
Understanding The Role of Adversarial Regularization in Supervised Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Despite numerous attempts sought to provide empirical evidence of adversarial
2 regularization outperforming sole supervision in various inverse problems, the
3 theoretical understanding of such phenomena remains elusive. In this study, we
4 aim to resolve whether adversarial regularization indeed performs better than sole
5 supervision at a fundamental level. To bring this insight into fruition, we study
6 vanishing gradient issue, asymptotic iteration complexity and gradient flow in the
7 context of sole supervision and adversarial regularization. The key ingredient is a
8 theoretical justification supported by empirical evidence of adversarial acceleration
9 in gradient descent. In addition, motivated by a recently introduced unit-wise
10 capacity based generalization bound, we analyze the generalization error in adversarial
11 framework. Guided by our observation, we cast doubts on the ability of this
12 measure to explain generalization. We therefore leave as open questions to explore
13 new measures that can explain generalization behavior in adversarial learning.

14 1 Introduction

15 At a fundamental level, we study the role of adversarial regularization in supervised learning through
16 the lens of theoretical justification. We intend to resolve the mystery of why supervised learning with
17 adversarial regularization accelerates gradient updates as compared to sole supervision. In light of
18 deeper understanding, we explore several crucial properties pertaining to adversarial acceleration in
19 gradient descent.

20 In recent years, the research community has witnessed pervasive use of Generative Adversarial
21 Networks (GANs) on a wide variety of complex tasks [1, 2, 3, 4]. Among many applications some
22 require generation of a particular sample subject to a conditional input. For this reason, there has
23 been a surge in designing conditional adversarial networks. In visual object tracking via adversarial
24 learning, Euclidean norm is used to regulate the generation process so that the generated mask
25 falls within a small neighborhood of actual mask [5]. In photo-realistic image super resolution,
26 Euclidean or supremum norm is used to minimize the distance between reconstructed and original
27 high resolution image [6, 7]. In medical image segmentation, multi-scale L_1 -loss with adversarial
28 regularization is shown to outperform sole supervision [8].

29 The authors of [1] use L_1 -loss as a supervision signal and adversarial regularization as a continuously
30 evolving loss function. Because GANs learn a loss that adapts to data, they fairly solve multitude of
31 tasks which would otherwise require hand-engineered loss. The authors of [9] use adversarial loss
32 on top of pixel, style, and feature loss to restrict the generated images on a manifold of real data.
33 Prior works on this operate under the synonym conditional GAN where a convex composition of
34 pixel and adversarial loss is primarily optimized [10, 11, 12]. Karacan et al. [13] use this technique
35 to efficiently generate images of outdoor scenes. The authors of [14] combined spatial and Laplacian
36 spectral channel attention in regularized adversarial learning to synthesize high resolution images.

37 Furthermore, Emami et al. [15] coalesce spatial attention with adversarial regularization and feature
38 map loss to achieve state-of-the-art image to image translation.

39 Additionally, the spectral and spatial super resolution based on adversarial regularization [16, 17]
40 is proven to achieve faster convergence and better empirical risk compared to purely supervised
41 learning [18]. Further, the authors of [6] showed improvement in perceptual quality of high resolution
42 images in adversarial setting [19]. Despite superior empirical performance of adversarial regulariza-
43 tion in diverse domains, the theoretical understanding of such phenomena remains elusive. So far the
44 theoretical analysis suggests that there is a constant that bounds the total empirical risk above [8].
45 As a consequence, this inhibits erroneous gradient estimation by the discriminator which apparently
46 improves perceptual quality. However, these benign properties of loss surface do not fully explain
47 this phenomenon at a fundamental level.

48 As per these prior works [17, 8, 16, 19, 20, 21], it is understandable that supervised learning with
49 adversarial regularization boosts empirical performance. In addition, this improvement is consistent
50 across a wide variety of inverse problems and network configurations. As much beneficial as this
51 regularization has been so far, to our knowledge, the theoretical understanding still remains relatively
52 less explored. Aiming to bridge this gap, we provide both theoretical and empirical evidence of faster
53 convergence due to adversarial regularization.

