

# Latent Preserving Generative Adversarial Network for Imbalance classification

Tanmoy Dam  
SEIT

University of New South Wales  
Canberra, Australia

Md Meftahul Ferdaus  
ATMRI

Nanyang Technological University  
Singapore

Mahardhika Pratama, *Senior Member, IEEE*  
STEM

University of South Australia  
Adelaide, Australia

Sreenatha G. Anavatti  
SEIT

University of New South Wales  
Canberra, Australia

Senthilnath Jayavelu, *Senior Member IEEE*  
Institute for Infocomm Research

A\*STAR  
Singapore

Hussein A. Abbass, *Fellow, IEEE*  
SEIT

University of New South Wales  
Canberra, Australia

**Abstract**—Many real-world classification problems have imbalanced frequency of class labels; a well-known issue known as the “class imbalance” problem. Classic classification algorithms tend to be biased towards the majority class, leaving the classifier vulnerable to misclassification of the minority class. While the literature is rich with methods to fix this problem, as the dimensionality of the problem increases, many of these methods do not scale-up and the cost of running them become prohibitive. In this paper, we present an end-to-end deep generative classifier. We propose a domain-constraint autoencoder to preserve the latent-space as prior for a generator, which is then used to play an adversarial game with two other deep networks, a discriminator and a classifier. Extensive experiments are carried out on three different multi-class imbalanced problems and a comparison with state-of-the-art methods. Experimental results confirmed the superiority of our method over popular algorithms in handling high-dimensional imbalanced classification problems. Our code is available on <https://github.com/TanmDL/SLPPL-GAN>.

**Index Terms**—class imbalance, adversarial learning, oversampling techniques

## I. INTRODUCTION

Class imbalance classification is an old-standing problem. Common methods to address the problem include cost-sensitive classification, undersampling, and oversampling techniques [1]. The former requires a domain expert to transform precision and recall rates into a utility function. Undersampling the majority class could lead to information loss. Oversampling the minority class has been commonly used [2], [3], but raises the main challenge solved in this paper: how to generate new meaningful samples.

A deep oversampling framework (DOS) was proposed in [4] to fulfill the end-to-end requirement of deep learning algorithms. A limitation of the DOS lies in its dependency on the class-wise neighborhood sizes, which are determined by costly parameter tuning. The generative adversarial network (GAN) [5] has gained popularity due to its unique capability in generating synthetic but realistic samples. GAN relies on random noises that may yield a highly entangled process and disruptions of feature orientations. A two-stage framework called BAGAN [6] is proposed by combining autoencoder

(AE) with conditional GAN (cGAN) [7]. The latent code learnt via AE is fed to the cGAN to replace random noises. Because of its power in generating realistic samples, GAN is applied to oversample minority class(es) [8]. This approach often leads to boundary distortion as witnessed in [9], [10]. This problem inspires the development of a discriminative feature-based sampling (DFBS) method, where the main goal is to produce discriminative latent features achieved via the use of the triplet loss to learn AE. This approach can cause intra-class instances to stay together, while inter-class instances are pushed apart.

Instances generated by their method are likely to be close to the boundaries of the minority class(es), calling for a reliable classifier [11], [12]. A deep generative classifier (DGC) is developed in [13] for solving unstable prediction problems in the imbalance classification problem. A model is perturbed by replacing the fixed values of latent variables with a probability distribution over possible values, whereas the data are perturbed by generating feature/label. To mitigate the majority class influence on the classifier, probabilistic latent codes are over-sampled at different fractions rates. However, the learned probabilistic latent codes can’t guarantee well separation at the encoded manifold that may adversely influence the classifier’s performance.

Generative adversarial minority oversampling (GAMO) is proposed in [9]. A mixture of convex generators is proposed to mitigate the generation of majority class samples at the generator side. The mixture of generators forcefully generates samples within a specific minority class distribution by utilizing real instances from that minority class. However, their generators are limited by generating class instances instead of generating the real data distribution. To determine class instances from the real data distribution, an adversarial game is played between the mixture of generators, a discriminator and a classifier. True data distributions could be far from such a convex hull of minority class(es) leading to the generation of instances with low information or overlapped instances.

