

JDGAN: Enhancing Generator on Extremely Limited Data via Joint Distribution

Wei Li

School of Artificial Intelligence and Computer Science, Jiangnan University,
Wuxi, Jiangsu, P.R.China.

Jiangsu Key Laboratory of Media Design and Software Technology, Wuxi,
Jiangsu, P.R.China.

Science Center for Future Foods, Jiangnan University, Wuxi, Jiangsu,
P.R.China. Email: cs_weili@jiangnan.edu.cn.

Linchuan Xu*

Corresponding author. Department of Computing, The Hong Kong
Polytechnic University, Hong Kong, P.R.China. Email:
linch.xu@polyu.edu.hk.

Jiannong Cao

IEEE Fellow. Department of Computing, The Hong Kong Polytechnic
University, Hong Kong, P.R.China. Email: csjcao@comp.polyu.edu.hk.

Zhixuan Liang

Department of Computing, The Hong Kong Polytechnic University, Hong
Kong, P.R.China. Email: zhixuan.liang@connect.polyu.edu.hk.

Senzhang Wang

Nanjing University of Aeronautics and Astronautics, Nanjing, Jiangsu,
P.R.China. Email: szwang@nuaa.edu.cn.

Thomas.C.Lam

Laboratory of Experimental Optometry, Centre for Myopia Research,
School of Optometry, The Hong Kong Polytechnic University, Hong Kong.
Centre for Eye and Vision Research, Hong Kong, P.R.China. Email:
thomas.c.lam@polyu.edu.hk.

Xiaohui Cui*

Corresponding author. School of Cyber Science and Engineering, Wuhan
University, Wuhan, Hubei, P.R.China. Email: xcui@whu.edu.cn.

Abstract

Generative Adversarial Network (GAN) is a thriving generative model and

Preprint submitted to Elsevier

June 13, 2022

considerable efforts have been made to enhance the generation capabilities via designing a different adversarial framework of GAN (e.g., the discriminator and the generator) or redesigning the penalty function. Although existing models have been demonstrated to be very effective, their generation capabilities have limitations. Existing GAN variants either result in identical generated instances or generate simulation data with low quality when the training data are diverse and extremely limited (a dataset consists of a set of classes but each class holds several or even one single sample) or extremely imbalanced (a category holds a set of samples and other categories hold one single sample). In this paper, we present an innovative approach to tackling this issue, which jointly employs joint distribution and reparameterization method to reparameterize the randomized space as a mixture model and learn the parameters of this mixture model along with that of GAN. In this way, we term our approach *Joint Distribution GAN* (JDGAN). We conduct extensive experiments using MNIST, CIFAR10, and Mass Spectrometry datasets, all using extremely limited amounts of data, to demonstrate significant performance improvement of our proposed JDGAN in both achieving the smallest *Fréchet Inception Distance* (FID) score and producing diverse generated data in different types.

Keywords: Mode Collapse, Joint Distribution, Reparameterization, GAN.

1. Introduction

Generative Adversarial Network (GAN) [13] has been demonstrated as the state-of-the-art generative model in various tasks of generating synthetic but realistic-like data [30] [46] [45] [44] [29]. A typical GAN model consists of two components, a generator and a discriminator. We usually view the generator as the forger who specializes in generating plausible data to fool the discriminator into accepting it as real data, while the discriminator could be regarded as a detective who can determine whether the current data are from the generator or the real dataset. The GAN model sidesteps the difficulty of approximating many intractable probabilistic computations, for it samples data from an easy-to-sample distribution so the Markov chains [36] are never needed. Its gradients are tuned by using back-propagation, which makes the training computationally inexpensive. Besides the elegant framework, the GAN model also enjoys the power of deep generative neural networks, which, in theory, have the capabilities of approximating any

complicated probability distributions with adequate data. In particular, the generator is a deep neural network, which is designed to transform a noise sampled from a fixed, easy-to-sample distribution (e.g., Uniform distribution with $(-1, 1)$) into the “realistic” data. Because of its huge potentials, it has become an active research field and many researchers try to enhance data-generation capability by modifying the framework and its two components or loss function (e.g., Wasserstein GAN (WGAN) [4], Least Squares GAN (LSGAN) [34], Mixture Generators GAN (MGAN) [18] and Relativistic average GAN (RaGAN) [20]). Although those GAN variants have been widely used in many applications [28] [31] [10] [19] and achieve impressively plausible simulated data, their generation capabilities are not viable as the training data samples are extremely limited and diverse, e.g., a dataset consists of a set of classes but each class just holds several or even one single sample. If the training dataset is extremely limited dataset, the generated instances either hold identical samples or low quality.

In deep learning, the overlapping area refers to such an area which is intersected by two distributions. Fig.1 shows a less formal but more pedagogical illustration for overlapping area. In Fig.1, there exists two distributions, which are generated distribution (p_G) and original data distribution (p_r). We also assume both of two distributions lying in a mapping space. From Fig.1, we can see that the two distributions just hold a small overlapping area in this mapping space, which is indicated by gray shade region. Under such a scenario, the generator may hold a part of data modes and missing other modes, resulting in mode collapse problem. This is because the Jensen Shannon Divergence is maxed out when the generated and real data distributions have disjoint support [4]. With this in place, the generated data distribution p_G is far away from the original data distribution p_r in the mapping space, which results in that the generator fails to correctly capture all data modes [21]. This becomes more challenging if the training data are extremely limited data samples.

To visually and effectively demonstrate our hypothesis, we apply recent GAN variants to a diverse and extremely limited training dataset. Here, we take the MNIST dataset as the example to form the extremely limited dataset. Only 10 samples are selected from MNIST (figures '0' to '9'), and each sample belongs to a specific category. They are diverse because each sample belongs to one category; they are extremely limited because each category only holds one sample. The generated images are shown in Fig.2.

From Fig.2, we can see that the mode collapse arises in most existing GAN

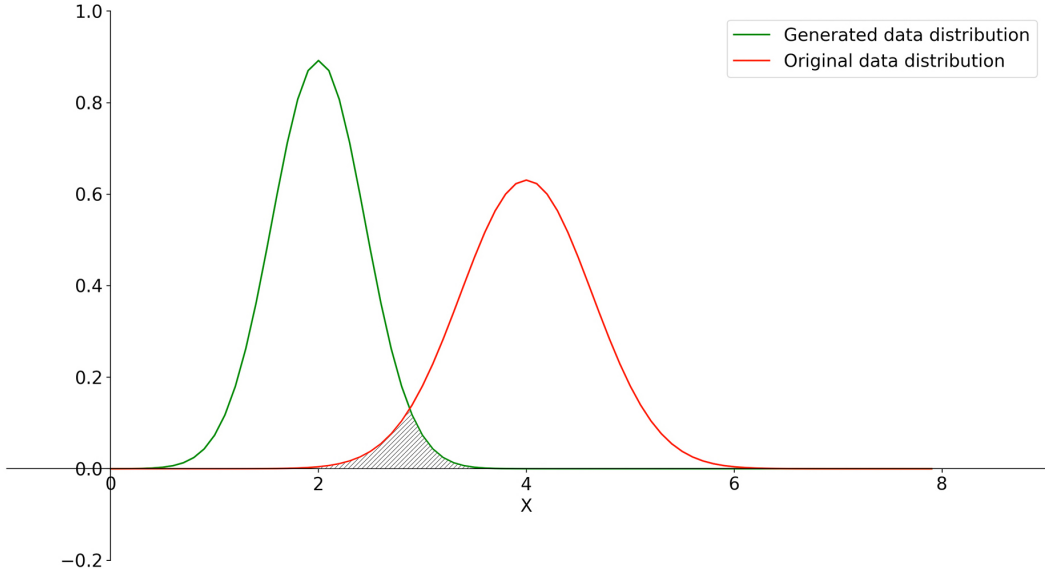


Figure 1: The gray shade region indicates the overlapping area between two distributions. Such a case also implicitly indicates low match degree of two distributions [16] [32].

variants, i.e., most data modes are missed and there exist identical generated images. Also, the generated data quality is not satisfactory in LSGAN. It is noticed that there are five GAN variants (sub-figures (a)-(d) and (f)) sample noise from Gaussian(0, 1). DeLiGAN (sub-figure (e)) [15] samples noise from the Mixture-of-Gaussian model [48]. Moreover, MGAN employs 10 generators and each generator only produces a single image. In addition, we demonstrate the performance of those GAN variants in imbalanced dataset (a category holds a set of samples while other categories just hold a single sample), and we can observe the details in experiment section.

Wasserstein GAN (WGAN) [4] utilizes Wasserstein distance to measure the dissimilarity between original data distribution p_r and generated data distribution p_G . However, only certain optimizers (e.g., RMSProp [41] or SGD [6]) are suitable for optimizing WGAN. Other momentum based ones (e.g., Adam [22]) may even turn the gradients negative, causing unstable training. RaGAN [20] argues that the minimax game should simultaneously decrease the probability that real data is real. The researchers [20] induce this property by using a “relativistic discriminator”, and prove that the relativistic discriminator makes training more stable. However, it is not suitable for

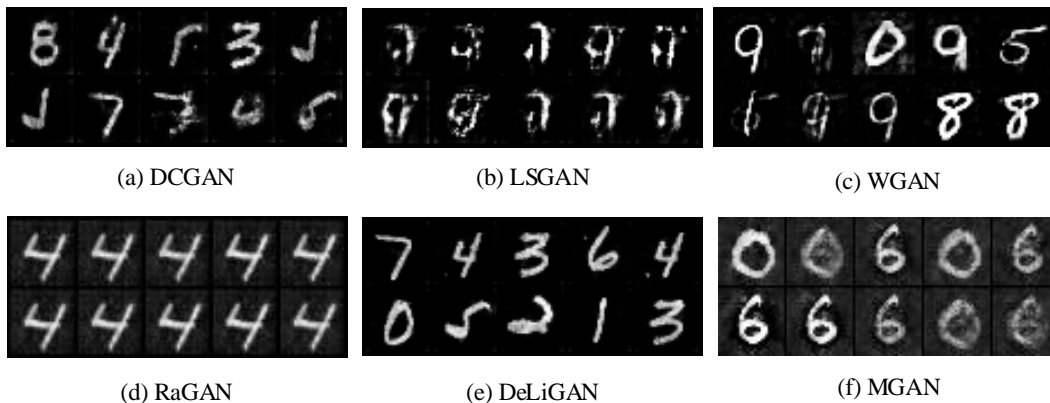


Figure 2: A case study on extremely limited MNIST samples. The generator samples noise from standard Gaussian distribution $(0, 1)$ for DCGAN, LSGAN, WGAN, RaGAN and MGAN, and samples noise from a Mixture-of-Gaussian model for DeLiGAN.

extremely limited training data. MGAN [18], on the other hand, employs a mixture of generators to learn disconnected manifolds. Since there are no restrictions enforcing generators to learn those manifolds mutually exclusively, generators may learn the same manifolds as each other, resulting in the generation of identical instances. DeLiGAN [15] modifies the latent space with Mixture-of-Gaussian model and samples noise from the MoG model. However, MoG model is a linear combination of multiple Gaussian distributions, which causes each Gaussian distribution having different weights. Therefore, the MoG model has been easily driven into the local optimal defect [7], and holding most contribution of a Gaussian to MoG model dominates other Gaussian models. The MoG model becomes a single Gaussian model under such a scenario. Therefore, the mode collapse arises.

