
ESTR3114

CHAN Yuk Kit
1155193356

HUI Man Hei
1155194718

Abstract

1 Generative Adversarial Networks (GANs) have emerged as a type powerful un-
2 supervised learning framework for generative modeling, being able to learn and
3 generate complex images, audio and other data. However, they are also known for
4 being notoriously difficult to train due to its adversarial nature between the genera-
5 tor and discriminator, which leads to high instability. This paper explores various
6 optimization strategies for GANs, first investigating the effectiveness of vanilla
7 GANs and Wasserstein GANs, along with three basic optimization algorithms:
8 Stochastic Gradient Descent, RMSProp, and Adam. We also study alternative
9 models, namely Consensus and Optimistic Optimizations. Rigorous experiments
10 are carried out to evaluate the performance difference between different algorithms,
11 highlighting their advantages and the importance of an appropriate optimization
12 method. This report contributes to the understanding of the optimization in GANs,
13 providing useful insights in the training process of adversarial generative models.

14 1 Introduction

15 Generative Adversarial Networks (GANs) can be conceptualized as addressing a statistical inference
16 problem, where the objective is to approximate the original data distribution using real-world
17 examples. In the image generation problem, this involves producing new images that are statistically
18 similar to a given dataset of real images. GAN is one of the most famous networks in the generation
19 field, known for its incredible ability to learn complex data distributions [1].

20 Despite its promising potentials, the training of GANs have remained challenging due to its unsuper-
21 vised and adversarial nature, which leads to problems like mode collapse and vanishing gradients,
22 making it challenging for the model to achieve convergence [1]. One common alternative, Wasserstein
23 GAN (WGAN), is introduced to improve upon the common problems in GANs by providing a more
24 meaningful learning curve [2].

25 In this paper, we will study three different optimizing algorithms: Stochastic Gradient Descent (SGD),
26 Root Mean Squared Propagation (RMSProp) and Adaptive Moment Estimation (Adam). We also
27 investigate two alternative algorithm for the framework: Consensus and Optimistic, and carry out
28 experiments to compare the performance of different optimization strategies with MNIST [3] and
29 CelabA [4] dataset. We hope this can provide insight on the optimization of GANs.

30 2 Basics of GAN

31 Goodfellow et al. propose the idea of GAN, an unsupervised adversarial generative network frame-
32 work [5]. The generative model and the discriminative model are in a zero-sum competition. The
33 generator tries to create samples that mimic the real data, such as counterfeiters trying to produce
34 fake currency. Meanwhile, the discriminator model seeks to distinguish these fakes samples from
35 genuine data, similar to how the police detect counterfeit bills [5]. The two models keep on competing
36 and improve their generation/discrimination overtime. The network is considered successful if the
37 generator can produce indistinguishable counterfeits.

38 **2.1 Numeric of GAN**

39 The GAN problem can be expressed as a min-max (minimax) problem with the following value
40 function $V(G, D)$:

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

41 The generator $G(z)$ generates data with a random noise input z over the distribution $p_z(z)$; while the
42 discriminator $D(x)$ returns the probability of its input x being real, with $p_{\text{data}}(x)$ representing the
43 real data distribution. The discriminator D tries to maximize the value of V , while the generator G
44 attempts to minimize it. Commonly, gradient descent¹ is performed in attempt to achieve optimality.

45 **3 Optimization in GAN**

46 In this section, we will study various methods of optimizing a GAN. We will first explain the
47 traditional and common algorithms for the optimization process: Stochastic Gradient Descent (SGD),
48 Root Mean Squared Propagation (RMSProp) and Adaptive Momentum Estimation (Adam). Then
49 we will analyze alternative approaches to the GAN problem, namely, Consensus and Optimistic
50 Optimization, and Wasserstein GAN. We also discuss the challenges of a vanilla GAN and compare
51 the different approaches to the standard GAN.

