

ClusterGAN : Latent Space Clustering in Generative Adversarial Networks

Sudipto Mukherjee¹, Himanshu Asnani¹, Eugene Lin¹, Sreeram Kannan¹,

1 University of Washington, Seattle.

{sudipm, asnani, lines, ksreeram}@uw.edu

Abstract

Generative Adversarial networks (GANs) have obtained remarkable success in many unsupervised learning tasks and unarguably, clustering is an important unsupervised learning problem. While one can potentially exploit the latent-space back-projection in GANs to cluster, we demonstrate that the cluster structure is not retained in the GAN latent space. In this paper, we propose ClusterGAN as a new mechanism for clustering using GANs. By sampling latent variables from a mixture of one-hot encoded variables and continuous latent variables, coupled with an inverse network (which projects the data to the latent space) trained jointly with a clustering specific loss, we are able to achieve clustering in the latent space. Our results show a remarkable phenomenon that GANs can preserve latent space interpolation across categories, even though the discriminator is never exposed to such vectors. We compare our results with various clustering baselines and demonstrate superior performance on both synthetic and real datasets.

1 Introduction

1.1 Motivation

Representation learning enables machine learning models to decipher underlying semantics in data and disentangle hidden factors of variation. These powerful representations have made it possible to transfer knowledge across various tasks. But what makes one representation better than another ? [2] mentioned several general-purpose priors that are not dependent on the downstream task, but appear as commonalities in good representations. One of the general-purpose priors of representation learning that is ubiquitous across data intensive domains is clustering. Clustering has been extensively studied in unsupervised learning with multifarious approaches seeking efficient algorithms [24], problem specific

distance metrics [29], validation [12] and the like. Even though the main focus of clustering has been to separate out the original data into classes, it would be even nicer if such clustering was obtained along with dimensionality reduction where the real data actually seems to come from a lower dimensional manifold.

In recent times, much of unsupervised learning is driven by deep generative approaches, the two most prominent being Variational Autoencoder (VAE) [17] and Generative Adversarial Network (GAN) [9]. The popularity of generative models themselves is hinged upon the ability of these models to capture high dimensional probability distributions, imputation of missing data and dealing with multimodal outputs. Both GAN and VAE aim to match the real data distribution (VAE using an explicit approximation of maximum likelihood and GAN through implicit sampling), and simultaneously provide a mapping from a latent space \mathcal{Z} to the input space \mathcal{X} . The latent space of GANs not only provides dimensionality reduction, but also gives rise to novel applications. Perturbations in the latent space could be used to determine adversarial examples that further help build robust classifiers [14]. Compressed sensing using GANs [3] relies on finding a latent vector that minimizes the reconstruction error for the measurements. Generative compression is yet another application involving \mathcal{Z} [26]. One of the most fascinating outcomes of the GAN training is the interpolation in the latent space. Simple vector arithmetic properties emerge which when manipulated lead to changes in the semantic qualities of the generated images [25]. This differentiates GANs from traditional dimensionality reduction techniques [22] [20] which lack interpretability. One potential application that demands such a property is clustering of cell types in genomics. GANs provide a means to understand the change in high-dimensional gene expression as one traverses from one cell type (i.e., cluster) to another in the latent space. Here, it is critical to have both clustering as well as good interpretability and interpolation ability. This brings us to the principal motivation of this work: ***Can we design a GAN training methodology that clusters in the latent space?***

1.2 Related Works

Deep learning approaches have been used for dimensionality reduction starting with variants of the autoencoder such as the stacked denoising autoencoders [28], sparse autoencoder [5] and deep CCA [1]. Architectures for deep unsupervised subspace clustering have also been built on the encoder-decoder framework [15]. Recent works have addressed this problem of joint clustering and dimensionality reduction in autoencoders. [30] solved this problem by initializing the cluster centroids and the embedding with a stacked autoencoder. Then they use alternating optimization to improve the clustering and report state-of-the-art results in both clustering accuracy and speed on real datasets. The clustering algorithm is referred to as DEC in their paper. Since K-means is often the most widely used algorithm for clustering, [31] improved upon DEC by introducing a modified cost function that incorporates the K-means loss. They optimized the non-convex objective using alternating SGD to obtain an embedding that is amenable to K-means clustering. Their algorithm DCN was shown to outperform

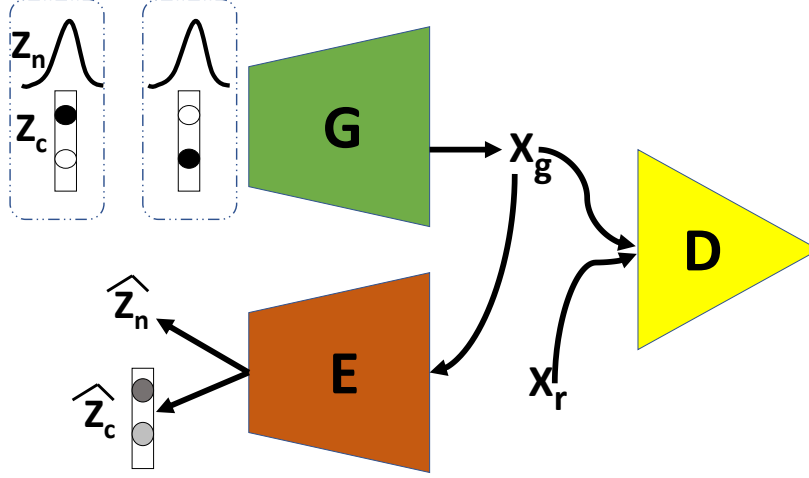


Figure 1: ClusterGAN Architecture

all standard clustering methods on a range of datasets. It is interesting to note that the vanilla autoencoder by itself did not explicitly have any clustering objective. But it could be improved to achieve this end by careful algorithmic design. Since GANs have outperformed autoencoders in generating high fidelity samples, we had a strong intuition in favour of the powerful latent representations of GAN providing improved clustering performance also.

Interpretable representation learning in the latent space has been investigated for GANs in the seminal work of [4]. The authors trained a GAN with an additional term in the loss that seeks to maximize the mutual information between a subset of the generator’s noise variables and the generated output. The key goal of InfoGAN is to create interpretable and disentangled latent variables. While InfoGAN does employ discrete latent variables, it is not specifically designed for clustering. In this paper, we show that our proposed architecture is superior to InfoGAN for clustering. The other prominent family of generative models, VAE, has the additional advantage of having an inference network, the encoder, which is jointly learnt during training. This enables mapping from \mathcal{X} to \mathcal{Z} that could potentially preserve cluster structure by suitable algorithmic design. Unfortunately, no such inference mechanism exists in GANs, let alone the possibility of clustering in the latent space. To bridge the gap between VAE and GAN, various methods such as Adversarially Learned Inference (ALI) [8], Bidirectional Generative Adversarial Networks (BiGAN) [7] have introduced an inference network which is trained to match the joint distributions of (x, z) learnt by the encoder \mathcal{E} and decoder \mathcal{G} networks. Typically, the reconstruction in ALI/BiGAN is poor as there is no deterministic pointwise matching between x and $\mathcal{G}(\mathcal{E}(x))$ involved in the training. Architectures such as Wasserstein Autoencoder [27], Adversarial Autoencoder [21], which depart from the traditional

GAN framework, also have an encoder as part of the network. So this led us to consider a formulation using an Encoder which could *both reduce the cycle loss as well as aid in clustering*.