54 2 Preliminaries

55 **Notations.** Let $X \subset \mathbb{R}^{d_x}$ and $Y \subset \mathbb{R}^{d_y}$ where d_x and d_y denote input and output dimensions,
56 respectively. The empirical distributions of X and Y are denoted by \mathcal{P}_X and \mathcal{P}_Y . Given an input
57 $x \in X$, $f(\theta; x) : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}$ represents a common neural network architecture with rectified linear
58 unit (ReLU) activation for both supervised and adversarial learning. Here, θ denotes the trainable
59 parameters of the generator $f(\theta; \cdot)$. The discriminator, $g(\psi; \cdot)$ has trainable parameters collected by
60 ψ . For $g : \mathbb{R}^{d_y} \rightarrow \mathbb{R}$, ∇g denotes its gradient and $\nabla^2 g$ denotes its Hessian. Given a vector x , $\|x\|$
61 represents the Euclidean norm. Given a matrix M , $\|M\|$ and $\|M\|_F$ represent spectral and Frobenius
62 norm, respectively.

63 **Definition 1.** (L -Lipschitz) A function f is L -Lipschitz if $\forall \theta$, $\|\nabla f(\theta)\| \leq L$.

64 **Definition 2.** (β -Smoothness) A function f is β -smooth if $\forall \theta$, $\|\nabla^2 f(\theta)\| \leq \beta$

65 **Problem Setup.** In Wasserstein GAN (WGAN) + Gradient Penalty (GP), the generator cost function
66 is given by

$$\arg \min_{\theta} -\mathbb{E}_{x \sim \mathcal{P}_X} [g(\psi; f(\theta; x))] \quad (1)$$

67 and the discriminator cost function,

$$\arg \min_{\psi} \mathbb{E}_{x \sim \mathcal{P}_X} [g(\psi; f(\theta; x))] - \mathbb{E}_{y \sim \mathcal{P}_Y} [g(\psi; y)] + \lambda_{GP} \mathbb{E}_{z \sim \mathcal{P}_Z} [(\|\nabla_z g(\psi; z)\| - 1)^2]. \quad (2)$$

68 Here, \mathcal{P}_Z represents the distribution over samples along the line joining samples from real and
69 generator distribution. Unlike sole supervision, the mapping function $f_{\theta}(\cdot)$ in augmented objective
70 has access to a feedback signal from the discriminator. Thus, the optimization in supervised learning
71 with adversarial regularization is carried out by

$$\arg \min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y) - g(\psi; f(\theta; x))]. \quad (3)$$

72 The discriminator cost function remains identical to Wasserstein discriminator as given by equation (2).
73 Here, \mathcal{P} denotes the joint empirical distribution over X and Y .

74 3 Theoretical Analysis

75 This section clearly states the assumptions and justifies their fidelity in the context of adversarial
76 regularization. The theoretical findings are intended to provide a reasonable justification to multitude
77 of tasks that owe the benefits to adversarial training. The technical overview begins with exploiting
78 vanishing gradient issue in the near optimal region. It then presents the main results by estimating the
79 iteration complexity and sub-optimality gap.

80 3.1 Warm-Up: Mitigating Vanishing Gradient in Near Optimal Region

81 **Assumption 1.** *The mapping function $f(\theta; x)$ is L -Lipschitz in θ .*

82 **Assumption 2.** *The loss function $l(p; y)$ where $p = f(\theta; x)$ is β -smooth in p .*

83 **Assumption 1** is a mild requirement that is easily satisfied in near optimal region. Different from
84 standard smoothness in optimization, it is trivial to justify **Assumption 2** by relating it to a quadratic
85 loss function as followed by most in practice.

