

Unsupervised Blind Source Separation with Variational Auto-Encoders - Slides

 ${\sf Julian \ Neri}^1 \quad {\sf Philippe \ Depalle}^1 \quad {\sf Roland \ Badeau}^2$

¹McGill University, CIRMMT, Montréal, Canada ²LTCI, Télécom Paris, Institut Polytechnique de Paris, France

August 26, 2021









Overview

Introduction Source Separation Problem Statement

Methodology

Architecture

Evaluation Datasets Images Audio

Conclusion

Source code and data are available on my website, http://www.music.mcgill.ca/~julian/vae-bss

EUSIPCO 2021 - Unsupervised Blind Source Separation with Variational Auto-Encoders - Neri et al. 2/17

Introduction

Unsupervised blind source separation: estimate the underlying sources from a single-channel mixture signal

- without knowing the true number of sources, and
- without clean target source signals for model training.



Applications: re-mixing, hearing aids, dictation, computer vision.

• Ex: Interact with an audio recording in real-time.

EUSIPCO 2021 - Unsupervised Blind Source Separation with Variational Auto-Encoders - Neri et al. 3/17

Introduction

Supervised

Requires the true underlying source files

 Mature enough to separate real musical recordings consisting of vocals, drum, bass, piano.

Unsupervised

Learns solely from the mixtures

- Harder than supervised, requiring strong prior information:
 - Make assumptions about (model) the sources.
 - Learn prior source information from data.

Examples: NMF [Fevotte et al., 2009], DNN [Halperin et al., 2019], [Wisdom et al., 2020].

Methodology - Overview

We address the problem with a new variational auto-encoder (VAE) [Kingma and Welling, 2014, Neri et al., 2021].

- Encoder disentangles (separates) latent sources
- Decoder independently generates each source signal



Methodology - VAE Decoder (Generative Network)

Gaussian Prior
$$p(\mathbf{Z}) = \prod_{k=1}^{K} p(\mathbf{z}_k) = \prod_{k=1}^{K} \mathcal{N}(\mathbf{z}_k | \mathbf{0}, \mathbf{I})$$

= number of sources



EUSIPCO 2021 - Unsupervised Blind Source Separation with Variational Auto-Encoders - Neri et al. 6/17

Methodology - VAE Encoder (Inference Network)



Approximate Posterior

$$q_{\phi}(\mathbf{Z}|\mathbf{x}) = \mathcal{N}\left(\mathbf{Z}|\mu_{\phi}(\mathbf{x}), \sigma_{\phi}^{2}(\mathbf{x})\mathbf{I}\right)$$

Methodology - VAE Training



$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{X}) = \sum_{n=1}^{N} \underbrace{\langle \ln p_{\boldsymbol{\theta}}(\mathbf{x}^{(n)} | \mathbf{Z}) \rangle_{q_{\boldsymbol{\phi}}}}_{\text{Expected log-likelihood}} - \underbrace{D_{KL}\left(q_{\boldsymbol{\phi}}(\mathbf{Z} | \mathbf{x}^{(n)}) \| p(\mathbf{Z})\right)}_{\text{KL divergence}}$$

Expected log-likelihood: reconstruction error between x and \hat{x} **KL divergence**: regulates posterior, helps source disentanglement

Architecture

Table: Encoder input, hidden, and output layer units (dimensions). MNIST: handwritten digit images. MUMS: audio spectrograms.

Detect	Input	Hidden Layers					Output
Dataset	D_x	L1	L2	L3	L4	L5	$2 \times D_Z$
MNIST	784	700	600	500	400	300	$2 \times 20K$
MUMS	32768	2560	2048	1536	1024	512	$2 \times 64K$

VAE Mask (VAEM) After model training (at run-time), makes sum of sources exactly equal to the input mixture.

$$\mathbf{\check{s}}_k = \mathbf{\widehat{s}}_k \odot (\mathbf{x} \oslash \mathbf{\widehat{x}})$$

Evaluation - Datasets

- Data was generated by randomly sampling and adding source signals.
- Sources were never mixed between testing and training datasets.
- For audio, each mixture was transformed into a spectrogram.
- Data normalized to values between 0 and 1.



