

Unsupervised source separation with deep priors

Maurice Frank

July 20, 2020

Contents

Method	5
Modeling the priors	7
Sampling from the posteriors	10
Modeling the posteriors	11
Difficulties of this method	14

Method

IN THIS CHAPTER, we introduce the theoretical idea of our probability modelling of musical source separation. First, we explicitly state our chosen model of the problem, next derive a suitable optimization objective from it and then explain how we can optimize towards that.

We propose the graphical model as shown in Figure 1 as the generative story of music tracks being generated from separate source tracks. For each source, a sample is taken from the latent source distribution. The observed mix is generated deterministically from the full set of sources. Without loss of generality, we fix this function to be the mean.

Our stated task in ?? is to retrieve the sources $\{s_1, \dots, s_N\}$ from a given mix m . Our model is trained without using sample tuples

$$(m, \{s_1, \dots, s_N\}) : f(\{s_1, \dots, s_N\}) = m$$

which would show the relation between separate sources and mixes. The general idea is visualized in Figure 2.

Looking back at the graphical model in Figure 1, it implies the following factorization:

$$p(m) = \int^N p(s_1, \dots, s_N, m) d^N s \quad (1)$$

$$= \int^N p(m|s_1, \dots, s_N) \cdot p(s_1, \dots, s_N) d^N s \quad (2)$$

$$= \int^N p(m|s_1, \dots, s_N) \cdot p(s_1) \cdots p(s_N) d^N s \quad (3)$$

While the conditional $p(m|s_1, \dots, s_N)$ is not even probabilistic, as the mix is generated deterministically with the mean of the sources, the model posterior $p(s_1, \dots, s_N|m)$ is intractable and precisely what we are interested in. Again we want to make it clear that this setup changes the typical optimization target, as seen in other source separation approaches. We factorize the distribution $p(m)$ of mixed songs, only then to extract the most likely *latent* components that generate this mix, the sources. Therefore we are explicitly modelling the generative process of the mixed songs, and only implicitly the separation task.

Already with the graphical model we introduce a fairly big as-

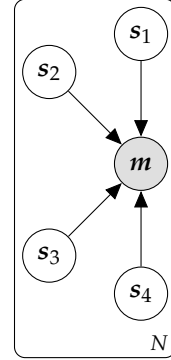


Figure 1: The used graphical model for the source separation task. We have the latent source channel variables s_k . Exemplary here, as in our data, we have four sources. The mix m is observed.

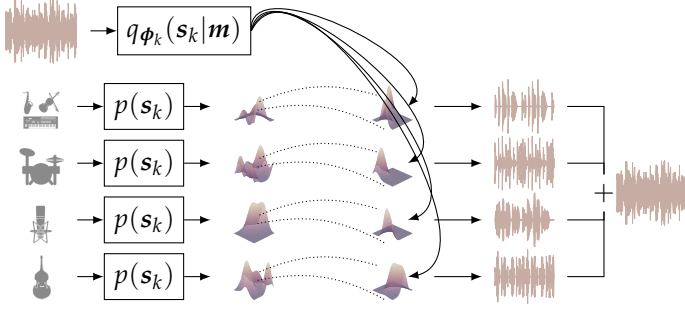


Figure 2: The method visualized, **not finished yet**

sumption, namely that we can model the source distributions independently from each other when modeling the joint:

$$p(\mathbf{m}, \mathbf{s}_1, \dots, \mathbf{s}_N) \equiv p(\mathbf{m} | \mathbf{s}_1, \dots, \mathbf{s}_N) \cdot \prod_k^N p(\mathbf{s}_k) \quad (4)$$

Intuitively this assumption does not, in general, hold for the musical domain. We can expect a different behaviour for a *guitar* in a song involving a second source like a *percussion set*. The joint of those two models is different than their independent densities ¹. Nevertheless, this assumption is crucial to be able to model the prior instrumental models without needing the tuples $(\mathbf{s}_1, \dots, \mathbf{s}_N)$ of co-occurring stems.

The general process is as follows:

First, because we assumed in our graphical model the latent sources to be independent, there exists a probability distribution $p(\mathbf{s}_k)$ for each source k . Thereout we need to choose a model which gives the parametrized approximated prior $p_{\theta}(\mathbf{s}_k)$ with the parameters θ which we optimize with the samples from \mathcal{D}_k .

Second, for each source, there exists a posterior given a mix $p(\mathbf{s}_k | \mathbf{m})$ from which we want to draw samples. Here two approaches are possible. Either we can propose an approximate posterior $q_{\phi_k}(\mathbf{s}_k | \mathbf{m})$, which is trained to minimize the divergence from the previously trained and fixed corresponding prior (VAE setting). Or, we can sample directly from the posterior using Stochastic Gradient Langevin Dynamics (SGLD)² without explicitly approximating the posterior distribution (sampling setting). Both optimization, either training or sampling, happen under the mixing constraint:

$$\sum_{k=1}^N \mathbf{s}_k = \mathbf{m} \quad (5)$$

When using the VAE setting we then can sample from each posterior to retrieve a set of (approximately) correct source tracks. In the case of using Langevin dynamics, the iterative sampling will directly give this set from the prior distributions.