86 **Lemma 1.** *Suppose **Assumption 1** and **Assumption 2** hold. If $\|\theta - \theta^*\| \leq \epsilon$, then*
87 $\|\nabla_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y)]\| \leq L^2 \beta \epsilon$.

88 *Proof.* Refer to Appendix D.1. **Lemma 1** provides an upper bound on the expected gradient over
89 empirical distribution \mathcal{P} in near optimal region. As the intermediate iterates (θ) move closer to the
90 optima (θ^*), i.e. $\epsilon \rightarrow 0$, the gradient norm vanishes in expectation. This essentially resonates with the
91 intuitive understanding of gradient descent. From another perspective, the issue of gradient descent
92 inherently resides in near optimal region. We therefore ask a fundamental question: can we attain
93 faster convergence without loosing any empirical risk benefits? The following sections are intended
94 to shed some light in this direction.

95 **Lemma 2.** *Suppose **Assumption 1** holds. For a differentiable discriminator $g(\psi; y)$, if $\|g - g^*\| \leq \delta$,*
96 *where $g^* \triangleq g(\psi^*)$ denote optimal discriminator, then $\|-\nabla_{\theta} \mathbb{E}_{x \sim \mathcal{P}_X} [g(\psi; f(\theta; x))]\| \leq L\delta$.*

97 *Proof.* Refer to Appendix D.2. **Lemma 2** indicates that the expected gradient of purely adversarial
98 generator does not produce erroneous gradients in the near optimal region, suggesting well behaved
99 composite empirical risk [8].

100 **Theorem 1.** *Let us suppose **Assumption 1** and **Assumption 2** hold. If $\|\theta - \theta^*\| \leq \epsilon$ and $\|g - g^*\| \leq$
101 δ , then $\|\nabla_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y) - g(\psi; f(\theta; x))]\| \leq (L^2 \beta \epsilon + L\delta)$.*

102 *Proof.* Refer to Appendix D.3. To focus more on the empirical success of adversarial regular-
103 ization, we study simple convex-concave minimax optimization. It will certainly be interest-
104 ing to borrow some ideas from the vast minimax optimization literature in various other set-
105 tings [22, 23, 24, 25]. According to **Theorem 1**, the expected gradient of augmented objec-
106 tive does not vanish in the near optimal region, i.e. $\|\Delta\theta\| \rightarrow L\delta$ as $\epsilon \rightarrow 0$. In the current
107 setting, the estimated gradients of $l(\theta)$ and $-g(\theta)$ at any instant during the optimization pro-
108 cess are positively correlated. Thus, the gradients of augmented objective is lower bounded by
109 $\|\nabla_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y) - g(\psi; f(\theta; x))]\| \geq \|\nabla_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y)]\|$. The upper and
110 lower bounds of the intermediate iterates justify non-vanishing gradient in near optimal region.
111 Having proven the contribution of discriminator in mitigating vanishing gradient, it seems natural to
112 wonder whether adversarial regularization improves the iteration complexity, which we discuss in the
113 following Section 3.2.

114 3.2 Main Results: Asymptotic Iteration Complexity

115 In this section, we analyze global iteration complexity of both sole supervision and adversarial
116 regularization [26, 27]. We restrict our analysis to a deterministic setting. For a deterministic
117 sequence of parameters $\{\theta_k\}_{k \in \mathbb{N}}$, the complexity of $\{\theta_k\}_{k \in \mathbb{N}}$ for a function $l(\theta)$ is defined as

$$\mathcal{T}_{\epsilon}(\{\theta_k\}_{k \in \mathbb{N}}, l) := \inf \{k \in \mathbb{N} \mid \|\nabla l(\theta_k)\| \leq \epsilon\}.$$