52 **3.1 Traditional Algorithms**

53 **3.1.1 Stochastic Gradient Descent**

54 In the original paper, Goodfellow et al. suggest using mini-batch SGD to be the training algorithm [5].
55 By using a small batch of sample for each update, basic mini-batch SGD increases the convergence
56 speed and introduce noise to improve performance. For parameters θ_d of discriminator D and θ_g
57 of generator G , the algorithm samples mini-batch for x and z , and updates the discriminator by
58 ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D \left(x^{(i)} \right) + \log \left(1 - D \left(G \left(z^{(i)} \right) \right) \right) \right],$$

59 then samples another mini-batch z to update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log \left(1 - D \left(G \left(z^{(i)} \right) \right) \right).$$

60 However, vanilla mini-batch descent does not guarantee a great performance as its gradient may be
61 trapped in saddle points, leading to a vanished gradient [6]. It is also difficult to choose a proper
62 learning rate, and different parameters may wish to have a different learning rate [7]. Different
63 learning rules are required to improve these problem, as we will introduce below.

64 **3.1.2 RMSProp**

65 RMSProp, or Root Mean Squared Propagation, is an algorithm with adaptive learning rate proposed
66 by Geoff Hinton in his lecture notes² [7]. The learning rate is calculated each update, by dividing
67 the default learning rate η with the exponentially decaying root mean square of the gradients g with
68 fraction β plus a smoothing term ϵ , which avoids division by zero (commonly set to 10^{-8}). The
69 formulas are as follows:

$$E[g^2]_t = \beta_1 E[g^2]_{t-1} + (1 - \beta_1) g_t^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

¹Note that descent is only done for generator G ; gradient ascent is performed for the discriminator D as it aims to maximize function V .

²http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf

71 Hinton suggests β to be 0.9, and Ruder suggests η to be 0.001 [7]. This adaptive method means that
 72 the learning rate will be adjusted inversely proportional to the gradient, so the update rate is lower
 73 when the gradient is high and vice versa, leading to a more stable training and improved convergence
 74 speed.

75 **3.1.3 Adam**

76 Before introducing Adaptive Momentum Estimation (Adam), we have to first understand a term
 77 called momentum[8]. The momentum term v_t is calculated by the sum of the current gradient g_t with
 78 a learning rate $-\eta$ plus a fraction β of the past momentum v_{t-1} . Carrying the exponentially decaying
 79 average of the past gradients, it is an analogy to the concept of momentum and potential energy in
 80 real life physics. It can be expressed with the following formula:

$$v_t = \beta v_{t-1} - \eta g_t$$

81 The parameters are then updated with the standard formula $\theta_{t+1} = \theta_t + v_t$. By carrying past gradients,
 82 the sum of vectors with different directions may lead to a smoother and quicker convergence. More
 83 importantly, the optimizer is able to maintain its velocity when approaching local minima and saddle
 84 points, which may allow it to escape these points to seek for the global optimum.

85 Adam [9] combines RMSProp and momentum, earning the benefits of both algorithms, and can be
 86 expressed via the following formulas, where the variables carry meaning similar to the previous cases:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t \quad v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

87 The first momentum m_t is calculated similar to momentum, an exponentially decaying average of past
 88 gradients; and the second momentum v_t is calculated similar to that in RMSProp, an exponentially
 89 decaying average of past squared gradients. The terms are then corrected with the following formulas
 90 to avoid bias towards zero, especially when the decay rates β_1 and β_2 are small (close to 1):

$$\hat{m}_t = \frac{m_t}{1 - \beta_1} \quad \hat{v}_t = \frac{v_t}{1 - \beta_2}$$

91 The parameters are then updated similar to that in RMSProp:

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t + \epsilon}} \hat{m}_t$$

92 The authors propose that $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$ [9]. It is obvious that Adam contains
 93 the advantages of both RMSProp and momentum, improving the common problems of vanilla SGD.

94 **3.2 Consensus Optimization**

95 Mescheder et al. demonstrated through convergence theory that an extremely small learning rate
 96 is necessary when the eigenvalue of the Jacobian of the updated gradient does not meet certain
 97 favorable conditions [10]. Here is the classical convergence he used to delve into future proof. Define
 98 $F : \Omega \rightarrow \Omega$ be a continuously differential function, for example $F(x) = x - \eta \nabla G(x)$ (A typical
 99 gradient descent with G as the lost function). Then, if there exists a \bar{x} satisfying following:

- 100 1. $F(\bar{x}) = \bar{x}$, and
 101 2. the absolute values of the eigenvalues of the Jacobian $F'(\bar{x})$ are all smaller than 1.

