

GENERATIVE ADVERSARIAL SOURCE SEPARATION

Y.Cem Subakan^b, Paris Smaragdis^{b,‡}

^bUIUC, [‡]Adobe Systems
{subakan2, paris}@illinois.edu

ABSTRACT

Generative source separation methods such as non-negative matrix factorization (NMF) or auto-encoders, rely on the assumption of an output probability density. Generative Adversarial Networks (GANs) can learn data distributions without needing a parametric assumption on the output density. We show on a speech source separation experiment that, a multi-layer perceptron trained with a Wasserstein-GAN formulation outperforms NMF, auto-encoders trained with maximum likelihood, and variational auto-encoders in terms of source to distortion ratio.

Index Terms— Generative Adversarial Networks, Source Separation, Generative Models

1. INTRODUCTION

Many popular audio modeling/source separation algorithms such as Non-Negative Matrix Factorization (NMF) [1], autoencoders [2], or tensor factorization models [3] are all generative models and they are trained with maximum likelihood (ML) which require the specification of output distributions. For instance, NMF models with different loss functions (e.g. KL-NMF, Euclidean NMF, IS-NMF) [4], actually have the same underlying mapping from latent space to observed space (same underlying network), but their performances typically differ on a given dataset. The output distribution/loss function therefore biases the model.

Generative Adversarial Networks (GANs) [5] offer a generative model learning framework, which does not require the specification of an output distribution. GANs are able to learn the processes which are implicitly defined via a transformation of a random variable. Namely, the generative process is defined such that a random latent variable is mapped to the data domain via getting transformed through a deterministic neural network. This removes the bias that comes from assuming a parametric output distribution and leads to more accurate modeling of distributions. [6]

GANs have been very popular in computer vision since their first introduction [6]. However, to the best of our knowledge, usage of GANs in the audio modeling domain has been

limited. In [7], authors train de-noising networks by using an adversarial framework. In this paper, we propose using GANs to learn a generative model over magnitude spectrogram frames, which are used in a speech source separation task.

Source separation is the task where the goal is to decompose a given signal into additive components which approximates the original sources as accurately as possible. In generative source separation, we train generative models to recover the sources from an observed mixture. We experimentally show on speech mixtures that an adversarially trained two layer perceptron outperforms NMF and ML-trained autoencoders in terms of source-to-distortion ratio [8]. In our experiments, we have observed that the original GAN formulation in [5] is hard to train. We therefore showed the performance improvement over standard audio models with the more recent Wasserstein-GAN formulation [9].

2. GENERATIVE SUPERVISED SOURCE SEPARATION

In generative source separation, an observed mixture signal $x \in \mathbb{R}^L$ is assumed to follow the generative process below:

$$\begin{aligned} h_1 &\sim p_{\text{latent}}(h_1), \quad s_1|h_1 \sim p_{\text{forward}_1}(s_1|h_1) \\ h_2 &\sim p_{\text{latent}}(h_2), \quad s_2|h_2 \sim p_{\text{forward}_2}(s_2|h_2) \\ x|s_1, s_2 &\sim p_{\text{mixture}}(x|s_1 + s_2) \end{aligned}$$

where the source $s_k \in \mathbb{R}^L$, $k \in \{1, 2\}$, follow the distribution $p_{\text{source}_k}(s_k) = \int p_{\text{forward}}(s_k|h_k)p(h_k)dh_k$, where $h_k \in \mathbb{R}^K$ is a latent variable with lower dimensionality, such that $K < L$, and $p_{\text{forward}}(s_k|h_k)$ is the forward model for the sources. Given the sources, the mixture x is assumed to be distributed according to the conditional distribution $p_{\text{mixture}}(x|s_1 + s_2)$, where x is conditioned on the sum of the sources. Note that we have not yet assumed parametric forms for the distributions above. Also note that, in our experiments we consider the case where there are only two sources, although methods discussed can be generalized to more sources. In the context of our audio application, mixture x here corresponds to a column of a magnitude spectrogram.