Thinking in the common terms of the VAE we need to choose:

¹ A practical counter-example: If learning the behaviour of drums, only from solo drum recordings, one will encounter complex rhythms and sounds. In many rock genres drums often exhibit simpler sounds, providing beat and tempo. These two distributions are not the same.

² Welling and Teh, “Bayesian Learning via Stochastic Gradient Langevin Dynamics”, 2011.

1. a parametrization $p_{\theta}(s_k)$ of the *prior* distribution
2. a parametrized approximate posterior distribution $q_{\phi_k}(s_k|\mathbf{m})$ with parameters η
3. a parametrized *encoder* which gives the parameters for the posterior distribution $\text{Encoder}_{\phi}(\mathbf{m}) = \eta$
4. a *decoder* which returns the input \mathbf{m} from the latents $\text{Decoder}(\{s_1, \dots, s_N\}) = \mathbf{m}$

As stated before the decoder in our case is the mean function, thus not probabilistic and without parameters. It is certainly possible to parametrize the decoder with trainable parameters. This would imply though that the latent samples are not in the *direct form* of the source sounds but under an (unknown) transformation.

Modeling the priors

The first step in the training procedure is the training of the independent source priors $p(s_k)$. We have to choose a model for estimating the density that results in smooth, tight densities but also is capable enough to capture the complex data space of audio data. We choose to estimate the priors, with flow models for which we gave an overview in ???. The invertibility of the flow gives us the ability to sample from the priors, which is necessary for the sampling-based approach. For the variational-based approach sampling is not necessary, nevertheless we utilize the same generative priors for both experiments. For different representations of the input, different network architectures are more suited. We experiment both with using directly the time-domain wave but also modelling on top of a spectral transformation of the time-domain data.

The two main variants of normalizing flows we build our priors from are the RealNVP and Glow. Their most important difference is the different formulation of the coupling layers. For the coupling layers, we have to split the layer input into two equally sized subsets, one of whom will be transformed using an invertible transformation parametrized by the other. The RealNVP masks the data spatially or over the channels. When done spatially a checkerboard binary mask is used (same for all channels), for the channel variant the channels are assigned in an even-odd switching. Glow, on the other hand, learns a 1×1 -convolution which permutes the data channel-wise and than just splits the data along the channel dimension into two. For both RealNVP and Glow prior work exists specifically adapting the architecture to the time-series domain by integrating WaveNet modules, FloWaveNet³ and WaveGlow⁴, respectively.

Remove Glow if no experiment will use it

For the case of the time-domain representation, the data has only one channel (mono-audio). Therefore a Glow-like shuffling of channels is not possible and we resort to using spatial masking

³ Kim et al., “FloWaveNet : A Generative Flow for Raw Audio”, 2019.

⁴ Prenger et al., “WaveGlow: A Flow-Based Generative Network for Speech Synthesis”, 2018.

for the coupling layers. A one-dimensional mask is used over the time dimension. In the case of using a spectral representation of the audio data, the input samples have a lower time resolution but instead the frequency distribution over a deeper channel dimension. Therefore we can also experiment with channel-shuffling.

We make it clear that the success of this method depends strongly on the expressiveness and smoothness of the learned prior models. In the case of learning the posterior, variational learning will fit the posterior distribution under the prior distribution. If the prior does not contain enough variation of the actual data distribution, the de-mixing model cannot fit a conditional posterior similar to the sample in the mixed song. In the case of the sampling approach, this constraint is even more stressed. The mixed is generated by iteratively sampling from the set of prior distributions, any possible source sample, therefore, has to be found on the manifold of the prior. If any of the priors does not capture the full variations of the data distributions, neither method will be able to extract the correct separated sources.

More practically we again point out the high difficulty of modelling audio data. As laid out in ?? audio consists of information at greatly different resolutions. A generative model of an instrument has to model the low-range components, e.g. pitch, timbre, tremolo and high-range components, e.g. notes and rhythm.

Prior architecture

See the visualization of the prior model’s architecture in Figure 3. We construct normalizing flow models, following the description in Kim et al.⁵.

The flow network, at the highest level, consists of n_b context blocks, each of whom contain one squeeze operation and n_f flows. The squeeze operation halves the length of the signal by doubling the number of channels. By doing so, we effectively double the receptive field of the WaveNet operators in the affine couplings. One flow contains one affine coupling layer after which the masking for the coupling (*even/odd*) is inverted. Before the coupling layer, we apply activation normalization⁶. The affine coupling is applying an affine transformation on the *odd* set of the input, with scale s and translation t coming from the *even* input put through a *non-causal* WaveNet.