When training a GAN, we usually learn a mapping from noise distribution $p_z(z)$ to original data distribution p_r . To guarantee learning successfully, it requires a lot of training samples such that the generator can disentangle the underlying factors of variation and make the generated data diverse. When training data are extremely limited, such a mapping mechanism becomes infeasible. To address this issue, this paper proposes to explore an alternative direction, which increases the power of the latent distribution. To this end, we propose **Joint Distribution GAN (JDGAN)**, which jointly employs multiple easy-to-sample distributions to construct the randomized space Z and learns these distributions together with the generator with reparamete-

terization, without modifying the framework. JDGAN addresses the mode collapse by increasing the dimension of $p_z(z)$ to increase that of p_G , because $p_z(z)$ is contained in p_G [3]. If the supports of p_r and p_G are not disjoint, the generator can correctly capture all the modes. To obtain a sample from the joint distribution, we employ the reparameterization method introduced by Kingma [23] to sample noise. Assuming the joint distribution is formed by a Gaussian distribution ($G(\mu, \theta)$) and a Uniform distribution ($U(a, b)$), we represent the noise from the joint distribution as a deterministic function of μ , θ , a and b .

In summary, the major innovations and contributions of this paper are as follows:

- This paper proposes a novel JDGAN model which modifies the input randomized space to enhance the generation capabilities of the GAN model, overcoming the problem of mode collapse.
- This paper demonstrates how to sample noise from the joint distribution and how to learn the parameters of the joint distribution together with the generator from theoretical and empirical perspectives, giving new insights into the success of JDGAN.
- Through comprehensive experiments on generating simulated data, we demonstrate the effectiveness of the proposed approach.

This paper is organized as follows. In section 2 we discuss some related work. We discuss the limitations of a single distribution in section 3 and present our main idea in section 4. In section 5 we will show our experimental results, and we conclude this work in section 6.

2. Related Work

The GAN model displays its powerful generative capabilities since its invention, however, there exist challenges (e.g., mode collapse and gradients vanishing). Hence, many modifications to the original GAN model have been proposed, and they can be mainly categorized into four types: modifying the components, modifying the penalty function, modifying the architecture and modifying the randomized space.

Modifying the components. The early researchers utilize some seemingly simple but powerful strategies to improve the performance of the GAN

model. One of the first major improvements is DCGAN [1], which hopes to bridge the gap between the success of CNNs for supervised learning and unsupervised learning. It modifies the two components of the GAN model by adding some functions such as BatchNorm [27], ReLu activation [12] for the generator and BatchNorm, LeakyReLu activation [47] for the discriminator, and replacing any pooling layers [24] with strided convolutions (discriminator) and fractional-strided convolutions (generator). These modifications are suitable for both the generator and the discriminator to learn good up-sampling and down-sampling operations, which could improve the quality of simulated data. Salimans et al. (Improved GAN) [39] propose heuristic approaches to stabilize the training of GAN. Specifically, they use the feature matching to address the instability of training a GAN model by changing the objective for the generator to prevent the problem of overtraining, and they use the mini-batch [27] discriminator to prevent the generator collapse, for the discriminator can easily tell whether the generator is producing same outputs. Mirza et al. (Conditional GAN) [35] train both the discriminator and the generator by using the new input that is conditioned on adding extra information y (e.g., class labels), and it can generate descriptive tags which are not part of training labels. This modification holds considerable flexibility for generation.

Modifying the penalty function. The GAN model generally adopts the JS divergence [33] to calculate the similarity between two distributions which are from different datasets, for the principle of GAN is to transform a distribution into another distribution. However, it is hard to achieve the transforming process, and always causes the mode collapse in practice. Thus, many researchers adopt different strategies to address this issue. WGAN [4] replaces the JS divergence with the Wasserstein distance [37]. The value of JS divergence for a GAN model could be a constant as the two distributions have no overlapping area (or the overlapping area can be neglected), which causes distribution transformation difficulty. The Wasserstein distance can reflect the dissimilarity between the two distributions without the overlapping area. Moreover, Gulrajani et al. (Improved WGAN) [14] find that the weight clipping adversely reduces the capability of the discriminator in WGAN, they then improve the WGAN by penalizing the norm of the discriminator gradients during training instead of performing parameter clipping. LSGAN [34] modifies the loss function with the Least Squares to generate samples that are closer to the real data. This study argues that the Sigmoid Cross Entropy loss function [50] for the discriminator would lead to the problem

of vanishing gradients as updating the generator using the fake samples that are on the correct side of the decision boundary and are far from the real data. The Least Squares loss function can penalize the fake samples that are lying in a long way on the correct side of the decision boundary moving toward the decision boundary even though they are classified correctly.

Modifying the Architecture. Since manifolds of original data are disconnected in the space and a single generator G only produces instances in certain regions of this space, the generated data focuses on several or even one single manifold. Researchers attempt to increase the quantity of generator to learn more about different manifolds. Tolstikhin et al. [42] added a new component to a mixture model by running a GAN algorithm on a re-weighted sample. Inspired by boosting techniques, this idea greedily aggregates many potentially weak individual predictors to form a strong composite predictor. This model is termed AdaGAN. Since AdaGAN utilizes a sequential training technique to train the model, the model is computationally expensive. Moreover, it is hard to search the *ChooseMixtureWeight* and the *UpdateTrainingWeight* functions for boosting techniques. Arora et al., [5], alternatively, trained a mixture of generators and discriminators to play the minimax game with the reward function being the weighted average reward function between any pair of generator and discriminator. This strategy is not only computationally expensive but also lacks a mechanism to enforce the divergence among generators. Ghosh et al. [11] employed many generators and trained them by using multi-class discriminators that, in addition to detecting whether a sample is fake or not, predict which generator produces this sample. The loss function in this study focuses on detecting whether a sample is fake and does not directly encourage generators to produce diverse instances. The recent GAN variant for increasing the quantity of generator is MGAN [18]. MGAN employs many generators $G_{1:K}$ and an extra classifier C to construct architecture. Although study [18] claimed such a design can help a model learn all manifolds, it is difficult to achieve this goal in practice. MGAN does not provide that G_i is mutually exclusive with G_j , $i \neq j$. In other words, G_i and G_j may learn the same manifold during training. In addition, MGAN adopts the shared parameters to save training cost. However, the shared parameters may cause all generators output the same instance, given that the parameters have an influence to the quality of data. RaGAN [20] argues that the Nash equilibrium should simultaneously decrease the probability that real data is real, which may improve the GAN performance. In this way, RaGAN utilizes a “relativistic discriminator” to

prove this property, and this discriminator estimates the probability that the given real data is more realistic than fake data, on average. However, it still suffers from the challenge in the case of extremely limited data. We can observe more details in section experiments.

Modifying randomized space. DeLiGAN [15] hopes to make \mathcal{Z} disconnected to learn the disconnected manifolds, with the mixture of Gaussian (MoG) model. However, it is hard to make \mathcal{Z} disconnected in practice. MoG model is a linear combination of multiple Gaussian distributions, and the combined distribution still belongs to a bounded continuous distribution. This is pedagogically shown in Fig.3. There still exists the uncovered manifolds such that some modes are missed in generated instances, given that $p_z(z)$ is contained in p_G [3]. When we sample the noise from the MoG model, it needs two steps. First, one of N Gaussian models is selected. Second, we draw the noise from the chosen Gaussian distribution. Considering an extreme scenario, the same Gaussian distribution has been chosen at each epoch. The MoG model becomes the single Gaussian model under such a scenario. The mode collapse is still not addressed. In addition, DeLiGAN does not give how to determine the number of Gaussian model. However, we would fully discuss how to determine the number of distributions in our study. We believe that the number of distributions would influence the diversity of generated data. More details are shown in experiment section.

3. Preliminaries

In this section, we discuss the limitations of sampling noises from a fixed, easy-to-sample distribution and give the explanation of reparameterization method.