102 Then there will have a neighborhood region that U such that any $x \in U$ converge to \bar{x} .

103 Therefore, the eigenvalues of the Jacobian influence the converge result. Under some unfavorable
 104 conditions, including eigenvalues of updated gradient having large complex values and zero-real
 105 parts, extremely small learning rate η is needed to project the eigenvalues into the required unit ball.

106 Under simple GAN setting, the problem tends to converge if³

$$\eta < \frac{1}{|\Re(\lambda)|} \frac{2}{1 + \left(\frac{\Im(\lambda)}{\Re(\lambda)}\right)^2}$$

³Although the definition $F(x) = x - \eta \nabla G(x)$ here is under a different setting compared to Mescheder's one, the negative direction and eigenvalues with positive real part will reach the same result as Mescheder's one.

107 It clearly shows that if $|\Im(\lambda)| >> |\Re(\lambda)|$ then $\frac{\Im(\lambda)}{\Re(\lambda)}$ would be extremely small. Detailed proposition
 108 and proof is in Appendix A. Therefore, Mescheder introduces a regularization term $L(w, \theta)$ to tackle
 109 this issue [10].

Algorithm 1: Consensus Algorithm

```

1  $\gamma \leftarrow$  regularization parameter
2  $\eta \leftarrow$  low learning rate
3  $G, D \leftarrow$  lost function for generator, reward function for discriminator
4 while not converged do
5    $L(w, \theta) = \frac{1}{2} \left\| \begin{pmatrix} \nabla_w G(w, \theta) \\ \nabla_\theta D(w, \theta) \end{pmatrix} \right\|^2$ 
6    $w_{t+1} \leftarrow w_t - \eta * \nabla_w G(w_t, \theta_t) + \eta * \gamma * \nabla_w L(w_t, \theta_t)$ 
7    $\theta_{t+1} \leftarrow \theta_t + \eta * \nabla_\theta D(w_t, \theta_t) - \eta * \gamma * \nabla_\theta L(w_t, \theta_t)$ 

```

110 Referring to Algorithm 1, this regularization term will push the eigenvalues to native real part.
 111 Therefore, it could maximize the maximum η value. The regularization parameter will be choose in
 112 [0.1, 10].

113 **3.3 Wasserstein GAN**

114 Wasserstein GAN (WGAN) [2] is a variation of GAN that makes use of the Wasserstein distance, also
 115 called earth mover's distance. The Wasserstein distance represents the minimum cost of transporting
 116 "earth" or mass, to transform one data distribution into another. The Wasserstein distance of the
 117 real data distribution \mathbb{P}_r and generated data distribution \mathbb{P}_g can be represented as the lower bound
 118 (infimum) of any transportation plan:

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|]$$

119 where $\Pi(\mathbb{P}_r, \mathbb{P}_g)$ represents the set of all joint distributions $\gamma(x, y)$ with marginals \mathbb{P}_r and \mathbb{P}_g .

120 The divergence used in the original GAN is called Jensen-Shanon (JS) divergence, which may not
 121 converge under some cases such as parallel lines [2]. Wasserstein distance however, do converge in
 122 the suggested cases. WGAN employs the Wasserstein distance to provide a continuous measure of
 123 the distance between the real and generated distributions, which facilitates smoother gradients and
 124 training [2]. The critic replaces the discriminator to provide an estimation of the Wasserstein distance
 125 between the real and generated samples. The loss function in Section 2.1 can be redefined as follows:

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

126 The authors also mention that the calculation of Wasserstein distance requires Lipschitz continuity
 127 [2], which ensures that small changes in input leads to controlled changes in the output, leading to
 128 stability to avoid exploding or vanishing gradients. They suggest weight clipping, which restrict the
 129 weights within a predefined range, to enforce this constraint, although admitting that it is a terrible
 130 way, where a parameter too large increases the difficulty to reach optimality, and a parameter too
 131 small can lead to vanishing gradient, but is still employed due to its simplicity and low cost [2].

132 **3.4 Optimistic Optimization**

133 Daskalakis et al. propose an alternative optimization method for solving GANs, or more specifically
 134 WGANs [11]. The proposed method, called Optimistic Mirror Descent (OMD), is a variant of
 135 gradient descent. By Section 3.1.1, we know that gradient descent is one of the method for solving
 136 the min-max optimization problem to train the GAN. To delve into the context of online learning, this
 137 descent algorithm can be adapted into a broader framework known as the no-regret learning algorithm
 138 and the gradient descent would be a special case for Follow-the-Regularized-Leader (FTRL) with a
 139 ℓ_2 regularizer. A detailed proof for the correlations between gradient descent and no-regret learning
 140 is in Appendix B.