The *even/odd* masking for the coupling layer is simply achieved by splitting the channel dimension into two equal parts. Flipping the masking exchanges the top half of channels with the bottom half. No channel transformation of the channels is learned as in Glow. Further observe that through the repeated squeeze operation the masking implicitly includes a checkerboard pattern on the temporal dimension, similar to the explicit pattern used in RealNVP.

As explained in ?? the invertible flow $f(\cdot)$ network directly maps the data space into a prior distribution $p(z)$. Through the invertibil-

⁵ Kim et al., “FloWaveNet : A Generative Flow for Raw Audio”, 2019.

⁶ Kingma and Dhariwal, “Glow: Generative Flow with Invertible 1x1 Convolutions”, 2018.

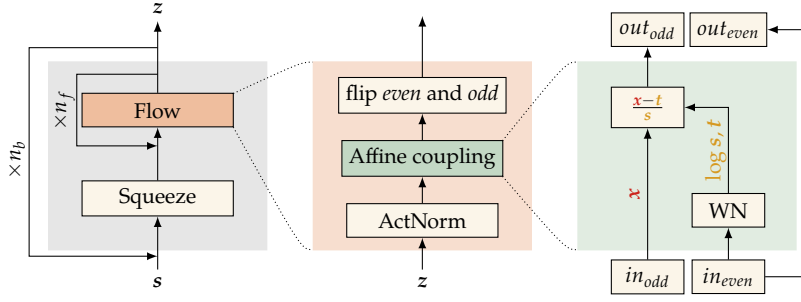


Figure 3: The building blocks for the prior model. The model consists of n_b blocks (left). Each block consists of n_f flows (middle). In each flow we apply activation normalization, followed by the affine coupling (right), after which the binary mask for the even/odd mapping is inverted. The affine coupling layer uses a WaveNet with the *even* set as the input to output scaling $\log s$ and translation t which the *odd* set is transformed.



Figure 4: Flip and squeeze operation for the coupling layers. The squeeze operation doubles the channels and halves the length while the flip operation inverts the tensor along the channel dimension. figure not finished yet

ity we can directly optimize the log probability with the change of variable formula:

$$\log p(x) = \log p(f(x)) + \log \det \left(\frac{\partial f(x)}{\partial x} \right) \quad (6)$$

Our network consists of n_b blocks with each n_f flows each with one activation normalization layer f_{AN} and an affine coupling layer f_{AC} . Therefore the log-determinant becomes:

$$\log \det \left(\frac{\partial f(x)}{\partial x} \right) = \sum_i^{n_b \cdot n_f} \log \det \frac{\partial f_{AN}(x)}{\partial x} + \log \det \frac{\partial f_{AC}(x)}{\partial x} \quad (7)$$

The activation normalization layer f_{AN} learns an affine transformation per channel i :

$$f_{AN}^i(x_i) = s_i \cdot x_i + t_i \quad (8)$$

The log-determinant is easily computed with⁶:

$$\log \det \frac{\partial f_{AN}(x)}{\partial x} = T \cdot \sum_i^C \log |s_i| \quad (9)$$

where C is the number of channels and T is the length of the signal.

For the affine coupling layer, the log-determinant is made tractable by only transforming half the elements in each application. The coupling transforms the *odd* part as follows:

$$\log s, t = \text{WN}(in_{odd}) \quad (10)$$

$$out_{odd} = f_{AC}(in_{odd}) = \frac{in_{odd} - t}{s} \quad (11)$$

The inverse is:
 $f_{AN}^{-1}(y_i) = \frac{y_i - t_i}{s_i}$

The inverse is:
 $f_{AC}^{-1}(out_{odd}) = out_{odd} \cdot s + t$

where $WN(\cdot)$ is a WaveNet. The Jacobian of the coupling function is a lower triangular matrix⁷ because of which the log-determinant becomes the product of the diagonal elements:

$$\log \det \frac{\partial f_{AC}(x)}{\partial x} = -\sum s \quad (12)$$

The purpose of the invertible flow is it to project the complex data distribution into a tractable defined prior distribution. In most cases, as in this, the prior $p(z)$ is set to be a factorized Gaussian distribution with zero mean and variance of one. As such the likelihood of $f(x) \sim p(z)$ can be evaluated with:

$$\log p(f(x)) = -\frac{1}{2}(\log(2 \cdot \pi)) + f(x)^2 \quad (13)$$

To train the model we maximize the likelihood in Equation (6) which we is the sum of the parts in Equations (9), (12) and (13). Optimization is done in mini-batches with the Adam optimizer⁸ and a decaying learning rate scheme.

It is noted that the likelihood under the prior is computed “pixel”-wise⁹. While for the optimization the actual loss on which the optimizers performs update steps is the mean likelihood, we can use the trained prior flow to assign a likelihood map to an input signal where each log-likelihood value is depending on this values receptive field.

