}_t) p(\mathbf{z}_t) p(\mathbf{n}_t; \boldsymbol{\Theta}_t). \quad (13)$$

We can then iteratively maximize $\mathcal{L}(r, \boldsymbol{\Theta}_t)$ with respect to the factorized distributions $r(\mathbf{s}_t, \mathbf{n}_t)$ and $r(\mathbf{z}_t)$, and the NMF parameters $\boldsymbol{\Theta}_t = \{\mathbf{W}, \mathbf{H}_t\}$. As given by the equation (10.9) in [4], the variational distributions $r(s_{ft}, n_{ft})$ and $r(z_{lt})$ can be updated using

$$\log r(s_{ft}, n_{ft}) \stackrel{c}{=} \mathbb{E}_{r(\mathbf{z}_t)} [\log p(x_{ft}, s_{ft}, n_{ft}, \mathbf{z}_t; \boldsymbol{\Theta}_t)], \quad (14)$$

$$\log r(z_{lt}) \stackrel{c}{=} \mathbb{E}_{r(\mathbf{s}_t, \mathbf{n}_t)} [\log p(\mathbf{x}_t, \mathbf{s}_t, \mathbf{n}_t, z_{lt}; \boldsymbol{\Theta}_t)]. \quad (15)$$

\mathbf{W} and \mathbf{H} can be updated by maximizing the following

$$\begin{aligned} \mathcal{Q}(\boldsymbol{\Theta}, \boldsymbol{\Theta}^{old}) &\stackrel{c}{=} \mathbb{E}_{r(\mathbf{s}, \mathbf{b}, \mathbf{z}; \boldsymbol{\Theta}^{old})} [\log p(\mathbf{x}, \mathbf{s}, \mathbf{b}, \mathbf{z}; \boldsymbol{\Theta})] \\ &\stackrel{c}{=} \mathbb{E}_{r(\mathbf{s}, \mathbf{b}; \boldsymbol{\Theta}^{old})} [\log p(\mathbf{b}; \boldsymbol{\Theta})], \end{aligned} \quad (16)$$

where $\boldsymbol{\Theta} = \{\mathbf{W}, \mathbf{H}\}$.

2.1 E-(s,n) step

We define $\sigma_{n,ft}^2 = (\mathbf{WH})_{ft}$ to make the notation less cluttered. Removing the terms independent from s_{ft} and n_{ft} in (14), we have

$$\begin{aligned} \log r(s_{ft}, n_{ft}) &\stackrel{c}{=} \log p(x_{ft}|s_{ft}, n_{ft}) + \mathbb{E}_{r(\mathbf{z}_t)}[\log p(s_{ft}|\mathbf{z}_t)] + \log p(n_{ft}; \sigma_{n,ft}^2) \\ &\stackrel{c}{=} \log p(x_{ft}|s_{ft}, n_{ft}) - \mathbb{E}_{r(\mathbf{z}_t)} \left[\frac{1}{\sigma_f^2(\mathbf{z}_t)} \right] |s_{ft}|^2 + \log p(n_{ft}; \sigma_{n,ft}^2). \end{aligned} \quad (17)$$

We can define γ_{ft}^2 as

$$\frac{1}{\gamma_{ft}^2} = \mathbb{E}_{r(\mathbf{z}_t)} \left[\frac{1}{\sigma_f^2(\mathbf{z}_t)} \right] \approx \sum_{d=1}^D \left[1/\sigma_f^2(\mathbf{z}_t^{(d)}) \right], \quad (18)$$

where $\{\mathbf{z}_t^{(d)}\}_{d=1,\dots,D}$ are randomly drawn from $r(\mathbf{z}_t)$.