1.3 Main Contributions

To the best of our knowledge, this is the first work that addresses the problem of clustering in the latent space of GAN. The main contributions of the paper can be summarized as follows:

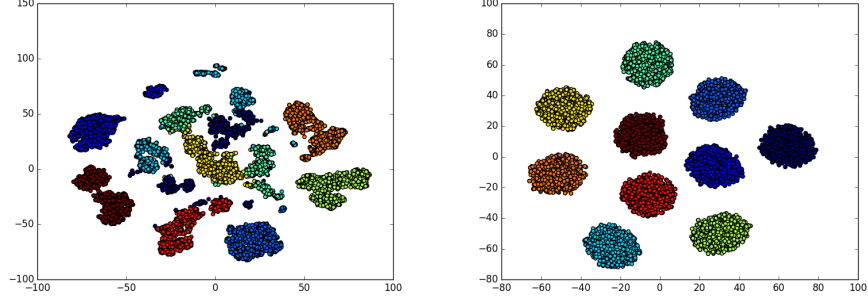
- We show that even though the GAN latent variable preserves information about the observed data, the latent points are smoothly scattered based on the latent distribution leading to no observable clusters.
- We propose three main algorithmic ideas in ClusterGAN in order to remedy this situation.
 1. We utilize a **mixture of discrete and continuous** latent variables in order to create a non-smooth geometry in the latent space.
 2. We propose a **novel backpropagation algorithm** accommodating the discrete-continuous mixture, as well as an **explicit inverse-mapping network** to obtain the latent variables given the data points, since the problem is non-convex.
 3. We propose to jointly train the GAN along with the inverse-mapping network with a **clustering-specific loss** so that the distance geometry in the projected space reflects the distance-geometry of the variables.
- We compare ClusterGAN and other possible GAN based clustering algorithms, such as InfoGAN, along with multiple clustering baselines on varied datasets. This demonstrates the superior performance of ClusterGAN for the clustering task.
- We demonstrate that ClusterGAN surprisingly retains good interpolation across the different classes (encoded using one-hot latent variables), even though the discriminator is never exposed to such samples.

The formulation is general enough to provide a *meta* framework that incorporates the additional desirable property of clustering in GAN training.

2 Discrete-Continuous Prior

2.1 Background

Generative adversarial networks consist of two components, the generator \mathcal{G} and the discriminator \mathcal{D} . Both \mathcal{G} and \mathcal{D} are usually implemented as neural



(a) Non-linear generator with $z \sim \mathcal{N}(0, I)$ (b) Linear generator with z one-hot encoded

Figure 2: TSNE visualization of latent vectors. Linear Generator recovers clusters, suggesting that representation power is not a bottleneck.

networks parameterized by Θ_G and Θ_D respectively. The generator can also be considered to be a mapping from latent space to the data space which we denote as $\mathcal{G} : \mathcal{Z} \mapsto \mathcal{X}$. The discriminator defines a mapping from the data space to a real value which can correspond to the probability of the sample being real, $\mathcal{D} : \mathcal{X} \mapsto \mathbb{R}$. The GAN training sets up a two player game between \mathcal{G} and \mathcal{D} , which is defined by the minimax objective : $\min_{\Theta_G} \max_{\Theta_D} \mathbf{E}_{x \sim \mathbb{P}_x^r} q(\mathcal{D}(x)) + \mathbf{E}_{z \sim \mathbb{P}_z} q(1 - \mathcal{D}(\mathcal{G}(z)))$, where \mathbb{P}_x^r is the distribution of real data samples, \mathbb{P}_z is the prior noise distribution on the latent space and $q(\cdot)$ is the quality function. For vanilla GAN, $q(x) = \log x$, and for Wasserstein GAN (WGAN) $q(x) = x$. We also denote the distribution of generated samples x_g as \mathbb{P}_x^g . The discriminator and the generator are optimized alternatively so that at the end of training \mathbb{P}_x^g matches \mathbb{P}_x^r .

2.2 Vanilla GAN does not cluster well in the latent space

One possible way to cluster using a GAN is to back-propagate the data into the latent space (using back-propagation decoding [19]) and cluster the latent space. However, this method usually leads to very bad results (see Fig. 3 for clustering results on MNIST). The key reason is that, if indeed, back-propagation succeeds, then the back-projected data distribution should look similar to the latent space distribution, which is typically chosen to be a Gaussian or uniform distribution, and we cannot expect to cluster in that space. Thus even though the latent space may contain full information about the data, the distance geometry in the latent space does not reflect the inherent clustering. In [11], the authors sampled from a Gaussian mixture prior and obtained diverse samples even in limited data regimes. However, even GANs with a Gaussian mixture failed to cluster, as shown in 3(c). As observed by the authors of DeLiGAN, Gaussian components tend to ‘crowd’ and become redundant. Lifting the space using categorical variables could solve the problem effectively. But continuity in latent space is traditionally viewed to be a pre-requisite for the objective of good

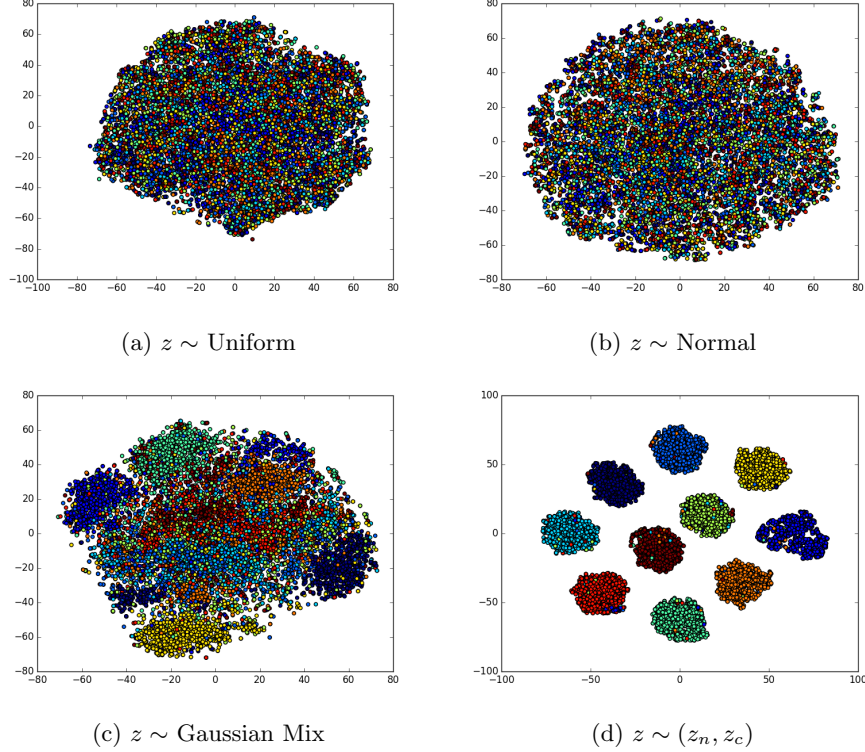


Figure 3: TSNE visualization of latent vectors for GANs trained with different priors on MNIST.

interpolation. In other words, interpolation seems to be at loggerheads with the clustering objective. We demonstrate in this paper how ClusterGAN can obtain **good interpolation and good clustering** simultaneously.

2.3 Sampling from Discrete-Continuous Mixtures

In ClusterGAN, we sample from a prior that consists of normal random variables cascaded with one-hot encoded vectors. To be more precise $z = (z_n, z_c)$, $z_n \sim \mathcal{N}(0, \sigma^2 I_{d_n})$, $z_c = e_k$, $k \sim \mathcal{U}\{1, 2, \dots, K\}$, e_k is the k^{th} elementary vector in \mathbb{R}^K and K is the number of clusters in the data. In addition, we need to choose σ in such a way that the one-hot vector provides sufficient signal to the GAN training that leads to each mode only generating samples from a corresponding class in the original data. To be more precise, we chose $\sigma = 0.10$ in all our experiments so that each dimension of the normal latent variables, $z_{n,j} \in (-0.6, 0.6) \ll 1.0 \forall j$ with high probability. Small variances σ are chosen to ensure the clusters in \mathcal{Z} space are separated. Hence this prior naturally enables us to design an algorithm that clusters in the latent space.