The goal in source separation is to compute accurate estimates for the sources given a mixture signal x . In supervised

generative source separation, the approach is to first train the forward models $p_{\text{forward}_1}(\cdot)$, $p_{\text{forward}_2}(\cdot)$ such that the source distributions $p_{\text{source}_1}(\cdot)$, $p_{\text{source}_2}(\cdot)$ are approximated as accurately as possible. Given the trained models for both sources, in testing we compute source estimates \hat{s}_1 , \hat{s}_2 such that the conditional likelihood $p_{\text{mixture}}(x|\hat{s}_1 + \hat{s}_2)$ is maximized (or equivalently the reconstruction error for the mixture is minimized). In the next section, we describe the specifics on how to go through supervised source separation with maximum likelihood training.

2.1. Maximum Likelihood Training for Sources

A common way to go about approximating the source distributions is through assuming that the sources are generated by transforming a K dimensional latent variable $h \sim p_{\text{latent}}(h)$, through a non-linear mapping (such as a neural network) $f_{\theta}(h)$ with parameters θ , and adding noise to the transformed variable. This corresponds to the following generative model:

$$h \sim p_{\text{latent}}(h), s|h \sim p_{\text{out}}(s; f_{\theta}(h)), \quad (1)$$

where $p_{\text{out}}(\cdot)$ is the output distribution which models the noise at the output of the mapping $f_{\theta}(h)$. E.g. When modeling spectrograms, $p_{\text{out}}(\cdot)$ is usually taken as Poisson distribution, which corresponds to the unnormalized KL divergence. Under these modeling assumptions, the optimization problem for approximating source distribution is written as follows:

$$\max_{\theta} \sum_t \log \int \mathcal{PO}(s^t; f_{\theta}(h)) p(h) dh, \quad (2)$$

where the integral over the latent variable h is intractable in the general case. Using variational auto-encoder framework in [10], the objective in expression (2) can be maximized by computing a variational lower bound.

In practice however, especially in audio modeling, the integral over the latent variable h is not computed, and only the conditional forward model $p_{\text{forward}}(s|h)$ is learnt, by simultaneously optimizing over the forward model parameters and the latent variables. This is written as the following optimization problem:

$$\max_{\theta, h} \sum_t \log \mathcal{PO}(s^t; f_{\theta}(h)), \quad (3)$$

If $f_{\theta}(h) = Wh$, where $W, h \geq 0$, then this formulation corresponds to the widely used Non-Negative Matrix Factorization (NMF) model [1, 11, 4]. It is also possible to include the latent variable estimation part in the model with an auto-encoder. This results in the following optimization problem:

$$\max_{\theta, \theta^{\text{enc}}} \sum_t \log \mathcal{PO}(s^t; f_{\theta}(f_{\theta^{\text{enc}}}^{\text{enc}}(s^t))), \quad (4)$$

where $f_{\theta^{\text{enc}}}^{\text{enc}}(\cdot) : \mathbb{R}^L \rightarrow \mathbb{R}^K$, is the encoder, and $f_{\theta}(\cdot) : \mathbb{R}^K \rightarrow \mathbb{R}^L$ is decoder part. In [2], $f_{\theta}(h) = \log(\exp(W_1 h) + 1)$, and $f_{\theta^{\text{enc}}}^{\text{enc}}(s) = \log(\exp(W_2 s) + 1)$ is used.

The conceptual problem with the training objectives discussed in this section is that by picking a specific output distribution, we are sacrificing from the generality of the approximated source distributions. To remove this assumption, in this paper we use generative adversarial networks, which is a neural network framework for learning generative models without explicitly specifying an output distribution when training the generator network.

