
Interesting properties of GAN samples

Anonymous Author(s)

Affiliation

Address

email

Abstract

In this paper we investigate numerical properties of samples produced with adversarial methods, specially Generative Adversarial Networks. We analyze summary statistics of real and fake data and compute distances based on the marginal distribution of perceptually significant features. We provide results on image, music and speech data and show that GAN generated samples have interesting signatures that can be used to detect adversarial attacks.

1 Introduction

Since the groundbreaking Generative Adversarial Networks paper [5] in 2014, GAN related publications use a grid of natural images to accompany theoretical and empirical results. Early GAN research focused on natural images and is expanding to other domains including language models [7] and music [15].

Unlike variational auto encoders and other methods [5], most of the evaluation of the output of Generators trained with the GAN framework is qualitative: authors normally list higher sample quality as one of the advantages of their method over other methods. Interestingly, unlike other optimization problems where analysis of the empirical risk is a strong indicator of progress, in GANs decrease in loss is not always correlated with increase in image quality [3], and thus authors still rely on visual inspection of generated images.

Based on visual inspection, authors confirm that they have not observed mode collapse or that their framework is robust to mode collapse if some criteria is met ([3],[7]). In practice, github issues where practitioners report mode collapse or not enough variety abound.

Verifiable Artificial Intelligence[14], more specifically verifying GAN samples is hard because it depends on the existence of perceptually meaningful features. Let's consider the generation of images of mammals: although it is possible to perform computation on color histograms, to compare fake and real samples, we do not yet have robust algorithms able to verify if an image follows specifications derived from anatomy.

This paper is related to this effort and focuses on understanding how GAN generators approximate modes in the real distribution and verifying if the generated samples violate specifications derived from the real distribution. We quantitatively evaluate GAN generated samples by marginalizing perceptually meaningful features and computing the distance between the distribution of these features in the real and the fake¹ data. We analyze real and fake data and offer the following contributions in this paper:

- We show that GAN samples have universal signatures.
- We show how GAN samples approximate modes of the real distribution.
- We show significant differences between the marginal distribution of features.

¹Data sampled from the generator

2 Related work

Despite its youth, several publications ([2],[13], [16], [12]) have investigated the use of the GAN framework for generation of samples and unsupervised feature learning. Following the procedure described in [1] and used in [5], earlier GAN papers evaluate the quality of the generator by fitting a Gaussian Parzen window² to the GAN samples and reporting the log-likelihood of the test set data under this distribution. It is known that this method has some drawbacks, including its high variance and bad performance in high dimensional spaces.

In their brilliant publications, [9], [3] and [7] propose alternative objective functions and algorithms that circumvent problems that are common when using the original GAN objective. The problems addressed include instability of learning, mode collapse and meaningful learning curves.

These alternatives do not eliminate the need or excitement³ of visually inspecting GAN samples during training. In [4], the authors propose a solution to the diversity problem by introducing a new hyper-parameter γ with a loss derived from the Wasserstein distance.

In the next section we describe the hypotheses evaluated in this paper.

3 Hypotheses

Hypothesis 1 (H1): *Generative models can approximate the distribution of real data and hallucinate fake data that has some variety and resembles the real data*

Although this hypothesis is trivial for experiments that have already been conducted, it is the first condition for our experiments with music and speech data. To our knowledge there are no publications where GANs are successful in hallucinating polyphonic music and speech data. During our experiments we prove that this hypothesis is true.

Hypothesis 2 (H2): *The real data has useful properties that can be extracted computationally.*

By useful we refer to properties that are closely related to the real data itself. For example, computing the distribution MNIST pixel values might be not useful for assessing drawing quality. However, it might be useful to evaluate if a random MNIST samples is real or fake data.

Hypothesis 3 (H3): *The fake data has properties that are hardly noticed with non-computational inspection.*

Visual inspection of generated samples has become the norm for the evaluation of samples generated using the GAN framework. We investigate if there are properties common to all GAN samples or properties that significantly differ between the real data and the fake data. This hypothesis supports the next hypothesis related to adversarial attacks.

Hypothesis 4 (H4): *The difference in properties can be used to identify the source (real or fake)*

The development of generative models for digital media announce the imminent rise of adversarial attacks. We investigate if these differences can be used to detect if the data was generated with the GAN framework or is an adversarial attack.

