

Conditional Variational Autoencoder with Balanced Pre-training for Generative Adversarial Networks

Yuchong Yao¹, Xiaohui Wang¹, Yuanbang Ma¹, Han Fang¹, Jiaying Wei¹, Liyuan Chen¹, Ali Anaissi¹ and Ali Braytee^{1,2}

¹School of Computer Science, The University of Sydney

²School of Computer Science, University of Technology Sydney

Abstract

Class imbalance occurs in many real-world applications, including image classification, where the number of images in each class differs significantly. With imbalanced data, the generative adversarial networks (GANs) leans to majority class samples. The two recent methods, Balancing GAN (BAGAN) and improved BAGAN (BAGAN-GP), are proposed as an augmentation tool to handle this problem and restore the balance to the data. The former pre-trains the autoencoder weights in an unsupervised manner. However, it is unstable when the images from different categories have similar features. The latter is improved based on BAGAN by facilitating supervised autoencoder training, but the pre-training is biased towards the majority classes. In this work, we propose a novel *Conditional Variational Autoencoder with Balanced Pre-training for Generative Adversarial Networks* (CAPGAN)¹ as an augmentation tool to generate realistic synthetic images. In particular, we utilize a conditional convolutional variational autoencoder with supervised and balanced pre-training for the GAN initialization and training with gradient penalty. Our proposed method presents a superior performance of other state-of-the-art methods on the highly imbalanced version of MNIST, Fashion-MNIST, CIFAR-10, and two medical imaging datasets. Our method can synthesize high-quality minority samples in terms of Fréchet inception distance, structural similarity index measure and perceptual quality.

1 Introduction

Computer vision contains many supervised learning problems, including image classification, image segmentation, and others [He *et al.*, 2016]. Modern image classifiers are generally deep learning models which need balanced imaging datasets such as MNIST [LeCun *et al.*, 2010], Fashion-MNIST [Xiao *et al.*, 2017], CIFAR-10 [Krizhevsky *et al.*, 2009], ImageNet [Deng *et al.*, 2009], and others. However,

class distribution in real-world datasets is often skewed, especially in the medical imaging domain, i.e. there are more normal images than cancerous images. The performance of the image classification models based on deep learning degrade significantly in the presence of class imbalance because these models will be biased towards the majority class samples and ignore the minority ones [Braytee *et al.*, 2019].

Interestingly, the literature shows that augmenting the minority classes with sample generation such as Generative Adversarial Networks (GANs) in image generation is a promising approach to deal with class imbalance [Rezaei *et al.*, 2020]. Briefly, GANs consist of a generator and a discriminator that adopt an adversarial training schema to allow the generator and discriminator to compete. GANs can generate synthetic minority samples to help restore balance to the data. However, GANs-based approaches have several limitations, including mode collapse, sub-optimal initialization, and training instability, which lead to unstable results. Further, bias can occur towards the majority class. Recent works combine GANs with other models, such as autoencoder, to borrow the reconstruction ability to enhance the initialization and training of the GANs-based models for generating minority samples. A recent powerful methods BAGAN [Mariani *et al.*, 2018], and BAGAN-GP [Huang and Jafari, 2021] are examples that have shown promising results to handle class imbalance on various benchmarks and datasets. Nevertheless, they still suffer from the following limitations. The pre-training in BAGAN-GP is created on imbalanced data which lead to be biased towards the majority classes. Further, BAGAN-GP used naive autoencoder model and objective function to obtain the pre-trained weights, which can be further optimized by more advanced architectures and objectives. Moreover, BAGAN and BAGAN-GP are lack of comprehensive evaluations, where they only evaluated the models under datasets with a small imbalance rate and they haven't evaluated on highly and extreme imbalance rates. To this end, we propose a new framework, namely, *Conditional Variational Autoencoder with Balanced Pre-training for Generative Adversarial Networks* (CAPGAN). The general objective of our framework is to synthesize high-quality samples for the majority and minority classes to overcome the class imbalance problem.

The major contributions of CAPGAN can be summarized as follows: (1) we facilitate a balanced pre-training stage to

¹We will open source the code upon acceptance

GAN components; (2) we utilize conditional variational autoencoder model in the pre-training stage for GAN initialization; (3) we propose a novel sophisticated objective function that encourages the model to capture the true distribution of the samples and generate high-quality samples; (4) we integrate the proposed balanced pre-training and the new objective function simultaneously to initialize and train the corresponding GAN components to enhance the training stability and generate more realistic and diverse samples.