2.2. Adversarial Training for Sources

We have seen in the previous section that maximum likelihood training involves a parametric assumption for the output distribution. As an alternative, in this paper we propose using an implicit generative model in training, which does not require an explicit loss function. An implicit generative model for the sources can be specified as follows:

$$h \sim p_{\text{latent}}(h), s = f_{\theta}(h), \quad (5)$$

where the source s is deterministically related to the latent variable h , unlike the source model in the previous section. This process implies an intractable density function $p_{\text{model}}(\cdot)$ for s in the general case where $f_{\theta}(h)$ is a complicated non-linear mapping such as a neural network. Learning under implicit generative models is a currently a very active field of research [12]. One way to attack this problem is to use discriminator function $D_{\xi}(\cdot)$ which aims to distinguish between the samples generated from the model and the training instances. The goal in training is then becomes to generate samples using the process in expression (5) so that, the discriminator $D_{\xi}(\cdot)$ becomes unable to distinguish between the generated samples and the training data. This described setup is known as a generative adversarial network (GAN) [5], and the corresponding minimax game is specified as follows:

$$\min_{\theta} \max_{\xi} \mathbb{E}_s \log D_{\xi}(s) + \mathbb{E}_h \log(1 - D_{\xi}(f_{\theta}(h))), \quad (6)$$

This expression can be recognized as the sum of Bernoulli log-likelihoods, where $D_{\xi}(\cdot)$ tries to maximize by outputting 1 for the training data s , and outputting 0 for the generated samples $f_{\theta}(h)$. The generator however tries to minimize the expression by *fooling* the discriminator. It can be shown that under some assumptions this scheme minimizes the Jensen-Shannon divergence between the actual source distribution $p_{\text{source}}(\cdot)$ and the model distribution $p_{\text{model}}(\cdot)$. However in practice, this scheme is unstable, and usually suffers from the mode collapse problem where the learnt distribution $p_{\text{model}}(\cdot)$ only captures a subset of the actual sample space. [6]. This is unfortunately not acceptable for our source separation application.

An alternate formulation known as the Wasserstein-GAN, alleviates the mode collapse problem by minimizing the Wasserstein-1 distance between the learnt and data distributions, which results in smooth gradients [9]. Authors show

that this can be achieved with following minimax game:

$$\min_{\theta} \max_{\xi \in \mathcal{W}} \mathbb{E}_s D_{\xi}(s) - \mathbb{E}_h D_{\xi}(f_{\theta}(h)), \quad (7)$$

where \mathcal{W} denotes the set for parameters ξ for which $D_{\xi}(\cdot)$ will be γ -Lipschitz continuous, for some γ . In the algorithm provided in the paper, this constraint is achieved by clipping the weights ξ . In our experiments, Wasserstein GANs showed significant improvement over the original GAN formulation. Note that $D_{\xi}(\cdot)$ is referred to as critic in this formulation.

2.3. Testing

After training the forward models for each source, given a test mixture, the estimates \hat{s}_1, \hat{s}_2 for the sources is obtained by minimizing the reconstruction error via finding the optimal latent variables as inputs to the forward mappings:

$$\hat{h}_1, \hat{h}_2 = \arg \max_{h_1, h_2} \log p_{\text{mixture}}(x; f_{\hat{\theta}_1}(h_1) + f_{\hat{\theta}_2}(h_2)), \quad (8)$$

and then we get the estimates for the sources by setting $\hat{s}_1, \hat{s}_2 = f_{\hat{\theta}_1}(\hat{h}_1), f_{\hat{\theta}_2}(\hat{h}_2)$, where $\hat{\theta}_1, \hat{\theta}_2$ denote the trained network parameters.

An extra benefit we get by training our generative models with GANs is that, in addition to the generator networks $f_{\hat{\theta}_k}(\cdot)$, we also get discriminators/critics $D_{\hat{\xi}_k}(\cdot)$. We can therefore use them in the separation stage to score how much the obtained source looks like the instances in training set. We also noticed that using a smoothing term to enforce smooth first difference across time improves the quality of the estimated sources for both GANs and Maximum likelihood based auto-encoders. Therefore, the optimization for separating the sources, given T mixture spectrogram columns $x^{1:T}$, becomes the following:

$$\begin{aligned} \max_{h_1^{1:T}, h_2^{1:T}} \frac{1}{T} \sum_{t=1}^T \log p_{\text{mixture}}(x^t; f_{\hat{\theta}_1}(h_1^t) + f_{\hat{\theta}_2}(h_2^t)) \\ + \frac{\alpha}{T} \sum_{t=1}^T \left(D_{\hat{\xi}_1}(f_{\hat{\theta}_1}(h_1^t)) + D_{\hat{\xi}_2}(f_{\hat{\theta}_2}(h_2^t)) \right) \\ + \frac{\beta}{T-1} \sum_{t=1}^{T-1} \left(\sum_{k=1}^2 \|f_{\hat{\theta}_k}(h_k^{t+1}) - f_{\hat{\theta}_k}(h_k^t)\|_1 \right), \end{aligned} \quad (9)$$

where α is a trade-off scalar between the reconstruction quality and discriminator/critic score. In our experiments we fixed $\alpha = 0.1$, but it can also potentially be optimized on a validation set. For the smoothing term, we used $\beta = 0.1$. Finally, note that for magnitude spectrograms it is very common to use Poisson distribution for $p_{\text{mixture}}(\cdot)$.

3. EMPIRICAL RESULTS

To show the validity of using GANs in source separation, we compare adversarially trained networks with auto-encoders

trained with maximum likelihood, variational auto-encoders and NMF.

The experiment set-up is as follows: We form mixtures of male and female speaker utterances and corresponding training data from the TIMIT speech corpus [13]. To form the training/test data pairs, we randomly pick male and female speaker pairs from the *train* folder of the TIMIT corpus. Each speaker has 10 available utterances. For both speakers, of the 10 available utterances, we use 9 for training and 1 for testing. The resulting training set for each source is around 30 seconds long. The selected test utterances are around 3 seconds, and the mixture signal is obtained by mixing the test utterances at 0 dB. We form 25 such mixtures/training sets and test each algorithm on these randomly selected sets (The speaker pairs are the same across algorithms). As the preprocessing step, we compute Fourier spectrograms of the utterances. We use 1024 point FFT, and a hop size of 256 samples. The learning and source separation are performed on the columns (Fourier magnitude vectors for each time window) of the magnitude spectrograms. When reconstructing the separated sources in the time domain, we use the Wiener filtering equation:

$$\hat{s}_k^{\text{time}} = \text{ISTFT}\left(\frac{\hat{s}_k}{\hat{s}_1 + \hat{s}_2} \odot x \odot x^{\text{phase}}\right), \text{ for } k \in \{1, 2\}, \quad (10)$$

where x and x^{phase} are respectively the magnitude and phase spectrograms of the mixture. The magnitude spectra for the estimated sources are denoted by \hat{s}_1 and \hat{s}_2 . The estimated time domain signal is denoted by \hat{s}_k^{time} . The division and the multiplication \odot are both element wise, and $\text{ISTFT}(\cdot)$ designates the inverse short time Fourier transform operation to get the time domain signal from a complex Fourier spectrogram. We obtain results with the following models:

- KL-NMF model.
- The auto-encoder model suggested in [2], trained with maximum likelihood using a Poisson likelihood (equivalently unnormalized KL divergence).
- Standard GAN with Gaussian random inputs.
- Wasserstein GAN with Gaussian random inputs.
- Autoencoding Wasserstein GAN, where instead of Gaussian random inputs, we feed the training samples to the generator network.
- Variational Autoencoder with Gaussian prior on the latent variable, as in [10], and Poisson likelihood at the output.

For all GANs, we used the following architecture for the generator:

$$f_{\theta}(h) = \text{SP}(W_2 \text{SP}(W_1 h)), \quad (11)$$

where $\text{SP}(\cdot)$ is the soft-plus nonlinearity, such that $\text{SP}(x) = \log(\exp(x) + 1)$. Note that we have omitted the bias terms from the equation to reduce clutter. For all GANs with Gaussian random inputs, we used 513 dimensional inputs h (This