2.4 Modified Backpropagation Based Decoding

Previous works [6] [19] have explored solving an optimization problem in z to recover the latent vectors, $z^* = \arg \min_z \mathcal{L}(\mathcal{G}(z), x) + \lambda \|z\|_p$, where \mathcal{L} is some suitable loss function and $\|\cdot\|_p$ denotes the norm. This approach is insufficient for clustering with traditional latent priors even if backpropagation was lossless and recovered accurate latent vectors. To make the situation worse, the optimization problem above is non-convex in z (\mathcal{G} being implemented as a neural network) and can obtain different embeddings in the \mathcal{Z} space based on initialization. Some of the approaches to address this issue could be multiple restarts with different initialiations to obtain z^* , or stochastic clipping of z at each iteration step. None of these lead to clustering, since they do not address the root problem of sampling from separated manifolds in \mathcal{Z} . But our sampling procedure naturally gives way to such an algorithm. We use $\mathcal{L}(\mathcal{G}(z), x) = \|\mathcal{G}(z) - x\|_1$. Since we sample from a normal distribution, we use the regularizer $\|z_n\|_2^2$, penalizing only the normal variables. We use K restarts, each sampling z_c from a different one-hot component and optimize with respect to only the normal variables, keeping z_c fixed. Adam [16] is used for the updates during Backprop decoding. Formally, Algorithm 1 summarizes the approach.

Input: Real sampler x , Generator function \mathcal{G} , Number of Clusters K ,
Regularization parameter λ , Adam iterations τ
Output: Latent embedding z^*
for $k \in \{1, 2, \dots, K\}$ **do**
 Sample $z_n^0 \sim \mathcal{N}(0, \sigma^2 I_{d_n})$
 Initialization $z_k^0 \leftarrow (z_n^0, e_k)$ (e_k is k^{th} elementary unit vector in K dimensions)
 for $t \in \{1, 2, \dots, \tau\}$ **do**
 Obtain the gradient of loss function
 $g \leftarrow \nabla_{z_n} (\|\mathcal{G}(z_k^{t-1}) - x\|_1 + \lambda \|z_n^{t-1}\|_2^2)$
 Update z_n^t using g with Adam iteration to minimize loss.
 Clipping of z_n^t , i.e., $z_n^t \leftarrow \mathcal{P}_{[-0.6, 0.6]}(z_n^t)$
 $z_k^t \leftarrow (z_n^t, e_k)$
 end
 Update z^* if z_k^τ has lowest loss obtained so far.
end
return z^*

Algorithm 1: DECODE_LATENT

2.5 Linear Generator clusters perfectly

The following lemma suggests that with discrete-continuous mixtures, we need only linear generation to generate mixture of Gaussians in the generated space.

Lemma 1. *Clustering with only z_n cannot recover a mixture of gaussian data in the linearly generated space. Further \exists a linear $G(\cdot)$ mapping discrete-continuous*

mixtures to a mixture of Gaussians.

Proof. If latent space has only the continuous part, $z_n \sim \mathcal{N}(0, \sigma^2 I_{d_n})$, then by the linearity property, any linear generation can only produce Gaussian in the generated space. Now we show there exists a $G(\cdot)$ mapping discrete-continuous mixtures to the generate data $X \sim \mathcal{N}(\mu_\omega, \sigma^2 I_{d_n})$, where $\omega \sim \mathcal{U}\{1, 2, \dots, K\}$ (K is the number of mixtures). This is possible if we let $z_n \sim \mathcal{N}(0, \sigma^2 I_{d_n})$, $z_c = e_k$, $k \sim \mathcal{U}\{1, 2, \dots, K\}$ and $G(z_n, z_c) = z_n + Az_c$, $A = \text{diag}[\mu_1, \dots, \mu_K]$ being a $K \times K$ diagonal matrix with diagonal entries as the means μ_i . \square

To illustrate this lemma, and hence the drawback of traditional priors \mathbb{P}_z for clustering, we performed a simple experiment. The real samples are drawn from a mixture of 10 Gaussians in \mathbb{R}^{100} . The means of the Gaussians are sampled from $\mathcal{U}(-0.3, 0.3)^{100}$ and the variance of each component is fixed at $\sigma = 0.12$. We trained a GAN with $z \sim \mathcal{N}(0, I_{10})$ where the generator is a multi-layer perceptron with two hidden layers of 256 units each. For comparison, we also trained a GAN with z sampled from one-hot encoded normal vectors, the dimension of categorical variable being 10. The generator for this GAN consisted of a linear mapping $W \in \mathbb{R}^{100 \times 10}$, such that $x = Wz$. After training, the latent vectors are recovered using Algorithm 1 for the linear generator, and 10 restarts with random initializations for the non-linear generator. Even for this toy setup, the linear generator perfectly clustered the latent vectors (Acc. = 1.0, NMI = 1.0, ARI = 1.0), but the non-linear generator performed poorly (Acc. = 0.73, NMI = 0.75, ARI = 0.60) (Figure 2). The situation becomes worse for real datasets such as MNIST when we trained a GAN using latent vectors drawn from uniform, normal or a mixture of Gaussians. None of these configurations succeeded in clustering in the latent space as shown in Figure 3.



Figure 4: Fashion items generated from distinct modes : Fashion-MNIST

2.6 Separate Modes for distinct classes in the data

It was surprising to find that trained in a purely unsupervised manner without additional loss terms, each one-hot encoded component generated points from a specific class in the original data. For instance, $z = (z_n, e_k)$ generated a particular digit $\pi(k)$ in MNIST, for multiple samplings of $z_n \sim \mathcal{N}(0, \sigma^2 I_{d_n})$ (π denotes a permutation). This was a necessary first step for the success of Algorithm 1. We also quantitatively evaluated the modes learnt by the GAN

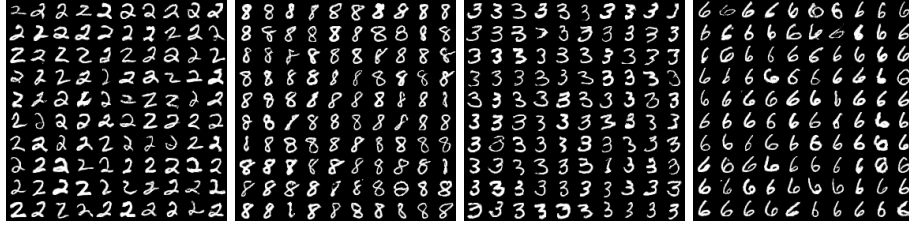


Figure 5: Digits generated from distinct modes : MNIST

by using a supervised classifier for MNIST. The supervised classifier had a test accuracy of 99.2%, so it had high reliability of distinguishing the digits. We sample from a mode k and generate a digit x_g . It is then classified by the classifier as \hat{y} . From this pair (k, \hat{y}) , we can map each mode to a digit and compute the accuracy of digit \hat{y} being generated from mode k . This is denoted as Mode Accuracy. Each digit sample x_r with label y can be decoded in the latent space by Algorithm 1 to obtain z . Now z can be used to generate x_g , which when passed through the classifier gives the label \hat{y} . The pair (y, \hat{y}) must be equal in the ideal case and this accuracy is denoted as Reconstruction Accuracy. Finally, all the mappings of points in the same class in \mathcal{X} space should have the same one-hot encoding when embedded in \mathcal{Z} space. This defines the Cluster Accuracy. This methodology can be extended to quantitatively evaluate mode generation for other datasets also, provided there is a reliable classifier. For MNIST, we obtained Mode Accuracy of 0.97, Reconstruction Accuracy of 0.96 and Cluster Accuracy of 0.95. Some of the modes in Fashion-MNIST and MNIST are shown in Figures 4 and 5, respectively. Supplementary materials contain the images from all modes in these two datasets.

2.7 Interpolation in latent space is preserved

The latent space in a traditional GAN with Gaussian latent distribution enforces that different classes are continuously scattered in the latent space, allowing nice inter-class interpolation, which is a key strength of GANs. In ClusterGAN, the latent vector z_c is sampled with a one-hot distribution and in order to interpolate across the classes, we will have to sample from a convex combination on the one-hot vector. While these vectors have never been sampled during the training process, we surprisingly observed very smooth inter-class interpolation in ClusterGAN. To demonstrate interpolation, we fixed the z_n in two latent vectors with different z_c components, say $z_c^{(1)}$ and $z_c^{(2)}$ and interpolated with the one-hot encoded part to give rise to new latent vectors $z = (z_n, \mu z_c^{(1)} + (1 - \mu) z_c^{(2)})$, $\mu \in [0, 1]$. As Figure 6 illustrates, we observed a nice transition from one digit to another as well as across different classes in FashionMNIST. This demonstrates that ClusterGAN learns a very smooth manifold even on the untrained directions of the discrete-continuous distribution. We also show interpolations from a vanilla GAN trained with Gaussian prior as reference.