2 Related Work

Generative Models on Class Imbalance. Autoencoder and generative adversarial networks (GANs) are two representative generative models proposed to handle class imbalance in imaging applications. Several studies suggest that acquiring more samples (especially the minority classes) to restore data balance is the most effective way to address the class imbalance. Few studies use autoencoder variants to handle class imbalance. For example, Taghanaki et al. (2020) state that variational autoencoder (VAE) can improve the performance on imbalanced data [Taghanaki *et al.*, 2020]. However, Li et al. (2018) find that the samples generated by VAE are not as diverse as the samples from GANs [Li *et al.*, 2018]. Hence, several studies find GANs variants are powerful to handle the imbalanced data. DCGAN [Shoohi and Saud, 2020] is proposed to synthesize samples for the minority classes. It leads to impressive results on various tasks (e.g. plant disease). WGAN is also widely used for data augmentation [Bhatia and Dahyot, 2019] and minority oversampling for CT images [Wang *et al.*, 2019]. CycleGAN applies image-to-image translation on the imbalanced data, which attempts to generate minority samples based on majority samples [Zhu *et al.*, 2017]. Many existing studies attempt to facilitate semi-supervised GANs (or conditional GANs) and unsupervised GANs together. For example, [Balasubramanian *et al.*, 2020] uses an unconditional GAN for diabetes image oversampling. Further, another study proposes a conditional GAN (CovidGAN) [Waheed *et al.*, 2020] to augment the minority Covid-19 CXR images. Although the generative models achieve impressive results for addressing the class imbalance, they suffer from mode collapse, training instability, and unstable results. Further, some studies argued that the generated minority samples would bias towards the majority classes and degrade the original performance for majority samples [Sampath *et al.*, 2021].

GAN-Autoencoder-based augmentation To overcome the limitation of GANs augmentation methods to handle the class imbalance, BAGAN combines the power of autoencoder and GANs [Mariani *et al.*, 2018]. It integrates an autoencoder with GANs to gain better reconstruction ability and produces a stable starting point for training. Particularly, it initializes the GAN model by integrating the pre-trained autoencoder (given that the GANs and the autoencoder have the same network architectures). BAGAN can only use an unsupervised autoencoder where it does not use the label information during the pre-training. However, label information is critical for imposing class conditioning on pre-trained weights. Recently, an improved version of BAGAN (i.e. BAGAN-

GP) [Huang and Jafari, 2021] states that BAGAN does not perform well on medical data and the results are not stable. BAGAN-GP introduces a supervised autoencoder, which utilizes the label information during pre-training. Besides, BAGAN-GP applies conditional GANs in its structure to improve the class-specific generative performance. The results show that BAGAN-GP is superior over BAGAN-GP on various datasets (MNIST, Fashion-MNIST, and CIFAR-10). Furthermore, BAGAN-GP is proved to be more effective than BAGAN on the medical imaging data. However, BAGAN-GP facilitates simple network structures and basic objective functions in its autoencoder. Also, although the pre-training takes class information into account, the pre-training is biased towards majority classes, leading to sub-optimal solutions.

3 Method

3.1 Supervised Initialization and Training

We initialize the discriminator and the generator in GAN with the weights of a pre-trained conditional variational autoencoder (CVAE). When the GANs are trained under the adversarial settings, the generator can produce class-specific samples with better representative ability given by the pre-training, especially for the minority classes. Also, the discriminator can identify whether the image belongs to one of the classes or is a fake sample. In this step, the discriminator and generator are updated in GAN by following the min-max adversarial settings [Goodfellow *et al.*, 2014] in Equation 1

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))] \quad (1)$$

where z denotes noise samples, and x is from data generating distribution. After we develop CVAE in the initialization step, the weights of CVAE’s components are transferred to GAN’s components. The generator and the decoder (with the embedding component) are designed to have the same network structures and topologies to allow the weights to be transferred from the pre-trained decoder to the generator. The discriminator is initialized to have the same network structures and topologies in the first few layers as the encoder, followed by dense layers that match the dimension in the final output. The weights of the final dense layers are randomly initialized. The initialization of the generator and discriminator in GAN is illustrated in Figure 1(a).

Training the CAPGAN follows the standard adversarial settings that the discriminator and the generator compete with each other. The loss function of the discriminator and the generator is inspired from DRAGAN [Kodali *et al.*, 2017]. The discriminator consists of three losses for fake images, real images, and wrong labels, while the generator has only a fake image loss. Moreover, we impose the gradient penalty term [Arjovsky *et al.*, 2017] in the discriminator loss, which aims to help the convergence of the discriminator as shown in Equation 2.

where $\hat{x} = \alpha x_r + (1 - \alpha)x_{noise}$, $\alpha \sim U(0, 1)$, x_r is the real input image, α is a normally distributed random number with