We call the reader's attention that approximating the distribution over features computed on the real data does not guarantee that the real data is being approximated. Formally speaking: consider $X \sim Z$, i.e. X distributed as Z , and $f(X) \sim W$, where $f : X \mapsto Y$. If $A \sim B$ and B approximates Z , then $f(A) \sim D$ must also approximate W . However, a distribution that approximates W is not guaranteed to approximate Z .

4 Method

In this section we describe our analysis method in detail, including briefly describing the datasets and features computed, as well as distance or divergence measures.

²Kernel Density Estimation

³Despite of authors promising on twitter to never touch GANs again.

78 4.1 Datasets

79 In our experiments we use the MNIST dataset, a MIDI dataset of 389 Bach Chorales downloaded
80 from the web and a subsample of the NIST 2004 telephone conversational speech dataset with 100
81 speakers, multiple languages and on average of 5 minutes per speaker.

82 4.2 Property extraction

83 The properties extracted from the datasets used on this paper can be perceptually meaningful or not.
84 We claim that both properties can be used to numerically identify the source of the sample. In the
85 context of this paper, samples are images of fixed size. Consider the single channel image I with
86 dimensions R by C , where $I_{r,c}$ is the pixel intensity of the pixel at row r and column c

87 4.2.1 Spectral Moments

88 The spectral centroid [11] is a feature commonly used in the audio domain, where it represents the
89 barycenter of the spectrum. This feature can be applied to other domains and we invite the reader
90 to visualize Figure 1 for examples on MNIST and Mel-Spectrograms [11]. For each column in an
91 image, we transform the pixel values into row probabilities by normalizing them by the column sum,
92 after which we take the expected row value.

93 Figure 1a shows the spectral centroid computed on sample of MNIST training data.

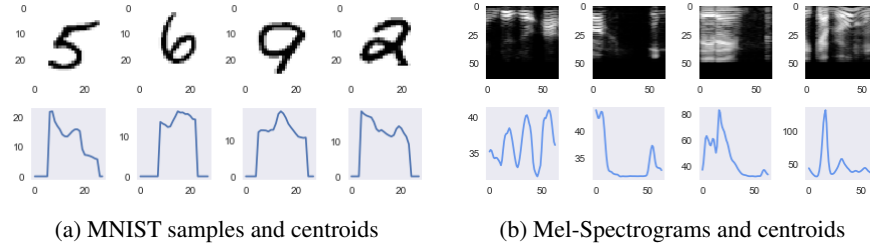


Figure 1: Spectral centroids on digits and Mel-Spectrograms

94 4.2.2 Spectral Slope

95 The spectral slope is computed by applying linear regression using a overlapping sliding window of
96 size 7. Figure 2 shows these features computed on MNIST and Mel-Spectrograms.

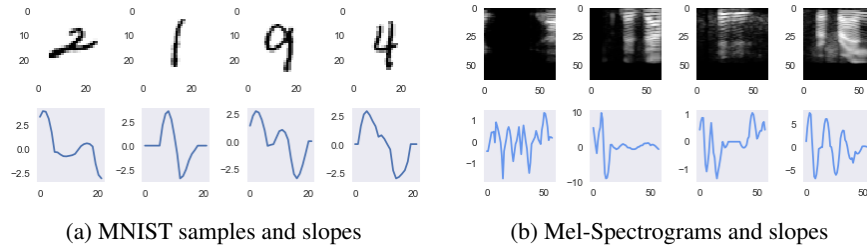


Figure 2: Spectral slopes on digits and Mel-Spectrograms

97 4.3 Generative Models

98 We investigate samples produced with the DCGAN architecture using the Least-Squares GAN
99 (LSGAN) [9] and the improved Wasserstein GAN (IWGAN) [7]. We also compare adversarial
100 MNIST samples produced with the fast gradient sign method (FGSM) [6].

5 Experiments

5.1 MNIST

We compare the distribution of features computed over the MNIST training set to the distribution of the same features computed over other datasets, including the MNIST test set, samples generated with GANs and adversarial samples computed using the FGSM. The training data is scaled to $[0, 1]$ and the random baseline is sampled from a Bernoulli distribution with probability equal to the normalized mean value of pixel intensities in the MNIST training data, 0.13. Each GAN model is trained until the loss plateaus and the generated samples look similar to the real samples.