Input: Functions \mathcal{G} , \mathcal{D} and \mathcal{E} , Regularization parameters β_n , β_c , learning rate η , parameters Θ_G^t , Θ_E^t

Output: $\Theta_G^{(t+1)}$, $\Theta_E^{(t+1)}$

Sample $z_{i=1}^m$ from \mathbb{P}^z , $z = (z_n, z_c)$

$g_{\Theta_G} \leftarrow$

$$\nabla_{\Theta_G} \left(- \sum_{i=1}^m q(\mathcal{D}(\mathcal{G}(z^{(i)}))) + \beta_n \sum_{i=1}^m \|z_n^{(i)} - \mathcal{E}(\mathcal{G}(z_n^{(i)}))\|_2^2 + \beta_c \sum_{i=1}^m \mathcal{H}(z_c^{(i)}, \mathcal{E}(\mathcal{G}(z_c^{(i)}))) \right)$$

$g_{\Theta_E} \leftarrow \nabla_{\Theta_E} \left(\beta_n \sum_{i=1}^m \|z_n^{(i)} - \mathcal{E}(\mathcal{G}(z_n^{(i)}))\|_2^2 + \beta_c \sum_{i=1}^m \mathcal{H}(z_c^{(i)}, \mathcal{E}(\mathcal{G}(z_c^{(i)}))) \right)$

Update Θ_G using $(g_{\Theta_G}, \Theta_G^t)$ with Adam ; similarly for Θ_E .

return Θ_G, Θ_E

Algorithm 2: UPDATE_PARAM

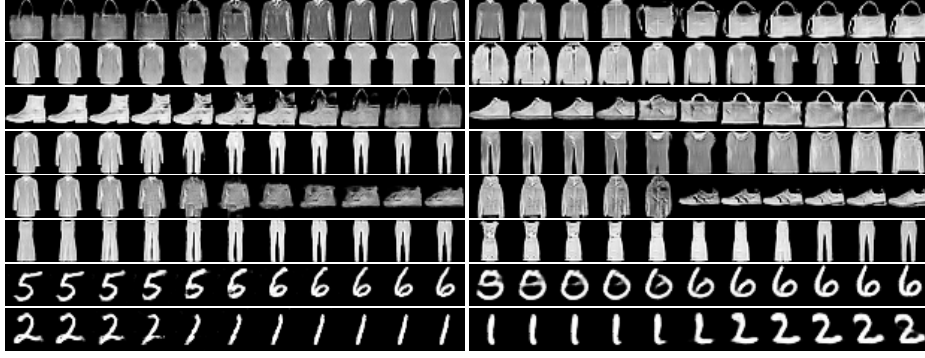


Figure 6: Comparison of Latent Space Interpolation : ClusterGAN (left) and vanilla WGAN (right)

3 ClusterGAN

Even though the above approach enables the GAN to cluster in the latent space, it may be able to perform even better if we had a clustering specific loss term in the minimax objective. For MNIST, digit strokes correspond well to the category in the data. But for more complicated datasets, we need to enforce structure in the GAN training. One way to ensure that is to enforce precise recovery of the latent vector. We therefore introduce an encoder $\mathcal{E} : \mathcal{X} \mapsto \mathcal{Z}$, a neural network parameterized by Θ_E . The GAN objective now takes the following form:

$$\begin{aligned} \min_{\Theta_G, \Theta_E} \max_{\Theta_D} \mathbf{E}_{x \sim \mathbb{P}_x^r} q(\mathcal{D}(x)) + \mathbf{E}_{z \sim \mathbb{P}_z} q(1 - \mathcal{D}(\mathcal{G}(z))) \\ + \beta_n \mathbf{E}_{z \sim \mathbb{P}_z} \|z_n - \mathcal{E}(\mathcal{G}(z_n))\|_2^2 + \beta_c \mathbf{E}_{z \sim \mathbb{P}_z} \mathcal{H}(z_c, \mathcal{E}(\mathcal{G}(z_c))) \quad (1) \end{aligned}$$

where $\mathcal{H}(\cdot, \cdot)$ is the cross-entropy loss. The relative magnitudes of the regularization coefficients β_n and β_c enable a flexible choice to vary the importance

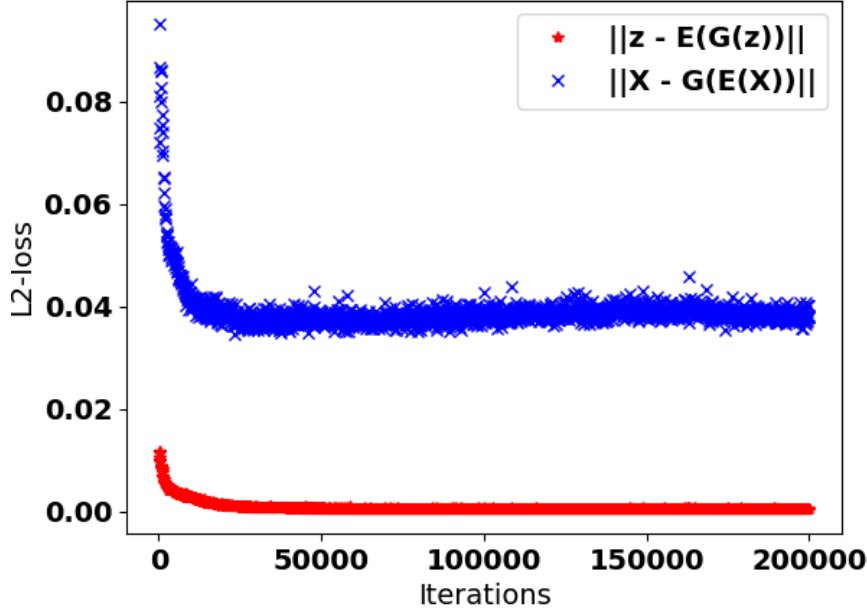


Figure 7: Decrease of Cycle Loss with iterations in MNIST. The mean square L2-distance $\|x - \mathcal{G}(\mathcal{E}(x))\|$ was 0.038 and $\|z - \mathcal{E}(\mathcal{G}(z))\|$ was 0.0004. Mean distances were computed on a test batch not used in training.

of preserving the discrete and continuous portions of the latent code. One could imagine other variations of the regularization that map $\mathcal{E}(\mathcal{G}(z))$ to be close to the centroid of the respective cluster, for instance $\|\mathcal{E}(\mathcal{G}(z^{(i)})) - \mu^{c(i)}\|_2^2$, in similar spirit as K-Means. The GAN training in this approach involves jointly updating the parameters of Θ_G and Θ_E (Algorithm 2).

As shown in Figure 7, in our architecture, both x is close to $\mathcal{G}(\mathcal{E}(x))$ and z is close to $\mathcal{E}(\mathcal{G}(z))$. Even though our architecture is optimizing for one type of cycle loss, both losses are small. The loss optimized for is even smaller.

4 Experiments

4.1 Datasets

Synthetic Data The data is generated from a mixture of Gaussians with 4 components in 2D, which constitutes the \mathcal{Z} space. We generated 2500 points from each Gaussian. The \mathcal{X} space is obtained by a non-linear transformation : $x = f(\mathbf{U} \cdot f(\mathbf{W}z))$, where $\mathbf{W} \in \mathbb{R}^{10 \times 2}$, $\mathbf{U} \in \mathbb{R}^{100 \times 10}$ with $W_{i,j} \sim \mathcal{N}(0, 1)$, $U_{i,j} \sim \mathcal{N}(0, 1)$. $f(\cdot)$ is the sigmoid function to introduce non-linearity.