3.1. The Limited Overlapping Area

Although the GAN model was introduced in Section 1, we formally present the GAN model as below to establish the continuity. The GAN model was introduced by Goodfellow [13] as a novel generative model to simultaneously train a generator and a discriminator by using Eq. (1).

$$\begin{aligned} \min_G \max_D V(G, D) = & \mathbb{E}_{x \sim p_r(x)} [\log D(x)] \\ & + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))] \end{aligned} \quad (1)$$

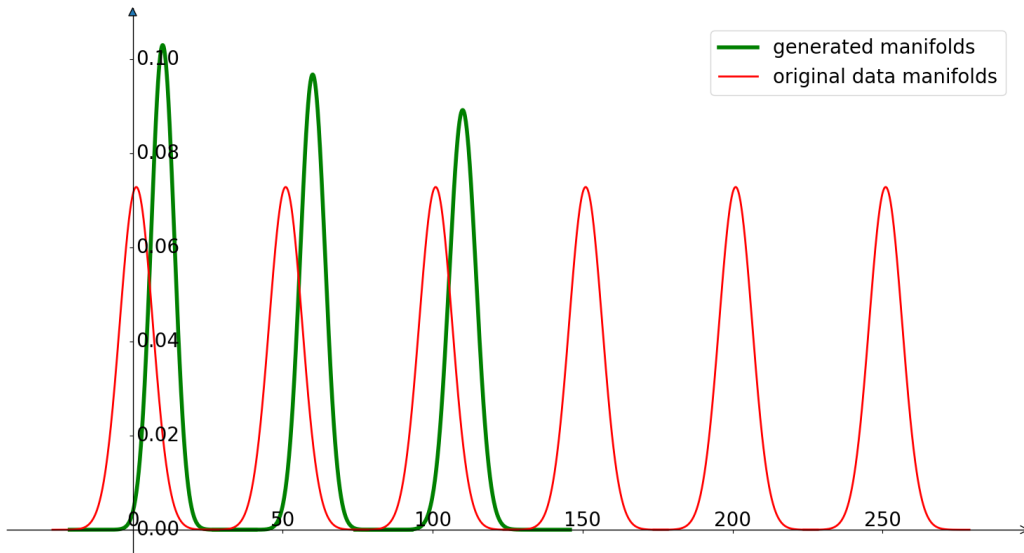


Figure 3: Assuming original data (red curve) hold 6 submanifolds, while MoG model (green curve) just covers part of submanifolds of original data.

where x comes from a distribution $p_r(x)$ (it is abbreviated as p_r) sampled from the original dataset and z comes from a fixed, easy-to-sample distribution $p_z(z)$ (e.g., Uniform with $(-1, 1)$ or Gaussian with $(0, 1)$). The generator builds a mapping function from $p_z(z)$ to original data space χ , and the discriminator would output a single score $\in [0, 1]$ to indicate whether the current data are from the generator or not. Generally, a small score indicates the current data generated by the generator, while a high score indicates the current data coming from the original dataset. We repeat this training process until both the discriminator and the generator reach to the Nash equilibrium [8] where $p_G = p_r = 0.5$.

Eq.(1) measures the difference between p_r and p_G by using the Jensen–Shannon (JS) divergence [9]. The closer the two distributions are, the smaller the JS divergence is. However, the two distributions are hard to match each other because they lie in low-dimensional manifolds [3] [4]. In other words, we cannot guarantee that the two probability distributions (p_G and p_r) have a large overlapping area. Under such a scenario, the matching degree of both p_G and p_r is not high [3]. With such, both p_r and p_G have a limited overlapping area. Note that p_G is defined via sampling from the simple prior $p_z(z)$ [3]. This implicitly indicates that we can utilize prior $p_z(z)$ to increase the

overlapping area.

3.2. The Limited Diversity

In general, the noise distribution $p_z(z)$ is a fixed, easy-to-sample distribution, i.e., Uniform with $(-1, 1)$ or Gaussian with $(0, 1)$. The model usually adopts independent identically distributed strategy to draw noise code z from such an distribution, and transform it into high-dimensional matrix ($G(z)$). Let G be a function composed by affine transformations and rectifiers. We get $G(z) = D_n W_n \dots D_1 W_1 z$ where W_i denotes the affine transformation and D_i indicates the rectifier. Note that the noise code z follows the distribution $p_z(z)$ and the generated data G_z follows distribution p_G . In this way, $p_z(z)$ is contained in p_G .

With this in place, the diversity of generated data G_z has been limited by such a low-dimensional noise code when GAN model has been trained successfully because the number of noise code z is limited. One may want to increase the dimension by increasing the number of noise samples (e.g., $z = 10 \rightarrow z = 100$). However, these noise samples could be linear correlation so the intrinsic dimension of those noise samples is far less than 100. It further increases the loss of diversity as we consider such a scenario where the neural network brings the dimension reduction mapping.

On the other hand, assuming the original data holds a set of disjoint manifolds and each one named submanifold $_i$, a part of manifolds of original data are covered by the support of the generated data distribution. This is because $p_z(z)$ is supported on a connected subspace of χ while $p_z(z)$ is contained in p_G . Therefore, generator G (a continuous function by design) can not correctly model a set of disjoint manifolds in χ [21]. Although there is a trade-off between covering all original data distribution and minimizing the volume of the off real-manifold space in the cover, such a trade-off also indicates that the generator may sacrifice certain submanifolds to learn a cover with less off real-manifold volume [21]. In this way, the diversity of the generated data is limited.

In addition, there is a lower bound for intrinsic dimension of a dataset [49], the diversity of data will be lost when the intrinsic dimension of generated data is smaller than such a lower bound. It is much worse on diversity when the training data samples are extremely limited and diverse (See Fig.2).

3.3. Reparameterization Method

In the traditional sampling method, the lower bound of the noise variable is fixed, e.g., Uniform with $(-1, 1)$ or Gaussian with $(0, 1)$. The reparameterization method can help us to determine the lower bound of the variable through back-propagation. For example, we want to take the gradient w.r.t. θ of the following expectation $\mathbb{E}_{p(z)}f_{\theta}(z)$ where p denotes the density. However, computing the gradient of this expectation (i.e., $\nabla_{\theta}\mathbb{E}_{p_{\theta}(z)}f_{\theta}(z)$) is often difficult because the integral is typically unknown and the parameters θ , with respect to which we are computing the gradient, are of the distribution $p_{\theta}(z)$. The reparameterisation method can “transform” this formula to another expression, i.e., $\nabla_{\theta}\mathbb{E}_{p_{\theta}(z)}f_{\theta}(z) = \nabla_{\theta}\mathbb{E}_{p_{\epsilon}}f(g(\epsilon, \theta))$. In this way, the gradient is now unrelated to the distribution with which we take the expectation, so easily passes through the integral. We take the Gaussian distribution as the example. The sample can be represented from the Gaussian as a deterministic function of μ and δ and an auxiliary noise variable ϵ , i.e., $z = \mu + \delta\epsilon$, and we can learn the parameters μ and δ along with the GAN parameters.

4. Joint Distribution GAN

In this section, we propose the Joint Distribution GAN (JDGAN). As mentioned in the introduction, the innovation of JDGAN lies in the input randomized space Z , which modifies Z with joint of multiple independent distributions. To this end, there are four important issues to be addressed, which are how to construct a new randomized space using joint distribution, how to draw samples from the joint distribution, how to choose the number of distributions and how to learn the parameters of JDGAN.

4.1. Constructing the New Randomized Space

No matter which GAN-based model we choose, we always hope the network to learn a mapping from the noise distribution $p_z(z)$ to the original data distribution p_r . In our study, we employ joint distribution to construct the randomized space for the generator. In joint distribution, the distribution of the random vector $Z = (Z_1, Z_2, \dots, Z_n)$ is named joint distribution of the random variable Z_i , and it is shown as follows.

$$P(Z) = P(Z_1 = z_1 \text{ and } Z_2 = z_2 \text{ and } \dots \text{ and } Z_n = z_n) \quad (2)$$

In Eq.4, each z_i indicates a distribution. The joint distribution of multiple variables can be generally represented as follows:

$$\begin{aligned} p_z(z) &= p_z(z_1, z_2, \dots, z_n) \\ &= p_z(z_1)p_z(z_2)\dots p_z(z_n) \end{aligned} \quad (3)$$

Here we assume that each distribution is independent of all the others, so Eq.(3) can be reformulated as follows:

$$\begin{aligned} p_z(z) &= p(z_1, z_2, \dots, z_N | \theta_1, \theta_2, \dots, \theta_N) \\ &= p(z_1 | \theta_{z_1}, z_2 | \theta_{z_2}, \dots, z_N | \theta_{z_N}) \\ &= p(z_1 | \theta_{z_1}) \cdot \dots \cdot p(z_i | \theta_{z_i}) \cdot \dots \cdot p(z_n | \theta_{z_n}) \\ &= \prod_i^n P(z_i | \theta_{z_i}) \end{aligned} \quad (4)$$

In this way, p_G which is the generating distribution parameterized by the generator should be a larger and more complicated distribution, so the overlapping area between p_r and p_G is increased. Moreover, the increasing distributions also increase the dimension of manifold, so the diversity of the generated data can also be increased. This is because adopting the joint distribution can construct a more complex $p_z(z)$. Since the supports of p_G and p_r lie on low-dimensional manifolds [3], we attempt to increase the dimension of manifold of p_G to overlap p_r , given that $p_z(z)$ has to be contained in p_G [3]. If the dimension of manifold of p_G is no longer less than that of p_r , it is certainly for p_r to be covered by p_G in the mapping space. Assuming an extreme scenario, p_G is complicated enough to fill up the mapping space. That is to say, p_r and p_G have non-disjoint supports. In this way, the diversity of the generated data can be guaranteed.

Since the JDGAN focuses on constructing the joint distribution to modify the randomized space Z , the generator and discriminator architectures of JDGAN are from the vanilla DCGAN. The details of components still adopt Conv-BatchNorm-ReLu (Generator G) and Conv-BatchNorm-LeakyReLu (Discriminator D). One may use other frameworks given that the designing of randomized space Z is independent of the architectures. Even though the common architectures are utilized, JDGAN is demonstrated to be more effective than state-of-the-art variants of GAN. More details are shown in the experiment section.

In the rest of this subsection, we discuss how to choose each distribution (i.e., $P(z_i|\theta_{z_i})$) to form the joint distribution for constructing randomized space Z . Theoretically, we can use any statistical distribution (even the same distribution) to jointly construct Z , and what we need to pay attention to is that those distributions have to be independent to each other when employing multiple distributions (we will consider the case of dependent joint distribution in future work). However, many distributions can be viewed as the combination of the Gaussian distributions and the Uniform distributions [17], and both of them are widely used in many applications. Here we take the two distributions as the example to construct the randomized space. More distributions are shown in the section experiments. We instantiate Eq.(4) via the joint distribution of a Gaussian distribution and a Uniform distribution as shown in the following equation:

$$p_z(z) = p(z_1, z_2) = p(z_1|G(\mu, \delta))p(z_2|U(a, b)), \quad (5)$$

By Eq.(5), the joining of two distributions belongs to non-linear combination, which increases the power of prior distribution. It is highly likely to produce a larger overlapping area (or a higher matching degree) than a single one and the joining of two distributions significantly increases the diversity. Next, we will introduce how to sample noise from the new randomized space.