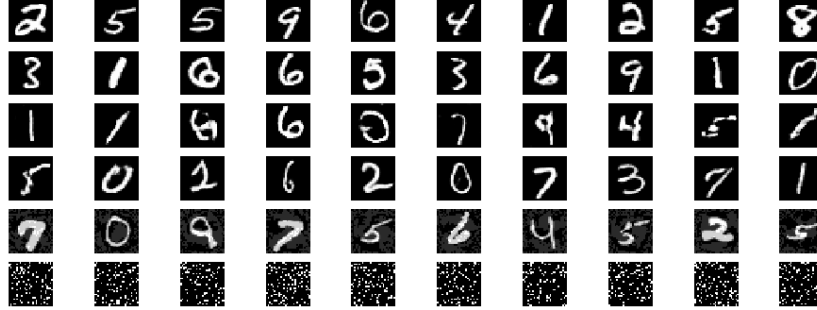


Figure 3: Samples drawn from MNIST train, test, LSGAN, IWGAN, FGSM and bernoulli respectively.

Visual inspection of these 10 generated samples in Figure 3 show that IWGAN seems to produce better samples than LSGAN. Below we use the training set as a reference and compare the distribution of pixel intensities. Table 1 reveals that although samples generated with LSGAN and IWGAN look similar to the training set, they are considerably different given the Kolmogorov-Smirnoff (KS) Two Sample Test and the Jensen Shannon-Divergence (JSD), specially if compared to the same statistics on the test data.

	KS Two Sample Test		JSD
	Statistic	P-Value	
mnist_train	0.0	1.0	0.0
mnist_test	0.003177	0.0	0.000029
mnist_lsgan	0.808119	0.0	0.013517
mnist_iwgan	0.701573	0.0	0.014662
mnist_adversarial	0.419338	0.0	0.581769
mnist_bernoulli	0.130855	0.0	0.0785009

Table 1: KS Two Sample Test and JSD over the distribution of pixel values for different samples

This phenomena can be understood by investigation of the empirical CDFs in Figure 4. The pixel values of the samples generated with the GAN framework are mainly bimodal and asymptotically approaching the modes of the distribution, 0 and 1. Such behavior will be present in any gradient descent method using an asymptotically converging non-linearity, such as sigmoid and tanh, immediately preceding the output of the generating function.

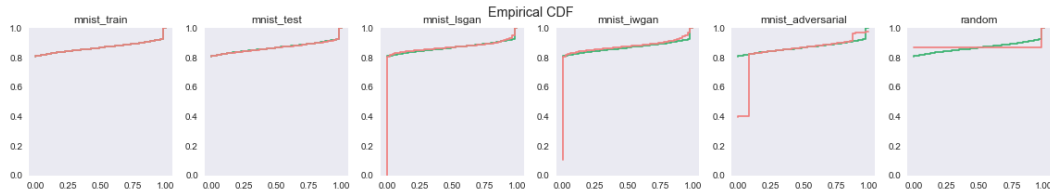


Figure 4: Pixel empirical CDF of training data as reference in green and other datasets

120 In addition, Figure 5 shows that the GAN generated samples smoothly approximate the modes of
 121 the distribution. This smooth approximation is considerably different from the training and test sets.
 122 These properties could be used to identify the source of the data, hence confirming hypotheses 2, 3
 123 and 4.

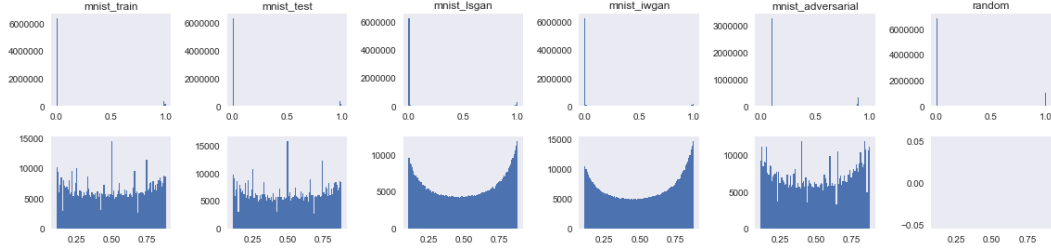


Figure 5: Histogram of pixel intensities for each dataset. First row shows histogram within the (0, 1) interval and 100 bins. Second row shows histogram between the (0.11, and 0.88) interval and 100 bins. Note the smooth curvature in GAN samples.