Dataset	Algorithm	ACC	NMI	ARI
Synthetic	ClusterGAN	0.99	0.99	0.99
	Info-GAN	0.88	0.75	0.74
	GAN with bp	0.95	0.85	0.88
	GAN with Disc. ϕ	0.99	0.98	0.98
	AGGLO.	0.99	0.99	0.99
	NMF	0.98	0.96	0.97
	SC	0.99	0.98	0.98
MNIST	ClusterGAN	0.95	0.89	0.89
	Info-GAN	0.89	0.86	0.82
	GAN with bp	0.95	0.90	0.89
	GAN with Disc. ϕ	0.70	0.62	0.52
	DCN	0.83	0.81	0.75
	AGGLO.	0.64	0.65	0.46
	NMF	0.56	0.45	0.36
Fashion-10	ClusterGAN	0.63	0.64	0.50
	Info-GAN	0.61	0.59	0.44
	GAN with bp	0.56	0.53	0.37
	GAN with Disc. ϕ	0.43	0.37	0.23
	AGGLO.	0.55	0.57	0.37
	NMF	0.50	0.51	0.34
Fashion-5	ClusterGAN	0.73	0.59	0.48
	Info-GAN	0.71	0.56	0.45
	GAN with bp	0.73	0.54	0.45
	GAN with Disc. ϕ	0.67	0.49	0.40
	AGGLO.	0.66	0.52	0.36
	NMF	0.67	0.48	0.40
10x_73k	ClusterGAN	0.81	0.73	0.67
	Info-GAN	0.62	0.58	0.43
	GAN with bp	0.65	0.59	0.45
	GAN with Disc. ϕ	0.33	0.17	0.07
	AGGLO.	0.63	0.58	0.40
	NMF	0.71	0.69	0.53
	SC	0.40	0.29	0.18
Pendigits	ClusterGAN	0.77	0.73	0.65
	Info-GAN	0.72	0.73	0.61
	GAN with bp	0.76	0.71	0.63
	GAN with Disc. ϕ	0.65	0.57	0.45
	DCN	0.72	0.69	0.56
	AGGLO.	0.70	0.69	0.52
	NMF	0.67	0.58	0.45
	SC	0.70	0.69	0.52

Table 1: Comparison of clustering metrics across datasets

MNIST It consists of 70k images of digits ranging from 0 to 9. Each data sample is a 28×28 greyscale image. We used the DCGAN with conv-deconv layers, batch normalization and leaky relu activations, the details of which are available in the Supplementary material.

Fashion-MNIST (10 and 5 classes) This dataset has the same number of images with the same image size as MNIST, but it is fairly more complicated. Instead of digits, it consists of various types of fashion products. Supervised methods achieve lower accuracy than MNIST on this dataset. For training a GAN, we used the same architecture as MNIST for this dataset. We also merged some categories which were similar to form a separate 5-class dataset. The five groups were as follows : {Tshirt/Top, Dress}, {Trouser}, {Pullover, Coat, Shirt}, {Bag}, {Sandal, Sneaker, Ankle Boot}.

10x.73k Even though GANs have achieved unprecedented success in generating realistic images, it is not clear whether they can be equally effective for other types of data. In this experiment, we trained a GAN to cluster cell types from a single cell RNA-seq counts matrix. Moreover, computer vision might have ample supply of labelled images, obtaining labels for some fields, for instance biology, is extremely costly and laborious. Thus, unsupervised clustering of data is truly a necessity for this domain. The dataset consists of RNA-transcript counts of 73233 data points belonging to 8 different cell types [33]. To reduce the dimension of the data, we selected 720 highest variance genes across the cells. The entries of the counts matrix \mathbf{C} are first transformed as $\log_2(1 + C_{ij})$ and then divided by the maximum entry of the transformation to obtain values in the range of $[0, 1]$. One of the major challenges in this data is sparsity. Even after subselection of genes based on variance, the data matrix was close to 40% zero entries.

Pendigits It is a very different dataset that consists of a time series of $\{x_t, y_t\}_{t=1}^T$ coordinates. The points are sampled as writers write digits on a pressure sensitive tablet. The total number of datapoints is 10992, and consists of 10 classes, each for a digit. It provided a unique challenge of training GANs for point cloud data.

In all our experiments in this paper, we used an improved variant (WGAN-GP) which includes a gradient penalty [10]. Using cross-validation for selecting hyperparameters is not an option in purely unsupervised problems due to absence of labels. We adapted standard architectures for the datasets [4] and avoided data specific tuning as much as possible. Some choices of regularization parameters $\lambda = 10$, $\beta_n = 10$, $\beta_r = 10$ worked well across all datasets.

4.2 Evaluation

Since clustering is an unsupervised problem, we ensured that all the algorithms are oblivious to the true labels unlike a supervised framework like conditional GAN [23]. We compared ClusterGAN with other possible GAN based clustering approaches we could conceive.

Algorithm 1 + K-Means is denoted as “GAN with bp”. For InfoGAN, we used $\arg \max_c \mathbb{P}(c | x)$ as an inferred cluster label for x . Further, the features $\phi(x)$ in

Dataset	Algorithm			
	Cluster GAN	WGAN (Normal)	WGAN (One-Hot)	Info GAN
MNIST	0.81	0.88	0.94	1.88
Fashion	0.91	0.95	6.14	11.04

Table 2: Comparison of Frechet Inception Distance (FID) (Lower distance is better)

Dataset : MNIST, Algorithm : ClusterGAN				
ACC				
K = 7	K = 9	K = 10	K = 11	K = 13
0.72	0.80	0.95	0.90	0.84

Table 3: Robustness to Cluster Number K

the last layer of the Discriminator could contain some class-specific discriminating features for clustering. So we used Kmeans on $\phi(x)$ to cluster, denoted as “GAN with Disc. ϕ ”. We also included clustering results from Non-negative matrix Factorization (NMF) [18], Agglomerative Clustering (AGGLO) [32] and Spectral Clustering (SC). AGGLO with Euclidean affinity score and ward linkage gave best results. NMF had both l-1 and l-2 regularization, initialized with Non-negative Double SVD and used KL-divergence loss. SC had rbf kernel for affinity. We reported normalized mutual information (NMI), adjusted Rand index (ARI), and clustering purity (ACC). Since DCN has been shown to outperform various deep-learning based clustering algorithms, we reported its metrics from the paper [31] for MNIST and Pendigits. We found DCN to be very sensitive to hyperparameter choice, architecture and learning rates and could not obtain reasonable results from it on the other datasets. But we outperformed DCN results on MNIST and Pendigits dataset ¹.

Since clustering metrics do not reveal the quality of generated samples from a GAN, we report the Frechet Inception Distance (FID) [13] for the image datasets. We found that ClusterGAN achieves good clustering without compromising sample quality as shown in Table 2.

In all datasets, we provided the true number of clusters to all algorithms. In addition, for MNIST, Table 3 provides the clustering performance of ClusterGAN as number of clusters is varied. Overestimates do not severely hurt ClusterGAN; but underestimate does.

¹For all baselines and GAN variants, Table 1 reports metrics for the model with best validation purity from 5 runs.

4.3 Scalability to Large Number of Clusters

We ran ClusterGAN on Coil-20 ($N = 1440$, $K = 20$) and Coil-100 ($N = 7200$, $K = 100$) datasets, where N is the number of Data points. ClusterGAN could obtain good clusters even with such high value of K . These data sets were particularly difficult for GAN training with only a few thousand data points. Yet, we found similar behavior as MNIST / Fashion-MNIST emerging here as well. Distinct modes generated distinct 3D-objects along with rotations as shown in Figure 8.

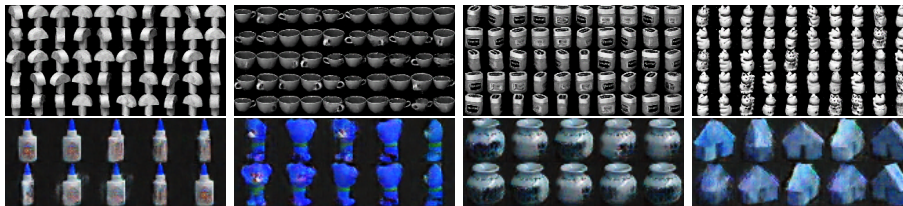


Figure 8: Scalability of ClusterGAN to large number of clusters : Modes of Coil-20 (above) and Coil-100 (below)

5 Discussion and Future Work

In this work, we discussed the drawback of training a GAN with traditional prior latent distributions for clustering and considered discrete-continuous mixtures for sampling noise variables. We proposed ClusterGAN, an architecture that enables clustering in the latent space. Comparison with clustering baselines on varied datasets using ClusterGAN illustrates that GANs can be suitably adapted for clustering. Future directions can explore better data-driven priors for the latent space. Another possibility is to improve results for problems that have a sparse generative structure such as compressed sensing.

