

COMP 551: Miniproject 4

Lucas Bennett, Bjorn Christensen
McGill University

Reproducibility Summary

2 Scope of Reproducibility

In this paper, we look to reproduce the findings of Kodali et al.'s (2017) paper "On Convergence and Stability of GANs", which demonstrated the ability for a regret minimization gradient penalty term introduced in their proposed DRAGAN model to provide improved stability as compared to vanilla GAN algorithms. The original paper proposed the idea of studying GAN training dynamics through the form of regret minimization, and how this process would eliminate the prevalence of mode collapse. We looked to reproduce the results of paper that demonstrate this finding, showing how regret minimization effects the stability of loss in generator and discriminator models, and effects inception score.

9 Methodology

10 The original code provided by the authors was not sufficient to fully reproduce the experiments and furthermore allowed
11 little expansion on their results. Therefore, we turned to a number of sources to create our own code that would replicate
12 the findings of the original paper. We decided to utilize PyTorch in order to create and evaluate the models. Although
13 the authors used TensorFlow to conduct their experiments, they directed readers to the PyTorch-GAN git repository
14 for a general implementation of DRAGAN using pytorch. The source code of this implementation highly influence
15 the architecture of our recreation. To study inception scoring, we implemented code influence by the PyTorch-Ignite
16 module documentation.

17 Results

18 While our results are indicative of the results demonstrated in the original paper, due to computational load, we were
19 not able to demonstrate complete corroboration. This is the case particularly for our findings of inception score. As we
20 were unable to run a model capable of reaching the maximum inception scores highlighted in the paper, it is difficult to
21 assess the added stability of a DRAGAN model solely from our inception score results. On the contrary, our results do
22 demonstrate additional stability in G and D loss in DRAGAN over vanilla GAN algorithms, supporting the results of the
23 paper. We observed the additional stability of DRAGAN in G and D loss on both MNIST and CelebA datasets. As such,
24 although our findings cannot completely validate the original study, they do provide strong suggestions in their support.

25 What was easy

Creating the discriminator and generator networks were quite simple using PyTorch and training on the MNIST dataset was able to be done using readily available code.

28 What was difficult

Interpreting the authors methodology in order to follow their pipeline using different packages was quite difficult. The DRAGAN and inception scoring implementations were not compatible, meaning we had to write a large portion of the code from scratch in order to integrate the steps. Furthermore, the code provided by the authors was only suited to analysing the MNIST data-set. As we were far more interested in datasets with RGB images, the need to translate the code to suit multiple channels and convolutional neural networks was also a challenge.

34 **1 Introduction**

35 Generative adversarial networks (GANs) are foundational models in image generation. Their standard structure consists
36 of two competing networks, a generator and a discriminator, engaging in a zero-sum game. However, the non-convex
37 nature of this classical architecture has produced numerous theoretical challenges which have proved difficult to
38 overcome. One such challenge is the prevalence of *mode collapses*, a phenomena in which the network converges to a
39 local-equilibria (modes), and the generator gets trapped without the discriminator being able to remove it. In their 2018
40 paper "On Convergence and Stability of GANs", Kodali et al. (2017) proposed a model possessing a novel training
41 method - which they called DRAGAN (Deep Regret Analytic Generative Adversarial Networks) - that could overcome
42 the challenge of mode collapse. It entails the addition of a gradient penalty during the learning process which aims to
43 pull the model out of local equilibria. The validity of their findings are investigated here.

44 **2 Scope of reproducibility**

45 In this paper, we look to reproduce the findings of Kodali et al.'s (2017) paper "On Convergence and Stability of GANs"
46 The two major claims from that paper that we are looking to investigate are,

- 47 • Claim 1: Introducing a gradient penalty into the classic GAN training procedure will reduce model instability.
48 • Calim 2: The inception score of a DRAGAN model will be more robust than the inception score of a vanilla
49 GAN model.

50 **3 Methodology**

51 For MNIST, we used a PyTorch DRAGAN implementation which was constructed based on the author's original
52 tensorflow implementation. We made some slight modifications to the code for ease of use but generally it is identical.
53 We did not conduct an inception score experiment on the MNIST dataset due to computational load, instead we modified
54 the external code in order to analyze the dynamics of the D and G loss. For experiments on the CelebA dataset, we
55 constructed our own code and DRAGAN implementations. This was done as the author had not provided code for
56 inception scoring convolutional models. Our code was highly influenced by implementations and examples in the
57 PyTorch documentation.