WaveNet

In the affine coupling layer we utilize a *non-causal* WaveNet for generating the scale $\log s$ and translation t values. The WaveNet is non-causal because the filters are centred, as such the receptive field of an output value is symmetrically in the past and future input space. The kernel size is fixed to 3. The WaveNet is using the original gated convolution scheme described in ???. Appended to the last skip connection output are two additional one-dimensional convolutional layers, each preceded by a ReLU. The last convolutional layer is initialized to zero. By that we initialize the affine transformation of the coupling layer to the identity at the beginning of training, helping to stabilize optimization. The last layer outputs both scaling and translation factors. Learning those factors with weight sharing does not change the formulation of the coupling layer, as the log-determinant does not depend on the network. Learning them in the same network might increase efficiency, as well as, increasing implementation simplicity. The WaveNet completely consists of convolutional operations. As such it can be applied to one-dimensional signals of variable length.

Sampling from the posteriors

After having estimated the N prior distributions the first possibility for retrieving sample estimates s_k is sampling from the posterior

⁷ Dinh, Sohl-Dickstein, et al., “Density Estimation Using Real NVP”, 2017.

⁸ Kingma and Ba, “Adam: A Method for Stochastic Optimization”, 2017.

⁹ time-points in our sound case

using Langevin dynamics¹⁰. With Stochastic Gradient Langevin Dynamics (SGLD) we can sample from $s_i \sim p(\cdot|\mathbf{m})$ without computing, explicitly, the posterior $p(s_i|\mathbf{m})$. For an overview of SGLD see ??.

Starting with the update step in SGLD we reformulate the update with our problem constraint:

$$\mathbf{s}_k^{(t+1)} = \mathbf{s}_k^{(t+1)} + \eta \cdot \nabla_{\mathbf{s}_k} \log p(\mathbf{s}_k^{(t)}|\mathbf{m}) + 2\sqrt{\eta}\epsilon_t \quad (14)$$

$$= \mathbf{s}_k^{(t+1)} + \eta \cdot \nabla_{\mathbf{s}_k} [\log p(\mathbf{m}|\mathbf{s}_k) + \log p(\mathbf{s}_k) - \log p(\mathbf{m})] + 2\sqrt{\eta}\epsilon_t \quad (15)$$

$$= \mathbf{s}_k^{(t+1)} + \eta \cdot \nabla_{\mathbf{s}_k} [\log p(\mathbf{m}|\mathbf{s}_k) + \log p(\mathbf{s}_k)] + 2\sqrt{\eta}\epsilon_t \quad (16)$$

$$= \mathbf{s}_k^{(t+1)} + \eta \cdot \left[\|\mathbf{m} - \frac{1}{N} \sum_j \mathbf{s}_j^{(t)}\| + \nabla_{\mathbf{s}_k} \log p(\mathbf{s}_k) \right] + 2\sqrt{\eta}\epsilon_t \quad (17)$$

η is the update step size and $\epsilon_t \sim \mathcal{N}(0, \mathbb{1})$.

Note that the final formulation of the update step is simply a constrained greedy hill-climb under the prior model with added Gaussian noise. Intuitively the artificially noised update makes it harder for the greedy optimization to get stuck in local minima of the prior surface. See Algorithm 1 for the algorithm.

```

1: for  $t = 1 \dots T$  do
2:   for  $k = 1 \dots N$  do
3:      $\epsilon_t \sim \mathcal{N}(0, \mathbb{1})$ 
4:      $\Delta \mathbf{s}_k^t \leftarrow \mathbf{s}^t + \eta \cdot \nabla \log p(\mathbf{s}^t) + 2\sqrt{\eta}\epsilon_t$ 
5:   end for
6:   for  $k = 1 \dots N$  do
7:      $\mathbf{s}_k^{t+1} \leftarrow \Delta \mathbf{s}_k^t - \frac{\eta}{\sigma^2} \cdot [\mathbf{m} - \frac{1}{N} \sum_i \mathbf{s}_i^t]$ 
8:   end for
9: end for

```

The idea of using SGLD in combination with deep parametrized priors for source separation was, concurrently to this work, introduced in Jayaram and Thickstun¹¹. The authors empirically show that

$$\mathbb{E} \left[\|\nabla_s \log p_\sigma(\mathbf{s})\|^2 \right] \propto 1/\sigma^2$$

They argue that this surprising proportionality is stemming from the severe non-smoothness of the prior model. If the prior model, in the extreme case, exhibits a Dirac delta peak as the probability mass, then the estimation of the gradient with added Gaussian noise will itself be proportional to the Gaussian. From there the authors argue, that the prior models have to be trained with noised samples, where the added noised is proportional to later used estimator noise.

Will add more discussion of their results

¹⁰ Welling and Teh, “Bayesian Learning via Stochastic Gradient Langevin Dynamics”, 2011.

$\log p(\mathbf{m})$ is independent from \mathbf{s}_k therefore no gradient.