118 For a given initialization θ_0 , risk function l and algorithm A_{ϕ} , where ϕ denotes hyperparameters
119 of training algorithm, such as learning rate and momentum coefficient, $A_{\phi}[l, \theta_0]$ denotes the se-
120 quence of iterates generated during training. We compute iteration complexity of an algorithm class
121 parameterized by p hyperparameters, $\mathcal{A} = \{A_{\phi}\}_{\phi \in \mathbb{R}^p}$ on a function class, \mathcal{L} as

$$\mathcal{N}(\mathcal{A}, \mathcal{L}, \epsilon) := \inf_{A_{\phi} \in \mathcal{A}} \sup_{\theta_0 \in \{\mathbb{R}^h \times d_x, \mathbb{R}^{d_y \times h}\}, l \in \mathcal{L}} \mathcal{T}_{\epsilon}(A_{\phi}[l, \theta_0], l).$$

122 We derive the asymptotic bounds under a less restrictive setting as introduced in [26]. The new
123 condition is weaker than commonly used Lipschitz smoothness assumption. Under this condition,
124 the authors of [26] aim to resolve the mystery of why adaptive gradient methods converge faster.
125 We use this theoretical tool to study the asymptotic convergence of sole supervision and adversarial

126 regularization in near optimal region. To circumvent the tractability issues in non-convex optimization,
 127 we follow the common practice of seeking an ϵ -stationary point, i.e. $\|\nabla l(\theta)\| < \epsilon$. We start by
 128 analyzing the iteration complexity of gradient descent with fixed step size. In this regard, we build on
 129 the assumptions made in [26]. To put more succinctly, let us recall the assumptions.

130 **Assumption 3.** The loss l is lower bounded by $l^* > -\infty$.

131 **Assumption 4.** The function is twice differentiable.

132 **Assumption 5.** ((L_0, L_1) -Smoothness). The function is (L_0, L_1) -smooth, i.e. there exist positive
 133 constants L_0 and L_1 such that $\|\nabla^2 l(\theta)\| \leq L_0 + L_1 \|\nabla l(\theta)\|$.

134 **Theorem 2.** Suppose the functions in \mathcal{L} satisfy **Assumption 3, 4 and 5**. Given $\epsilon > 0$, the iteration
 135 complexity in sole supervision is upper bounded by $\mathcal{O}\left(\frac{(l(\theta_0)-l^*)(L_0+L_1L^2\beta\epsilon)}{\epsilon^2}\right)$.

136 *Proof.* Refer to Appendix D.4.

137 **Corollary 1.** Using first order Taylor series, the upper bound in **Theorem 2** becomes $\mathcal{O}\left(\frac{l(\theta_0)-l^*}{h\epsilon^2}\right)$.

138 *Proof.* Refer to Appendix D.5.

139 **Assumption 6.** (Existence of useful gradients) For arbitrarily small $\zeta > 0$, the norm of the gradients
 140 provided by discriminator is lower bounded by ζ , i.e. $\|\nabla g(\psi; f(\theta; x))\| \geq \zeta$.

141 **Assumption 6** requires discriminator to provide useful gradients until convergence. This is a valid
 142 assumption in convex-concave minimax optimization problems. It is trivial to prove this in the inner
 143 maximization loop under concave setting. In other words, the stated assumptions are mild and derived
 144 from prior analysis for the sole pupose of mathematical simplicity. Keeping this in mind, we analyze
 145 the global iteration complexity in adversarial setting.

146 **Theorem 3.** Suppose the functions in \mathcal{L} satisfy **Assumption 3, 4 and 5**. Given **Assumption 6** holds,
 147 $\epsilon > 0$ and $\delta \leq \frac{\sqrt{2\epsilon\zeta}}{L}$, the iteration complexity in adversarial regularization is upper bounded by
 148 $\mathcal{O}\left(\frac{(l(\theta_0)-l^*)(L_0+L_1L^2\beta\epsilon)}{\epsilon^2+2\epsilon\zeta-L^2\delta^2}\right)$.

149 *Proof.* Refer to Appendix D.6.