In deep generative classifiers, cross entropy loss [13] or

mean-squared loss [9] are commonly used to update the classifier. These loss functions may create hard partition(s) between majority and minority classes, where the decision boundary is influenced by the majority class samples leading to over-fitting. To mitigate the influence of majority class samples on the decision boundary, a classifier is updated twice in [9], [13] based on concepts borrowed from [14], [15]. First, it is updated through the real feature distribution. Afterward, generated minority class samples are used to update the classifier. A similar approach is also followed in this paper. In addition, mixtures of minority class-specific generators are replaced with a single generator only. This is done by deploying a domain-constrained AE to learn the class-specific latent code, preserved and used as a prior for the generative network. Besides, in GAMO, the minority samples are generated only from the feature space while ignoring the possibility of an adversarial data space oversampling approach. Thus, the only three-player strategy of adversarial minority oversampling of GAMO is extended here to two different three-player strategies: 1) adversarial minority oversampling (AMO); 2) adversarial data space oversampling (ADSO).

The main contributions of this paper are summarized:

- We define a joint learning framework to preserve the latent space in a low dimensional manifold by utilizing a supervised autoencoder (AE). The learned latent space works as a prior for a single generator, which engages in an adversarial game with a classifier and a discriminator. Replacement of minority class(es)-convex generators with only one generator makes our proposed three player adversarial architecture more scalable than GAMO [9].
- We leveraged all two possible strategies among three players game to improve the classifier’s performance. First strategy is related to ADSO and the last strategy is based on AMO.
- A set of experiments have been carried out under different imbalance ratios, where the experimental results support that an ADSO-based classifier performs better than the AMO-based approach and state-of-the-art baselines, and it does so with significant margins.

## II. PROBLEM FORMULATIONS

Let us assume a set of  $N$  samples  $\{x_i, y_i\}_{i=1}^N \in (X_{org}, Y_{org})$ , a multi-class imbalance dataset distributed over  $C$  classes. Here,  $x_i \in R^d$  is an  $i$ -th input image with its corresponding target class  $y_i$ . Without loss of generalisation, the class distribution follows  $p_{maj} \geq p_2 \geq \dots p_l \geq \dots \geq p_{min}$ ,  $N = \sum_{l=c}^C p_l$ , in which  $p_{maj}$  and  $p_{min}$  denote the majority class and the minority class, respectively. The relationship between the majority class and minority class(es) is set as  $p_{maj} \geq 50 * p_{min}$  following the problem setting of [9], [16], [17]. The main objective is to design a deep neural network estimating the underlying data distribution of  $C$  classes, thereby producing robust decision boundaries.

The network structure of our approach consists of three parts: supervised latent preserving prior learning (SLPPL), ADSO and AMO. The adversarial game is played among

preserved latent prior generator( $G_\omega$ ), a discriminator( $Dis_\xi$ ) and a classifier( $Q_\rho$ ) to improve the classification performance. In SLPPL, the  $Enc_\theta$  encodes the original class instances ( $x_i \in X_{org}$ ) into a lower-dimensional latent space ( $z_i \in Z$ ) and  $Dec_\phi$  takes the encoded latent space ( $z_i$ ) to produce the reconstructed original data ( $\hat{x}_i \in X_{org}$ ). The  $Enc_\theta$  learns class distribution ( $Enc_\theta(y_i|x_i)$ ) for  $i$ -th class instance. After learning the latent space ( $z_i$ ) from data directly, a Gaussian multivariate normal distribution is constructed as a prior for  $G_\omega$ . In the adversarial game, we introduce mainly two strategies, ADSO and AMO, among the pretrained prior  $G_\omega$ , a  $Dis_\xi$  and a  $Q_\rho$ .