Algorithm 1: The Langevin sampling procedure for source separation is fairly straight forward. For a fixed number of steps T we sample we take a step into the direction of the gradient under the priors and the gradient of the mixing constraint while adding Gaussian noise ϵ_t .

¹¹ Jayaram and Thickstun, “Source Separation with Deep Generative Priors”, 2020.

Modeling the posteriors

Instead of formulating a sampling regime for approaching s_k we can also use the prior models to variationally train a model to estimate the parameters of an approximate posterior $q_{\phi_k}(s_k|\mathbf{m})$. While estimating the true posterior $p(s_k|\mathbf{m})$ is intractable, we choose a certain parametrized distribution as an approximation of the posterior and optimize the parameters to align with the true one.

We follow the same steps as previously shown generally for the latent variable models. First, we introduce an approximate posterior $q_{\phi_k}(s_k|\mathbf{m})$ for each source channel. Next, we express the mix density as the expectation over those posteriors:

$$\log p(\mathbf{m}) = \mathbb{E}_{q_{\phi_k}(s_k|\mathbf{m})}^N [\log p(\mathbf{m})] \quad (18)$$

From here we introduce the approximate posteriors in the expectation as before, just now for N separate priors:

$$\mathbb{E}_{q_{\phi_k}(s_k|\mathbf{m})}^N [\log p(\mathbf{m})] = \mathbb{E}_{q_{\phi_k}(s_k|\mathbf{m})}^N \left[\log \frac{p(\mathbf{m}, s_1, \dots, s_N)}{p(s_1, \dots, s_N|\mathbf{m})} \right] \quad (19)$$

$$= \mathbb{E}_{q_{\phi_k}(s_k|\mathbf{m})}^N \left[\log \frac{p(\mathbf{m}|s_1, \dots, s_N) \cdot \prod_k^N p(s_k)}{\prod_k^N q_{\phi_k}(s_k|\mathbf{m})} + \log \frac{\prod_k^N q_{\phi_k}(s_k|\mathbf{m})}{p(s_1, \dots, s_N|\mathbf{m})} \right] \quad (20)$$

Note that again we use the assumption in Equation (4) to factorize the joint. This assumption is not being used for the independent approximate posteriors. While the true posterior certainly is the posterior of the joint source model, we choose our approximate posteriors to be independent. The expectation in (18) over those is still correct. The error arising from the independence assumption is captured in the tightness of the ELBO.

We arrive at the ELBO by leaving out the divergence of the approximate posterior to the true posterior:

$$\mathbb{E}_{q_{\phi_k}(s_k|\mathbf{m})}^N [\log p(\mathbf{m})] \geq \sum_k^N \mathbb{E}_{q_{\phi_k}(s_k|\mathbf{m})} \left[\log \frac{p(s_k)}{q_{\phi_k}(s_k|\mathbf{m})} \right] + \mathbb{E}_{q_{\phi_k}(s_k|\mathbf{m})} [p(\mathbf{m}|s_1, \dots, s_N)] \quad (21)$$

The lower bound, as in the normal VAE, will be our optimization target for training the inference models that give the variational parameters ϕ_k . We can also formulate the objective for just one source, as the Encoders are trained independently¹². Thus we come to:

¹² While the optimization of the Encoders is independent, the training is, of course, concurrent as the mixing condition depends on all N source samples.

$$\mathbb{E}_{q_{\phi_k}(s_k|\mathbf{m})} [\log p(\mathbf{m})] \geq \mathbb{E}_{q_{\phi_k}(s_k|\mathbf{m})} \left[\log \frac{p(s_k)}{q_{\phi_k}(s_k|\mathbf{m})} \right] + \mathbb{E}_{q_{\phi_k}(s_k|\mathbf{m})} [p(\mathbf{m}|s_1, \dots, s_N)] \quad (22)$$

$$= \mathbb{E}_{q_{\phi_k}(s_k|\mathbf{m})} \left[\log \frac{p(s_k)}{q_{\phi_k}(s_k|\mathbf{m})} \right] + \|\mathbf{m} - \frac{1}{N} \sum_j^N \mathbf{s}_j\| \quad (23)$$

The Kullback-Leibler divergence term is computationally expensive, therefore the divergence is the stochastic estimate using just one mini-batch.

As both the prior $p(s_k)$ and the approximate posterior $q_{\phi_k}(s_k|\mathbf{m})$ are parametrized by deep neural network, the gradient of Equation (23) is tractable with respect to the mini-batch.

Choosing a posterior

With the optimization regime set up, we have to choose a fitting posterior distribution and the deep neural network which computes its parameters from data. In many typical VAE settings, the prior density is a fixed probability distribution, typically a multivariate normal distribution and accordingly the approximate posterior also is set to be a normal distribution. As our prior is a greatly complex distribution estimated by the flow, we are not incentivized by this to choose the posterior accordingly. Further, we have the benefit of knowing our posteriors to model sound in the time-domain. Therefore we take a look at the distributions induced by our data samples, in the two cases of the Toy Data and the real musical recordings. For a more in-depth discussion of the datasets see ??.