124 5.2 Bach Chorales

125 We investigate the properties of Bach Chorales generated with the GAN framework and verify if they
 126 satisfy the musical specifications of Bach chorales. Bach Chorales are polyphonic pieces of music,
 127 normally written for 4 or 5 voices, that follow a set of specifications/rules⁴. A global specification
 128 could enforce that durations be chosen from a set of valid durations; a local specification could assert
 129 that only certain transitions between states (notes) are valid. For this experiment, we convert the
 130 dataset of Bach chorales to piano rolls. The piano roll is a representation in which the rows represent
 131 note numbers, the columns represent time steps and the cell values represent note intensities. We
 132 compare the distribution of features computed over the training set, test set, gan generated samples
 133 and a random baseline sampled from a Bernoulli distribution with probability equal to the normalized
 134 mean value of intensities in the training data. After scaling, the intensities in the training and test data
 135 are strictly bimodal and equal to 1 or 0. Figure 6 below shows training, test, IWGAN and Bernoulli
 136 samples.

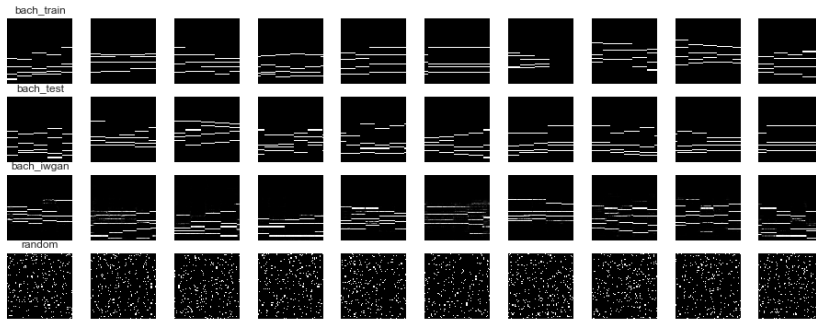


Figure 6: Samples drawn from Bach Chorales train, test, IWGAN, and Bernoulli respectively.

137 Figure 7 shows a behavior similar to our previous MNIST experiments: the IWGAN asymptotically
 138 approximates the modes of the distribution of the distribution of intensity values. In the interest of
 139 space, we refer the reader to the appendix for statistics and other relevant information.

140 Following, we investigate if the samples violate the specifications of Bach chorales. For doing so, we
 141 first convert the all data to boolean by thresholding at 0.5 such that values above the threshold are set
 142 to 1. We then compute boolean Chroma [11] features from the boolean piano rolls and compute a
 143 empirical transition matrix, where the positive entries represent existing transitions. The transition
 144 matrix built on the training data is taken as the reference specification, i.e. anything that is not

⁴specifications define the characteristics of the style

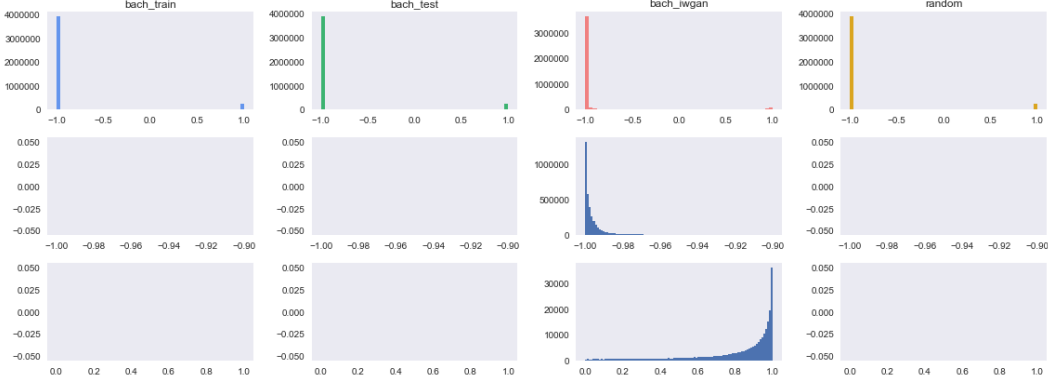


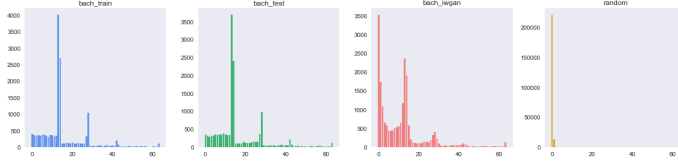
Figure 7

145 included is a violation of the specification. Table 2 shows the number of violations given each dataset.
 146 Compared to the test sample that has only 429 violations, the improved WGAN samples has over
 147 5000 violations, 10 times more than the test set! We use these facts to confirm hypotheses 2, 3 and 4.

	bach_train	bach_test	bach_iwgan	bach_bernoulli
Number of Violations	0	429	5029	58284

Table 2: Number of specification violation with respect to training data as reference.