References

- [1] Galen Andrew, Raman Arora, Jeff Bilmes, and Karen Livescu. Deep canonical correlation analysis. In *International Conference on Machine Learning*, pages 1247–1255, 2013.
- [2] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- [3] Ashish Bora, Ajil Jalal, Eric Price, and Alexandros G. Dimakis. Compressed sensing using generative models. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 537–546, 2017.

- [4] Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2172–2180, 2016.
- [5] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 215–223, 2011.
- [6] Antonia Creswell and Anil A Bharath. Inverting the generator of a generative adversarial network (ii). *arXiv preprint arXiv:1802.05701*, 2018.
- [7] Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*, 2016.
- [8] Vincent Dumoulin, Ishmael Belghazi, Ben Poole, Olivier Mastropietro, Alex Lamb, Martin Arjovsky, and Aaron Courville. Adversarially learned inference. *arXiv preprint arXiv:1606.00704*, 2016.
- [9] 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.
- [10] 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.
- [11] Swaminathan Gurumurthy, Ravi Kiran Sarvadevabhatla, and R Venkatesh Babu. Deligan: Generative adversarial networks for diverse and limited data.
- [12] Maria Halkidi, Yannis Batistakis, and Michalis Vazirgiannis. On clustering validation techniques. *Journal of intelligent information systems*, 17(2-3):107–145, 2001.
- [13] 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.
- [14] Andrew Ilyas, Ajil Jalal, Eirini Asteri, Constantinos Daskalakis, and Alexandros G Dimakis. The robust manifold defense: Adversarial training using generative models. *arXiv preprint arXiv:1712.09196*, 2017.
- [15] Pan Ji, Tong Zhang, Hongdong Li, Mathieu Salzmann, and Ian Reid. Deep subspace clustering networks. In *Advances in Neural Information Processing Systems*, pages 23–32, 2017.

- [16] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [17] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [18] Daniel D Lee and H Sebastian Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788, 1999.
- [19] Zachary C Lipton and Subarna Tripathi. Precise recovery of latent vectors from generative adversarial networks. *arXiv preprint arXiv:1702.04782*, 2017.
- [20] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.
- [21] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*, 2015.
- [22] Sebastian Mika, Bernhard Schölkopf, Alex J Smola, Klaus-Robert Müller, Matthias Scholz, and Gunnar Rätsch. Kernel pca and de-noising in feature spaces. In *Advances in neural information processing systems*, pages 536–542, 1999.
- [23] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- [24] Andrew Y Ng, Michael I Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. In *Advances in neural information processing systems*, pages 849–856, 2002.
- [25] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [26] Shibani Santurkar, David Budden, and Nir Shavit. Generative compression. *arXiv preprint arXiv:1703.01467*, 2017.
- [27] Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schoelkopf. Wasserstein auto-encoders. *arXiv preprint arXiv:1711.01558*, 2017.
- [28] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, and Pierre-Antoine Manzagol. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of Machine Learning Research*, 11(Dec):3371–3408, 2010.
- [29] Shiming Xiang, Feiping Nie, and Changshui Zhang. Learning a mahalanobis distance metric for data clustering and classification. *Pattern Recognition*, 41(12):3600–3612, 2008.

- [30] Junyuan Xie, Ross Girshick, and Ali Farhadi. Unsupervised deep embedding for clustering analysis. In *International conference on machine learning*, pages 478–487, 2016.
- [31] Bo Yang, Xiao Fu, Nicholas D. Sidiropoulos, and Mingyi Hong. Towards k-means-friendly spaces: Simultaneous deep learning and clustering. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 3861–3870, 2017.
- [32] Wei Zhang, Xiaogang Wang, Deli Zhao, and Xiaoou Tang. Graph degree linkage: Agglomerative clustering on a directed graph. In *European Conference on Computer Vision*, pages 428–441. Springer, 2012.
- [33] Grace XY Zheng, Jessica M Terry, Phillip Belgrader, Paul Ryvkin, Zachary W Bent, Ryan Wilson, Solongo B Ziraldo, Tobias D Wheeler, Geoff P McDermott, Junjie Zhu, et al. Massively parallel digital transcriptional profiling of single cells. *Nature communications*, 8:14049, 2017.

6 Supplementary Material

6.1 Additional Results

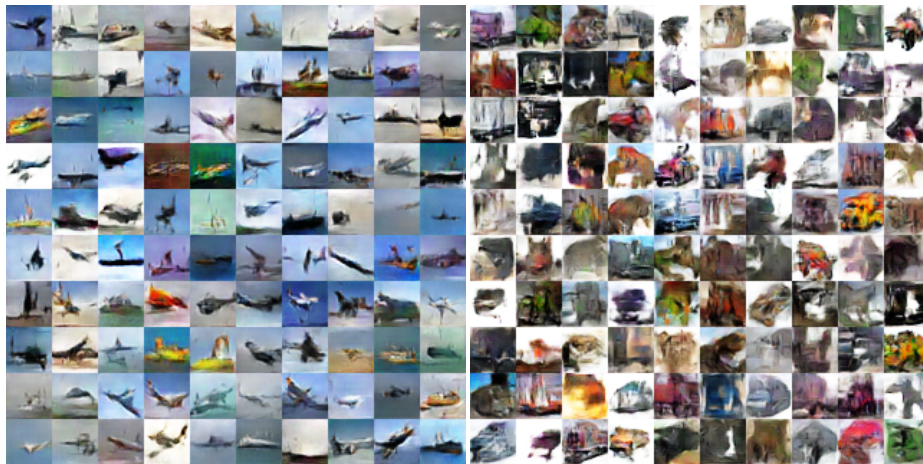


Figure 9: Modes of CIFAR-10 pick up features that easily segregate categories, but may not correspond to dataset labels. Blue background (left) and predominant white background (right).

We also ran ClusterGAN on CIFAR-10, which is a dataset with considerable intra-class variability. For CIFAR-10, the modes of ClusterGAN generate images based on a commonality (like blue background or white background), which may not be in correspondence with the labels. Blue background is present across images in airplane, bird, ship categories. Left to itself, a purely unsupervised model learns a representation that most easily separates the images into categories; but there is no impulsion to enforce higher order semantics of the labels. We believe that maximizing mutual information between an intermediate convolutional layer and the generated image would preserve the high-level semantics. We leave this exploration as future work.

6.2 Hyperparameter and Architecture Details

The networks were trained with Adam Optimizer (learning rate $\eta = 1e-04$, $\beta_1 = 0.5$, $\beta_2 = 0.9$) for all datasets. The number of discriminator updates was 5 for each generator update in training. Gradient penalty coefficient for WGAN-GP was set to 10 for all experiments. The dimension of z_c is the same as the number of classes in the dataset. Most networks used Leaky ReLU activations and Batch Normalization (BN), details for each dataset are provided below. (In the architecture without encoder, Algorithm 1 used Adam optimizer to minimize the objective for 5000 iterations per point.)

Synthetic Data

We used batch size = 64, z_n of 6 dimensions. LReLU activation with leak = 0.2 was used. $\beta_n = 10$, $\beta_c = 10$.

Generator	Encoder	Discriminator
Input $z = (z_n, z_c) \in \mathbb{R}^{10}$	Input $X \in \mathbb{R}^{16}$	Input $X \in \mathbb{R}^{16}$
FC 256 LReLU BN	FC 256 LReLU BN	FC 256 LReLU BN
FC 256 LReLU BN	FC 256 LReLU BN	FC 256 LReLU BN
FC 16 Sigmoid	FC 10 linear for \hat{z} Softmax on last 4 to obtain \hat{z}_c	FC 1 linear

MNIST and Fashion-MNIST

We used batch size = 64, z_n of 30 dimensions. LReLU activation with leak = 0.2 was used. $\beta_n = 10$ for MNIST and $\beta_n = 0$ for Fashion-MNIST, $\beta_c = 10$ for both .

