

Maurice Frank

Unsupervised music source separation with deep generative priors



# Contents

<i>Proposal</i>	5
<i>Abstract</i>	5
<i>Research Question</i>	5
<i>Related works</i>	5
<i>Methodology</i>	9
<i>Planning</i>	9



# Proposal

## Abstract

### Research Question

Source separation is the task of finding a set of latent sources  $\mathbf{s} = [s_1, \dots, s_n]$  to an observed mix of those sources  $\mathbf{m}$ . The induced model proposes a mixing function  $\mathbf{m} = f(\mathbf{s})$  which might just be a linear mixing  $\mathbf{m} = \mathbf{A} \cdot \mathbf{s}$ . The task is to find the inverse model  $f^{-1}(\cdot)$  which retrieves  $\mathbf{s}$ .

Can we learn an sound source separation model in an unsupervised manner. Unsupervised relating to missing pairing of sources to mixes.

$$p(s'|m) \cdot p(s) \quad (1)$$

### Related works

In this chapter we discuss previous research in supervised and semi-supervised source separation.

### Deep Latent-Variable Models

For our process we have observations from the data space  $\mathbf{x} \in \mathcal{D}$  for which there exists an unknown data probability distribution  $p^*(\mathcal{D})$ . We collect a data set  $\{\mathbf{x}_1 \dots \mathbf{x}_N\}$  with  $N$  samples. We introduce an approximate model with density<sup>1</sup>  $p_\theta(\mathcal{D})$  and model parameters  $\theta$ . Learning or modelling means finding the values for  $\theta$  which will give the closest approximation of the true underlying process:

$$p_\theta(\mathcal{D}) \approx p^*(\mathcal{D}) \quad (2)$$

The model  $p_\theta$  has to be complex enough to be able to fit the data density while little enough parameters to be learnable. Every choice for the form of the model will *induce* biases<sup>2</sup> about what density we

<sup>1</sup> We write density and distribution interchangeably to denote a probability function.

<sup>2</sup> called *inductive biases*

can model, even before we maximize a learning objective using the parameters  $\theta$ .

In the following described models we assume the sampled data points  $\mathbf{x}$  to be drawn from  $\mathcal{D}$  *independent and identically distributed*<sup>3</sup>. Therefore we can write the data log-likelihood as:

$$p_{\theta}(\mathcal{D}) = \prod_{\mathbf{x} \in \mathcal{D}} p_{\theta}(\mathbf{x}) \quad (3)$$

$$\log p_{\theta}(\mathcal{D}) = \sum_{\mathbf{x} \in \mathcal{D}} \log p_{\theta}(\mathbf{x}) \quad (4)$$

The maximum likelihood estimation of our model parameters maximizes this objective.

To form a latent-variable model we introduce a *latent variable*<sup>4</sup>. The data likelihood now is the marginal density of the joint latent density:

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} \quad (5)$$

Typically we introduce a factorization of the joint. Most commonly and simplest:

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) p(\mathbf{z}) d\mathbf{z} \quad (6)$$

This corresponds to the graphical model in which  $\mathbf{z}$  is generative parent node of the observed  $\mathbf{x}$ , see Figure 1. The density  $p(\mathbf{z})$  is called the *prior distribution*.

If the latent is small, discrete, it might be possible to directly marginalize over it. If for example  $\mathbf{z}$  is a discrete random variable and the conditional  $p_{\theta}(\mathbf{x}|\mathbf{z})$  is a Gaussian distribution than the data model density  $p_{\theta}(\mathbf{x})$  becomes a mixture-of-Gaussians, which we can directly estimate by maximum likelihood estimation of the data likelihood.

For more complicated models the data likelihood  $p_{\theta}(\mathbf{x})$  as well as the model posterior  $p_{\theta}(\mathbf{z}|\mathbf{x})$  are intractable because of the integration over the latent  $\mathbf{z}$  in Equation (6).

To formalize the search for an intractable posterior into a tractable optimization problem we follow the *variational principle*<sup>5</sup> which introduces an approximate posterior distribution  $q_{\phi}(\mathbf{z}|\mathbf{x})$ , also called the *inference model*. Again the choice of model here carries inductive biases as such that even in asymptotic expectation we can not obtain the true posterior.