## III. OUR APPROACH

### A. SLPPL

High dimensional data always retain their characteristics in a low dimensional encoded manifold that inspired us to design a latent prior for each class distribution [18]. A low dimensional manifold learning approach attains optimum performance by freely moving the latent space when the distribution of data is uniform in nature [18]. If data are not distributed uniformly, the majority class-driven latent space always dominates minority classes. To mitigate a biased prediction problem favouring only the majority class, we utilize jointly  $(X_{org}, Y_{org})$  a class distribution learning approach as well as a reconstruction learning approach to deliver the latent space bounded as much as possible. The bounded latent space is obtained by considering label information in SLPPL under the deterministic autoencoder (AE) framework. Unlike, arbitrarily chosen prior based on VAE [19], this approach reduces the stochasticity in the latent space. Finally, the obtained latent space is derived as a significant smooth manifold. The AE objective is to minimize the reconstruction loss ( $L_{rec}$ ) between input data ( $X_{org}$ ) and decoded output ( $Dec_\phi(Enc_\theta(X_{org}))$ ):

$$L_{rec}(X_{org}, Dec_\phi(Enc_\theta(X_{org}))) = \mathbb{E}_{x_i \in X_{org}} \|x_i - Dec_\phi(Enc_\theta(x_i))\|_2 \quad (1)$$

Similarly, for the bounded latent space, the encoder network ( $Enc_\theta$ ) estimates the class distribution ( $Y_{org}$ ).

$$L_{bce}(Y_{org}, (Enc_\theta(X_{org}))) = \mathbb{E}_{y_i \in Y_{org}} y_i \log(Enc_\theta(x_i)) \quad (2)$$

Finally, merging two losses, the final loss function for SLPPL can be minimized with respect to two networks parameter.

$$\min_{\theta, \phi} L_{SLPPL} = L_{rec} + L_{bce} \quad (3)$$

The proposed SLPPL is able to preserve stable network parameters ( $\phi, \theta$ ), and to maintain stable class distributions in the low-dimensional manifold where (3) is solved using the ADAM optimiser [16]. The  $i$ -th sample is easily encoded in the latent space afterward as  $z_i \in Z = Enc_\theta(x_i)$  and preserves a multivariate normal distribution (MND) [18] for the  $i$ -th class. The MND definition for the  $i$ -th class is  $z_i \in \mathcal{N}(\mu_i, \sigma_i^2)$  where  $\mu_i \in R^q$  and  $\sigma_i \in R^{q \times q}$  are the mean and variance of the latent space respectively. The prior ( $z_i$ ) preserving latent space

is applied to improve image generation through adversarial game between  $Q_\rho$  and  $Dis_\xi$  through two adversarial strategies: ADSO and AMO. In both the oversampling cases, the  $G_\omega$  network structure is same as  $Dec_\phi$ . Hence, the initialisation of  $G_\omega$  network weights is taken from pretrained  $Dec_\phi$ , whereas the structure of  $Dis_\xi$  and  $Q_\rho$  are almost similar to  $Enc_\theta$  network but the last layer of  $Dis_\xi$  gives a single output. Besides, the learned feature layers of  $Enc_\theta$  are used to initialise the weights of both the  $Dis_\xi$  and  $Q_\rho$ . However, to reduce the over-fitting effect at  $Q_\rho$ , we apply dropout after activation function [20].

## B. ADSO & AMO

In ADSO, the minority class(es) is repeated to form a balanced data distribution  $(X_{bal}, Y_{bal})$  in the data space. An adversarial game is played among  $G_\omega$ ,  $Dis_\xi$  and  $Q_\rho$  afterward, where the  $G_\omega$  network is updated by fooling only the discriminator  $Dis_\xi$  but by favouring the classifier  $Q_\rho$ . The generator network  $G_\omega$  aims to generate samples to be classified by  $Q_\rho$  as the same class, and thus, the generated samples assign real scores while updating the generator  $G_\omega$  through the classifier  $Q_\rho$ . The discriminator  $Dis_\xi$  enforces the generator  $G_\omega$  to follow the real data distribution. While updating the classifier  $Q_\rho$  through the generator  $G_\omega$ , generated samples assign a fake score to the classifier  $Q_\rho$ . In other words, the classifier  $Q_\rho$  assigns a high probability in such a way that the generated samples are classified as other classes. We can formulate the following optimisation problem for three players adversarial game among  $G_\omega$ ,  $Dis_\xi$  and  $Q_\rho$ :