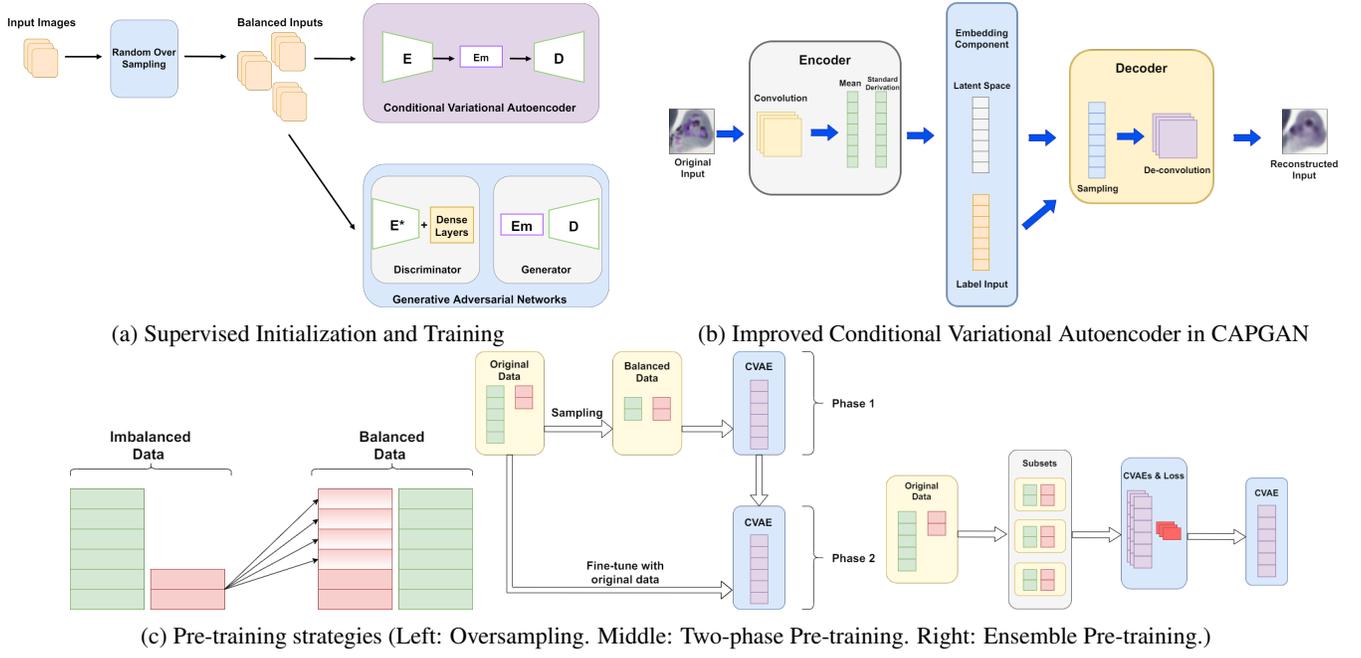


Figure 1: Our proposed CAPGAN framework

$$GP = \lambda \mathbb{E}_{\hat{x} \sim \hat{X}} [(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2] \quad (2) \quad \text{tion 6}$$

values in uniform distribution $U(0, 1)$, λ denotes gradient penalty weight and $\|\nabla_x D(x)\|_2$ is the norm of gradients. Both the discriminator and the generator attempt to minimize their losses to obtain better performance.

3.2 Improved Conditional Variational Autoencoder in CAPGAN

The simple autoencoder utilized by BAGAN-GP is replaced with a more powerful Conditional Convolutional Variational Autoencoder in CAPGAN. Variational autoencoder has an advanced architecture compared to the autoencoder. Autoencoder may create some samples in the latent space with no valid meaning or hard to interpret after decoding, leading to a poor generative performance for the decoder in creating new samples from the latent space. However, variational autoencoder overcomes latent space irregularity by encoding the inputs into mean and standard derivation (i.e. learns a distribution over the latent space). Hence, the latent space is continuous, regularized, and enables easier sampling and interpolation. Reparametrization is applied to integrate the learned mean and standard derivation from the encoder in CVAE as described in equation 3, 4, 5. Convolutional layers and transposed convolutional layers are utilized in the encoder and the decoder of CVAE. There is an embedding component (a shallow sub-network) in our proposed CVAE that takes the class labels of the input images and encodes them into class-specific information (i.e. same size as the latent space). The outputs from the encoder and the embedding component are fed into the decoder together as the input, where the latent output and the class-specific information are multiplied Equa-

$$\begin{aligned} (\mathbf{g}^*, \mathbf{h}^*) &= \arg \min KL(q_x(y), p(y | x)) \\ &= \arg \min (\mathbb{E}_{y \sim q_x} (\log q_x(y)) - \mathbb{E}_{y \sim q_x} (\log \frac{p(x | y)p(y)}{p(x)})) \\ &= \arg \max (\mathbb{E}_{y \sim q_x} (\log p(x | y)) - KL(q_x(y), p(y))) \end{aligned} \quad (3)$$

$$\begin{aligned} \mu_x &= \mathbf{g}(x) = \mathbf{g}_2(\mathbf{g}_1(x)), \\ \sigma_x &= \mathbf{h}(x) = \mathbf{h}_2(\mathbf{h}_1(x)) \end{aligned} \quad (4)$$

$$\mathbf{g}_1(x) = \mathbf{h}_1(x) \quad (5)$$

$$z = \mu + \exp(0.5 \times \sigma) \times \alpha, \quad \alpha \sim U(0, 1) \quad (5)$$

$$\mathbf{O} = z \odot \epsilon(z, y) \quad (6)$$

where z is the reparameterized variable in the latent dimension with a mean μ and standard derivation σ , α is a normally distributed random number with values in uniform distribution $U(0, 1)$, e and y denotes the embedding component and class label respectively, \mathbf{O} is the embedded output. In this way, we manage to pass the class information into the CVAE and train it in a supervised fashion. This is critical for the GAN initialization as the training of the GAN is supervised. Transferring weights from an unsupervised CVAE could mislead the GAN and result in sub-optimal solutions. Therefore, it is essential to facilitate the embedding component to make CVAE training conditioned on the class labels. The illustration of the proposed CVAE architecture component in CAP-

GAN is shown in Figure 1b.