We recognize $r(s_{ft}, n_{ft})$ to be a product of Gaussian, up to the normalization constant. We have

$$r(s_{ft}, n_{ft}) \sim \mathcal{N}_c(x_{ft}; s_{ft} + n_{ft}, \sigma_\epsilon^2) \mathcal{N}_c(n_{ft}; 0, \sigma_{n,ft}^2) \mathcal{N}_c(s_{ft}; 0, \gamma_{ft}^2) \quad (19)$$

$$\sim \mathcal{N}_c(x_{ft}; s_{ft} + n_{ft}, \sigma_\epsilon^2) \mathcal{N}_c([s_{ft}, n_{ft}]; \mathbf{0}, \mathbf{\Psi}), \quad (20)$$

with $\mathbf{\Psi}$ defined as

$$\mathbf{\Psi}_{ft} = \begin{bmatrix} \gamma_{ft}^2 & 0 \\ 0 & \sigma_{n,ft}^2 \end{bmatrix}. \quad (21)$$

Finally, with $\mathbf{S}_{ft} = [s_{ft}, n_{ft}]^T$ and $(\cdot)^H$ denoting the Hermitian transpose, we can rewrite $r(s_{ft}, n_{ft})$ using Bayes' theorem, as in [5]

$$r(\mathbf{S}_{ft}) = \frac{1}{|\pi \mathbf{\Sigma}|} \exp \left(-(\mathbf{S}_{ft} - \boldsymbol{\mu}_{ft})^H \mathbf{\Sigma}_{ft}^{-1} (\mathbf{S}_{ft} - \boldsymbol{\mu}_{ft}) \right), \quad (22)$$

where, using $\mathbf{M} = [1, 1]$, $\boldsymbol{\mu}_{ft}$ and $\mathbf{\Sigma}_{ft}$ are defined as

$$\mathbf{\Sigma}_{ft} = \left(\mathbf{\Psi}_{ft}^{-1} + \frac{\mathbf{M}^T \mathbf{M}}{\sigma_\epsilon^2} \right)^{-1} \quad (23)$$

$$\boldsymbol{\mu}_{ft} = \frac{\mathbf{\Sigma}_{ft} \mathbf{R}^T}{\sigma_\epsilon^2} x_{ft}. \quad (24)$$

With $A = \mathbf{\Psi}_{ft}^{-1}$, $U = \mathbf{M}^T$, $C = 1/\sigma_\epsilon^2$ and $V = \mathbf{M}$, we can use the Woodbury matrix identity :

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1} \quad (25)$$

to rewrite $\mathbf{\Sigma}_{ft}$ as

$$\mathbf{\Sigma}_{ft} = \frac{\gamma_{ft}^2 \sigma_{n,ft}^2}{\gamma_{ft}^2 + \sigma_{n,ft}^2 + \sigma_\epsilon^2} \begin{bmatrix} \frac{\sigma_{n,ft}^2}{\sigma_{n,ft}^2 + \sigma_\epsilon^2} & -1 \\ -1 & \frac{\gamma_{ft}^2}{\gamma_{ft}^2 + \sigma_\epsilon^2} \end{bmatrix} \xrightarrow{\sigma_\epsilon^2 \rightarrow 0} \frac{\gamma_{ft}^2 \sigma_{n,ft}^2}{\gamma_{ft}^2 + \sigma_{n,ft}^2} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}. \quad (26)$$

Using (24), we have

$$\boldsymbol{\mu}_{ft} = \begin{bmatrix} \frac{\gamma_{ft}^2}{\gamma_{ft}^2 + \sigma_{n,ft}^2 + \sigma_\epsilon^2} \\ \frac{\sigma_{n,ft}^2}{\gamma_{ft}^2 + \sigma_{n,ft}^2 + \sigma_\epsilon^2} \end{bmatrix} x_{ft} \xrightarrow{\sigma_\epsilon^2 \rightarrow 0} \begin{bmatrix} \frac{\gamma_{ft}^2}{\gamma_{ft}^2 + \sigma_{n,ft}^2} \\ \frac{\sigma_{n,ft}^2}{\gamma_{ft}^2 + \sigma_{n,ft}^2} \end{bmatrix} x_{ft}, \quad (27)$$

in which we recognize the Wiener filter applied to x_{ft} .