$$\min_{\omega, \rho} \max_{\xi} L_{ADSO}(G_\omega, Dis_\xi, Q_\rho) \quad (4)$$

The total loss can be expressed as  $L_{ADSO} = L_{ADSO}^{G_\omega} + L_{ADSO}^{Dis_\xi} + L_{ADSO}^{Q_\rho}$ , where

$$L_{ADSO}^{G_\omega} = \mathbb{E}_{G(z_i) \in Z} (f(1 - Dis_\xi(G(z_i))) + \mathbb{E}_{G(z_i) \in Z} (y_i \log(Q_\rho(G_\omega(z_i)))) \quad (5)$$

$$L_{ADSO}^{Q_\rho} = \mathbb{E}_{y_i \in Y_{bal}} y_i \log(Q_\rho(x_i)) + \mathbb{E}_{G(z_i) \in Z} (y_i \log(1 - Q_\rho(G_\omega(z_i)))) \quad (6)$$

$$L_{ADSO}^{Dis_\xi} = \mathbb{E}_{x_i \in X_{bal}} (f(Dis_\xi(x_i)) + \mathbb{E}_{G(z_i) \in Z} f(1 - Dis_\xi(G_\omega(z_i)))) \quad (7)$$

In AMO, generator network  $G_\omega$  aims to generate samples to be classified by  $Q_\rho$  as the same class, and thus, the generated samples assign real scores while updating the generator  $G_\omega$  through the classifier  $Q_\rho$ . To mitigate, majority class biases at  $Q_\rho$ , is updated through the real samples as well as minority-class generated samples.  $Dis_\xi$  network forces  $G_\omega$  to learn the real-data distributions. The following optimisation can be formulated by playing three players game among  $G_\omega$ ,  $Dis_\xi$  and  $Q_\rho$ :

$$\min_{\omega, \rho} \max_{\xi} L_{AMO}(G_\omega, Dis_\xi, Q_\rho) \quad (8)$$

The total loss can be expressed as  $L_{AMO} = L_{AMO}^{G_\omega} + L_{AMO}^{Dis_\xi} + L_{AMO}^{Q_\rho}$ , where

$$L_{AMO}^{G_\omega} = \mathbb{E}_{G(z_i) \in Z} (f(1 - Dis_\xi(G(z_i))) + \mathbb{E}_{G(z_i) \in Z} (y_i \log(1 - Q_\rho(G_\omega(z_i)))) \quad (9)$$

$$L_{AMO}^{Q_\rho} = \mathbb{E}_{y_i \in Y_{orig}} y_i \log(Q_\rho(x_i)) + \mathbb{E}_{G(z_j) \setminus P_{maj} \in Z} (y_j \log(Q_\rho(G_\omega(z_j)))) \quad (10)$$

$$L_{AMO}^{Dis_\xi} = \mathbb{E}_{x_i \in X_{orig}} (f(Dis_\xi(x_i)) + \mathbb{E}_{G(z_i) \in Z} f(1 - Dis_\xi(G_\omega(z_i)))) \quad (11)$$

Cross-entropy (CE) loss and complementary CE (CCE) loss are represented by the expressions  $\log Q(\cdot)$  and  $\log(1 - Q(\cdot))$ , respectively. For both the cases (AMO, ADSO), the functional operator  $f(\cdot)$  selects different GANs types. For vanilla [5] and Wasserstein GANs (WGANs) [21], ( $f$ ) is represented as  $f(x) = \log x$  and  $f(x) = x$  respectively. We follow WGAN's zero center (0)-gradient penalty (0-gp) for all the GAN strategies [21].

## IV. EXPERIMENTS

### A. Datasets

Two single-channel (MNIST [22], and Fashion-MNIST [23]) and a three-channel (CelebA [24]) image sets are used here. Properties of these datasets are tabulated in Table I where IR indicates the imbalance ratio.