4.2. Sampling from Joint Distribution

We take Eq.(5) as the example. In Eq.(5), the new randomized space has been determined by two different distributions ($p(z_1)$ follows the Gaussian distribution and $p(z_2)$ follows the Uniform distribution). In order to draw noise code z from Eq.(5), we employ the “reparameterization method” [23] to sample noise code z_1 from the Gaussian distribution and use the inverse transformation method [43] with the cumulative distribution function $(\frac{x-a}{b-a})$ to sample noise code z_2 from the Uniform distribution. Since we repeatedly sample noise code from the two distributions and both of the two distributions are univariate Gaussian distribution and uniform distribution, all the elements of z are not the same and each element within noise code is a scalar. In this way, the noise code drawn from the joint distribution can be denoted as follows:

$$z = (\mu + \delta\varepsilon) [\zeta(b - a) + a] \quad (6)$$

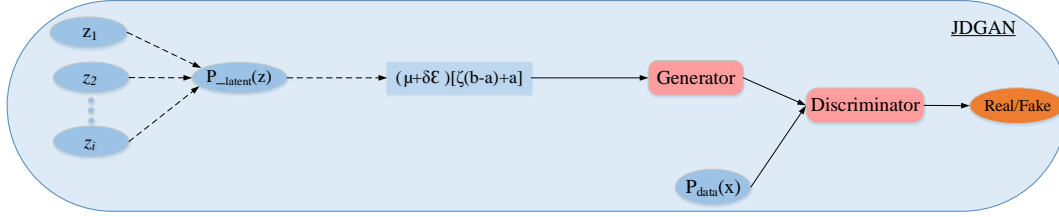


Figure 4: JDGAN architecture. Different from the regular GAN model, the generator samples noise from the joint distribution. The $P_{\text{latent}}(z)$ is made up of multiple distributions (z_i). Those distributions are independent to each other.

In Eq.(6), the first term follows the standard Gaussian and the second term follows the standard Uniform distributions respectively, which means that the data we sampled from the randomized space Z are a deterministic function of the parameters of μ , δ , b and a . Substituting Eq.(5) into Eq.(6), we get:

$$P_r(G(z)) = \int p_r[G(\mu + \delta\epsilon)|\epsilon]p(\epsilon)d\epsilon \cdot \int p_r[G(\zeta(b-a) + a)|\zeta]p(\zeta)d\zeta \quad (7)$$

In this way, our new objective is to learn μ , δ , b and a together with the generator to minimize $\mathbb{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))]$, and the architecture of JDGAN is shown in Fig.4.

4.3. Learning Parameters of JDGAN

For each distribution, we first need to initialize their parameters. We still take the Eq.(6) as the example. We initialize μ_i and δ_i with 0 and 1 for $p(z_1)$ and assign (-1, 1) to $p(z_2)$ (a and b). In order to train our model, generator samples a latent vector z from the joint distribution (Eq.(6)). After feeding z into the generator, we obtain the generated data and then put these data into discriminator. The discriminator would output a signal to update generator's parameters by using vanilla GAN training procedure (Eq.(1)). Also, μ , δ , a and b are trained simultaneously along with the generator's parameters, with gradient methods. The derivative with respect to each parameter is easy to obtain using the back-propagation method (e.g., $\frac{\partial G}{\partial \mu} = \frac{\partial G}{\partial z} \cdot \frac{\partial z}{\partial \mu}$).

4.4. Determining the Number of Distributions

The goal of jointly employing multiple distributions in the randomized space is to increase the generating capabilities. However, we need to know how to choose the optimal number of multiple distributions. In this paper, the *Intrinsic Dimension* [2] [26] is utilized to determine the number of distributions. A dataset is usually projected into a low-dimensional manifold, and the minimal dimension of such a manifold is called *intrinsic dimension* [2]. Assuming there exist independent identically distributed (IID) sample observations X_1, \dots, X_n from a high-dimensional space \mathbb{R}^D , and those observations represent an embedding of lower-dimensional samples, i.e., $X_i = g(Y_i)$. Y_i is from an unknown smooth density f on \mathbb{R}^m with $m < D$. The unknown term m is the expected *intrinsic dimension*. Here, we utilize the widely used maximum likelihood method [26] to estimate *intrinsic dimension*, which is shown in Eq.(8).

$$\hat{m}(x) = \left[\frac{1}{n} \sum_{j=1}^n \log \frac{x_j}{T_j(x)} \right]^{-1} \quad (8)$$

where $T_j(x)$ is the Euclidean distance from the point x to its j th nearest neighbor. More deductions are shown in [26]. We take the extremely limited MNIST as the example to demonstrate the optimal number of distributions. We feed the extremely limited training samples and generated instances from DCGAN (See Fig.2 (a)) into Eq.(8), and get $\hat{m}(x)_{DCGAN} = 3.30$, $\hat{m}(x)_{MNIST} = 5.79$. In this way, two distributions may be the optimal choice in our cases. A detailed discussion is shown in experiment section.

Besides, we observe in our experiments that a larger number not always brings better performance. The reason behind it could be as follows. If the number of distributions is large, it can cause that the noise distribution is more complicated than the raw data distribution. In other words, a GAN model transforms a complicated distribution into a relatively simpler distribution, and such a transformation process could abandon some valuable information for mapping the former one into the latter one, which would cause the loss of diversity.

5. Empirical Evaluation

5.1. Experiment Settings

For the experiments, the implementation details of JDGAN are shown in Fig.5. Six recent GAN models are employed as baselines, which are **DCGAN** [1], **WGAN** [4], **MGAN** [18], **LSGAN** [34], **RaGAN** [20] and **DeLiGAN** [15], respectively.

Two commonly used public image datasets, the MNIST dataset and the CIFAR-10 dataset, and a real-word medical Mass Spectrometry dataset, are studied. JDGAN is proposed to enhance the capabilities of capturing the diverse modes of data. We thus mainly compare the diversity of generated images with these GAN-based models. Since deep neural networks can theoretically approximate any kind of distributions with adequate training data, we reduce the amount of data so as to test how many benefits the designed JDGAN can bring. However, it is not easy to justify how limited the data should be. In this way, we study diverse and extremely limited data which means that there are many categories of images but each category holds only one image because no more data can further be reduced. To make a fair comparison, the quantity of generated data for all models is the same as that of training data (10 samples for image dataset and 5 samples for medical dataset).

5.2. MNIST Dataset

We first generate image simulation data using JDGAN on MNIST dataset where each sample is a gray image with 1*28*28 size. We sample only one image from each category as the training data, and the hyperparameters are shown in Fig.5. In Fig.5, we use the Adam optimizer [22] and Binary Cross Entropy [25] to update the generator and the discriminator and the multiple distribution’s parameters (e.g., μ , δ , b and a). We set the parameter of LeakyRelu [47] as 0.02 and that of Dropout [40] as 0.5. The activation of last layer for the generator and the discriminator is *Sigmoid* for MNIST and MS datasets and *Tanh* for CIFAR-10 dataset.

In this case, we test the number of distributions in Z from 2 to 4 and the drawing noise samples from only multiple Gaussian distributions ($z = (\mu + \delta\varepsilon_1) * \dots * (\mu + \delta\varepsilon_n)$) and only multiple Uniform distributions ($z = (\zeta_1(b - a) + a) * \dots * (\zeta_n(b - a) + a)$) which are used to form the joint distribution in Z , and the generated images are shown in Table.1. Note that the parameters a , b have been initialized to $(-1, 1)$ and the parameters μ and δ have been

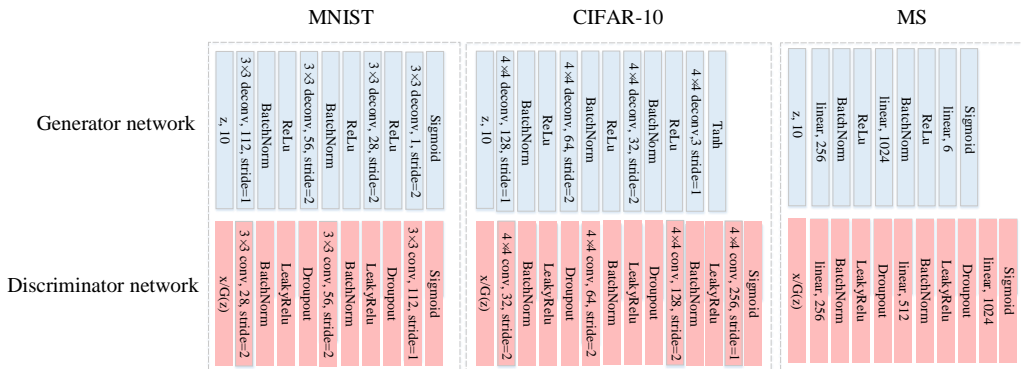


Figure 5: Architectural details of JDGAN model. “ $K \times K$, conv/deconv, C , stride = S ” denotes a convolutional/deconvolutional layer with $K \times K$ kernel, C output filters and stride = S . BN indicates a batch normalization layer. 10 indicates the dimension of z . In Mass Spectrometry (MS) dataset, the out channel at the last layer for generator is set to 6, for 6 features (Acq Time, Intensity, Precursor Intensity Acquisition, Apex Time, Elution Peak Width, MS2Counts) can determine whether the eyes are diseased or not.

initialized to $(0, 1)$. After that, we incorporate the two distributions (Eq.(7)) into a joint distribution, and the results are shown in Table.2. Note that we have already studied all the baselines on this dataset, and the results shown in Fig.2 demonstrate poor performance of all the baselines. We now report the *Fréchet Inception Distance* (FID) [16] [32] scores obtained by our JDGAN (Row 1 and column1 in Table.2) and baselines (Fig.2) in Table.3 on MNIST dataset. FID measures the Fréchet distance between the two distributions and it is the 2-Wasserstein distance. In other words, FID can loyally reflect the matching degree (it is also called overlapping area) of two distributions (p_r and p_G). The smaller the FID score is, the better. From Table.3, we can see that the JDGAN achieves the smallest score, which proves the effectiveness of our proposed JDGAN on generating diverse simulation data and more closer distribution as original data distribution. Also, we visualize the generated data distributions for all models to further validate the promising performance of our proposed JDGAN on overlapping area, which is shown in Fig.6. It shows that the distribution of simulation data produced by JDAGN is closer to the original data distribution than that produced by baselines, and thus suffers the least from the mode collapse.