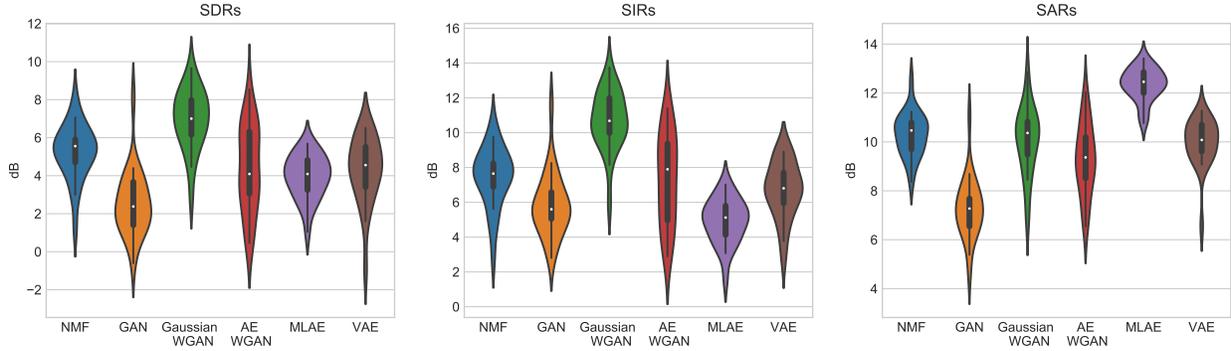


Fig. 1. The distributions of the BSS eval scores for our speech source separation experiment. Acronyms for each algorithm are indicated below each violin plot. In order from left to right the algorithms are ordered as: NMF, Standard GAN, Wasserstein GAN with Gaussian inputs, Autoencoding Wasserstein GAN, Autoencoder trained with maximum likelihood, and Variational Autoencoder. Each violin shows the distribution of the corresponding score in dB over different speaker pairs. The subplots are organized as follows: **Left:** SDR scores, **Middle:** SIR scores, **Right:** SAR scores.

is the dimensionality of the data items since we use 1024 point fft), and 100 hidden units. Therefore W_1 was of size 100×513 , and W_2 was of size 513×100 . For the auto-encoding GAN, and the auto-encoder trained with maximum likelihood, the network architecture of generator/forward model are exactly the same, except that the inputs are the data items s , instead of random variable h . For the VAE, we used the encoder $f_{\theta}^{\text{enc}}(s) = W_2 \text{ReLU}(W_1 s)$, both for the mean and the variance terms of the latent variable, where W_1 was of size 100×513 , and W_2 was of size 20×100 . For the decoder of the VAE, we used $f_{\theta}(h) = \text{SP}(W_3 h)$, where W_3 was of size 513×20 .

For the discriminator/critic networks of GANs, we use the following architecture:

$$D_{\xi}(s) = \sigma(V_2 \tanh(V_1 s)), \quad (12)$$

where V_1 is of size 90×513 (we use 90 hidden units for the discriminator), and V_2 is of size 1×90 . In standard GANs, $\sigma(\cdot)$ is the sigmoid function. In Wasserstein GAN, we do not use a non-linearity at the end of the network, and therefore $\sigma(\cdot)$ is the identity function. This gives smoother gradients.

In training and testing, for all neural network models we use the RMSprop algorithm [14] with a learning rate of 0.001. During the training of GANs, we do 5 iterations of discriminator/critic updates per generator update. For all neural network models, we do 4000 training iterations, and 20000 test iterations. For Wasserstein-GAN we clip the critic parameters at -0.01 and 0.01 for lower and upper limits respectively.

We report the BSS-eval [8] scores obtained after recovering the sources from the mixture signals. The BSS-eval scores are Source to Distortion Ratio (SDR), Source to Interference Ratio (SIR), and Source to Artifacts Ratio (SAR), where SDR being the summary measure on how good the separation is. For each speaker pair, we have averaged the BSS-eval scores of the recovered sources, and in Figure 1, with violin plots

we show the distribution of the averages of the two BSS-eval scores over all speaker pairs.