 $\begin{array}{l} \textbf{MNIST}: 32\,\times\,32 \text{ images of handwritten digits} \\ \textbf{MUMS}: 256\,\times\,128 \text{ spectrograms of audio} \ (5.5 \text{ kHz}\,\times\,1.5 \text{ s}) \end{array}$

EUSIPCO 2021 - Unsupervised Blind Source Separation with Variational Auto-Encoders - Neri et al. 10/17

Evaluation - Baseline Methods, and Ideal Masks

Baseline / State-of-the art

- NMF nonnegative matrix factorization [Fevotte et al., 2009]
- GLO generative latent optimization [Halperin et al., 2019]
- MixIT mixture invariant training [Wisdom et al., 2020]

Ideal masks (for audio)

- IRM ideal ratio mask [Stöter et al., 2018]
- IBM ideal binary mask [Stöter et al., 2018]

Ours

- AE auto-encoder
- VAE variational auto-encoder
- VAEM variational auto-encoder with masking

Evaluation - MNIST Handwritten Digits



Evaluation - MNIST Handwritten Digits



EUSIPCO 2021 - Unsupervised Blind Source Separation with Variational Auto-Encoders - Neri et al. 12/17

Evaluation - MNIST Handwritten Digits

Three Model Sources



Figure: Learned Latent Space Mapped to 2D.

Evaluation - MUMS Audio Spectrograms

MixI I	VAE	VAEM	IBM	IRM
6.64 17.59	14.33 29.92	17.10 29.55	23.97 48.89	22.66 34.39
	6.64 17.59 8.26	6.64 14.33 17.59 29.92 8.26 14.87	6.64 14.33 17.10 17.59 29.92 29.55 8.26 14.87 18.20	6.64 14.33 17.10 23.97 17.59 29.92 29.55 48.89 8.26 14.87 18.20 24.06

Evaluated with bss eval measures [Vincent et al., 2006] [Roux et al., 2019]:

- SI-SDR: scale-invariant signal-to-distortion ratio
- SIR: signal-to-interference ratio
- SAR: signal-to-artifact ratio

Evaluation - MUMS Audio Spectrograms

Mixture of Bassoon (GT-1) & Violin (GT-2)



Listen to more examples at http://www.music.mcgill.ca/~julian/vae-bss

Conclusion

VAE framework offers a viable solution to unsupervised blind separation of both image and audio mixtures.

- Unlike previous methods, it automatically avoids over-separation.
- Generating sources independently using the same decoder is memory efficient and crucial for separation.
- Disentangling low-dimensional encoded sources is practically effective and aligns with human perception.
- \rightarrow **Future research**: model temporal evolution of sources to separate waveforms and videos of arbitrary duration.

References I

Fevotte, C., Bertin, N., and Durrieu, J. (2009).

Nonnegative matrix factorization with the Itakura-Saito divergence: with application to music analysis.

Neural Computation, 21:793-830.



Halperin, T., Ephrat, A., and Hoshen, Y. (2019). Neural separation of observed and unobserved distributions.

In Proc. 36th Int. Conf. Machine Learning (ICML).



Kingma, D. P. and Welling, M. (2014).

Auto-encoding variational Bayes.

In Int. Conf. Learning Representations (ICLR).



Neri, J., Badeau, R., and Depalle, P. (2021).

Unsupervised blind source separation with variational auto-encoders. In Proc. 29th European Signal Processing Confernce (EUSIPCO).

References II

Roux, J. L., Wisdom, S., Erdogan, H., and Hershey, J. R. (2019). SDR – half-baked or well done?

In *IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, pages 626–630.

- Stöter, F.-R., Liutkus, A., and Ito, N. (2018).

The 2018 signal separation evaluation campaign.

In Latent Variable Analysis and Signal Separation, pages 293–305. Springer Int. Publishing.

Vincent, E., Gribonval, R., and Févotte, C. (2006).

Performance measurement in blind audio source separation.

IEEE Trans. Audio, Speech, Lang. Process., 14(4):1462–1469.



Wisdom, S., Tzinis, E., Erdogan, H., Weiss, J., Wilson, K., and Hershey, J. (2020).

Unsupervised sound separation using mixture invariant training. Advances in Neural Information Processing Systems.