The three sets of results demonstrate the promising results produced by our proposed JDGAN. However, there is an interesting scenario where the diversity is decreasing when the number of distributions is increasing. Note

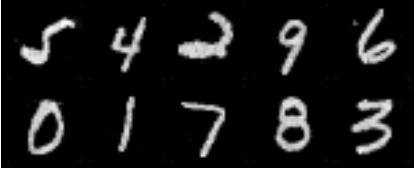
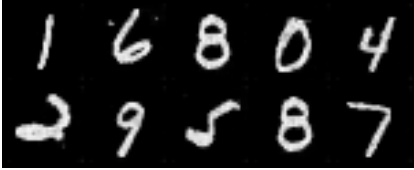
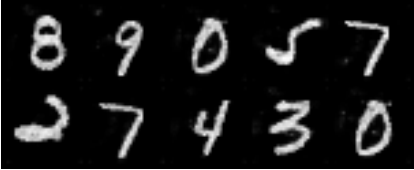
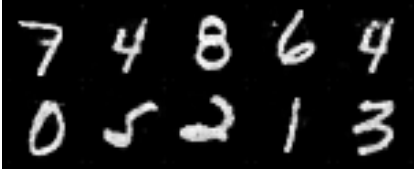
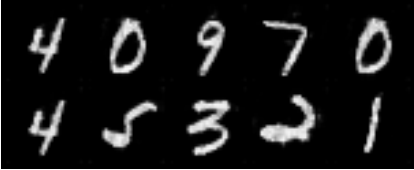
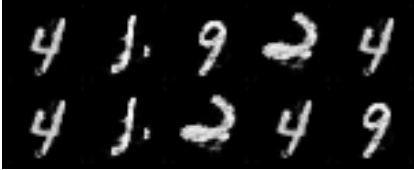
Number	Simulation Data1	Simulation Data2
2		
3		
4		

Table 1: We jointly employ multiple Gaussian distributions (Simulation Data1) or Uniform distributions (Simulation Data2). The number=2 indicates that we employ two same distributions to construct Z . The best performance is achieved with the number=2 as the randomized space is constructed by 2 Gaussian distributions, for the diversity outperforms others.

that when employing a large number of distributions to construct the randomized space, the noise distribution could be more complicated than original data distribution. We think it may be because the noise’s diversity is richer under this scenario, and the generator may abandon some values during training when the GAN model transforms a complicated distribution into an relatively simpler distribution. The abandoned values would cause incomplete information, which may be a part of a handwritten figure. Thus, to fool the discriminator, the generator easily tends to generate the “realistic” but repeated simulated data rather than the “unrealistic” but diverse instances. Thus, a suitable number is 2 for the MNIST dataset in this scenario.

Notice that all the experiments shown above are conducted on the scenario with diverse and extremely limited data. To validate the generalization of our proposed JDGAN, we test its performance on other scenarios.

1. We apply our proposed JDGAN to the adequate training data (60000 figure samples) and compare it with the state-of-the-art GAN vari-

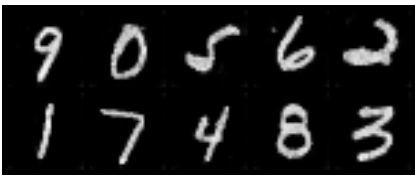
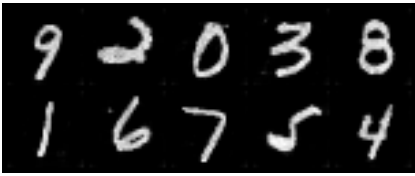
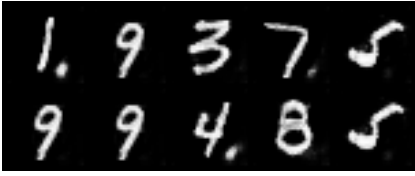
Number	Simulation Data
1	
2	
3	

Table 2: We jointly employ multiple Uniform and Gaussian distributions. The number=1 indicates that the joint distribution is determined by a Uniform and a Gaussian (Eq.(6)), and number=2 indicates that the joint distribution is determined by two Uniform distributions and two Gaussian distributions. The best performance is achieved with the number=1 and number=2, for they show the promising diversity.

ants. The generated images are shown in Fig.7. It shows that the proposed JDGAN can also achieve good performance with adequate training samples.

2. The training dataset is imbalanced dataset. On this scenario, the category '0' holds 8 samples, and the categories '1' and '2' hold 1 sample. The results are shown in Fig.8. The generated data produced by JDGAN are more diverse than other two models, because the three categories are arisen while DeLiGAN lacks category '2'.

Those experimental results show the effectiveness of our proposed JDGAN on capturing diverse modes in different scenarios.

5.3. CIFAR-10 Dataset

The CIFAR-10 dataset is a colorful image dataset with $3*32*32$ size. We continue to apply all models to this dataset, and the architectural details of

Models	FID Score
DCGAN	202.1
DeLiGAN	106.3
LSGAN	266.5
WGAN	179.4
RaGAN	290.3
MGAN	401.8
JDGAN	66.2

Table 3: FID scores of generated data shown in Fig.2 and Row 1, column1 in Table.2.

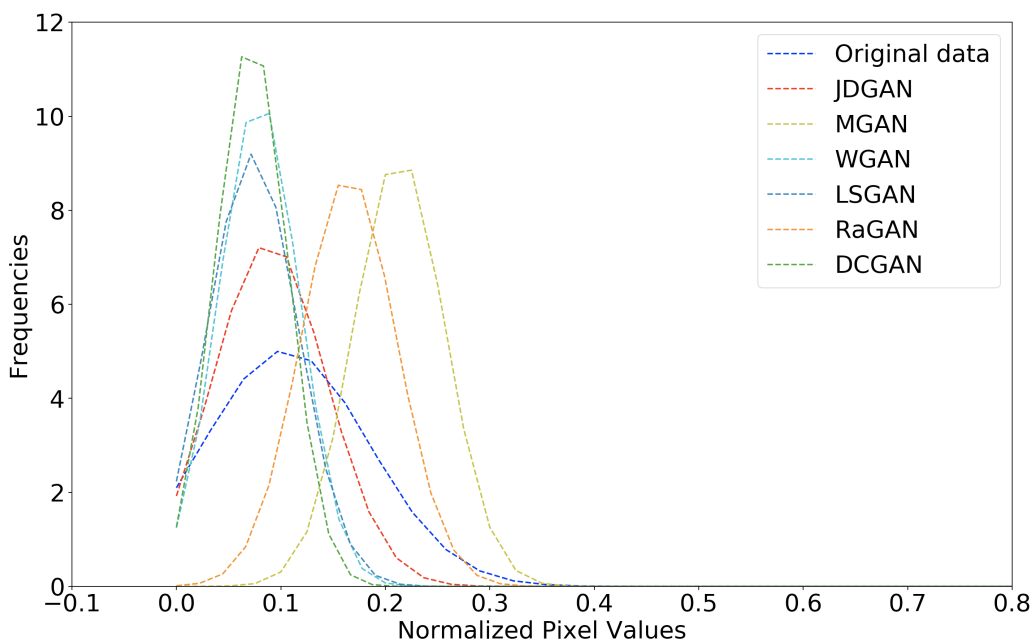


Figure 6: The overlapping area on the MNIST dataset. The blue line indicates the original data distribution, and other lines indicate the generated data distribution, which are produced by JDGAN, MGAN, WGAN, LSGAN, RaGAN and DCGAN respectively. Obviously, JDGAN is closer to original data distribution than baselines.

JDGAN are the same as the Fig.5 except for the last layer of the generator, which replaces *Sigmoid* with *Tanh*, given that the CIFAR-10 holds colorful images and the *Tanh* function can cover the color space of the training distribution. We set the noise distribution as the Gaussian distribution (0,

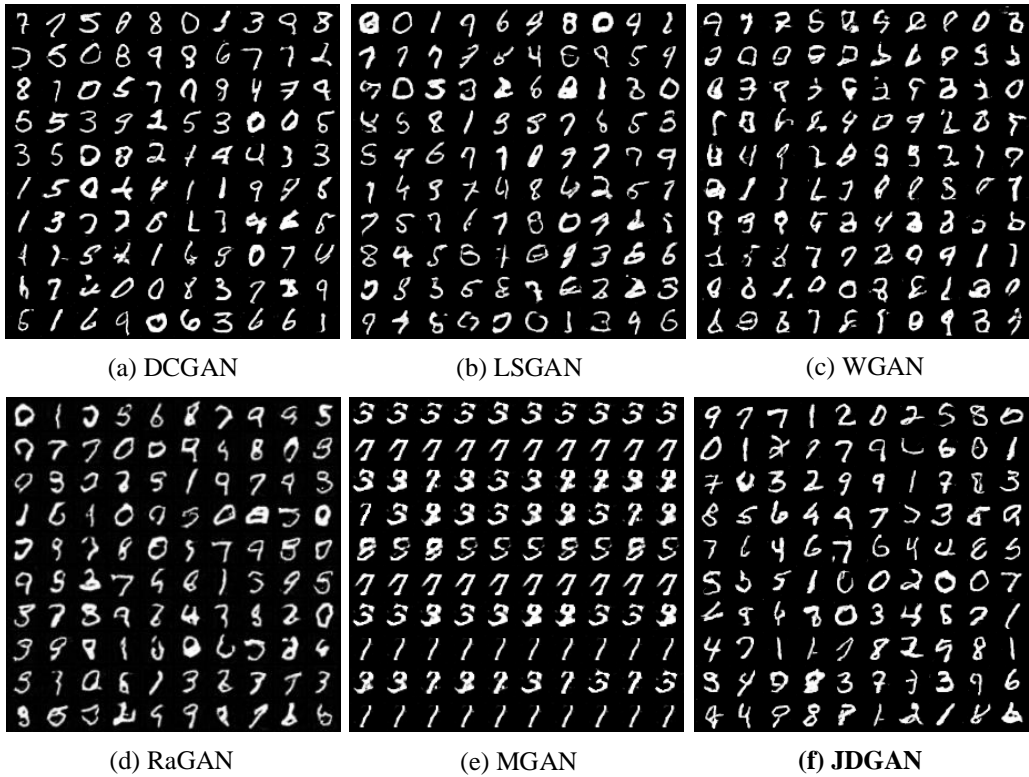


Figure 7: Generated images produced by baselines and JDGAN on adequate MNIST training samples.