3.3 Improved Objective Function

The BAGAN-GP method applies L2 minimization to train autoencoder, but this may lead to two drawbacks: firstly, mean absolute error (MAE) only enforces pixel-to-pixel similarity, which fails to capture the class-wise distributions of the input samples. Secondly, MAE is not suitable for training more advanced and sophisticated CVAE [Kingma and Welling, 2013]. The new objective function of our proposed CVAE model is composed of three components: (1) the Kullback–Leibler (KL) divergence; (2) the cross-entropy loss; (3) and the mean squared error. KL-divergence measures the difference between two probability distributions which is considered critical for training the CVAE because it encourages the model to learn a distribution in the latent space. By minimizing the KL-divergence, the learned mean and standard derivation for the target distribution of latent space are optimized, which allows the decoder to sample and generate better results. Further, the KL-divergence is denoted as the latent loss in the objective function. The remaining two components are related to the reconstruction loss. Cross entropy loss is more suitable for Bernoulli distribution as it expresses the negative Bernoulli log-likelihood, while the mean squared error assumes a Gaussian distribution. Incorporating these two components is significant to optimize the model’s performance on more complex distributions. By minimizing both cross-entropy loss and mean squared loss, the CVAE can learn better reconstruction ability on more sophisticated distributions and gain better generative performance. The objective function is presented in Equation 7 as follows

$$Objective = D_{KL}(p \parallel q) + H(p, q) + MSE \quad (7)$$

$$D_{KL}(p \parallel q) = \sum_x p(x) \ln \frac{p(x)}{q(x)} \quad (8)$$

$$H(p, q) = -\sum_x p(x) \log q(x) \quad (9)$$

$$MSE = \frac{1}{m} \sum_{i=1}^m (x_i - \hat{x}_i)^2 \quad (10)$$

3.4 Random Oversampling Pre-Training Strategy

The discriminator and the generator in the GAN require a good initialization point to synthesize balanced samples towards all classes. The original autoencoder pre-training strategy in the BAGAN-GP method is imbalanced towards the majority classes, which introduces burdens for the GAN training. In CAPGAN, we redesigned the pre-training strategy to allow the pre-trained weights to be balanced for the GAN initialization. Three different pre-training strategies are explored as show in Figure 1c. Two strategies that are implemented but not adopted called two-phase pre-training and ensemble pre-training. The two-phase pre-training resembles the ideas from two-phase learning, where the CVAE is first trained on balanced data, then the CVAE is fine-tuned on the

original imbalanced dataset. However, this approach suffers from overfitting and high cost on tuning. The ensemble strategy attempts to fit multiple CVAEs with different subsets of the majority classes and combines them with the entire minority samples. The final weights would be a weighted average of all weights from the CVAEs according to their training loss. This method has drawbacks such as computationally expensive and overfitting and it may be infeasible in practice.

In CAPGAN, we adopt random oversampling (ROS) for pre-training strategy. ROS has been proved to be effective in many applications for addressing class imbalance. The simplicity and compatibility of the method make it a popular choice for many class imbalance applications. The imbalanced data is randomly oversampled to make samples in each class are balanced before they are fed into the CVAE. Although there is a potential risk for overfitting the minority samples as they are replicated multiple times, the ROS pre-training shows an improved performance with minor computational costs. The reason is due to transferring the weights from the CVAE to the GAN only during the initialization step. These weights from the ROS pre-training strategy manage to produce a balanced and good enough starting for the GAN to achieve great results. Furthermore, ROS pre-training can be easily scaled to large and complex datasets due to its simplicity and computational efficiency.

4 Results and Discussion

4.1 Datasets

The experiments are conducted on general vision and medical imaging datasets. For general vision datasets, we consider MNIST, Fashion-MNIST, and CIFAR-10. All samples in the three datasets are resized into a uniform size, which is the same as the original sample size in CIFAR-10 (i.e. 32×32). For medical imaging datasets, we used small-scale blood cells data (Cells [Huang and Jafari, 2021]) and a breast cancer cell data (BreakHis [Spanhol *et al.*, 2015]). We scaled down the sample size in both datasets due to the computational and time constraints. The samples in Cells and BreakHis datasets are reshaped to 64×64 and 32×32 , respectively. The summary of datasets is shown in Table 2.

4.2 Imbalance Rate

The general vision datasets are balanced initially. We impose imbalance on MNIST, Fashion-MNIST, and CIFAR-10 by choosing one class as the majority class and treating the other classes as minority classes and sampling subsets for those classes. The imbalance rate is defined as the number of samples between the largest majority class and the smallest minority class. For general vision datasets, we construct imbalanced datasets using the following imbalance rates: 5, 10, 20, 50, and 100. Cells and BreakHis are originally imbalanced, we conduct the experiments on those datasets using the original imbalance rate, unless for BreakHis, we force imbalance in addition to the original imbalance rate.