Algorithm 2: Optimistic Algorithm

```
1  $\eta \leftarrow$  low learning rate  
2  $G, D \leftarrow$  lost function for generator, reward function for critic  
3 while not converged do  
4    $w_{t+1} \leftarrow w_t - (2\eta * \nabla_w G(w_t, \theta_t) - \eta * \nabla_w G(w_{t-1}, \theta_{t-1}))$   
5    $\theta_{t+1} \leftarrow \theta_t + 2\eta * \nabla_\theta D(w_t, \theta_t) - \eta * \nabla_\theta D((w_{t-1}, \theta_{t-1}))$ 
```

141 To understand the intuition of OMD, we can start from FTRL perspective. Since gradient descent is
142 equivalent with the conjunction of FTRL and ℓ_2 regularizer,

$$w_{t+1} = \operatorname{argmin}_w \left[\left(\sum_{i=1}^t \nabla_w^T w \right) + \frac{1}{2\eta} \|w\|_2^2 \right]$$

143 It is noted that if the summation here include the gradient of next iteration, then the optimization here
144 will lead to a constant term that arises only from the regularization. Therefore, it is done by adding a
145 item called predictor M_{t+1}

$$w_{t+1} = \operatorname{argmin}_w \left[\left(\sum_{i=1}^t \nabla_w^T w \right) + M_{t+1}^T w + \frac{1}{2\eta} \|w\|_2^2 \right]$$

146 There are group of choices for choosing the predictor suggested by Daskalakis et al. and we choose
147 to use the simple predictor, past gradient to test the ability of this algorithm.

148 4 Experiment

149 4.1 Setting

150 All experiments were conducted using Python 3.9 with TensorFlow 2.6 and CUDA 11.2 on a machine
151 equipped with an NVIDIA GeForce GTX 1650. The models were implemented based on the
152 Deep Convolution GAN (DCGAN) architecture, with generator and discriminator using both 4
153 convolution layers, while the datasets used are MNIST [3] and CelabA [4]. For DCGAN, we further
154 experiment two set of models structure (with Batch normalization and without Batch normalization)
155 to see whether these optimization algorithm could still perform in a poor condition. We repeat the
156 experiment five times to ensure there is a fair evaluation⁴.

157 4.2 Results

158 We did three major sets of comparison between the optimization algorithms, including Vanilla GAN
159 Optimization *versus* Wasserstein GAN Optimization, Vanilla GAN Optimization *versus* Consensus
160 Optimization, Wasserstein GAN Optimization *versus* Optimistic Optimization. Both Vanilla GAN
161 and Consensus Optimization will use the lost function stated in Section 2.1. While, Wasserstein GAN
162 and Optimistic Optimization will use the the lost function stated in Section 3.3.

163 4.2.1 Vanilla vs Wasserstein

164 In comparison between Vanilla and Wasserstein in DCGAN with Batch Normalization (BN), we has
165 done in total 12 types of experiment, including three types of optimizer explained in Section 3.1 and
166 two learning rate chosen (1e-4, 1e-5)⁵.

167 A complete sheet of generated samples is provided in Appendix C. It is important to emphasize that
168 the samples were not selectively chosen or cherry-picked. The FID score evaluates test performance,
169 with lower scores indicating better results. Based on Figure 1 (and Appendix C.1, Appendix C.2),
170 both GAN and WGAN can be trained to generate images using RMSProp across various learning

⁴Github Link:https://github.com/Mysterchan/estr3114_project/tree/main

⁵Both Vanilla and Wasserstein, we have tired all three optimizer under two learning rate. So, In total
2 × 3 × 2 = 12 types of experiment is conducted