58 **3.1 Model descriptions**

59 Both the DRAGAN and vanilla GAN models used identical discriminator(D) and generator(G) networks. For tests with
60 the MNIST dataset, both the D and G networks consisted of a single hidden linear layers. For tests with the CIFAR10
61 and CelebA datasets, the D and G networks consisted of 4 hidden convolutional layers with ReLu activation, resulting
62 in 10,949,504 parameters in the G network and 2,765,568 parameters in the D network (all trainable). The models
63 were not pretrained; train was done using an ADAM optimizer. Training of the DRAGAN model different due to the
64 inclusion of a regularization term applying a gradient penalty.

65 **3.2 Datasets**

66 The Datasets used were MNIST and CelebA. The MNIST dataset contains black and white images of hand-drawn
67 numbers ranging from 0-9. The images are composed of 28X28 pixels. The distribution of each category in the MNIST
68 training set ranges from 5421 to 6742 occurrences, totaling 60,000 total images, which means we have a good spread
69 to learn each number from. Furthermore the numbers in the train set were derived from a separate group of writers
70 from the test set, about 250 for each. The CelebA set includes 202,599 photos selected from 10,177 celebrities. These
71 images are described by 3 64x64 matrices (pixel dimensions), with each matrix representing an RGB channel.

72 For the MNIST dataset we used a 60,000-10,000 train-test split. CelebA uses a 162,770-19,962 train-test split.

73 No preprocessing was done on the MNIST dataset. In order to preprocess the CelebA dataset, we normalized the pixel
74 values of our images from 0-255 to 0-1 and centered the images using torchvisions transform utility.

75 MNIST: <https://www.tensorflow.org/datasets/catalog/mnist>

76 CelebA: <https://www.kaggle.com/datasets/jessicali9530/celeba-dataset>

77 **3.3 Hyperparameters**

78 We chose the hyperparameters using manual search to find those that would provide sufficient qualitative results without
79 requiring extensive run times. For training of the CelebA dataset, we found smaller learning rates (0.0001) were
80 required to generate sufficient results. Larger learning rates blew up d-loss and resulted in the extinction of g-loss.
81 Conversely, smaller learning rates were deemed too slow. We also manually chose the size of the latent space to be 1000.
82 This was done by finding a balance between the decreasing speed of training and increasing model performance as
83 we increased the size of the latent space. We inherited the lambda gradient penalty parameter set in the original paper
84 (lambda = 10). We found a batch size of 32, produced good results while maintaining fast experimental speed.

85 **3.4 Experimental setup and code**

86 Experiments on the MNIST dataset consisted of a standard mini-batch stochastic gradient descent algorithm to train
87 the vanilla GAN and DRAGAN models, which was modified to record the G and D loss of every iteration during
88 training and visualized the fake image that had been generated. We then plotted the results of the G and D losses using
89 Matplotlib.
90 The setup of our experiments on the CelebA dataset were more sophisticated. We define two training step functions
91 for vanilla GAN and DRAGAN models. These consisted of a gradient descent optimization step on both the G and
92 D networks, with and without gradient penalty respectively. We then sequentially initialized a PyTorch-Ignite engine
93 with each training step, initializing new G and D networks and optimizers. We defined a series of operations that
94 would be conducted as the engine trains a dataset, generation of fake images on a set of 64 fixed noise samples every 500
95 iterations, and recording D and G loss each iteration. This allowed us to analyze the dynamics and the performance of
96 the model quantitatively and qualitatively throughout training. After every epoch of training, we recorded the FID and
97 inception score of the model.

98 **3.5 Computational requirements**

99 Our experiment on the MNIST dataset was done across 1000 iterations of training conducted on the apple M2 CPU,
100 which required roughly 2 hours per model. We do not recommend using a CPU to run these experiments, and would
101 suggest using a device with CUDA support.
102 Experiments on the CelebA dataset were done using a cuda0 GPU with a nccl parallel backend. Each epoch of training
103 on all batches required under 4 minutes (training batch size of 64), although it varies with batch size. Evaluation of
104 FID and inception score after each epoch took roughly 2.5 minutes when using a test-batch size of 256. Increasing the
105 size of the test batch greatly increased the speed of our experiments without altering results. We do not suggest greatly
106 altering the batch size of the training set as it minimally reduces experimental time, but has a large influence on results
107 (see additional results 2).