148 In addition to experiments with Chroma features, we computed the distribution of note durations on
 149 the boolean piano roll described above. Figure 8a shows the distribution of note durations within
 150 each dataset. The train and test data are approximately bimodal and, again, the improved WGAN
 151 smoothly approximates the dominating modes of the distribution. Table 8b provides a numerical
 152 comparisson between datasets.



(a) Duration distribution

	KS Two Sample Test		JSD
	Statistic	P-Value	
train	0.0	1.0	0.0
test	0.09375	0.929	0.002
iwgan	0.21875	0.080	0.084
bernoulli	0.93750	0.0	0.604

(b) Test statistics

153 5.3 Speech

154 Within the speech domain, we investigate dynamic compressed Mel-Spectrogram samples produced
 155 with GANs trained on a subset of the NIST 2004 dataset, with 100 speakers. The generated samples
 156 can be seen in Figure 9, thus confirming hypothesis 1. We obtain the Mel-Spectrogram is obtained
 157 by projecting a spectrogram onto a mel scale, which we do by using the python library librosa [10]
 158 to project the spectrogram onto 64 mel bands, with window size equal to 1024 samples and hop
 159 size equal to 160 samples, i.e. frames of 100ms long. Dynamic range compression is computed as
 160 described in [8], with $\log(1 + C * M)$, where C is the compression constant scalar set to 1000 and
 161 M is the matrix representing the Mel-Spectrogram. The GAN models are trained using DCGAN
 162 with the improved Wasserstein GAN.

163 In Figure 10a we show the empirical CDFs of intensity values. Unlike our previous experiments
 164 where intensity (Bach Chorales) or pixel value (MNIST) was linear, in this experiment intensities are
 165 compressed using the log function. This considerably reduces the distance between the empirical
 166 CDFs of the training data and GAN samples, specially around the saturating points of the non-
 167 linearity, -1 and 1 in this case. In Table 10b we show numerical analysis of the differences and
 168 confirm hypotheses 2 and 3.

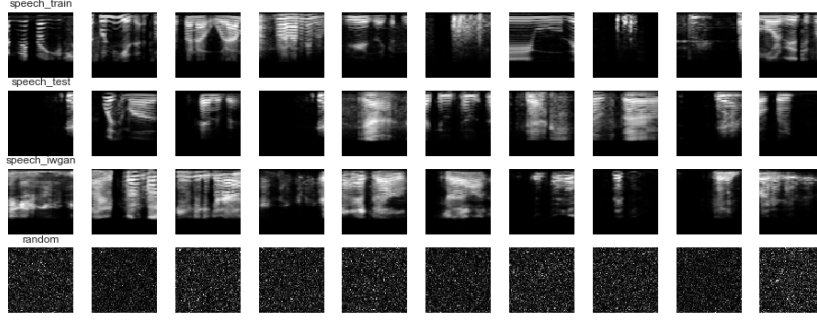
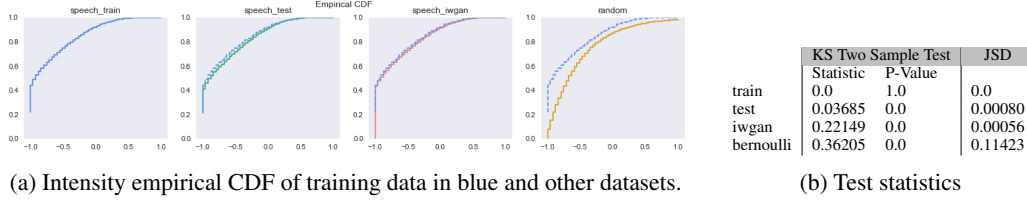


Figure 9: Samples drawn from Mel-Spectrogram Speech train, test, IWGAN, and exponential respectively.