150 **Corollary 2.** Using first order Taylor series, the upper bound in **Theorem 3** becomes $\mathcal{O}\left(\frac{l(\theta_0)-l^*}{h\epsilon^2+h\zeta\epsilon}\right)$.

151 *Proof.* Refer to Appendix D.7. Since $2\epsilon\zeta - L^2\delta^2 \geq 0$, the supervised learning with adversarial
 152 regularization has a *tighter* global iteration complexity compared to sole supervision. In a simplified
 153 setup, one can easily verify this hypothesis by following the proof using first order Taylor's approx-
 154 imation as given by **Corollary 1 and 2**. In this case, $h\zeta\epsilon > 0$ ensures *tighter* iteration complexity
 155 bound. This result is significant because it improves the convergence rates from $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$ to $\mathcal{O}\left(\frac{1}{\epsilon^2+\epsilon\zeta}\right)$.

156 Notice that for a too strong discriminator, **Assumption 6** does not hold. For a too weak discriminator,
 157 $\|g - g^*\| \leq \delta$ does not hold when δ is arbitrarily small. In these cases, the generator does not receive
 158 useful gradients from the discriminator to undergo accelerated training. However, for a sufficiently
 159 trained discriminator, i.e. $\|g - g^*\| \leq \delta \leq \frac{\sqrt{2\epsilon\zeta}}{L}$, the adversarial regularization accelerates gradient
 160 updates. Notably, both the empirical risk and iteration complexity benefit from this provided the
 161 discriminator and generator are trained alternatively as typically followed in practice¹.

162 4 Discussion

163 In this study, we investigated the reason behind slow convergence of purely supervised learning in
 164 near optimal region, and how adversarial regularization circumvents this issue. Further, we explored
 165 several crucial properties at this juncture of understanding the role of adversarial regularization
 166 in supervised learning. Particularly intriguing was the genericness of these theorems around the
 167 central theme. To make a fair assessment, standard theoretic tools were employed in all the theorems.
 168 In theoretical analysis, the asymptotic iteration complexity, gradient flow, provable convergence
 169 guarantee and the analysis of generalization error provided further insights to the empirical findings
 170 of adversarial regularization as reported in copious literature.

¹Refer to Appendix for further analysis and experimental results.