4.3 Evaluation Metrics and Compared methods

The Fréchet Inception Distance (FID) and Structural Similarity Index Measure (SSIM) are the metrics used for the eval-

Imbalance Rate	Model	MNIST				Fashion-MNIST				CIFAR-10			
		avg(Minority)		Majority		avg(Minority)		Majority		avg(Minority)		Majority	
		FID	SSIM	FID	SSIM	FID	SSIM	FID	SSIM	FID	SSIM	FID	SSIM
5	DCGAN	251.33	2.57×10^{-1}	207.93	2.96×10^{-1}	379.96	2.60×10^{-1}	316.61	2.79×10^{-1}	495.27	5.67×10^{-2}	335.43	9.30×10^{-2}
	BAGAN-GP	176.15	2.56×10^{-1}	157.39	2.86×10^{-1}	267.05	2.67×10^{-1}	240.21	2.82×10^{-1}	463.85	5.78×10^{-2}	364.04	8.94×10^{-2}
	CAP-GAN	165.59	2.69×10^{-1}	157.16	2.79×10^{-1}	264.65	2.67×10^{-1}	232.47	2.88×10^{-1}	366.06	6.00×10^{-2}	288.58	8.00×10^{-2}
10	DCGAN	213.39	2.58×10^{-1}	528.09	3.36×10^{-1}	472.96	2.53×10^{-1}	293.46	2.96×10^{-1}	495.27	5.67×10^{-2}	335.43	9.30×10^{-2}
	BAGAN-GP	187.44	2.62×10^{-1}	165.08	2.70×10^{-1}	328.19	2.66×10^{-1}	254.57	3.03×10^{-1}	486.36	5.76×10^{-2}	354.85	9.06×10^{-2}
	CAP-GAN	177.22	2.64×10^{-1}	160.12	2.90×10^{-1}	271.89	2.75×10^{-1}	241.99	2.94×10^{-1}	399.18	6.00×10^{-2}	291.54	8.00×10^{-2}
20	DCGAN	278.80	2.68×10^{-1}	210.22	2.85×10^{-1}	509.82	2.67×10^{-1}	468.58	2.25×10^{-1}	659.34	5.48×10^{-2}	460.03	8.69×10^{-2}
	BAGAN-GP	193.40	2.66×10^{-1}	174.21	2.78×10^{-1}	325.61	2.69×10^{-1}	288.96	2.87×10^{-1}	529.59	5.48×10^{-2}	480.33	7.40×10^{-2}
	CAP-GAN	173.59	2.70×10^{-1}	149.04	2.86×10^{-1}	260.61	2.64×10^{-1}	230.71	2.98×10^{-1}	368.46	6.00×10^{-2}	271.10	8.00×10^{-2}
50	DCGAN	479.63	2.55×10^{-1}	722.07	3.13×10^{-1}	540.12	2.74×10^{-1}	514.94	2.65×10^{-1}	742.54	5.71×10^{-2}	453.33	8.32×10^{-2}
	BAGAN-GP	204.70	2.58×10^{-1}	149.28	2.79×10^{-1}	388.60	2.54×10^{-1}	317.47	2.81×10^{-1}	529.69	5.17×10^{-2}	455.16	7.23×10^{-2}
	CAP-GAN	176.47	2.59×10^{-1}	163.90	2.83×10^{-1}	268.31	2.70×10^{-1}	225.08	2.95×10^{-1}	411.23	7.00×10^{-2}	346.43	8.00×10^{-2}
100	DCGAN	511.31	1.83×10^{-1}	711.49	1.68×10^{-1}	703.26	2.78×10^{-1}	811.17	2.69×10^{-1}	701.54	5.26×10^{-2}	520.11	6.32×10^{-2}
	BAGAN-GP	228.19	2.54×10^{-1}	167.74	2.79×10^{-1}	416.67	2.65×10^{-1}	316.33	2.70×10^{-1}	554.70	4.72×10^{-2}	456.06	6.33×10^{-2}
	CAP-GAN	168.00	2.66×10^{-1}	160.88	2.81×10^{-1}	286.45	2.70×10^{-1}	236.04	2.94×10^{-1}	370.82	7.00×10^{-2}	339.45	8.00×10^{-2}
p-value		FID		SSIM		FID		SSIM		FID		SSIM	
	CAP-GAN vs BAGAN-GP	1.3×10^{-3}		7.1×10^{-1}		5.0×10^{-6}		8.0×10^{-1}		5.0×10^{-16}		3.3×10^{-2}	
	CAP-GAN vs DCGAN	2.8×10^{-12}		2.4×10^{-1}		8.9×10^{-15}		8.1×10^{-1}		4.5×10^{-11}		1.5×10^{-1}	

Table 1: Averaged FID and SSIM for General Vision Benchmarks

Dataset	Resolution	Class	Training Samples Per Class			
			Min	Median	Mean	Max
MNIST	28×28	10	6000	6000	6,000	6000
Fashion-MNIST	28×28	10	6000	6000	6000	6000
CIFAR-10	32×32	10	5,000	5,000	5,000	5,000
<i>BreakHis*</i>	700×460	2	2,480	N/A	N/A	5,429
Cells	100×101	4	106	887	1,721	5,600

Table 2: Datasets characteristics. BreakHis only has two classes so the median and mean are N/A's.

uations. Lower FID or higher SSIM indicates better performance. For each class in each dataset, the model under evaluation generated 1,000 samples. Those generated samples are compared with the test samples to compute the corresponding FID and SSIM with the test set. Two baseline models are used for comparison: Conditional Deep Convolutional Generative Adversarial Networks (DCGAN)[Radford *et al.*, 2015] and the BAGAN-GP[Huang and Jafari, 2021].