We define $\boldsymbol{\mu}_t$ and $\boldsymbol{\Sigma}_t$ with

$$r(\mathbf{s}_t, \mathbf{n}_t) \sim \prod_f \mathcal{N}_c(\boldsymbol{\mu}_{ft}, \boldsymbol{\Sigma}_{ft}) = \mathcal{N}_c(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t). \quad (28)$$

We note $[\boldsymbol{\mu}_{t,s}, \boldsymbol{\mu}_{t,n}] = \boldsymbol{\mu}_t$ and define $\boldsymbol{\Sigma}_{t,ss}$ and $\boldsymbol{\Sigma}_{t,nn}$ to be the diagonal terms of $\boldsymbol{\Sigma}_t$.

2.2 E-z step

After removing the terms independent from \mathbf{z}_t from (15), we can update $r(\mathbf{z}_t)$ using

$$\log r(\mathbf{z}_t) \stackrel{c}{=} \log p(\mathbf{z}_t) + \mathbb{E}_{r(\mathbf{s}_t, \mathbf{n}_t)} [\log p(\mathbf{s}_t | \mathbf{z}_t)] \quad (29)$$

$$\stackrel{c}{=} \log p(\mathbf{z}_t) - \frac{1}{2} \sum_f \left[\log(2\pi\sigma_f^2(\mathbf{z}_t)) + \frac{|\boldsymbol{\mu}_{ft,ss}|^2 + \boldsymbol{\Sigma}_{ft,ss}}{\sigma_f^2(\mathbf{z}_t)} \right] \quad (30)$$

$$\stackrel{c}{=} \log p(\mathbf{z}_t) + p\left(\mathbf{s}_t = \sqrt{|\boldsymbol{\mu}_{t,s}|^2 + \boldsymbol{\Sigma}_{t,ss}} | \mathbf{z}_t\right), \quad (31)$$

where $p(\mathbf{s}_t | \mathbf{z}_t) = \prod_f p(s_{ft} | \mathbf{z}_t)$.

We can invert (31) using Bayes' theorem to obtain

$$r(\mathbf{z}_t) = p\left(\mathbf{z}_t | \mathbf{s}_t = \sqrt{|\boldsymbol{\mu}_{t,s}|^2 + \boldsymbol{\Sigma}_{t,ss}}\right). \quad (32)$$

As maximizing (8) minimizes $\mathcal{D}_{KL}(q_\phi(\mathbf{z} | \mathbf{s}) || p_\theta(\mathbf{z} | \mathbf{s}))$, we can use the probabilistic encoder $q_\phi(\mathbf{z} | \mathbf{s})$ as an approximation of the true posterior $p_\theta(\mathbf{z} | \mathbf{s})$. We further assume that this still holds for \mathbf{s} of the form $\mathbf{s}_t = \sqrt{|\boldsymbol{\mu}_{t,s}|^2 + \boldsymbol{\Sigma}_{t,ss}}$. Finally, we have

$$r(\mathbf{z}_t) \approx q_\phi(\mathbf{z}_t | \mathbf{s}_t = \sqrt{|\boldsymbol{\mu}_{t,s}|^2 + \boldsymbol{\Sigma}_{t,ss}}) \quad (33)$$

$$\sim \mathcal{N}(\tilde{\boldsymbol{\mu}}_l(|\boldsymbol{\mu}_{t,s}|^2 + \boldsymbol{\Sigma}_{t,ss}), \tilde{\sigma}_l^2(|\boldsymbol{\mu}_{t,s}|^2 + \boldsymbol{\Sigma}_{t,ss})). \quad (34)$$