References

- [1] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, “Image-to-image translation with conditional adversarial networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1125–1134, 2017.
- [2] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, “Unpaired image-to-image translation using cycle-consistent adversarial networks,” in *Proceedings of the IEEE international conference on computer vision*, pp. 2223–2232, 2017.
- [3] T. Park, M.-Y. Liu, T.-C. Wang, and J.-Y. Zhu, “Semantic image synthesis with spatially-adaptive normalization,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2337–2346, 2019.
- [4] T. Karras, S. Laine, and T. Aila, “A style-based generator architecture for generative adversarial networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4401–4410, 2019.
- [5] Y. Song, C. Ma, X. Wu, L. Gong, L. Bao, W. Zuo, C. Shen, R. W. Lau, and M.-H. Yang, “Vital: Visual tracking via adversarial learning,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8990–8999, 2018.
- [6] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, *et al.*, “Photo-realistic single image super-resolution using a generative adversarial network,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4681–4690, 2017.
- [7] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, and C. Change Loy, “Esrgan: Enhanced super-resolution generative adversarial networks,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 0–0, 2018.
- [8] Y. Xue, T. Xu, H. Zhang, L. R. Long, and X. Huang, “Segan: Adversarial network with multi-scale l-1 loss for medical image segmentation,” *Neuroinformatics*, vol. 16, no. 3-4, pp. 383–392, 2018.
- [9] W. Xian, P. Sangkloy, V. Agrawal, A. Raj, J. Lu, C. Fang, F. Yu, and J. Hays, “Texturegan: Controlling deep image synthesis with texture patches,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8456–8465, 2018.
- [10] M. Mirza and S. Osindero, “Conditional generative adversarial nets,” *arXiv preprint arXiv:1411.1784*, 2014.
- [11] E. L. Denton, S. Chintala, R. Fergus, *et al.*, “Deep generative image models using laplacian pyramid of adversarial networks,” in *Advances in neural information processing systems*, pp. 1486–1494, 2015.
- [12] X. Wang and A. Gupta, “Generative image modeling using style and structure adversarial networks,” in *European Conference on Computer Vision*, pp. 318–335, Springer, 2016.
- [13] L. Karacan, Z. Akata, A. Erdem, and E. Erdem, “Learning to generate images of outdoor scenes from attributes and semantic layouts,” *arXiv preprint arXiv:1612.00215*, 2016.
- [14] L. Rout, I. Misra, S. M. Moorthi, and D. Dhar, “S2a: Wasserstein gan with spatio-spectral laplacian attention for multi-spectral band synthesis,” in *Proceedings of the IEEE conference on computer vision and pattern recognition workshop*, 2020.
- [15] H. Emami, M. M. Aliabadi, M. Dong, and R. B. Chinnam, “Spa-gan: Spatial attention gan for image-to-image translation,” *arXiv preprint arXiv:1908.06616*, 2019.
- [16] L. Rout, “Alert: Adversarial learning with expert regularization using tikhonov operator for missing band reconstruction,” *IEEE Transactions on Geoscience and Remote Sensing*, 2020.
- [17] A. Rangnekar, N. Mokashi, E. Ientilucci, C. Kanan, and M. Hoffman, “Aerial spectral super-resolution using conditional adversarial networks,” *arXiv preprint arXiv:1712.08690*, 2017.
- [18] C. Lanaras, J. Bioucas-Dias, S. Galliani, E. Baltsavias, and K. Schindler, “Super-resolution of sentinel-2 images: Learning a globally applicable deep neural network,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 146, pp. 305–319, 2018.
- [19] C. Dong, C. C. Loy, K. He, and X. Tang, “Image super-resolution using deep convolutional networks,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 2, pp. 295–307, 2015.