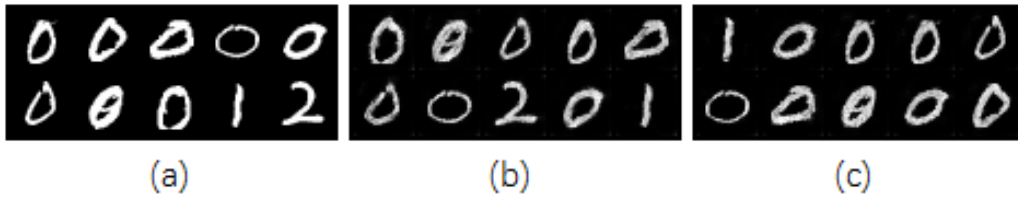


Figure 8: Sub-figure (a) indicates the training data, and sub-figure (b) indicates the generated data produced by JDGAN. Sub-figure (c) shows the generated data produced by DeLiGAN.

1) for DCGAN [1], WGAN [4], RaGAN [20], MGAN [18] and LSGAN [34], and sample noise samples from the Mixture-of-Gaussians model for DeLiGAN [15]. For MGAN, we employ 10 generators and each one is expected to produce one instance. Similar to MNIST dataset, we still sample 10 im-

ages from CIFAR10 dataset, and each image represents a specific category. The generated images of different models are shown in Fig.9. From these generated data, we can observe that there are many identical images in sub-figures (a), (b), (d) and (f), and we even observe the noise generated images in sub-figures (c), (d) and (f).

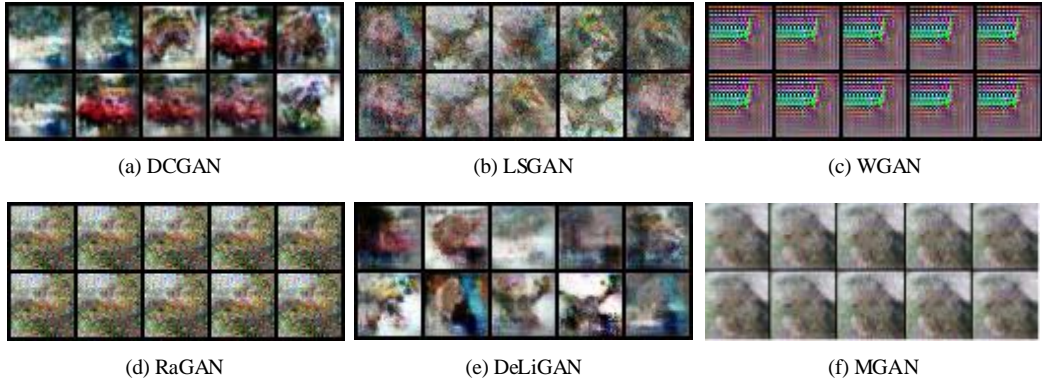


Figure 9: The generated images are produced by DCGAN, LSGAN, WGAN, RaGAN, DeLiGAN and MGAN, respectively. In this case, the training data only contain 10 samples, and each one represents a specific category. The generator samples noise from standard Gaussian distribution $(0, 1)$ for DCGAN, LSGAN, WGAN, RaGAN and MGAN, and samples noise from a Mixture-of-Gaussians model for DeLiGAN.

Similar to MNIST, we continue to construct the randomized space Z by using only multiple Gaussian distributions, only multiple Uniform distributions and the joint distribution of these two distributions respectively. The generated results are shown in Table.4 and Table.5, respectively.

From Table.4 and Table.5, We can observe that the best performance is achieved with the number=3 in Table.4 when the randomized space is determined by three same distributions, and that is achieved with the number=2 in Table.5 if the randomized space is determined by two different distributions.

We think that the information possessed by CIFAR-10 dataset is richer than MNIST dataset, and there is a need to sample noise from more complicated distributions for the generator to capture more diverse modes. Similar to MNIST dataset, the more distributions could also cause generating identical simulated images, for transforming complicated distributions into a relatively simpler distribution could abandon some information of diversity. As for less distributions (e.g., number=2 in Table.4), such a combination cannot


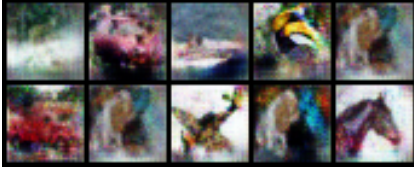


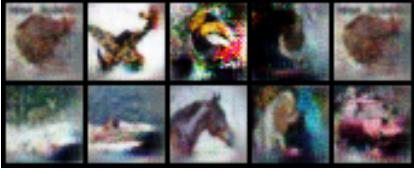
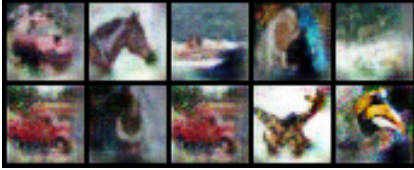
Number	Simulation Data1	Simulation Data2
2		
3		
4		

Table 4: We jointly employ multiple Gaussian distributions (Simulation Data1) and multiple Uniform distributions (Simulation Data2). The best performance is achieved with the number=3, which means we employing three same distributions to construct Z .

satisfy the rich diversity and the overlapping area could not be enough, thus, the generator tends to generate “realistic” but identical simulation images. In conclusion, we recommend the number of distributions as two different distributions. We now report the *Fréchet Inception Distance* (FID) scores and overlapping areas obtained by our JDGAN (Row 1 in Table.5) and baselines (Fig.9) in Table.6 and Fig.10. From Table.6, we can observe that the JDGAN still achieves the smallest score; from Fig.10, we can observe that the JDGAN holds the most similar distribution as original data distribution than other GAN variants. Both results further prove the effectiveness of our proposed JDGAN.

We continue to test its performance on adequate CIFAR10 training samples (50000 images), and the generated images are shown in Fig.11. The experimental results also further illustrate that the JDGAN is not only good at generating diverse simulation data in the extremely limited dataset but also suitable for the adequate data samples.

Number	Simulation Data
1	
2	
3	

Table 5: We jointly employ multiple Uniform and Gaussian distributions, and the best performance is achieved with the number=1.

5.4. Other distributions

The Uniform distribution and Gaussian distribution are the most common distributions in statistics, so they are widely used in many GAN variants (e.g, Uniform for the vanilla GAN and DCGAN, Gaussian for the LSGAN and WGAN). To validate the performance of proposed JDGAN on other distributions, we employ the Cauchy distribution, Exponential distribution, Uniform distribution and Gaussian distribution to form the joint distribution, and samples noise code Z from this joint distribution. According to the previous experimental results, we form the joint distributions of Cauchy and Exponential, Cauchy and Exponential and Uniform, Cauchy and Exponential and Uniform and Gaussian to construct the latent space. In this case, the training data is still the extremely limited samples, and the generated results are shown in Table.7. From Table.7, we can see that the best performance is still in the case of number=2, larger number resulting in identical instances. We compare the FID score of combination of Cauchy and Exponential (FID_{CE}) with that of combination of Gaussian and Uniform (FID_{GU}), we get FID_{CE} on MNIST and CIFAR10 are 101.4 and 352.1 respectively, which are larger than FID_{GU} 66.2 and 310.8. In this way, the

Models	FID Score
DCGAN	389.6
DeLiGAN	381.9
LSGAN	436.7
WGAN	444.3
RaGAN	453.3
MGAN	456.7
JDGAN	310.8

Table 6: FID scores of generated data shown in Fig.9.

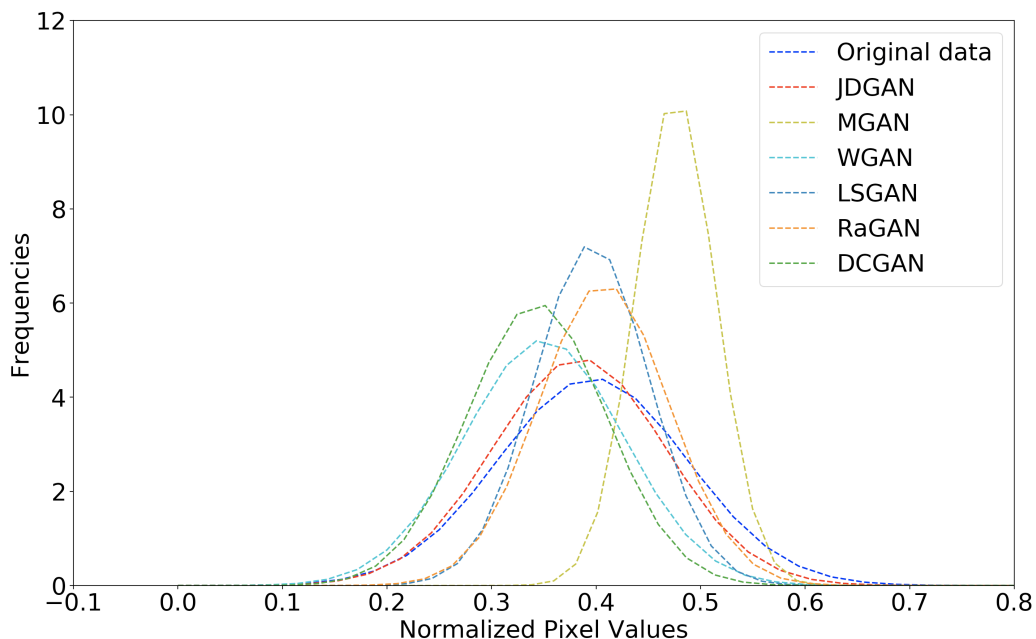


Figure 10: The overlapping area on CIFAR-10 dataset. The blue line indicates the original data distribution, other lines indicate the simulation data distribution, which are generated by JDGAN, MGAN, WGAN, LSGAN, RaGAN and DCGAN respectively. Obviously, JDGAN captures the distribution with the highest similarity to that of the original data.

combination of Uniform and Gaussian could be the better choice in our cases.