In the case of the **Toy data**, the waveforms are regular, simple patterns. If we plot the histograms of an example wave for each source class (Figure 5), the probability distribution induced by those sources is again simply formulated. Both the *triangle* and *sawtooth* waves induce a uniform distribution, the *square* wave collapses on the only two values contained in the signal (the maxima) and the sine wave also gives a symmetrical distribution with the values around the maxima greatly more likely. We propose modelling these posteriors with a **Beta distribution** scaled onto the interval $[-1, 1]$.

To learn the variational parameters ϕ of the encoder network $\text{Encoder}_{\phi}(\mathbf{m}) = \eta$ which infers the parameters η of the approximate posterior we need the gradient of samples taken from the posterior with respect to the parameters. The Beta distribution is parametrized by the two concentrations $\eta = \{\alpha, \beta\}$. Therefore we need:

$$\frac{\partial}{\partial \alpha} z, \quad \frac{\partial}{\partial \beta} z \quad z \sim q_{\phi_k}(s_k|\mathbf{m}) \quad (24)$$

As described in ?? with choosing a Gaussian posterior, one can use the reparametrization trick to get low-variance gradients, simply and efficiently. In the case of using a Beta posterior, the reparametrization trick is not applicable. There exist multiple approaches to compute implicit gradients through approximate differentiation of different probability distribution. Figurnov et al.¹³ provides an new fast general gradient estimator. We employ this estimator in the form that it is implemented in the PyTorch framework.

In the case of the **real musical data** we can not construct any informed posterior distribution, as the signals do not follow any constrained prior form. For the lack of an informed decision and the computational ease of the reparametrization trick, we resort to a Gaussian posterior. More specifically we propose employing a **truncated normal distribution**. The latent sound values are restricted

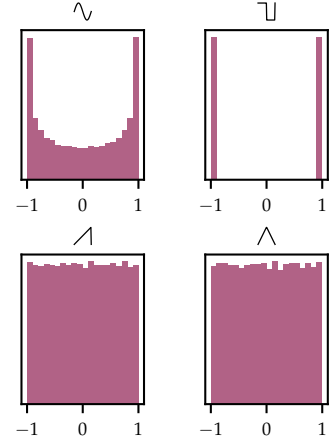


Figure 5: The logs of different classes of the different outputs

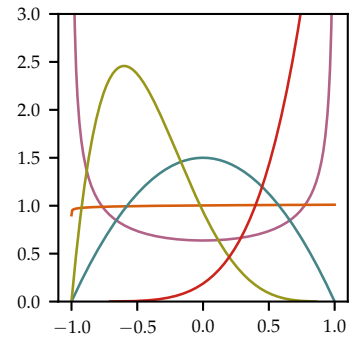


Figure 6: Beta distributions with different concentrations α and β

¹³ Figurnov et al., “Implicit Reparameterization Gradients”, 2019.

to $[-1, 1]$ and a with a truncated Gaussian we can restrict the support to this interval. The reparametrization trick is still applicable, for a concrete description of the sampling refer to ???. The normal distribution is parametrized by mean and variance $\eta = \{\sigma, \mu\}$.

Choosing an Encoder

The inference model $\text{Encoder}_\phi(\mathbf{m})$ infers the parameters for the chosen posterior distribution from data. The model takes the time-domain signals as an input, therefore we use the architectural design choices as with the prior networks and employ a non-causal WaveNet architecture. As laid out before this model is fully convolutional, not autoregressive, as the prediction of the parameters is a one-to-one mapping, where each time-point is mapped to a set of parameters. As it is practice we use the same network head for both parameters and append after the last skip-connection output of the WaveNet two convolutional layers with the final one mapping to the two separate parameters. The WaveNet used here is fairly similar to the modules used in the prior flows. They employ the gated convolution and have fixed kernel sizes of 3.

We note that we fix the encoder to work in the time-domain. We propose this to have access to the phase information of the original signals during the mapping $\mathbf{m} \rightarrow \mathbf{s}_k$. This does not make is necessary though that the prior models are also working in the time-domain. The samples taken from the approximate posterior during the training of the inference model can be transformed with the same spectral transformation used to train the prior models. We can compute the gradient of the likelihood from the prior under the spectral transformation. Because we are not using the generative priors to generate in this setting, the phase information lost through the transformation is not disadvantageous and still is accessible to the inference model.

As we have N different sources with each a trained prior network we will also train N inference networks each with a separate approximate posterior.

Difficulties of this method

We want to make the difficulties with this method clear and list them in the following explicitly:

1. *All sources known and trained with prior*
2. The separation capabilities in both approaches completely rely on the quality of the learned priors. The priors are learned unsupervised and unconditional from random cuts of the instruments stems. The learned manifolds are extremely complex and hard to train. There is little prior work on successfully building generative models for sound, without further conditional information.