108 **4 Results**

109 Our results found faster rates of convergence in G loss for the DRAGAN model in both the MNIST and CelebA datasets
110 as compared to the vanilla GAN model. Furthermore, there was far less spread in both G and D loss when using
111 DRAGAN, supporting the claim of the original paper that DRAGAN would increase the stability of the model. However,
112 we found the inception score of the DRAGAN to be generally reduced as compared to the inception score of vanilla
113 GAN. It is of note that the inception score of DRAGAN did show a consistent increase whereas the inception score of the
114 vanilla GAN showed a sharp decrease followed by a recovery. These results do suggest the results of the original paper,
115 however, we would need more iterations to sufficiently reproduce them.

116 **4.1 Results reproducing original paper**

117 **4.1.1 Result 1**

118 In this experiment, we found analyzed the dynamics of G and D loss of the vanilla GAN and DRAGAN models on the
119 MNIST and CelebA datasets. The results supports the first claim of the paper, suggesting that DRAGAN improves
120 model stability. We were able to successfully reproduce these results on both the MNIST dataset and CelebA dataset as
shown in figure 1 and 2 respectively.

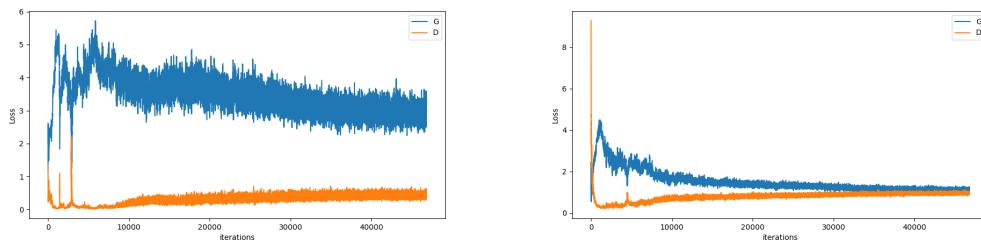


Figure 1: (MNIST) G and D loss for vanilla GAN (left) and DRAGAN (right) models.

121

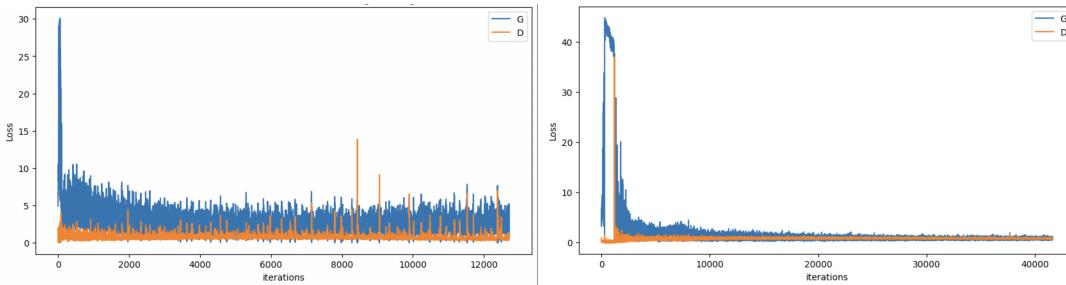


Figure 2: (CelebA) G and D loss for vanilla DSGAN (left) and DRAGAN (right) models.

122 We see faster convergence and limited spred of G loss when we introduce the regularization term to the model.

123 **4.1.2 Result 2**

124 These results demonstrate the inception score of vanilla GAN and DRAGAN at various epochs. The number of epochs
125 we were able to study is limited by computational need, however, the robustness of the inception score for the DRAGAN
model suggests support of claim 2.