Following the derivation in<sup>6</sup> we introduce the inference model into the data likelihood:

<sup>3</sup> meaning the sample of one datum does not depend on the other data points

<sup>4</sup> Latent variables are part of the directed graphical model but not observed.

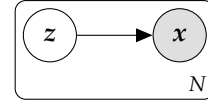


Figure 1: The graphical model with a introduced latent variable  $\mathbf{z}$ . Observed variables are shaded.

<sup>5</sup> Michael I. Jordan et al. "An Introduction to Variational Methods for Graphical Models". In: *Machine Learning* 37.2 (1999), pp. 183–233.

<sup>6</sup> Diederik P. Kingma and Max Welling. "An Introduction to Variational Autoencoders". In: (2019). arXiv: [1906.02691](https://arxiv.org/abs/1906.02691) [cs, stat], p. 20.

$$\log p_{\theta}(\mathbf{x}) = \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x})] \quad (7)$$

$$= \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x})} \right] \quad (8)$$

$$= \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p_{\theta}(\mathbf{z}|\mathbf{x})} \right] \quad (9)$$

$$= \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] + \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p_{\theta}(\mathbf{z}|\mathbf{x})} \right] \quad (10)$$

$$= \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] + \mathbb{D}_{\text{KL}}[q_{\phi}(\mathbf{z}|\mathbf{x}) \| p_{\theta}(\mathbf{z}|\mathbf{x})] \quad (11)$$

Note that we separated the likelihood into two parts. The second part is the (positive) Kullback-Leibler divergence of the approximate posterior from the true intractable posterior. This unknown divergence states the ‘correctness’ of our approximation <sup>7</sup>.

The first term is the *variational free energy* <sup>8</sup> or *evidence lower bound* (ELBO):

$$\text{ELBO}_{\theta, \phi}(\mathbf{x}, \mathbf{z}) = \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \quad (12)$$

We can introduce the same factorization as in Equation (6):

$$\text{ELBO}_{\theta, \phi}(\mathbf{x}, \mathbf{z}) = \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \quad (13)$$

$$= \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] + \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] \quad (14)$$

$$= -\mathbb{D}_{\text{KL}}[q_{\phi}(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z})] + \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] \quad (15)$$

$$(16)$$

Under this factorization we separated the lower bound into two parts. First the divergence of the approximate posterior from the latent prior distribution and second the data posterior likelihood from the latent <sup>9</sup>.

The optimization of the  $\text{ELBO}_{\theta, \phi}$  allows us to jointly optimize the parameters  $\theta$  and  $\phi$ . Using data samples we can compute the gradient of the ELBO with respect to  $\theta$  <sup>10</sup> allowing for an unbiased Monte Carlo estimate of the gradient. We can *not* though do the same for the variational parameters  $\phi$ , as the expectation of the ELBO is over the approximate posterior which depends on  $\phi$ .

Through a change of variable of the latent random variable we can make the gradient with respect to the variational parameters tractable. With this reparameterization trick.<sup>11</sup>

<sup>7</sup> More specifically the divergence marries two errors of our approximate model. First it gives the error of our posterior estimation from the true posterior, by definition of divergence. Second it specifies the error of our complete model likelihood from the marginal likelihood. This is called the *tightness* of the bound.

<sup>8</sup> TODO

<sup>9</sup> this will later be the reconstruction error. How well can we return to the data density from latent space

<sup>10</sup>  $\nabla_{\theta} \text{ELBO}_{\theta, \phi} \approx \nabla_{\theta} \log p_{\theta}(\mathbf{x}, \mathbf{z})$

<sup>11</sup> Diederik P. Kingma and Max Welling. “Auto-Encoding Variational Bayes”. In: (2014). arXiv: 1312.6114 [cs, stat].

## The VAE framework

VAE<sup>12,13</sup>

$\beta$ -VAE<sup>14</sup> - introduces  $\beta$  as controlling hyperparameter in the VAE objective - constraint that controls the capacity of the latent space - gives trade off between reconstruction quality and representation simplicity - similar to information bottleneck<sup>15</sup>

VQ-VAE<sup>16</sup>

## Flow based models

Another class of common deep latent models are based on *normalizing flows*.<sup>17</sup> They use a flow for the approximate posterior  $q_\phi(z|x)$ . A normalizing flow is a function  $f(x)$  that maps the input density to a fixed, prescribed density  $p(\epsilon) = p(f(x))$ , in that normalizing the density<sup>18</sup>. Again this is commonly set to be a factorized Gaussian distribution.