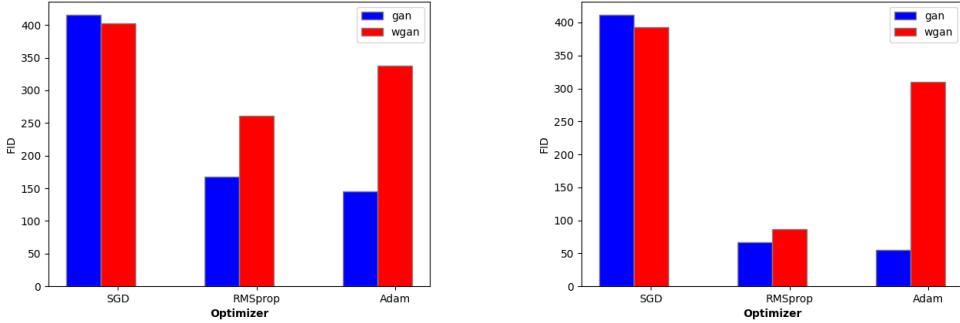
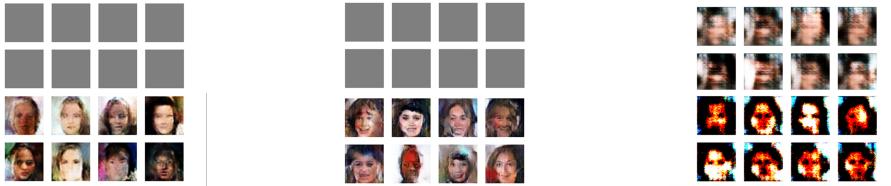


Figure 1: The Comparison between Optimizer, Model and Learning rate on CelebA

171 rates. However, both models fail to generate plausible results when trained with SGD. While GANs
 172 can be successfully trained using Adam, WGANs are unable to replicate this result. This observation
 173 aligns with findings reported by Arjovsky et al. [2] regarding challenges in training WGANs.

174 In the experiments conducted without batch normalization (BN), we set the learning rate to 1e-5. As
 175 shown in Figure 2(a), we observed that WGANs are capable of generating images using RMSProp in
 176 the absence of batch normalization.



(a) GAN (Upper 2) and WGAN (Lower 2) when without BN (b) GAN (Upper 2) and Consensus (Lower 2) when without BN (c) WGAN (Upper 2) and Optimistic (Lower 2) under Adam

Figure 2: The Comparison of different algorithm

177 4.2.2 Vanilla vs Consensus

178 In the "no batch normalization" setting with a learning rate of 1e-5, we compared the performance
 179 of vanilla GAN and consensus GAN under three different optimizers: SGD, Adam, and RMSProp.
 180 Both models failed to generate meaningful images when trained with SGD and Adam. However,
 181 consensus optimization successfully enabled image generation when using RMSProp.

182 4.2.3 Wasserstein vs Optimistic

183 We followed the learning rate of 1e-4 as suggested by Daskalakis et al. [11], with batch normalization
 184 enabled. As stated, WGANs produced noisy and blurred images, while we did not replicate the
 185 exact results in optimistic GAN reported in the original paper (likely due to the difference in model
 186 size, which could not be resolved due to memory limitations), our findings nonetheless suggest that
 187 optimistic GANs are capable of generating higher-quality images compared to WGANs.

188 5 Conclusion

189 In this project, we explore several optimization algorithms like Consensus and Optimistic for GAN,
 190 and gain a broader understanding of how different variants of gradient descent, such as SGD,
 191 RMSProp and Adam, contribute to solving the min-max optimization problem in GANs. Through
 192 the completion of the experiments, we also discover how different algorithms perform under various
 193 scenarios, providing insights regarding the optimization of GANs.