3. We have to train a prior model for each source. In the sampling case, we concurrently sample from each of these models, in the case of the inference model we have to train an additional model for all sources. Because the separation happens under the mixing condition, at the time of training, inference or sampling all parts are used.
4. *Unstable sampling / training because of many objectives*

Bibliography

- Burgess, Christopher P., Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner (2018). “Understanding Disentangling in Beta-VAE”. In: arXiv: [1804.03599 \[cs, stat\]](#).
- Chorowski, Jan, Ron J. Weiss, Samy Bengio, and Aaron van den Oord (2019). “Unsupervised Speech Representation Learning Using WaveNet Autoencoders”. In: *IEEE/ACM Trans. Audio Speech Lang. Process.* 27.12, pp. 2041–2053. arXiv: [1901.08810](#).
- Cohen-Hadria, Alice, Axel Roebel, and Geoffroy Peeters (2019). “Improving Singing Voice Separation Using Deep U-Net and Wave-U-Net with Data Augmentation”. In: arXiv: [1903.01415 \[cs, eess\]](#).
- Dauphin, Yann N., Angela Fan, Michael Auli, and David Grangier (2017). “Language Modeling with Gated Convolutional Networks”. In: arXiv: [1612.08083 \[cs\]](#).
- Défossez, Alexandre, Nicolas Usunier, Léon Bottou, and Francis Bach (2019). “Demucs: Deep Extractor for Music Sources with Extra Unlabeled Data Remixed”. In: arXiv: [1909.01174 \[cs, eess, stat\]](#).
- Défossez, Alexandre, Neil Zeghidour, Nicolas Usunier, Leon Bottou, and Francis Bach (2018). “SING: Symbol-to-Instrument Neural Generator”. In: *Advances in Neural Information Processing Systems 31*. Ed. by S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett. Curran Associates, Inc., pp. 9041–9051.
- Dinh, Laurent, David Krueger, and Yoshua Bengio (2015). “NICE: Non-Linear Independent Components Estimation”. In: arXiv: [1410.8516 \[cs\]](#).
- Dinh, Laurent, Jascha Sohl-Dickstein, and Samy Bengio (2017). “Density Estimation Using Real NVP”. In: arXiv: [1605.08803 \[cs, stat\]](#) (cit. on p. 10).
- Figurnov, Michael, Shakir Mohamed, and Andriy Mnih (2019). “Implicit Reparameterization Gradients”. In: arXiv: [1805.08498 \[cs, stat\]](#) (cit. on p. 13).
- Grais, Emad M., Dominic Ward, and Mark D. Plumbley (2018). “Raw Multi-Channel Audio Source Separation Using Multi-Resolution Convolutional Auto-Encoders”. In: arXiv: [1803.00702 \[cs\]](#).
- Higgins, Irina, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner (2016). “Beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework”. In:
- Hochreiter, Sepp and Jürgen Schmidhuber (1997). “Long Short-Term Memory”. In: *Neural Computation* 9.8, pp. 1735–1780.
- ISO/TC 43 Acoustics (1975). *ISO 16 Acoustics Standard Tuning Frequency*.
- Jansson, Andreas, Eric J. Humphrey, Nicola Montecchio, Rachel M. Bittner, Aparna Kumar, and Tillman Weyde (2017). “Singing Voice Separation with Deep U-Net Convolutional Networks”. In: *ISMIR*.
- Jayaram, Vivek and John Thickstun (2020). “Source Separation with Deep Generative Priors”. In: arXiv: [2002.07942 \[cs, stat\]](#) (cit. on p. 11).