For a finite normalizing flow we consider a chain of invertible, smooth mappings.

NICE<sup>19</sup> - volume preserving transformations - coupling layer - triangular shape

Normalizing Flow<sup>20</sup>

RealNVP<sup>21</sup> - non-volume preserving

Glow<sup>22</sup> - invertible 1x1 convs - ActNorm - zero init

<sup>23</sup> introduced WaveNet an autoregressive generative model for raw (*time-domain*) audio. WaveNet closely similar to the earlier PixelCNN<sup>24</sup> but adapted for the audio domain. Unmodified Cnns are unsuitable to the application to raw audio because of the form of data. as digital audio is sampled at an extremely high sample rate commonly 16kHz up to 44kHz the features of interest lie at scale of stringly different magnitudes. On the one hand recognizing phase, frequency of a wave might require features at those ms scales on the other hand the modelling of speech or music audio happens at the scale of seconds or minutes. As such a generative model for this domain has to capture those different time scales. The wavenet accomplishes this by using dilated convolutions a common tool in signal processing.<sup>25</sup> A dilated convolutions uses a kernel with an inner stride. Using a stack of dilated convolutions increases the receptive field of the features without increasing the computational complexity.

- gated convs -pixelcnn -lstm<sup>26</sup> - dilated convs - global conditioning -  $\mu$ -law encoding<sup>27</sup> - slow cause autoreg (better with<sup>28</sup>) -

PixelCNN++<sup>29</sup>

<sup>12</sup> Kingma and Welling, "Auto-Encoding Variational Bayes".

<sup>13</sup> Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. "Stochastic Backpropagation and Approximate Inference in Deep Generative Models". In: (2014). arXiv: [1401.4082 \[cs, stat\]](#).

<sup>14</sup> Irina Higgins et al. "Beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework". In: (2016).

<sup>15</sup> Christopher P. Burgess et al. "Understanding Disentangling in Beta-VAE". In: (2018). arXiv: [1804.03599 \[cs, stat\]](#).

<sup>16</sup> Aäron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. "Neural Discrete Representation Learning". In: *Advances in Neural Information Processing Systems* 30. Ed. by I. Guyon et al. Curran Associates, Inc., 2017, pp. 6306–6315.

<sup>17</sup> Esteban Tabak and Cristina V. Turner. "A family of nonparametric density estimation algorithms". In: *Communications on Pure and Applied Mathematics* 66.2 (2013), pp. 145–164.

<sup>18</sup> The extreme of this idea is, of course, an infinitesimal, continuous-time flow with a velocity field.

<sup>19</sup> Laurent Dinh, David Krueger, and Yoshua Bengio. "NICE: Non-Linear Independent Components Estimation". In: (2015). arXiv: [1410.8516 \[cs\]](#).

<sup>20</sup> Danilo Jimenez Rezende and Shakir Mohamed. "Variational Inference with Normalizing Flows". In: (2016). arXiv: [1505.05770 \[cs, stat\]](#).

<sup>21</sup> Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. "Density Estimation Using Real NVP". In: (2017). arXiv: [1605.08803 \[cs, stat\]](#).

<sup>22</sup> Diederik P. Kingma and Prafulla Dhariwal. "Glow: Generative Flow with Invertible 1x1 Convolutions". In: (2018). arXiv: [1807.03039 \[cs, stat\]](#).

<sup>23</sup> Aäron van den Oord et al. "WaveNet: A Generative Model for Raw Audio". In: (2016). arXiv: [1609.03499 \[cs\]](#).

<sup>24</sup> Aäron van den Oord et al. "Conditional Image Generation with PixelCNN Decoders". In: (2016). arXiv: [1606.05328 \[cs\]](#).

<sup>25</sup> P. Dutilleul. "An Implementation of the Algorithm à Trous to Compute the Wavelet Transform". In: *Wavelets*. Ed. by Jean-Michel Combes, Alexander Grossmann, and Philippe Tchamitchian. Inverse Problems and Theoretical Imaging. Berlin, Heidelberg: Springer, 1990, pp. 298–304.