5.5. Mass Spectrometry Dataset

Protein is highly complex biochemical entity, and is present in all living organisms. Proteins are the downstream products from genes which

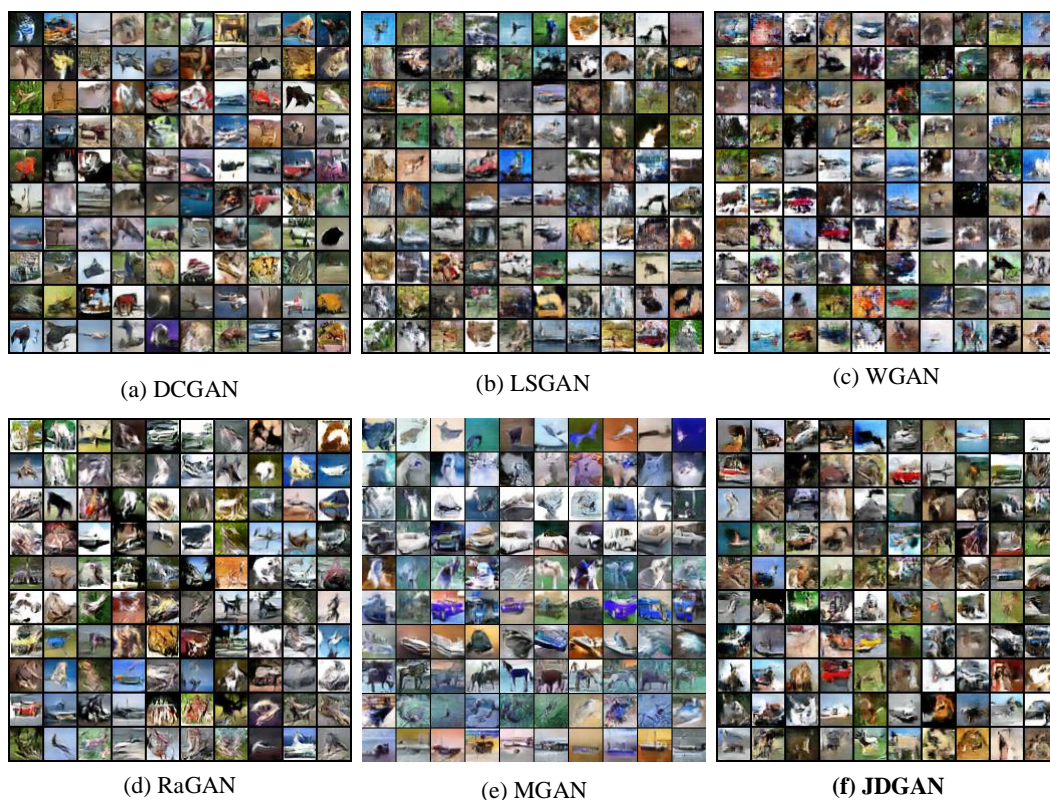


Figure 11: Generated images produced by GAN variants and JDGAN on adequate CIFAR10 training samples.

carry out almost all essential biological and chemical functions in a living body. Primary structure of a protein is determined by the sequence of specific amino acids (peptide). In recent years, mass spectrometry has become the core analytical technique in studying protein identification and quantification [38]. Mass spectrometry works by ionizing chemical compounds to generate charged molecules or molecule fragments and measuring their mass to charge ratios (m/z). The generated peak spectra reflect the identification and abundance of a specific protein.

In a typical tandem mass spectrometry (MS/MS) experiment, protein mixture is digested by an enzyme into smaller peptides. The sample peptides are then ionized and subjected to MS. As a medical dataset, the main problem is that enough available data collection could be expensive and unrealistic (here we just hold 5 patients and the number of peptide in each patient

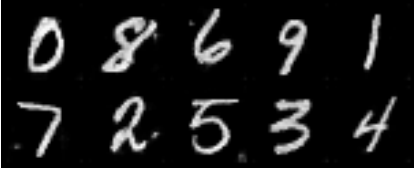
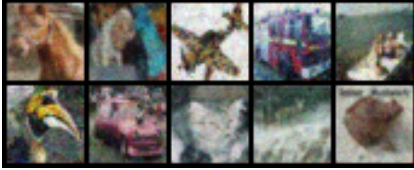
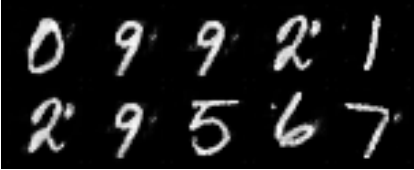
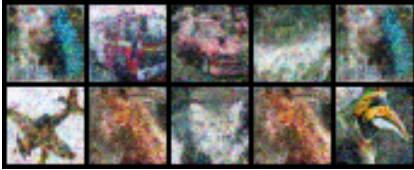
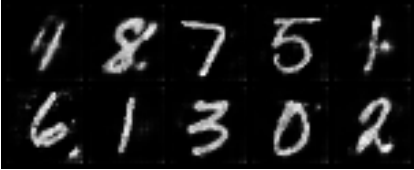
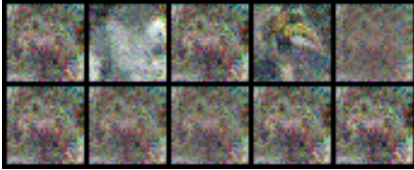
Number	Simulation Data on MNIST	Simulation Data on CIFAR10
2		
3		
4		

Table 7: Number=2 indicates that one single Cauchy and one single Exponential are employed; number=3 indicates that one single Cauchy, one single Exponential and one single Uniform are employed; number=4 indicates that one single Cauchy, one single Exponential, one single Uniform and one single Gaussian are employed.

is 74552). Thus, we hope to use the generative model to generate more simulation data to supplement the original data. We then apply our proposed JDGAN (Eq.(6)) to MS dataset, and the hyper parameters are shown in the right part of Fig.5. Although there are many features for each patient, we just pick up 6 features (Acq Time, Intensity, Precursor Intensity Acquisition, Apex Time, Elution Peak Width, MS2Counts) because the 6 features correlate to the differential protein expressions under different biological conditions. In most of the cases, between healthy and disease conditions.

The results are shown in Fig.12. In Fig.12, the blue line indicates the threshold, which is a group of mean values calculated by the healthy proteomic values and diseased proteomic values. In most cases, a proteomic value is below the threshold if this proteomic value belongs to the healthy set, while a proteomic value is larger than the threshold as the proteomic value belongs to the diseased set. The red points and green points are simulation data generated by JDGAN, we can see that most simulation data are below the threshold as they belong to healthy set and that are above the

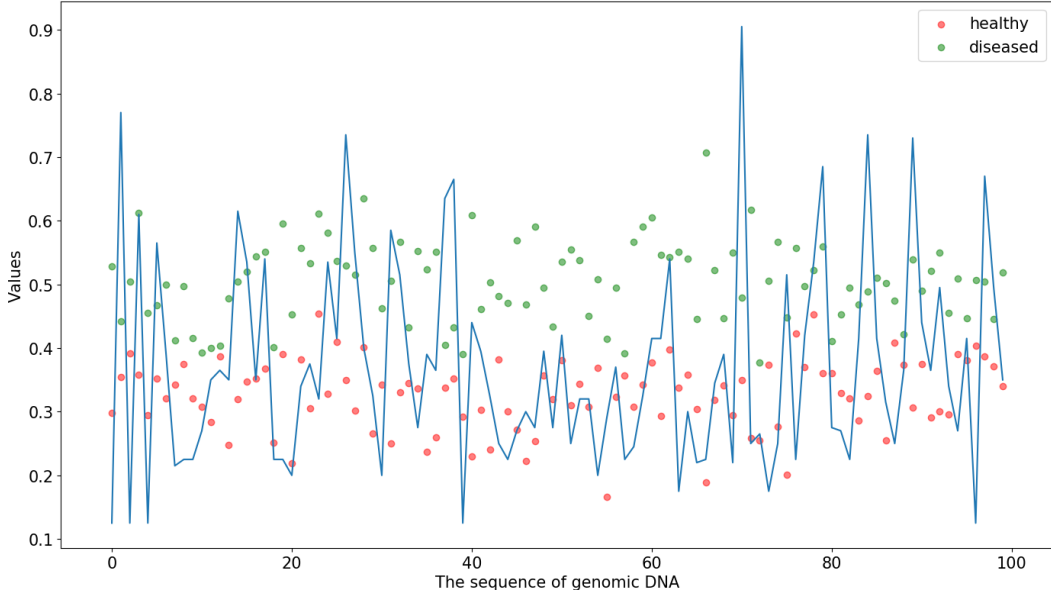


Figure 12: The green points and red points are belonging to the simulated data generated by JDGAN. The blue line indicates a threshold calculated by original proteomic values included by peptide, which can distinguish whether the peptides within a biological tissue is healthy or not. The tissue is healthy if the corresponding proteomic value is below the threshold, or diseased as the proteomic value is larger than the threshold.

threshold as they belong to the diseased set.