- Jordan, Michael I., Zoubin Ghahramani, Tommi S. Jaakkola, and Lawrence K. Saul (1999). "An Introduction to Variational Methods for Graphical Models". In: *Machine Learning* 37.2, pp. 183–233.
- Kalchbrenner, Nal, Erich Elsen, Karen Simonyan, Seb Noury, Norman Casagrande, Edward Lockhart, Florian Stimberg, Aäron van den Oord, Sander Dieleman, and Koray Kavukcuoglu (2018). "Efficient Neural Audio Synthesis". In: arXiv: [1802.08435 \[cs, eess\]](#).
- Kaspersen, Esbern Torgard (2019). "HydraNet: A Network For Singing Voice Separation". In:
- Kim, Sungwon, Sang-gil Lee, Jongyoon Song, Jaehyeon Kim, and Sungroh Yoon (2019). "FloWaveNet : A Generative Flow for Raw Audio". In: arXiv: [1811.02155 \[cs, eess\]](#) (cit. on pp. 7, 8).
- Kingma, Diederik P. and Jimmy Ba (2017). "Adam: A Method for Stochastic Optimization". In: arXiv: [1412.6980 \[cs\]](#) (cit. on p. 10).
- Kingma, Diederik P. and Prafulla Dhariwal (2018). "Glow: Generative Flow with Invertible 1x1 Convolutions". In: arXiv: [1807.03039 \[cs, stat\]](#) (cit. on p. 8).
- Kingma, Diederik P. and Max Welling (2014). "Auto-Encoding Variational Bayes". In: arXiv: [1312.6114 \[cs, stat\]](#).
- Kotelnikov, Vladimir (1933). "On the Carrying Capacity of the Ether and Wire in Telecommunications". In: *Material for the First All-Union Conference on Questions of Communication*. Moscow: Izd. Red. Upr. Svyazi RKKA.
- Lluís, Francesc, Jordi Pons, and Xavier Serra (2019). "End-to-End Music Source Separation: Is It Possible in the Waveform Domain?" In: arXiv: [1810.12187 \[cs, eess\]](#).
- Narayanaswamy, Vivek Sivaraman, Sameeksha Katoch, Jayaraman J. Thiagarajan, Huan Song, and Andreas Spanias (2019). "Audio Source Separation via Multi-Scale Learning with Dilated Dense U-Nets". In: arXiv: [1904.04161 \[cs, eess, stat\]](#).
- Native Instruments (n.d.). *Stem Audio Format*.
- Neal, Radford M. (2012). "MCMC Using Hamiltonian Dynamics". In: arXiv: [1206.1901 \[physics, stat\]](#).
- Paine, Tom Le, Pooya Khorrami, Shiyu Chang, Yang Zhang, Prajit Ramachandran, Mark A. Hasegawa-Johnson, and Thomas S. Huang (2016). "Fast Wavenet Generation Algorithm". In: arXiv: [1611.09482 \[cs\]](#).
- Prenger, Ryan, Rafael Valle, and Bryan Catanzaro (2018). "WaveGlow: A Flow-Based Generative Network for Speech Synthesis". In: arXiv: [1811.00002 \[cs, eess, stat\]](#) (cit. on p. 7).
- Pulse Code Modulation (PCM) of Voice Frequencies* (1972). G.711. ITU-T.
- Rafii, Zafar, Antoine Liutkus, Fabian-Robert Stöter, Stylianos Ioannis Mimitis, and Rachel Bittner (2017). *MUSDB18 - a Corpus for Music Separation*. Version 1.0.0.
- Rethage, Dario, Jordi Pons, and Xavier Serra (2018). "A Wavenet for Speech Denoising". In: arXiv: [1706.07162 \[cs\]](#).
- Rezende, Danilo Jimenez and Shakir Mohamed (2016). "Variational Inference with Normalizing Flows". In: arXiv: [1505.05770 \[cs, stat\]](#).
- Rezende, Danilo Jimenez, Shakir Mohamed, and Daan Wierstra (2014). "Stochastic Backpropagation and Approximate Inference in Deep Generative Models". In: arXiv: [1401.4082 \[cs, stat\]](#).
- Slizovskaia, Olga, Leo Kim, Gloria Haro, and Emilia Gomez (2019). "End-to-End Sound Source Separation Conditioned On Instrument Labels". Version 2. In: arXiv: [1811.01850 \[cs, eess\]](#).
- Stoller, Daniel, Sebastian Ewert, and Simon Dixon (2018). "Wave-U-Net: A Multi-Scale Neural Network for End-to-End Audio Source Separation" (Paris, France).

- Stöter, Fabian-Robert, Antoine Liutkus, and Nobutaka Ito (2018). “The 2018 Signal Separation Evaluation Campaign”. In: arXiv: [1804.06267 \[cs, eess\]](#).
- Szegedy, Christian, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich (2014). “Going Deeper with Convolutions”. In: arXiv: [1409.4842 \[cs\]](#).
- Tabak, Esteban and Cristina V. Turner (2013). “A family of nonparametric density estimation algorithms”. In: *COMMUN. PURE & APPL. MATHS.* 66.2, pp. 145–164.
- Van den Oord, Aäron, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu (2016). “WaveNet: A Generative Model for Raw Audio”. In: arXiv: [1609.03499 \[cs\]](#).
- Van den Oord, Aäron, Nal Kalchbrenner, Oriol Vinyals, Lasse Espeholt, Alex Graves, and Koray Kavukcuoglu (2016). “Conditional Image Generation with PixelCNN Decoders”. In: arXiv: [1606.05328 \[cs\]](#).
- Van den Oord, Aäron, Oriol Vinyals, and Koray Kavukcuoglu (2017). “Neural Discrete Representation Learning”. In: *Advances in Neural Information Processing Systems 30*. Ed. by I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett. Curran Associates, Inc., pp. 6306–6315.
- Welling, Max and Yee Whye Teh (2011). “Bayesian Learning via Stochastic Gradient Langevin Dynamics”. In: *Proceedings of the 28th International Conference on International Conference on Machine Learning*. ICML’11. Bellevue, Washington, USA: Omnipress, pp. 681–688 (cit. on pp. [6](#), [11](#)).