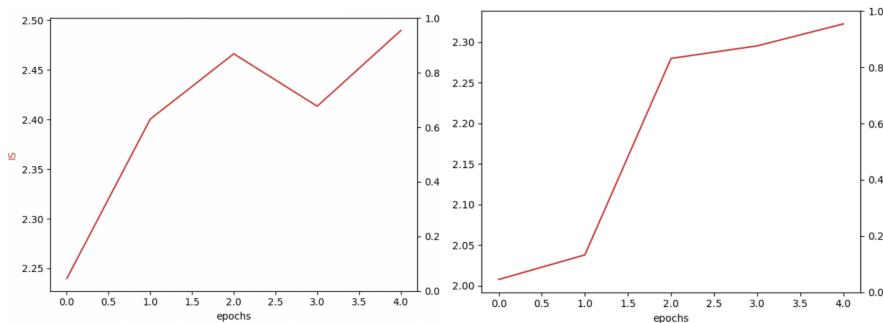


Figure 3: (CelebA) Inception scores for vanilla DSGAN (left) and DRAGAN (right) models across various epochs of training.

126

127 **4.2 Results beyond original paper**

128 **4.2.1 Additional Result 1**

129 We wanted to provide additional qualitative results to those presented in the paper. To do so, we generated a set of 64
 130 fixed noise latent space samples and visualized the generated images of the models after every 500 iterations. A GIF of
 131 the models progress is saved with our code.

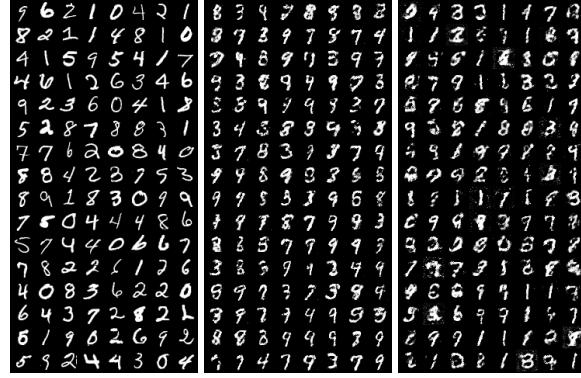


Figure 4: (MNIST) Real(left) and generated samples from vanilla GAN (middle) and DRAGAN (right)



Figure 5: (CelebA) Comparison of real and generated images for DRAGAN (top) and vanilla GAN (bottom)

132 In figure 4 and figure 5, we see that the vanilla GAN model is producing more accurate depictions of the numbers and
 133 faces. This is an interesting finding which may not be expecting with the results of the original paper. It is of note, that

134 the vanilla GAN appears to be overproducing images of one type. For example, there is a strong over representation of
135 the number nine in the MNIST generated images. This could be indicative of mode collapse.

136 **5 Discussion**

137 While our results are indicative of the original papers claims, due to computational load, we were not able to demonstrate
138 complete corroboration. This is the case particularly for our findings of inception score. As we were unable to run
139 a model capable of reaching the maximum inception scores highlighted in the paper. Furthermore, we were limited
140 in the number of epochs we were able to record the inception score, therefore, we cannot say with certainty that the
141 DRAGAN algorithm prevents mode collapse. However, the variance observed in the inception score of vanilla GAN
142 opposed to the steady increasing inception score of DRAGAN is promising. On the contrary, our results do demonstrate
143 additional stability in G and D loss in DRAGAN over vanilla GAN algorithms, supporting the first claim of the paper.
144 We observed the additional stability of DRAGAN in G and D loss on both MNIST and CelebA datasets. As such,
145 although our findings cannot completely validate the original study, they do provide strong suggestions in their support.

146 **5.1 What was easy**

147 Creating the discriminator and generator networks were quite simple using PyTorch and training on the MNIST dataset
148 was able to be done using readily available code.

149 **5.2 What was difficult**

150 Interpreting the authors methodology in order to follow their pipeline using different packages was quite difficult. The
151 DRAGAN and inception scoring implementations were not compatible, meaning we had to write a large portion of the
152 code from scratch in order to integrate the steps. Furthermore, the code provided by the authors was only suited to
153 analysing the MNIST data-set. As we were far more interested in datasets with RGB images, the need to translate the
154 code to suit multiple channels and convolutional neural networks was also a challenge.

155 **References**

156 Kodali,N. Abernethy, J. Hays, J. Kira, Z. (2017). On Convergence and Stability of GANs.