2.3 M-step

We now maximize $\mathcal{L}(r, \boldsymbol{\Theta}_t)$ with respect to $\boldsymbol{\Theta}_t$ with fixed $r(\mathbf{y}_t)$ using (16) :

$$\begin{aligned} \mathcal{Q}(\boldsymbol{\Theta}, \boldsymbol{\Theta}^{old}) &\stackrel{c}{=} \mathbb{E}_{r(\mathbf{s}, \mathbf{b}, \boldsymbol{\Theta}^{old})} [\log p(\mathbf{b}; \boldsymbol{\Theta})] \\ &\stackrel{c}{=} \sum_{ft} d_{IS} \left(|\boldsymbol{\mu}_{ft,n}|^2 + \boldsymbol{\Sigma}_{ft,nn}, (\mathbf{W}\mathbf{H})_{ft} \right) \end{aligned} \quad (35)$$

We obtain classic NMF's multiplicative updates as in [6] :

$$\mathbf{H} \leftarrow \mathbf{H} \odot \frac{\mathbf{W}^T \left((\mathbf{W}\mathbf{H})^{\odot -2} \odot \mathbf{V} \right)}{\mathbf{W}^T (\mathbf{W}\mathbf{H})^{\odot -1}}, \quad (36)$$

$$\mathbf{W} \leftarrow \mathbf{W} \odot \frac{\left((\mathbf{W}\mathbf{H})^{\odot -2} \odot \mathbf{V} \right) \mathbf{H}^T}{(\mathbf{W}\mathbf{H})^{\odot -1} \mathbf{H}^T}, \quad (37)$$

where $\mathbf{V} = \{|\boldsymbol{\mu}_{ft,n}|^2 + \boldsymbol{\Sigma}_{ft,nn}\}_{(ft)}$.

3 Distribution definitions

Proper complex Gaussian: The complex proper Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$, noted $\mathcal{N}_c(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is defined as:

$$\mathcal{N}_c(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\pi \det \boldsymbol{\Sigma}} \exp \left(-(\mathbf{x} - \boldsymbol{\mu})^H \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right), \quad (38)$$

where H denotes the Hermitian transpose. In the univariate case, it simplifies to

$$\mathcal{N}_c(x; \mu, \sigma) = \frac{1}{\pi \sigma^2} \exp \left(-\frac{|x - \mu|^2}{\sigma^2} \right). \quad (39)$$

Real-valued Gaussian The real-valued multivariate distribution with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$, noted $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is defined as:

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{\det 2\pi \boldsymbol{\Sigma}}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right), \quad (40)$$

which, in the univariate case, simplifies to

$$\mathcal{N}(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(x - \mu)^2}{2\sigma^2} \right). \quad (41)$$

References

- [1] M. Pariente, A. Deleforge, and E. Vincent, “A statistically principled and computationally efficient approach to speech enhancement using variational autoencoders,” *Submitted to INTERSPEECH*, 2019.
- [2] A. Ozerov and C. Fevotte, “Multichannel nonnegative matrix factorization in convolutive mixtures for audio source separation,” *IEEE Transactions on Audio, Speech and Language Processing*, vol. 18, no. 3, pp. 550–563, 2010.
- [3] D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” in *ICLR*, 2014.
- [4] C. M. Bishop, *Pattern Recognition and Machine Learning*. Springer, 2006.
- [5] D. J. C. MacKay, *Information Theory, Inference & Learning Algorithms*. New York, NY, USA: Cambridge University Press, 2002.
- [6] C. Fevotte, N. Bertin, and J. Durrieu, “Nonnegative Matrix Factorization with the Itakura-Saito Divergence: With Application to Music Analysis,” *Neural Computation*, vol. 21, no. 3, pp. 793–830, 2009.



**RESEARCH CENTRE
NANCY – GRAND EST**

615 rue du Jardin Botanique
CS20101
54603 Villers-lès-Nancy Cedex

Publisher
Inria
Domaine de Voluceau - Rocquencourt
BP 105 - 78153 Le Chesnay Cedex
inria.fr

ISSN 0249-6399