<sup>26</sup> Sepp Hochreiter and Jürgen Schmidhuber. "Long Short-Term Memory". In: *Neural Computation* 9.8 (1997), pp. 1735–



*Sound*NSynth<sup>30</sup>In<sup>31</sup>FloWaveNet<sup>32</sup>*Source separation*WaveNet for Speech denoising<sup>33</sup>WaveNet-VAE unsupervised speech rep learning<sup>34</sup>Wave-U-Net<sup>35</sup>DeMucs<sup>36</sup>Source Sep in Time Domain<sup>37</sup>*Methodology**Datasets**ToyData**MusDB**Planning*

<sup>30</sup> Nal Kalchbrenner et al. "Efficient Neural Audio Synthesis". In: (2018). arXiv: [1802.08435 \[cs, eess\]](#).

<sup>31</sup> Ryan Prenger, Rafael Valle, and Bryan Catanzaro. "WaveGlow: A Flow-Based Generative Network for Speech Synthesis". In: (2018). arXiv: [1811.00002 \[cs, eess, stat\]](#).

<sup>32</sup> Sungwon Kim et al. "FloWaveNet : A Generative Flow for Raw Audio". In: (2019). arXiv: [1811.02155 \[cs, eess\]](#).

<sup>33</sup> Dario Rethage, Jordi Pons, and Xavier Serra. "A Wavenet for Speech Denoising". In: (2018). arXiv: [1706.07162 \[cs\]](#).

<sup>34</sup> Jan Chorowski et al. "Unsupervised Speech Representation Learning Using WaveNet Autoencoders". In: *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 27.12 (2019), pp. 2041–2053. arXiv: [1901.08810](#).

<sup>35</sup> Daniel Stoller, Sebastian Ewert, and Simon Dixon. "Wave-U-Net: A Multi-Scale Neural Network for End-to-End Audio Source Separation" (Paris, France). 2018.

<sup>36</sup> Alexandre Défossez et al. "Demucs: Deep Extractor for Music Sources with Extra Unlabeled Data Remixed". In: (2019). arXiv: [1909.01174 \[cs, eess, stat\]](#).

<sup>37</sup> Francesc Lluís, Jordi Pons, and Xavier Serra. "End-to-End Music Source Separation: Is It Possible in the Waveform Domain?" In: (2019). arXiv: [1810.12187 \[cs, eess\]](#).



# Bibliography

- Burgess, Christopher P. et al. "Understanding Disentangling in Beta-VAE". In: (2018). arXiv: [1804.03599 \[cs, stat\]](#) (cit. on p. 8).
- Chorowski, Jan et al. "Unsupervised Speech Representation Learning Using WaveNet Autoencoders". In: *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 27.12 (2019), pp. 2041–2053. arXiv: [1901.08810](#) (cit. on p. 9).
- Défossez, Alexandre et al. "Demucs: Deep Extractor for Music Sources with Extra Unlabeled Data Remixed". In: (2019). arXiv: [1909.01174 \[cs, eess, stat\]](#) (cit. on p. 9).
- Dinh, Laurent, David Krueger, and Yoshua Bengio. "NICE: Non-Linear Independent Components Estimation". In: (2015). arXiv: [1410.8516 \[cs\]](#) (cit. on p. 8).
- Dinh, Laurent, Jascha Sohl-Dickstein, and Samy Bengio. "Density Estimation Using Real NVP". In: (2017). arXiv: [1605.08803 \[cs, stat\]](#) (cit. on p. 8).
- Dutilleul, P. "An Implementation of the Algorithme à Trous to Compute the Wavelet Transform". In: *Wavelets*. Ed. by Jean-Michel Combes, Alexander Grossmann, and Philippe Tchamitchian. Inverse Problems and Theoretical Imaging. Berlin, Heidelberg: Springer, 1990, pp. 298–304 (cit. on p. 8).
- Higgins, Irina et al. "Beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework". In: (2016) (cit. on p. 8).
- Hochreiter, Sepp and Jürgen Schmidhuber. "Long Short-Term Memory". In: *Neural Computation* 9.8 (1997), pp. 1735–1780 (cit. on p. 8).
- Jordan, Michael I. et al. "An Introduction to Variational Methods for Graphical Models". In: *Machine Learning* 37.2 (1999), pp. 183–233 (cit. on p. 6).
- Kalchbrenner, Nal et al. "Efficient Neural Audio Synthesis". In: (2018). arXiv: [1802.08435 \[cs, eess\]](#) (cit. on p. 8).
- Kim, Sungwon et al. "FloWaveNet : A Generative Flow for Raw Audio". In: (2019). arXiv: [1811.02155 \[cs, eess\]](#) (cit. on p. 9).
- Kingma, Diederik P. and Prafulla Dhariwal. "Glow: Generative Flow with Invertible 1x1 Convolutions". In: (2018). arXiv: [1807.03039 \[cs, stat\]](#) (cit. on p. 8).
- Kingma, Diederik P. and Max Welling. "Auto-Encoding Variational Bayes". In: (2014). arXiv: [1312.6114 \[cs, stat\]](#) (cit. on pp. 7, 8).
- Kingma, Diederik P. and Max Welling. "An Introduction to Variational Autoencoders". In: (2019). arXiv: [1906.02691 \[cs, stat\]](#) (cit. on p. 6).
- Lluís, Francesc, Jordi Pons, and Xavier Serra. "End-to-End Music Source Separation: Is It Possible in the Waveform Domain?" In: (2019). arXiv: [1810.12187 \[cs, eess\]](#) (cit. on p. 9).