194 **References**

- 195 [1] D. Saxena and J. Cao, “Generative adversarial networks (gans) challenges, solutions, and future
196 directions,” *ACM Computing Surveys (CSUR)*, vol. 54, no. 3, pp. 1–42, 2021.
- 197 [2] M. Arjovsky, S. Chintala, and L. Bottou, “Wasserstein generative adversarial networks,” in
198 *International conference on machine learning*, pp. 214–223, PMLR, 2017.
- 199 [3] L. Deng, “The mnist database of handwritten digit images for machine learning research,” *IEEE
200 Signal Processing Magazine*, vol. 29, no. 6, pp. 141–142, 2012.
- 201 [4] Z. Liu, P. Luo, X. Wang, and X. Tang, “Large-scale celebfaces attributes (celeba) dataset,”
202 *Retrieved August*, vol. 15, no. 2018, p. 11, 2018.
- 203 [5] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and
204 Y. Bengio, “Generative adversarial networks,” *Communications of the ACM*, vol. 63, no. 11,
205 pp. 139–144, 2020.
- 206 [6] Y. N. Dauphin, R. Pascanu, C. Gulcehre, K. Cho, S. Ganguli, and Y. Bengio, “Identifying and
207 attacking the saddle point problem in high-dimensional non-convex optimization,” *Advances in
208 neural information processing systems*, vol. 27, 2014.
- 209 [7] S. Ruder, “An overview of gradient descent optimization algorithms,” *arXiv preprint arXiv:1609.04747*, 2016.
- 211 [8] N. Qian, “On the momentum term in gradient descent learning algorithms,” *Neural networks*,
212 vol. 12, no. 1, pp. 145–151, 1999.
- 213 [9] D. P. Kingma, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*,
214 2014.
- 215 [10] L. Mescheder, S. Nowozin, and A. Geiger, “The numerics of gans,” *Advances in neural
216 information processing systems*, vol. 30, 2017.
- 217 [11] C. Daskalakis, A. Ilyas, V. Syrgkanis, and H. Zeng, “Training gans with optimism,” *arXiv
218 preprint arXiv:1711.00141*, 2017.

219 **A Propositions and proof related to Consensus Optimization**

220 Noted that the proof here is for gradient descent while [10] gives for gradient ascent only.

221 The value of the Jacobian $F'(x)$ would be equal to $I - \eta A$ where A is the jacobian of the gradient.

222 **Lemma 1.** *If λ is one of the eigenvalues of the A , $1 - \eta\lambda$ would be the a eigenvalue of the Jacobian
223 $F'(x)$*

224 *Proof.* Suppose $A\vec{v} = \lambda\vec{v}$ where \vec{v} is one of the corresponding eigenvectors. Then we have,

$$F'(x)\vec{v} = (I - \eta A)\vec{v} = \vec{v} - \eta A\vec{v} = \vec{v} - \eta\lambda\vec{v} = (1 - \eta\lambda)\vec{v}$$

225 Therefore, $(1 - \eta\lambda)$ would be the eigenvalue for matrix $F'(x)$. \square

226 **Lemma 2.** *Assume that $A \in \mathbb{R}^{n \times n}$ only has eigenvalues with positive real-part and let $\eta > 0$. Then
227 the eigenvalues of Jacobian $F'(x) = I - \eta A$ lie in the unit ball if and only if*

228 *Proof.* For $\lambda = a + bi$ with $a > 0$ we have eigenvalue $(1 - \eta\lambda)$ by lemma 1.

$$|1 - \eta\lambda|^2 = (1 - \eta a)^2 + \eta^2 b^2 = 1 + \eta^2 b^2 + \eta^2 a^2 - 2\eta a.$$

229 Then, we have,

$$\eta(\eta b^2 + \eta a^2 - 2a) < 0$$

230 Since $\eta > 0$, the absolute value of the eigenvalue is smaller than 1 if and only if

$$\eta < \frac{2a}{b^2 + a^2} = \frac{2a^{-1}}{(b/a)^2 + 1}$$

²³¹ By putting back $a = \Re(\lambda)$ and $b = \Im(\lambda)$, we have

$$\eta < \frac{1}{|\Re(\lambda)|} \frac{2}{1 + \left(\frac{\Im(\lambda)}{\Re(\lambda)}\right)^2}$$

²³²

□

²³³ B No-regret learning

²³⁴ The section aims to show Follow-the-Regularized-Leader (FTRL) with a ℓ_2 regularizer is equivalent
²³⁵ to gradient descent.

²³⁶ For the FTRL framework with ℓ_2 regularizer at iteration t , we have the following,

$$x_t = \operatorname{argmin}_x \left[\left(\sum_{i=1}^{t-1} \ell_i \right) + \frac{1}{2\eta} \|x\|_2^2 \right] \quad \text{where } \ell_i \text{ is the lost function at time } i$$

²³⁷ Consider the first-order approximation of ℓ_i , we have $\ell_i \approx \ell_i(x_i) + \nabla \ell_i(x_i)^T x$. Then the cumulative
²³⁸ loss term becomes

$$\sum_{i=1}^{t-1} \ell_i \approx \left(\sum_{i=1}^{t-1} \ell_i(x_i) \right) + \mathbf{g}_{t-1}^T x \quad \text{where } \mathbf{g}_{t-1} = \sum_{i=1}^{t-1} \nabla \ell_i(x_i)$$