### B. Evaluation metrics and baselines

Five metrics are used here to measure the imbalance classification performance: 1) average class specific accuracy (ACSA); 2) macro-averaged F-measure ( $F_{macro}$ ); 3) macro-averaged geometric mean ( $G_{mean}$ ); 4) precision of majority class ( $P_{maj}$ ); and 5) recall of minority class ( $R_{min}$ ).

Our proposed method is compared against seven different baselines namely BAGAN [6], DFBS [25], GAMO [9], mmDGMs [26], BayesCNN [27], DGC [13], and data space oversampling (DSO) plus baseline  $Q_\rho$ . Results of all above-mentioned baselines are adopted from [13]. For our two proposed models, the size of latent space ( $z_i \in R^q$ ) is set to  $q = 64$  for the MNIST and Fashion-MNIST datasets [13], assigned as  $q = 128$  for the CelebA dataset [13].

### C. Numerical Results

The overall classification performances in terms of average of ACSA,  $F_{macro}$ , and  $G_{macro}$  are reported in Table II. The ADSO-based strategy we suggest is the one that performs the best overall for handling imbalance classification. In contrast with the ADSO, some limitations of baselines for getting comparative poorer performance are as follows: DFBS can not create sufficient margins among classes. GAMO utilizes computationally expensive MSE loss that requires real samples

TABLE I  
THE DETAILED DESCRIPTION OF EXPERIMENTAL DATASETS

Datasets	Data Dimensions	IR	Classes	Training Set	Testing Set
MNIST	$28 \times 28 \times 1$	100	10	[4000, 2000, 1000, 750, 500, 350, 200, 100, 60, 40]	[980, 1135, 1032, 1010, 982, 892, 985, 1028, 974, 1009]
Fashion-MNIST	$28 \times 28 \times 1$	100	10	[4000, 2000, 1000, 750, 500, 350, 200, 100, 60, 40]	[1000, 1000, 1000, 1000, 1000, 1000, 1000, 1000, 1000, 1000]
CelebA	$64 \times 64 \times 3$	100	5	[15000, 1500, 750, 300, 150]	[2660, 5422, 412, 3428, 535]

TABLE II  
OVERALL CLASSIFICATION PERFORMANCE ON VARIOUS DATASETS

Methods	MNIST			FMNIST			CelebA		
	$ACSA$	$F_{macro}$	$G_{macro}$	$ACSA$	$F_{macro}$	$G_{macro}$	$ACSA$	$F_{macro}$	$G_{macro}$
BAGAN	0.8848	0.8785	0.9295	0.7814	0.7610	0.8546	0.5972	0.5152	0.6554
DFBS	0.7812	0.7838	0.8683	0.5135	0.4620	0.6382	0.2109	0.1335	0.2664
GAMO	0.8826	0.8794	0.9308	0.7929	0.7880	0.8740	0.6409	0.5903	0.7472
BayesCNN	0.9158	0.9141	0.9512	0.7934	0.7835	0.8701	0.5517	0.4936	0.6534
mmDGMs	0.9066	0.9039	0.9449	0.8091	0.7984	0.8796	0.3760	0.0618	0.3754
DGC	0.9480	0.9474	0.9704	0.8364	0.8314	0.9010	0.6755	0.6454	0.7779
DSO+ $Q_\theta$	0.9339	0.9325	0.9619	0.8450	0.8436	0.9089	0.7078	0.6663	0.7997
AMO	0.9403	0.9386	0.9656	0.8428	0.8378	0.9046	0.6702	0.6210	0.7628
ADSO	<b>0.9613</b>	<b>0.9609</b>	<b>0.9781</b>	<b>0.8675</b>	<b>0.8648</b>	<b>0.9221</b>	<b>0.7595</b>	<b>0.7217</b>	<b>0.8359</b>

TABLE III  
CLASSIFICATION PERFORMANCE ON THE LARGEST ( $P_{maj}$ ) AND SMALLEST CLASS ( $R_{min}$ )