- Paine, Tom Le et al. “Fast Wavenet Generation Algorithm”. In: (2016). arXiv: [1611.09482 \[cs\]](#) (cit. on p. 8).
- Prenger, Ryan, Rafael Valle, and Bryan Catanzaro. “WaveGlow: A Flow-Based Generative Network for Speech Synthesis”. In: (2018). arXiv: [1811.00002 \[cs, eess, stat\]](#) (cit. on p. 9).
- Recommendation G. 711. Pulse Code Modulation (PCM) of Voice Frequencies*. 1988 (cit. on p. 8).
- Rethage, Dario, Jordi Pons, and Xavier Serra. “A Wavenet for Speech Denoising”. In: (2018). arXiv: [1706.07162 \[cs\]](#) (cit. on p. 9).
- Rezende, Danilo Jimenez and Shakir Mohamed. “Variational Inference with Normalizing Flows”. In: (2016). arXiv: [1505.05770 \[cs, stat\]](#) (cit. on p. 8).
- Rezende, Danilo Jimenez, Shakir Mohamed, and Daan Wierstra. “Stochastic Backpropagation and Approximate Inference in Deep Generative Models”. In: (2014). arXiv: [1401.4082 \[cs, stat\]](#) (cit. on p. 8).
- Salimans, Tim et al. “PixelCNN++: Improving the PixelCNN with Discretized Logistic Mixture Likelihood and Other Modifications”. In: (2017). arXiv: [1701.05517 \[cs, stat\]](#) (cit. on p. 8).
- Stoller, Daniel, Sebastian Ewert, and Simon Dixon. “Wave-U-Net: A Multi-Scale Neural Network for End-to-End Audio Source Separation” (Paris, France). 2018 (cit. on p. 9).
- Tabak, Esteban and Cristina V. Turner. “A family of nonparametric density estimation algorithms”. In: *Communications on Pure and Applied Mathematics* 66.2 (2013), pp. 145–164 (cit. on p. 8).
- Van den Oord, Aäron, Oriol Vinyals, and Koray Kavukcuoglu. “Neural Discrete Representation Learning”. In: *Advances in Neural Information Processing Systems 30*. Ed. by I. Guyon et al. Curran Associates, Inc., 2017, pp. 6306–6315 (cit. on p. 8).
- Van den Oord, Aäron et al. “Conditional Image Generation with PixelCNN Decoders”. In: (2016). arXiv: [1606.05328 \[cs\]](#) (cit. on p. 8).
- Van den Oord, Aäron et al. “WaveNet: A Generative Model for Raw Audio”. In: (2016). arXiv: [1609.03499 \[cs\]](#) (cit. on p. 8).