Experiments indicate that, The Wasserstein GAN with a Gaussian noise input outperforms NMF, ML auto-encoder and Variational auto-encoder in terms of source to distortion ratio. Note that we are obtaining these results with very similar underlying networks. For all models except VAE, we have kept the exact generator architecture defined in Equation (11). We have also observed that the standard GAN formulation is not very reliable. Although occasionally we have seen good SDRs with it, we have observed through inspecting its outputs that it is not able to capture the variety in the source spectrogram distribution as good as the Wasserstein GAN, and therefore the source separation performance of the standard GAN is not as good. We have also experimented with training an auto-encoder with adversarial training, and have seen that although it is less reliable than the Wasserstein GAN with Gaussian inputs, it is sometimes able to give great SDRs. In general, adversarial methods give great SIRs, by losing a bit from SAR, especially compared to the ML-autoencoder. Finally, note that the code for our experiments is available at https://github.com/ycemsubakan/source_separation_misc.

4. CONCLUSIONS

In this paper we have experimentally shown that Wasserstein GANs can obtain good performance in generative source separation. In addition to not requiring the specification of an output distribution, GANs fit into the source separation task nicely since the discriminator/critic functions help in source separation. We believe that there exists many research opportunities to use GANs in the audio domain. One natural next step from this paper is to extend the results showed in this paper with an end-to-end generative adversarial audio model.

5. REFERENCES

- [1] Paris Smaragdis and Judith C. Brown, “Non-negative matrix factorization for polyphonic music transcription,” in *In IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, 2003, pp. 177–180.
- [2] Paris Smaragdis and Shrikant Venkataramani, “A neural network alternative to non-negative audio models,” in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2017, New Orleans, LA, USA, March 5-9, 2017*, 2017, pp. 86–90.
- [3] Ali Taylan Cemgil, Umut Simsekli, and Yusuf Cem Sübakan, “Probabilistic latent tensor factorization framework for audio modeling,” in *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics, WASPAA 2011, New Paltz, NY, USA, October 16-19, 2011*, 2011, pp. 137–140.
- [4] Cédric Févotte and Jérôme Idier, “Algorithms for non-negative matrix factorization with the beta-divergence,” *CoRR*, vol. abs/1010.1763, 2010.
- [5] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio, “Generative adversarial nets,” in *NIPS*, 2014.
- [6] Ian J. Goodfellow, “NIPS 2016 tutorial: Generative adversarial networks,” *CoRR*, vol. abs/1701.00160, 2017.
- [7] Santiago Pascual, Antonio Bonafonte, and Joan Serra, “SEGAN: speech enhancement generative adversarial network,” *CoRR*, vol. abs/1703.09452, 2017.
- [8] Emmanuel Vincent, Rmi Gribonval, and Cdric Fvotte, “Performance measurement in blind audio source separation,” *IEEE Trans. Audio, Speech and Language Processing*, vol. 14, no. 4, pp. 1462–1469, 2006.
- [9] Martín Arjovsky, Soumith Chintala, and Léon Bottou, “Wasserstein GAN,” *CoRR*, vol. abs/1701.07875, 2017.
- [10] Diederik P. Kingma and Max Welling, “Auto-encoding variational bayes.,” *CoRR*, vol. abs/1312.6114, 2013.
- [11] Paris Smaragdis, Cédric Févotte, Gautham J. Mysore, Nasser Mohammadiha, and Matthew D. Hoffman, “Static and dynamic source separation using nonnegative factorizations: A unified view,” *IEEE Signal Process. Mag.*, vol. 31, no. 3, pp. 66–75, 2014.
- [12] S. Mohamed and B. Lakshminarayanan, “Learning in Implicit Generative Models,” *ArXiv e-prints*, Oct. 2016.
- [13] J. S. Garofolo, L. F. Lamel, W. M. Fisher, J. G. Fiscus, D. S. Pallett, and N. L. Dahlgren, “TIMIT acoustic phonetic continuous speech corpus,” 1993.
- [14] T. Tieleman and G Hinton, “Root mean square propagation (rmsprop),” Lecture 6.5 - rmsprop, COURSERA: Neural Networks for Machine Learning, Accessed: 2017-April-12.