²³⁹ Therefore, the reformed x_t would be $\arg \min_x [(\sum_{i=1}^{t-1} \ell_i(x_i)) + \mathbf{g}_{t-1}^T x + \frac{1}{2\eta} \|x\|_2^2]$. Because we hope
²⁴⁰ to find the minimum point of this item, we would apply differentiation on it.

$$\nabla_x \left(\left(\sum_{i=1}^{t-1} \ell_i(x_i) \right) + \mathbf{g}_{t-1}^T x + \frac{1}{2\eta} \|x\|_2^2 \right) = \mathbf{g}_{t-1} + \frac{1}{\eta} x = 0$$

$$x_t = -\eta \mathbf{g}_{t-1} = x_{t-1} - \eta \nabla \ell_{t-1}(x_{t-1})$$

²⁴¹ Hence, we get the gradient descent from FTRL with ℓ_2 regularizer.

242 **C Sheets of samples**

243 **C.1 CelebA**



Figure 3: Samples from GAN using Adam Optimizer and Under 1e-5 rate in CelabA

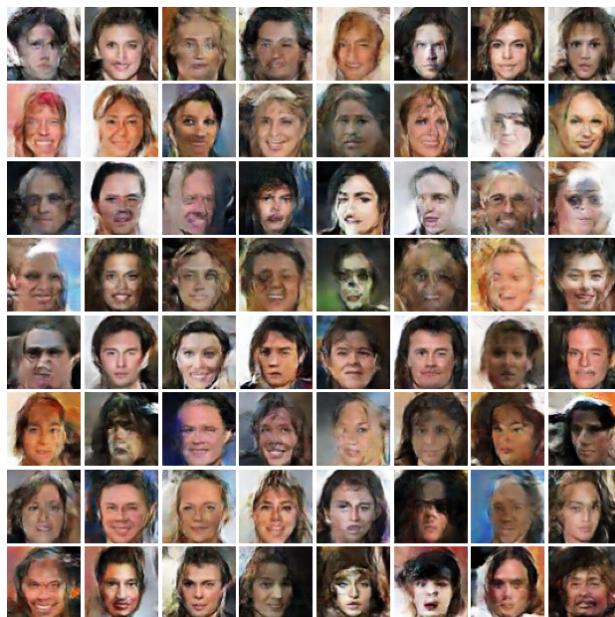


Figure 4: Samples from GAN using RMSProp Optimizer and Under 1e-5 rate in CelabA



Figure 5: Samples from WGAN using Adam Optimizer and Under 1e-5 rate in CelabA

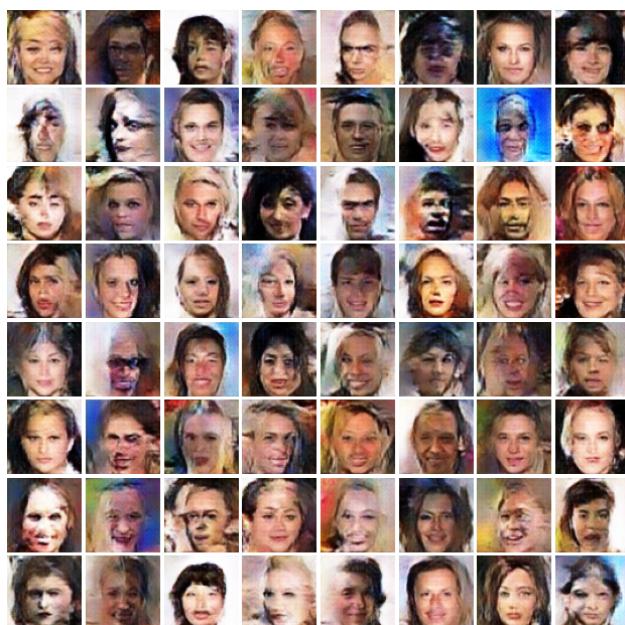


Figure 6: Samples from WGAN using RMSProp Optimizer and Under 1e-5 rate in CelabA