4.4 Experiment Setting

The hyperparameter values during the training are summarized in Table 3. The compared methods are evaluated under the same experiment settings and implemented using TensorFlow 2.0 framework. We utilize NVIDIA Tesla P100 to train the models.

Hyperparameter	Values
Learning Rate (CVAE)	0.0006, 0.0007, 0.0008, 0.001, 0.0005
CVAE Epoch	30, 40, 50
Adam beta1 (CVAE)	0.5, 0.6, 0.7, 0.8
Learning Rate (Generator)	0.00005, 0.0001, 0.0002, 0.0005, 0.0008, 0.0013, 0.0015, 0.001, 0.002
Learning Rate (Discriminator)	0.00005, 0.0001, 0.0008, 0.0013, 0.0015, 0.0002, 0.002
Gradient Penalty Weight	5, 10
Train Ratio	2, 3, 4, 5, 6, 7, 8, 10
Batch Size	32, 64, 128, 256
Latent Dimension	64, 128, 256, 512

Table 3: Hyperparameter Optimization for CAP-GAN

4.5 Results for General Vision Datasets

We first conduct experiments on the low, moderate, high and extreme imbalanced versions of MNIST, Fashion-MNIST, and CIFAR-10 balanced datasets using different imbalance rates. For each dataset, there is one majority class and nine minority classes. The average values of FIDs and SSIMs for the majority and the minority class are presented in Table 1. It is clearly shown from the results that BAGAN-GP and CAPGAN outperform the DCGAN consistently across all datasets, suggesting that the countermeasures for class imbalance are effective. It also indicates that generative models which are proved to be successful on balanced datasets could not handle class imbalance. We further observe that CAPGAN achieves superior results than BAGAN-GP in almost all experiment settings, which indicates that CAPGAN is much more powerful than BAGAN-GP as a generative model for imbalanced data. We also note, as shown in Figures 2 and 3, that the compared methods produce unstable results when the imbalance rate increase. In particular, the FID of BAGAN-GP and DCGAN increase significantly as the imbalance rate becomes higher in all three datasets. In contrast, the FID of CAPGAN remains at a steady level or increases much lower as the imbalance rate increases. Therefore, CAPGAN is proven as a powerful generative model that is able to maintain a low FID in the presence of high and extreme imbalance rates (i.e., 50 and 100). As for SSIM, although the improvements are not as evident as the FID, CAPGAN managed to score higher SSIM scores under most experiment configurations than DCGAN and BAGAN-GP. In terms of SSIM evaluation metric, although the improvements are not as evident as the FID, CAPGAN is able to attain higher scores under most experiment configurations compared to the state-of-the-arts. We believe the CVAE initialization and the proposed pre-training help CAPGAN to achieve such results. Furthermore, using the Student's t-test, we investigate whether the results that are produced by CAPGAN are significantly different to the state-of-the-art. The statistical results show that

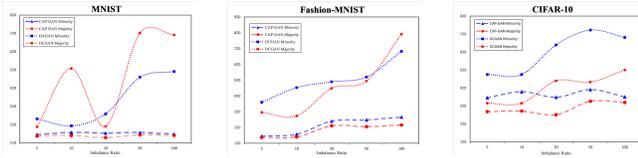


Figure 2: Comparison of DCGAN and CAPGAN

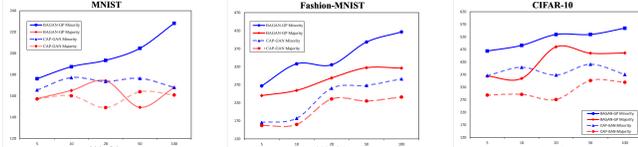


Figure 3: Comparison of BAGAN-GP and CAPGAN

the p-value in all tests for FID is less than 0.05, which rejects the null hypothesis that CAPGAN and the compared methods have equal performance.

We investigate the generated images from DCGAN, BAGAN-GP, and CAPGAN, along with the original images on CIFAR-10. As shown in Figure 4, the generated images are all from the minority classes. Both DCGAN and BAGAN-GP produce very blurring images when the imbalance rate is extreme (i.e. 100). Most of the generated samples using these two models lose details and textures, and several samples look like unrecognizable objects (i.e., noises). Further, they are lack of diversity, which indicates that DCGAN and BAGAN-GP suffers from mode collapse. On the contrary, CAPGAN generates more realistic samples that capture the details and the textures of the objects. Each sample is easily recognized regarding its original class. Furthermore, CAPGAN produces samples of high diversity, which is crucial for reducing the likelihood of overfitting when using the CAPGAN for oversampling and achieving a better performance at high and extreme imbalance rates.