Generator	Encoder	Discriminator
Input $z = (z_n, z_c) \in \mathbb{R}^{40}$	Input $X \in \mathbb{R}^{28 \times 28}$	Input $X \in \mathbb{R}^{28 \times 28}$
FC 1024 ReLU BN	4×4 conv, 64 stride 2 LReLU	4×4 conv, 64 stride 2 LReLU
FC $7 \times 7 \times 128$ ReLU BN	4×4 conv, 128 stride 2 LReLU	4×4 conv, 128 stride 2 LReLU
4×4 upconv, 64 stride 2 ReLU BN	FC 1024 LReLU	FC 1024 LReLU
4×4 upconv, 1 stride 2, Sigmoid	FC 40 linear for \hat{z} Softmax on last 10 to obtain \hat{z}_c	FC 1 linear

For Fashion-MNIST, we used $z_n = 40$. Rest of the architecture remained identical.

10x_73k

We used batch size = 64, z_n of 30 dimensions. LReLU activation with leak = 0.2 was used. $\beta_n = 10$, $\beta_c = 10$.

Generator	Encoder	Discriminator
Input $z = (z_n, z_c) \in \mathbb{R}^{38}$	Input $X \in \mathbb{R}^{720}$	Input $X \in \mathbb{R}^{720}$
FC 256 LReLU	FC 256 LReLU	FC 256 LReLU
FC 256 LReLU	FC 256 LReLU	FC 256 LReLU
FC 720 Linear	FC 38 linear for \hat{z} Softmax on last 8 to obtain \hat{z}_c	FC 1 linear

Pendigits

We used batch size = 64, z_n of 5 dimensions. LReLU activation with leak = 0.2 was used. $\beta_n = 10$, $\beta_c = 10$.

Generator	Encoder	Discriminator
Input $z = (z_n, z_c) \in \mathbb{R}^{15}$	Input $X \in \mathbb{R}^{16}$	Input $X \in \mathbb{R}^{16}$
FC 256 LReLU BN	FC 256 LReLU BN	FC 256 LReLU BN
FC 256 LReLU BN	FC 256 LReLU BN	FC 256 LReLU BN
FC 16 Sigmoid	FC 15 linear for \hat{z} Softmax on last 10 to obtain \hat{z}_c	FC 1 linear

Coil-20, Coil-100 and CIFAR-10

For Coil-20, we used batch size = 64, z_n of 20 dimensions. For Coil-100, we used batch size = 512, z_n of 20 dimensions. For CIFAR-10, we used batch size = 64, z_n of 50 dimensions.

$\beta_n = 10$, $\beta_c = 10$, LReLU activation with leak = 0.2 for all datasets.

Generator	Encoder	Discriminator
Input $z = (z_n, z_c) \in \mathbb{R}^{d_z}$	Input $X \in \mathbb{R}^{32 \times 32 \times 3}$	Input $X \in \mathbb{R}^{32 \times 32 \times 3}$
FC $2 \times 2 \times 448$ ReLU BN	4×4 conv, 64 stride 2 LReLU	4×4 conv, 64 stride 2 LReLU
4×4 upconv, 256 stride 2 ReLU BN	4×4 conv, 128 stride 2 LReLU BN	4×4 conv, 128 stride 2 LReLU BN
4×4 upconv, 128 stride 2 ReLU BN	4×4 conv, 256 stride 2 LReLU BN	4×4 conv, 256 stride 2 LReLU BN
4×4 upconv, 64 stride 2 ReLU BN	4×4 conv, 512 stride 2 LReLU BN	4×4 conv, 512 stride 2 LReLU BN
4×4 upconv, 3 stride 2, Sigmoid	FC d_z linear for \hat{z} Softmax on last K to obtain \hat{z}_c	FC 1 linear

For InfoGAN, we used the implementation of the authors <https://github.com/openai/InfoGAN> for MNIST and Fashion-MNIST. For the other datasets, we used our hyperparameters for Generator and Discriminator and added the Q network (FC 128-BN-LReLU-FC dim z_c). For “GAN with bp”, we used the same Generator and Discriminator hyperparameters as ClusterGAN. Features for “GAN with Disc. ϕ ” was obtained from the trained Discriminator of experiments “GAN with bp”.

6.3 Reporting Clustering Performance

In [30], the authors ran all algorithms multiple times with a single hyperparameter change and reported the best accuracy. For fair comparison, we used 5 runs to determine the best model using validation. To be more precise, we first split

our datasets into Train, Validation and Test portions. The GAN was trained only on the Train split of the data in an unsupervised manner, sometimes with a single hyper-parameter change such as β_n, z_n , leak or batch-norm . For each dataset, we saved the model with the best purity on the Validation split from the 5 runs. Table 1 reports the metrics on the Test split for the saved model. The metrics on the entire dataset for the saved model was either identical upto 2 decimals or slightly better than Test. But we report only for the Test split, since it has neither been used for training nor for validation runs.