Methods	MNIST		FMNIST		CelebA	
	$R_{min}$	$P_{maj}$	$R_{min}$	$P_{maj}$	$R_{min}$	$P_{maj}$
BAGAN	0.5354	0.8541	0.7306	0.5709	0.0192	0.5064
DFBS	0.5946	0.5118	0.4412	0.3395	0.0522	0.2174
GAMO	0.6394	0.8812	0.7928	0.6165	0.2302	0.6687
BayesCNN	0.7578	0.8896	0.8474	0.6022	0.1063	0.5225
mmDGMs	0.6525	0.8459	0.8160	0.5942	0.0006	0.4110
DGC	0.8276	0.9270	0.8864	0.6900	0.2987	0.7603
DSO+ $Q_\theta$	0.7393	0.8761	0.9130	0.7598	0.2783	0.7483
AMO	0.7363	0.9055	0.9170	0.6752	0.2405	0.5785
ADSO	<b>0.8840</b>	<b>0.9348</b>	<b>0.9510</b>	<b>0.7885</b>	<b>0.3223</b>	<b>0.8132</b>

in each generator. The mode collapse problem may appear in BAGAN due to the initialization mechanism of subsequent GAN. Though Bayes CNN has adopted a model perturbation strategy, very few instances in minority class(es) may not be sufficient to train the complicated model. mmDGMs determines class boundaries by adopting the discriminative classifier, limiting their performance in imbalanced datasets. The perturbation mechanism of both data and model supports the DGC to outperform the above-mentioned baselines. However, well separation at the encoded manifold can not be guaranteed by the learned probabilistic latent codes. A high recall on minority class is expected from the  $Q_\theta$  while maintaining a high precision on majority class. Recall of the smallest class ( $R_{min}$ ) and precision of the largest class ( $P_{maj}$ ) are listed in Table III. For Fashion-MNIST, it is observed that the performance of all six baselines have improved in the minority class but has not been significant in the majority class. These results confirm that class boundaries are not determined clearly by these existing methods. In the three-channel CelebA image dataset, poor performance is seen in both majority and

minority classes from the first five baselines. Since the learned probabilistic latent in DGC does not guarantee well separation at the encoded manifold. Similar phenomena are also observed in minority generative samples for our proposed AMO method because of the majority class influences at  $Q_\theta$ . Even if we preserve the prior for the  $G_\omega$ ,  $Q_\theta$  has influenced it to generate majority classes, i.e., AMO can't beat ADSO. In contrast, under three-player GAN settings, the proposed ADSO-based  $G_\omega$  tries to fool itself by generating a subset of each class sample that is relevant to  $Q_\theta$ .

## V. CONCLUSIONS

We propose a latent preserving based deep generative model for handling imbalanced classification problems. The class constraints AE is used to preserve the latent space, utilized as a prior for  $G_\omega$ . By playing two adversarial games among the latent space preserved across  $G_\omega$ , a  $Dis_\xi$ , and a  $Q_\theta$ , improvement in the  $Q_\theta$ 's performance is witnessed. From the experimental results on all three datasets, our ADSO-based strategy performed better than all the baselines. Our future work includes real-world applications like semiconductor fault detection, medical diagnosis, etc.

## ACKNOWLEDGEMENTS

S. Jayavelu and Md M. Ferdous acknowledges funding from the Accelerated Materials Development for Manufacturing Program at A\*STAR via the AME Programmatic Fund by the Agency for Science, Technology and Research under Grant No. A1898b0043. T. Dam acknowledges UIPA funding from UNSW Canberra.

## REFERENCES

- [1] P. Branco, L. Torgo, and R. P. Ribeiro, "A survey of predictive modeling on imbalanced domains," *ACM Computing Surveys (CSUR)*, vol. 49, no. 2, pp. 1–50, 2016.