4.6 Results for Medical Imaging Benchmarks

Medical imaging data are often imbalanced due to the high cost of generating real images. The generative models play an important role to produce synthetic images at medical applications. The results for medical imaging data are presented in Table 4 and Table 5. Cells is a small-scaled medical dataset with a high imbalance rate (i.e. 52.83). Similar to the results of general vision datasets. Firstly, we investigate the impact of the class imbalance on the baseline DCGAN. We find that BAGAN-GP and CAPGAN achieve better performance than DCGAN on almost all metrics, especially FID. It

		Cells			
		avg(Minority)		Majority	
Imbalance Rate	Model	FID	SSIM	FID	SSIM
52.83	DCGAN	438.48	3.75×10^{-1}	266.99	4.46×10^{-1}
	BAGAN-GP	260.79	3.51×10^{-1}	247.34	4.27×10^{-1}
	CAP-GAN	228.76	3.48×10^{-1}	160.00	4.35×10^{-1}
Improvement on FID (%)		CAP-GAN vs BAGAN-GP		CAP-GAN vs DCGAN	
		12.28 ↑		35.31 ↑	
		47.83 ↑		40.07 ↑	

Table 4: Averaged FID and SSIM for Cells



(a) Original Image (b) DCGAN (c) BAGAN-GP (d) CAPGAN age

Figure 4: Generated Images for CIFAR-10 with an Imbalanced Rate of 100

		BreakHis			
		Minority		Majority	
Imbalance Rate	Model	FID	SSIM	FID	SSIM
2.19	DCGAN	456.39	5.66×10^{-2}	294.23	6.34×10^{-2}
	BAGAN-GP	251.00	6.76×10^{-2}	231.76	6.47×10^{-2}
	CAP-GAN	237.47	5.91×10^{-2}	216.01	6.06×10^{-2}
10	DCGAN	565.67	5.25×10^{-2}	416.73	5.86×10^{-2}
	BAGAN-GP	211.28	8.95×10^{-2}	210.61	9.30×10^{-2}
	CAP-GAN	203.60	5.94×10^{-2}	181.82	5.95×10^{-2}
Improvement on FID (%) of 2.19		CAP-GAN vs BAGAN-GP		CAP-GAN vs DCGAN	
		5.39 ↑		6.79 ↑	
		47.97 ↑		26.58 ↑	
Improvement on FID (%) of 10		CAP-GAN vs BAGAN-GP		CAP-GAN vs DCGAN	
		3.63 ↑		13.67 ↑	
		64.01 ↑		56.37 ↑	

Table 5: FID and SSIM for BreakHis

shows that the powerful generative architecture on traditional and balanced datasets could not adapt to more challenging and high imbalanced medical imaging data, especially for the minority classes, which is considered as the class of interest. Secondly, we test the performance of the proposed CAPGAN against the-state-of-arts BAGAN-GP. CAPGAN outperforms the BAGAN-GP in FID. The improvements on FID are around 12.78% and 35.31% for minority and majority classes, respectively. The boost on SSIM is not as significant as FID, where CAPGAN and BAGAN-GP achieved comparable results.

For BreakHis dataset, although the original imbalance rate is lower than Cells, it is challenging in other aspects because it contains images of different scales, such as the objects in different images might be collected at a different magnification rate. As presented in Table 5, BAGAN-GP and CAPGAN consistently outperform the DCGAN, especially under a high imbalance rate. The results suggest the effectiveness of class imbalance countermeasures. Furthermore, CAPGAN achieves better FID than BAGAN-GP under both imbalance rates. As for SSIM, all three models obtained very low SSIM, so it would be pointless to analyze the statistics regarding the SSIM. We believe that the variety of scales in the datasets lead to unsatisfying performance for all three models since they do not have any technique to deal with samples of different scales.

We illustrate the generated images from DCGAN, BAGAN-GP, and CAPGAN in the Cells. As shown in Figure 5, the generated samples from DCGAN and BAGAN-GP are poor in their quality due to a lack of textures and details.

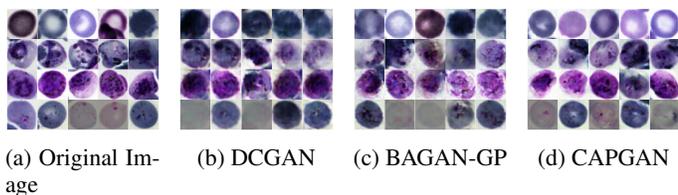


Figure 5: Generated Images for Cells

Besides, we can observe many irregular dark or grey shapes on those generated images, which cover a large portion of the cell and make the cell unrecognizable. However, the generated images by CAPGAN are more realistic and diverse cell images. The details and textures are clear with small noises covering cell's body. Furthermore, the cell images have richer and more diverse colours than those generated by DCGAN and BAGAN-GP, which are dark and have low contrast.