6.4 Generated Modes

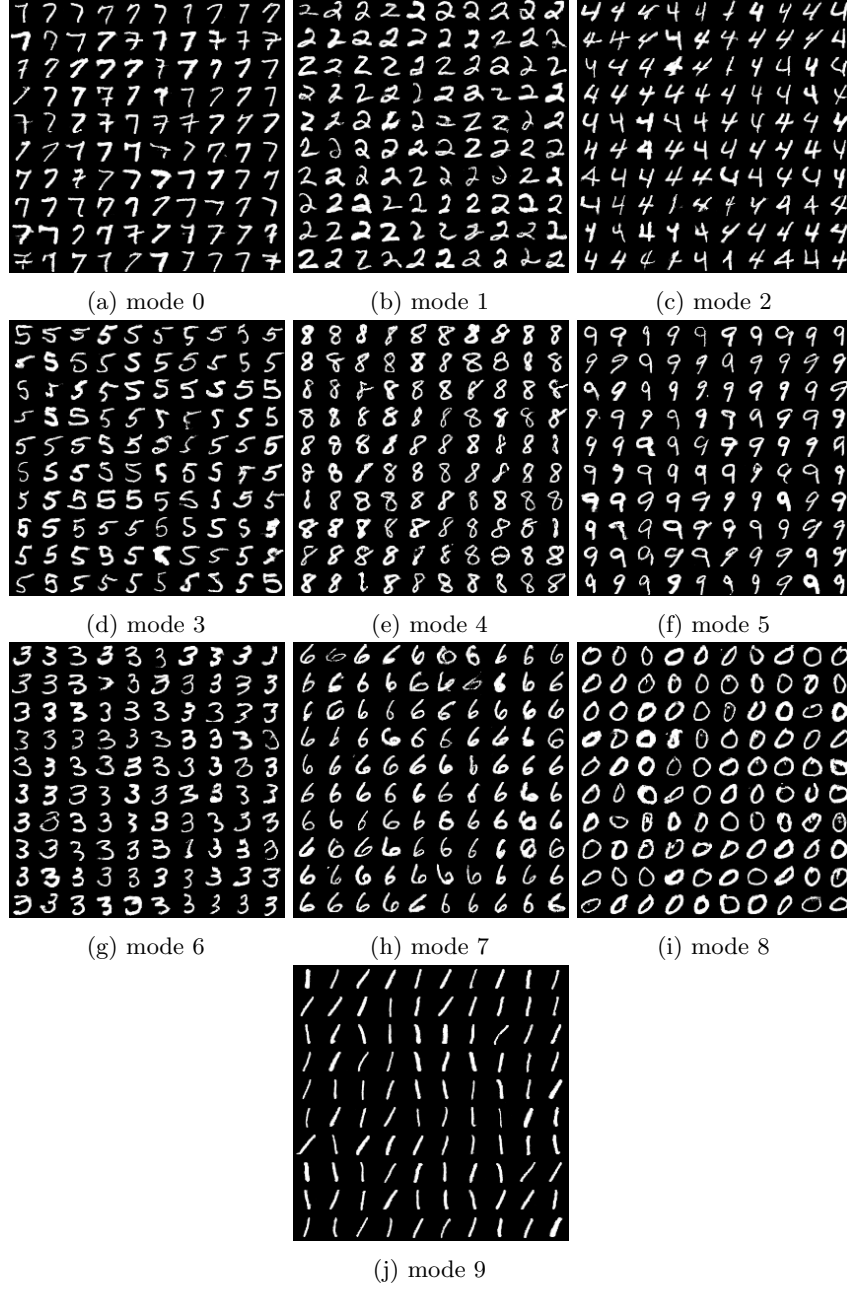


Figure 10: Generated digits from distinct modes

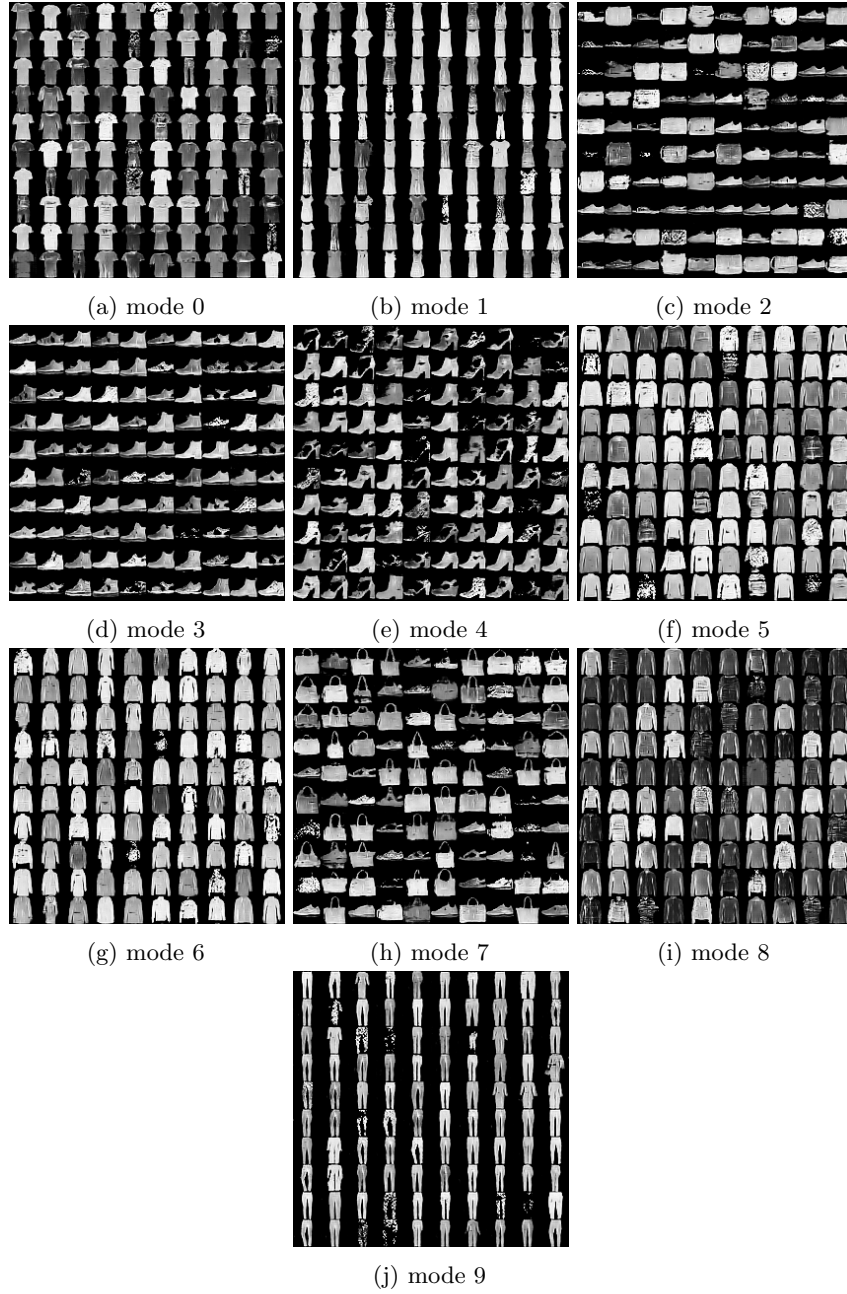


Figure 11: Generated fashion items from distinct modes

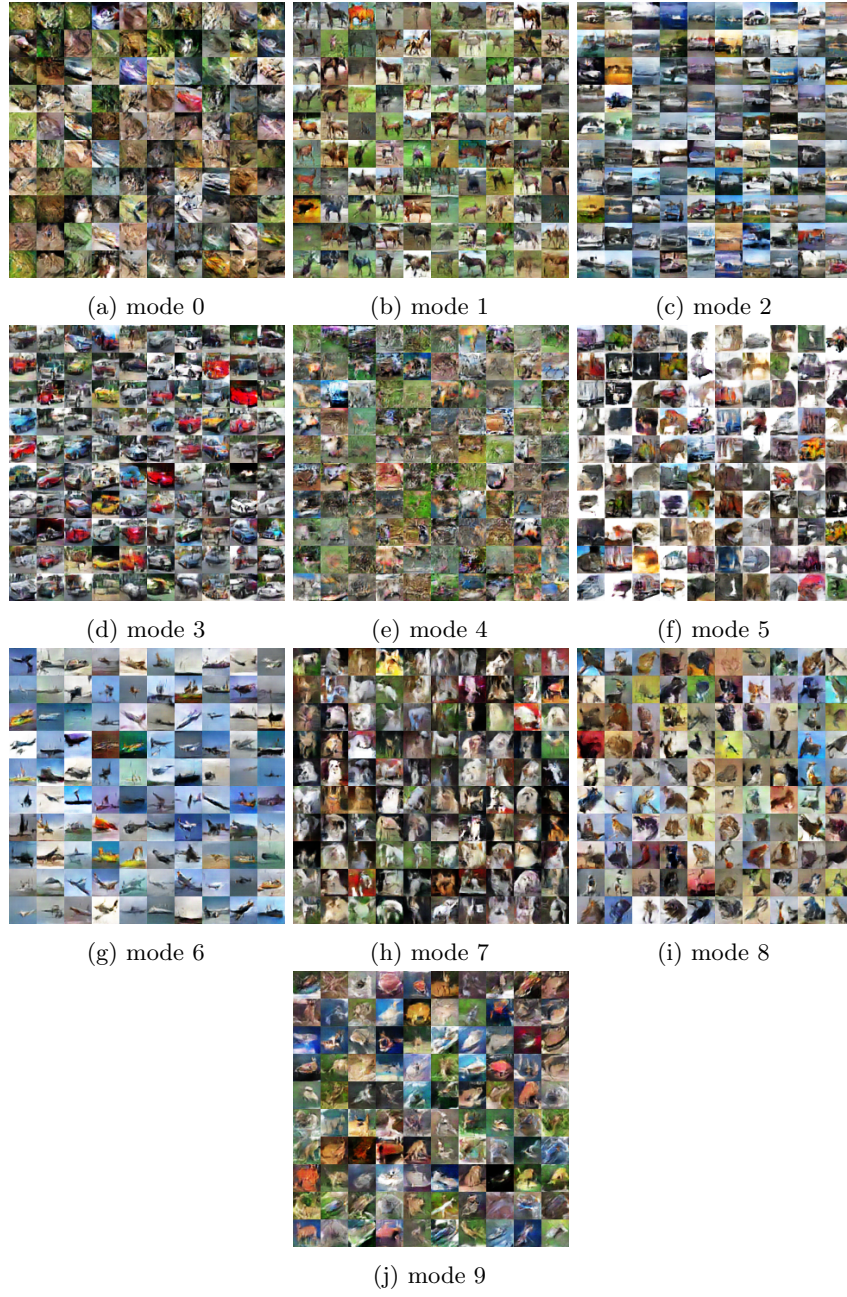


Figure 12: Generated categories from distinct modes of CIFAR-10