5.6. Discussion

In fact, it is hard to obtain the original data distribution. It becomes even harder for extremely limited data with rich information (e.g., many categories, color or shape information, pixel details). Let $p_z \dashrightarrow z$ and $p_r \dashrightarrow x$ where the dash line indicates model drawing samples from a certain distribution. In the traditional GAN generating mechanism, G always specializes in $z \mapsto x$. Considering $p_z(z)$ is contained in p_G , the generated data distribution cannot be complicated to match the original data distribution completely. Under such a scenario, the generated data diversity cannot be guaranteed. See Fig.2 and Fig.9. JDGAN increases the diversity of generated instances by employing the joint distribution to construct the randomized space. See Table.1, Table.2, Table.4 and Table.5.

Note that the crucial point of JDGAN is how to determine the number of distributions. In this paper, we utilize the *Intrinsic Dimension* strategy

to evaluate the *Intrinsic Dimensions* of original data and generated data, because the *Intrinsic Dimension* can loyally reflect the minimal dimension of data manifold [26]. By calculating the *Intrinsic Dimensions* of original data and generated data (Fig.2 (a)) with Eq.(8), we get $\hat{m}(x)_{raw\ data} = 5.79$ and $\hat{m}(x)_{generated\ data} = 3.30$. This suggests that two distributions may be the optimal choice in our cases, and greater or less than this number would decrease the performance. We still use Eq.(8) to calculate the *Intrinsic Dimension* of synthetic data produced by JDGAN (Table.5 row1), and get $\hat{m}(x)_{synthetic\ data} = 5.73$, which is very close to that of raw data. Under such a scenario, the generated distribution is closer to original data distribution (See Table.3 and Table.6) and the diversity of generated data can also be guaranteed (See Table.1, Table.2, Table.4, Table.5 and Table.7).

6. Conclusion

In this paper, instead of modifying the framework of a GAN model, we explore another direction, which is to modify the randomized space by jointly employing multiple distributions and learn these distributions together with the generator with reparameterization method, to increase the diversity of simulated data and the overlapping area between the raw data distribution and the generating distribution. Also, we explore how to combine these distributions and display the different performance on different numbers of distributions. Specifically, we jointly employ only multiple Gaussian distributions, only multiple Uniform distributions and the mixture of two distributions respectively, and test each performance. We conducted extensive experiments with MNIST and CIFAR-10 as well as mass spectrometry datasets to validate our proposed JDGAN. The results have shown that our approach is effective and better than other GAN-based models.

Acknowledgment

This work was supported in part by the National Key R&D Program of China (No.2018YFC1604000), in part by the RGC Collaborative Research Fund (CRF) 2018/19 - Group Research Grant (No.:C5026-18G), and in part by the RGC General Research Fund GRF (No.15102015/15M).

- [1] Luke Metz Alec Radford and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. 2015.

- [2] Laurent Amsaleg, Oussama Chelly, Teddy Furon, Stéphane Girard, Michael E Houle, Ken-ichi Kawarabayashi, and Michael Nett. Estimating local intrinsic dimensionality. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 29–38, 2015.
- [3] Martin Arjovsky and Léon Bottou. Towards principled methods for training generative adversarial networks. *arXiv preprint arXiv:1701.04862*, 2017.
- [4] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017.
- [5] Sanjeev Arora, Rong Ge, Yingyu Liang, Tengyu Ma, and Yi Zhang. Generalization and equilibrium in generative adversarial nets (gans). In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 224–232. JMLR. org, 2017.
- [6] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010*, pages 177–186. Springer, 2010.
- [7] Miguel A. Carreira-Perpinan. Mode-finding for mixtures of gaussian distributions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(11):1318–1323, 2000.
- [8] Constantinos Daskalakis, Paul W. Goldberg, and Christos H. Papadimitriou. *The complexity of computing a Nash equilibrium*. ACM, 2009.
- [9] Dominik Maria Endres and Johannes E Schindelin. A new metric for probability distributions. *IEEE Transactions on Information theory*, 2003.
- [10] Jesse Engel, Kumar Krishna Agrawal, Shuo Chen, Ishaan Gulrajani, Chris Donahue, and Adam Roberts. Gansynth: Adversarial neural audio synthesis. *The International Conference on Learning Representations*, 2019.
- [11] Arnab Ghosh, Viveka Kulharia, Vinay P Namboodiri, Philip HS Torr, and Puneet K Dokania. Multi-agent diverse generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8513–8521, 2018.

- [12] Bordes A. Glort, X. and Y. Bengio. Deep sparse rectifier neural networks. *International Conference on Artificial Intelligence and Statistics*, pages 315–323, 2011.
- [13] 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.
- [14] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *Advances in Neural Information Processing Systems*, pages 5769–5779, 2017.
- [15] Swaminathan Gurumurthy, Ravi Kiran Sarvadevabhatla, and R Venkatesh Babu. Deligan: Generative adversarial networks for diverse and limited data. In *CVPR*, pages 4941–4949, 2017.
- [16] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in neural information processing systems*, pages 6626–6637, 2017.
- [17] CC Heyde. Central limit theorem. *Encyclopedia of Actuarial Science*, 2006.
- [18] Quan Hoang, Tu Dinh Nguyen, Trung Le, and Dinh Phung. Mgan: Training generative adversarial nets with multiple generators. 2018.
- [19] Sicong Huang, Qiyang Li, Cem Anil, Xuchan Bao, Sageev Oore, and Roger B Grosse. Timbretron: A wavenet (cyclegan (cqt (audio))) pipeline for musical timbre transfer. *The International Conference on Learning Representations*, 2019.
- [20] Alexia Jolicoeur-Martineau. The relativistic discriminator: a key element missing from standard gan. *The International Conference on Learning Representations*, 2019.
- [21] Mahyar Khayatkhoei, Maneesh K Singh, and Ahmed Elgammal. Disconnected manifold learning for generative adversarial networks. In

- Advances in Neural Information Processing Systems*, pages 7343–7353, 2018.
- [22] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
 - [23] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
 - [24] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
 - [25] Dirk P Kroese, Reuven Y Rubinstein, Izack Cohen, Sergey Porotsky, and Thomas Taimre. Cross-entropy method. In *Encyclopedia of Operations Research and Management Science*, pages 326–333. Springer, 2013.
 - [26] Elizaveta Levina and Peter J Bickel. Maximum likelihood estimation of intrinsic dimension. In *Advances in neural information processing systems*, pages 777–784, 2005.
 - [27] Mu Li, Tong Zhang, Yuqiang Chen, and Alexander J Smola. Efficient mini-batch training for stochastic optimization. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 661–670. ACM, 2014.
 - [28] Wei Li, Wei Ding, Rajani Sadasivam, Xiaohui Cui, and Ping Chen. Hisgan: A histogram-based gan model to improve data generation quality. *Neural Networks*, 119:31–45, 2019.
 - [29] Wei Li, Li Fan, Zhenyu Wang, Chao Ma, and Xiaohui Cui. Tackling mode collapse in multi-generator gans with orthogonal vectors. *Pattern Recognition*, 110:107646.
 - [30] Wei Li, Linchuan Xu, Zhixuan Liang, Senzhang Wang, Jiannong Cao, Chao Ma, and Xiaohui Cui. Sketch-then-edit generative adversarial network. *Knowledge-Based Systems*, page 106102, 2020.
 - [31] Yijun Li, Sifei Liu, Jimei Yang, and Ming-Hsuan Yang. Generative face completion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, volume 1, page 6, 2017.

- [32] Mario Lucic, Karol Kurach, Marcin Michalski, Sylvain Gelly, and Olivier Bousquet. Are gans created equal? a large-scale study. In *Advances in neural information processing systems*, pages 700–709, 2018.
- [33] Christopher D Manning and Hinrich Schütze. *Foundations of statistical natural language processing*. MIT press, 1999.
- [34] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2813–2821. IEEE, 2017.
- [35] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- [36] James R Norris. *Markov chains*. Number 2. Cambridge university press, 1998.
- [37] Ludger Rüschendorf. The wasserstein distance and approximation theorems. *Probability Theory and Related Fields*, 70(1):117–129, 1985.
- [38] Eduard Sabidó, Nathalie Selevsek, and Ruedi Aebersold. Mass spectrometry-based proteomics for systems biology. *Current opinion in biotechnology*, 23(4):591–597, 2012.
- [39] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *Advances in Neural Information Processing Systems*, pages 2234–2242, 2016.
- [40] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [41] T Tieleman and G Hinton. Rmsprop: Divide the gradient by a running average of its recent magnitude. coursera: Neural networks for machine learning. *Tech. Rep., Technical report*, page 31, 2012.
- [42] Ilya O Tolstikhin, Sylvain Gelly, Olivier Bousquet, Carl-Johann Simon-Gabriel, and Bernhard Schölkopf. Adagan: Boosting generative models.

- In *Advances in Neural Information Processing Systems*, pages 5424–5433, 2017.
- [43] Curtis R Vogel. *Computational methods for inverse problems*, volume 23. Siam, 2002.
- [44] Jia Wang, Jiannong Cao, Wei Li, and Senzhang Wang. Cane: community-aware network embedding via adversarial training. *Knowledge and Information Systems*, pages 1–28, 2020.
- [45] Jia Wang, Jiannong Cao, Senzhang Wang, Zhongyu Yao, and Wengen Li. Irda: Incremental reinforcement learning for dynamic resource allocation. *IEEE Transactions on Big Data*, 2020.
- [46] Senzhang Wang, Jiannong Cao, and Philip Yu. Deep learning for spatio-temporal data mining: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 2020.
- [47] Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. Empirical evaluation of rectified activations in convolutional network. *Computer Science*, 2015.
- [48] Guoshen Yu, Guillermo Sapiro, and Stéphane Mallat. Solving inverse problems with piecewise linear estimators: From gaussian mixture models to structured sparsity. *IEEE Transactions on Image Processing*, 21(5):2481–2499, 2012.
- [49] Kai Yu, Tong Zhang, and Yihong Gong. Nonlinear learning using local coordinate coding. In *Advances in neural information processing systems*, pages 2223–2231, 2009.
- [50] Ning Zhang, Evan Shelhamer, Yang Gao, and Trevor Darrell. Fine-grained pose prediction, normalization, and recognition. *arXiv preprint arXiv:1511.07063*, 2015.