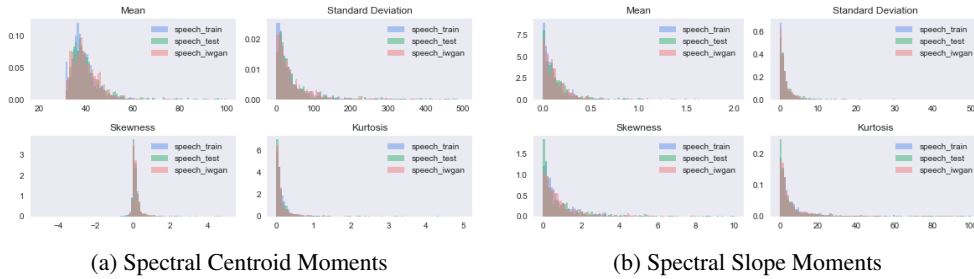


(a) Intensity empirical CDF of training data in blue and other datasets.

(b) Test statistics

Figure 10: empirical CDF and statistical tests of speech intensity

Figure 11 shows the distribution of moments computed on spectral centroids and slope. The distributions from different sources considerably overlap, indicating that the generator has efficiently approximated the real distribution of these features.



(a) Spectral Centroid Moments

(b) Spectral Slope Moments

Figure 11: Moments of spectral centroid (left) and slope(right)

Figure 12 shows statistics used to compare the reference (training data) and other datasets. The difference between KS-Statistics and JSD of the test data and generated samples are negligible. Interestingly, the spectral slope p-values of the improved WGAN are considerably higher than the test data. For these reasons and although Table 10b shows a significant difference between the KS-Statistic of test data and generated data with respect to the training data, we refrain from confirming hypothesis 4). An adversary can easily manipulate the generated data to considerably decrease the difference in this difference and keeping the high similarity in features higher to simulate such moments of spectral centroid or slope.

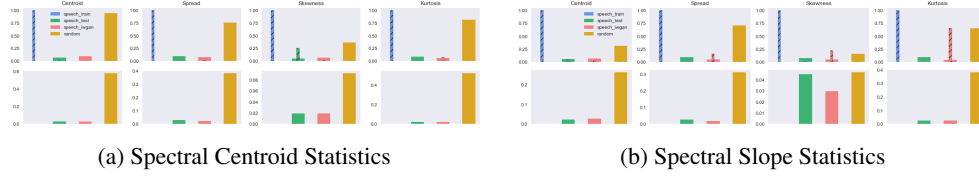


Figure 12: Statistics of spectral centroid (left) and slope(right)

6 Conclusions

Acknowledgments

References

References

- [1] Quickly generating representative samples from an rbm-derived process.
- [2] Martin Arjovsky and Léon Bottou. Towards principled methods for training generative adversarial networks. In *NIPS 2016 Workshop on Adversarial Training. In review for ICLR*, volume 2016, 2017.
- [3] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017.
- [4] David Berthelot, Tom Schumm, and Luke Metz. Began: Boundary equilibrium generative adversarial networks. *arXiv preprint arXiv:1703.10717*, 2017.
- [5] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [6] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- [7] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. Improved training of wasserstein gans. *arXiv preprint arXiv:1704.00028*, 2017.
- [8] Yanick Lukic, Carlo Vogt, Oliver Dürr, and Thilo Stadelmann. Speaker identification and clustering using convolutional neural networks. In *Machine Learning for Signal Processing (MLSP), 2016 IEEE 26th International Workshop on*, pages 1–6. IEEE, 2016.
- [9] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. *arXiv preprint ArXiv:1611.04076*, 2016.
- [10] Brian McFee, Colin Raffel, Dawen Liang, Daniel PW Ellis, Matt McVicar, Eric Battenberg, and Oriol Nieto. librosa: Audio and music signal analysis in python. In *Proceedings of the 14th python in science conference*, 2015.
- [11] Geoffroy Peeters. A large set of audio features for sound description (similarity and classification) in the cuidado project. 2004.
- [12] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [13] Tim Salimans, Ian J. Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *CoRR*, abs/1606.03498, 2016.
- [14] Sanjit A. Seshia and Dorsa Sadigh. Towards verified artificial intelligence. *CoRR*, abs/1606.08514, 2016.

- 215 [15] Li-Chia Yang, Szu-Yu Chou, and Yi-Hsuan Yang. Midinet: A convolutional generative ad-
216 versarial network for symbolic-domain music generation using 1d and 2d conditions. *CoRR*,
217 abs/1703.10847, 2017.
- 218 [16] Junbo Zhao, Michael Mathieu, and Yann LeCun. Energy-based generative adversarial network.
219 *arXiv preprint arXiv:1609.03126*, 2016.