Figure 7: Samples from GAN using Adam Optimizer and Under 1e-5 rate in MNIST

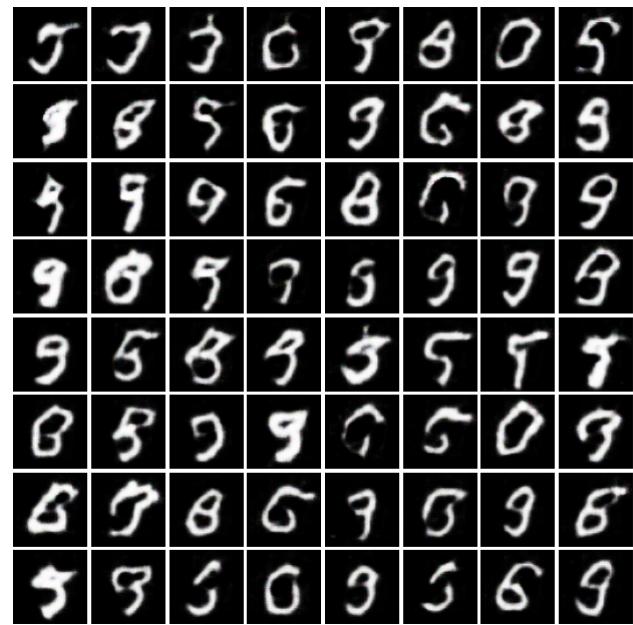


Figure 8: Samples from GAN using RMSProp Optimizer and Under 1e-5 rate in MNIST

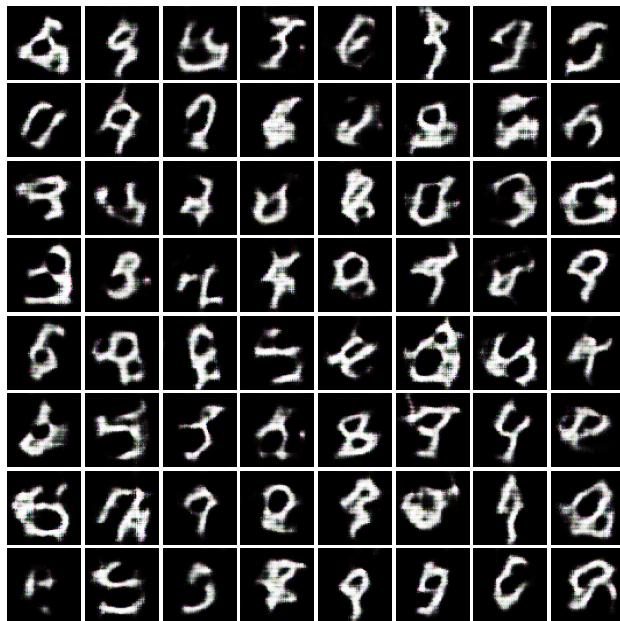


Figure 9: Samples from WGAN using Adam Optimizer and Under 1e-5 rate in MNIST

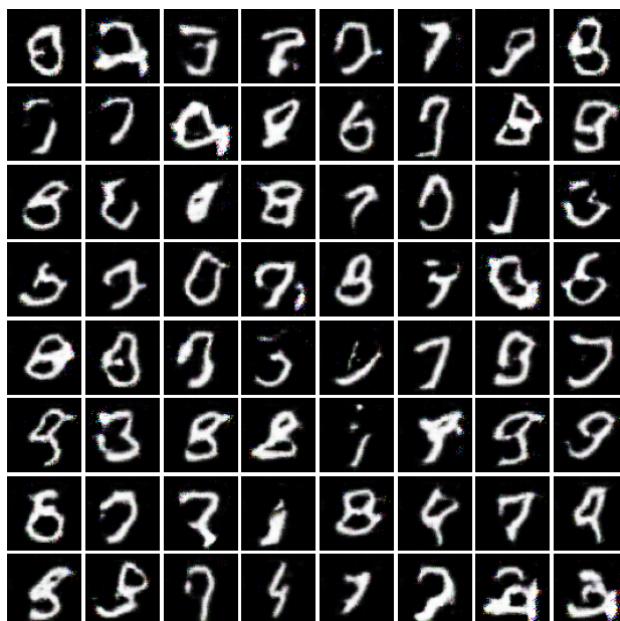


Figure 10: Samples from WGAN using RMSProp Optimizer and Under 1e-5 rate in MNIST