- [2] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: synthetic minority over-sampling technique," *Journal of artificial intelligence research*, vol. 16, pp. 321–357, 2002.
- [3] H. He, Y. Bai, E. A. Garcia, and S. Li, "Adasyn: Adaptive synthetic sampling approach for imbalanced learning," in *2008 IEEE international joint conference on neural networks (IEEE world congress on computational intelligence)*. IEEE, 2008, pp. 1322–1328.
- [4] S. Ando and C. Y. Huang, "Deep over-sampling framework for classifying imbalanced data," in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 2017, pp. 770–785.
- [5] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," *Advances in neural information processing systems*, vol. 27, 2014.
- [6] G. Mariani, F. Scheidegger, R. Istrate, C. Bekas, and C. Malossi, "Bagan: Data augmentation with balancing gan," *arXiv preprint arXiv:1803.09655*, 2018.
- [7] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014.
- [8] G. Douzas and F. Bacao, "Effective data generation for imbalanced learning using conditional generative adversarial networks," *Expert Systems with applications*, vol. 91, pp. 464–471, 2018.
- [9] S. S. Mullick, S. Datta, and S. Das, "Generative adversarial minority oversampling," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 1695–1704.
- [10] S. Santurkar, L. Schmidt, and A. Madry, "A classification-based study of covariate shift in gan distributions," in *International Conference on Machine Learning*. PMLR, 2018, pp. 4480–4489.
- [11] C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra, "Weight uncertainty in neural network," in *International Conference on Machine Learning*. PMLR, 2015, pp. 1613–1622.
- [12] A. Srivastava, L. Valkov, C. Russell, M. U. Gutmann, and C. Sutton, "Veegan: Reducing mode collapse in gans using implicit variational learning," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, pp. 3310–3320.
- [13] X. Wang, Y. Lyu, and L. Jing, "Deep generative model for robust imbalance classification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14 124–14 133.
- [14] R. Statistics, "The approach based on influence functions," *New York*, 1986.
- [15] P. W. Koh and P. Liang, "Understanding black-box predictions via influence functions," in *International Conference on Machine Learning*. PMLR, 2017, pp. 1885–1894.
- [16] T. Dam, M. M. Ferdous, S. G. Anavatti, S. Jayavelu, and H. A. Abbass, "Does adversarial oversampling help us?" in *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 2021, pp. 2970–2973.
- [17] T. Dam, S. G. Anavatti, and H. A. Abbass, "Mixture of spectral generative adversarial networks for imbalanced hyperspectral image classification," *IEEE Geoscience and Remote Sensing Letters*, 2020.
- [18] P. Ghosh, M. S. Sajjadi, A. Vergari, M. Black, and B. Scholkopf, "From variational to deterministic autoencoders," in *International Conference on Learning Representations*, 2019.
- [19] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [20] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov, "Improving neural networks by preventing co-adaptation of feature detectors," *arXiv preprint arXiv:1207.0580*, 2012.
- [21] L. Mescheder, A. Geiger, and S. Nowozin, "Which training methods for gans do actually converge?" in *International conference on machine learning*. PMLR, 2018, pp. 3481–3490.
- [22] Y. LeCun, "The mnist database of handwritten digits," <http://yann.lecun.com/exdb/mnist/>, 1998.
- [23] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017.
- [24] Z. Liu, P. Luo, X. Wang, and X. Tang, "Deep learning face attributes in the wild," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 3730–3738.
- [25] Y.-H. Liu, C.-L. Liu, and S.-M. Tseng, "Deep discriminative features learning and sampling for imbalanced data problem," in *2018 IEEE International Conference on Data Mining (ICDM)*. IEEE, 2018, pp. 1146–1151.
- [26] C. Li, J. Zhu, and B. Zhang, "Max-margin deep generative models for (semi-) supervised learning," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 11, pp. 2762–2775, 2017.
- [27] K. Shridhar, F. Laumann, and M. Liwicki, "A comprehensive guide to bayesian convolutional neural network with variational inference," *arXiv preprint arXiv:1901.02731*, 2019.

This figure "fig1.png" is available in "png" format from:

<http://arxiv.org/ps/2209.01555v1>