5 Conclusion

In this work, we propose a method to mitigate the problem of class imbalance in the imaging domain by utilizing generative models. The proposed method CAPGAN facilitates a conditional convolutional variational autoencoder (CVAE), which has an embedding component to perform training in a supervised fashion. A new objective function is applied to the CVAE training to improve the generative and reconstruction ability. Moreover, we present several pre-training strategies which could lead to produce balanced weights for the generative model. The generator and discriminator in CAPGAN are initialized by the pre-trained components in the CVAE and are trained in an adversarial setting. A gradient penalty term is added to the loss function of the discriminator to help stabilize the GAN training. We demonstrate the efficiency of CAPGAN on various datasets, including hand-crafted imbalanced datasets from general vision datasets and two imbalanced medical imaging datasets. We compare our proposed model with DCGAN and BAGAN-GP. The results show that CAPGAN outperforms these two methods by generating higher quality images given imbalanced datasets. Empirical results indicate that CAPGAN can retain high performance as the imbalance rate increases and can deliver acceptable results even under extreme imbalanced situations.

References

- [Arjovsky *et al.*, 2017] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.
- [Balasubramanian *et al.*, 2020] Ramji Balasubramanian, V Sowmya, EA Gopalakrishnan, Vijay Krishna Menon, VV Sajith Variyar, and KP Soman. Analysis of adversarial based augmentation for diabetic retinopathy disease grading. In *2020 11th International Conference on Computing, Communication and Networking Technologies (ICCCNT)*, pages 1–5. IEEE, 2020.
- [Bhatia and Dahyot, 2019] Snehal Bhatia and Rozenn Dahyot. Using wgan for improving imbalanced classification performance. In *AICS*, pages 365–375, 2019.
- [Braytee *et al.*, 2019] Ali Braytee, Wei Liu, Ali Anaissi, and Paul J Kennedy. Correlated multi-label classification with incomplete label space and class imbalance. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 10(5):1–26, 2019.
- [Deng *et al.*, 2009] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [Goodfellow *et al.*, 2014] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.
- [He *et al.*, 2016] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [Huang and Jafari, 2021] Gaofeng Huang and Amir Hossein Jafari. Enhanced balancing gan: minority-class image generation. *Neural Computing and Applications*, pages 1–10, 2021.
- [Kingma and Welling, 2013] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [Kodali *et al.*, 2017] Naveen Kodali, Jacob Abernethy, James Hays, and Zsolt Kira. On convergence and stability of gans. *arXiv preprint arXiv:1705.07215*, 2017.
- [Krizhevsky *et al.*, 2009] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [LeCun *et al.*, 2010] Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. *ATT Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- [Li *et al.*, 2018] Zhe Li, Yi Jin, Yidong Li, Zhiping Lin, and Shan Wang. Imbalanced adversarial learning for weather image generation and classification. In *2018 14th IEEE International Conference on Signal Processing (ICSP)*, pages 1093–1097. IEEE, 2018.
- [Mariani *et al.*, 2018] Giovanni Mariani, Florian Scheidegger, Roxana Istrate, Costas Bekas, and Cristiano Malossi. Bagan: Data augmentation with balancing gan. *arXiv preprint arXiv:1803.09655*, 2018.
- [Radford *et al.*, 2015] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [Rezaei *et al.*, 2020] Mina Rezaei, Tomoki Uemura, Janne Näppi, Hiroyuki Yoshida, Christoph Lippert, and Christoph Meinel. Generative synthetic adversarial network for internal bias correction and handling class imbalance problem in medical image diagnosis. In *Medical Imaging 2020: Computer-Aided Diagnosis*, volume 11314, page 113140E. International Society for Optics and Photonics, 2020.
- [Sampath *et al.*, 2021] Vignesh Sampath, Iñaki Maurtua, Juan José Aguilar Martín, and Aitor Gutierrez. A survey on generative adversarial networks for imbalance problems in computer vision tasks. *Journal of big Data*, 8(1):1–59, 2021.
- [Shoohi and Saud, 2020] Liqaa M Shoohi and Jamila H Saud. Dcgan for handling imbalanced malaria dataset based on over-sampling technique and using cnn. *Medico Legal Update*, 20(1):1079–1085, 2020.
- [Spanhol *et al.*, 2015] Fabio A Spanhol, Luiz S Oliveira, Caroline Petitjean, and Laurent Heutte. A dataset for breast cancer histopathological image classification. *Ieee transactions on biomedical engineering*, 63(7):1455–1462, 2015.
- [Taghanaki *et al.*, 2020] Saeid Asgari Taghanaki, Mohammad Havaei, Alex Lamb, Aditya Sanghi, Ara Danielyan, and Tonya Custis. Jigsaw-vae: Towards balancing features in variational autoencoders. *arXiv preprint arXiv:2005.05496*, 2020.
- [Waheed *et al.*, 2020] Abdul Waheed, Muskan Goyal, Deepak Gupta, Ashish Khanna, Fadi Al-Turjman, and Plácido Rogerio Pinheiro. Covidgan: data augmentation using auxiliary classifier gan for improved covid-19 detection. *Ieee Access*, 8:91916–91923, 2020.
- [Wang *et al.*, 2019] Qingfeng Wang, Xuehai Zhou, Chao Wang, Zhiqin Liu, Jun Huang, Ying Zhou, Changlong Li, Hang Zhuang, and Jie-Zhi Cheng. Wgan-based synthetic minority over-sampling technique: Improving semantic fine-grained classification for lung nodules in ct images. *IEEE Access*, 7:18450–18463, 2019.
- [Xiao *et al.*, 2017] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017.
- [Zhu *et al.*, 2017] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.