- 224 [20] M. Henaff, A. Canziani, and Y. LeCun, “Model-predictive policy learning with uncertainty
225 regularization for driving in dense traffic,” *arXiv preprint arXiv:1901.02705*, 2019.
- 226 [21] M. Sarmad, H. J. Lee, and Y. M. Kim, “Rl-gan-net: A reinforcement learning agent controlled
227 gan network for real-time point cloud shape completion,” in *Proceedings of the IEEE Conference*
228 *on Computer Vision and Pattern Recognition*, pp. 5898–5907, 2019.
- 229 [22] T. Lin, C. Jin, and M. I. Jordan, “On gradient descent ascent for nonconvex-concave minimax
230 problems,” *arXiv preprint arXiv:1906.00331*, 2019.
- 231 [23] T. Lin, C. Jin, M. Jordan, *et al.*, “Near-optimal algorithms for minimax optimization,” *arXiv*
232 *preprint arXiv:2002.02417*, 2020.
- 233 [24] C. Jin, P. Netrapalli, and M. I. Jordan, “What is local optimality in nonconvex-nonconcave
234 minimax optimization?,” *arXiv preprint arXiv:1902.00618*, 2019.
- 235 [25] P. Mertikopoulos, C. Papadimitriou, and G. Piliouras, “Cycles in adversarial regularized learn-
236 ing,” in *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*,
237 pp. 2703–2717, SIAM, 2018.
- 238 [26] J. Zhang, T. He, S. Sra, and A. Jadbabaie, “Why gradient clipping accelerates training: A
239 theoretical justification for adaptivity,” in *International Conference on Learning Representations*,
240 2019.
- 241 [27] Y. Carmon, J. C. Duchi, O. Hinder, and A. Sidford, “Lower bounds for finding stationary points
242 i,” *Mathematical Programming*, pp. 1–50, 2019.
- 243 [28] A. Kudo, Y. Kitamura, Y. Li, S. Iizuka, and E. Simo-Serra, “Virtual thin slice: 3d conditional
244 gan-based super-resolution for ct slice interval,” in *International Workshop on Machine Learning*
245 *for Medical Image Reconstruction*, pp. 91–100, Springer, 2019.
- 246 [29] M. Schmidt, N. Le Roux, and F. Bach, “Minimizing finite sums with the stochastic average
247 gradient,” *Mathematical Programming*, vol. 162, no. 1-2, pp. 83–112, 2017.
- 248 [30] D. Zhou, P. Xu, and Q. Gu, “Stochastic nested variance reduction for nonconvex optimization,”
249 in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*,
250 pp. 3925–3936, Curran Associates Inc., 2018.
- 251 [31] Y. Nesterov, “Efficiency of coordinate descent methods on huge-scale optimization problems,”
252 *SIAM Journal on Optimization*, vol. 22, no. 2, pp. 341–362, 2012.
- 253 [32] Y. Carmon, J. C. Duchi, O. Hinder, and A. Sidford, “Accelerated methods for nonconvex
254 optimization,” *SIAM Journal on Optimization*, vol. 28, no. 2, pp. 1751–1772, 2018.
- 255 [33] J. Duchi, E. Hazan, and Y. Singer, “Adaptive subgradient methods for online learning and
256 stochastic optimization,” *Journal of Machine Learning Research*, vol. 12, no. Jul, pp. 2121–
257 2159, 2011.
- 258 [34] M. Staib, S. J. Reddi, S. Kale, S. Kumar, and S. Sra, “Escaping saddle points with adaptive
259 gradient methods,” *arXiv preprint arXiv:1901.09149*, 2019.
- 260 [35] D. Zhou, Y. Tang, Z. Yang, Y. Cao, and Q. Gu, “On the convergence of adaptive gradient
261 methods for nonconvex optimization,” *arXiv preprint arXiv:1808.05671*, 2018.
- 262 [36] J. Zhang, S. P. Karimireddy, A. Veit, S. Kim, S. J. Reddi, S. Kumar, and S. Sra, “Why adam
263 beats sgd for attention models,” *arXiv preprint arXiv:1912.03194*, 2019.
- 264 [37] S. Lacoste-Julien, M. Schmidt, and F. Bach, “A simpler approach to obtaining an $o(1/t)$ conver-
265 gence rate for the projected stochastic subgradient method,” *arXiv preprint arXiv:1212.2002*,
266 2012.
- 267 [38] B. Neyshabur, S. Bhojanapalli, D. McAllester, and N. Srebro, “Exploring generalization in deep
268 learning,” in *Advances in Neural Information Processing Systems*, pp. 5947–5956, 2017.
- 269 [39] V. Nagarajan and J. Z. Kolter, “Generalization in deep networks: The role of distance from
270 initialization,” in *Neural Information Processing Systems (NeurIPS) Workshop, Deep Learning:*
271 *Bridging Theory and Practice*, 2017.
- 272 [40] B. Neyshabur, Z. Li, S. Bhojanapalli, Y. LeCun, and N. Srebro, “The role of over-parametrization
273 in generalization of neural networks,” in *Proceedings of Intenational Conference on Learning*
274 *Representations (ICLR)*, 2019.

275 **A More Related Works**

276 **A.1 Adversarial Regularization**

277 The notion of adversarial regularization has also been studied in Reinforcement Learning (RL). The
278 authors of [20] use adversarial learning with expert regularization to learn a predictive policy that
279 allows to drive in simulated dense traffic. In [21], the authors use RL agent controlled GAN along
280 with L_2 -distance between global feature vectors to convert noisy, partial point cloud into high-fidelity
281 data. In medical image analysis, a 3d conditional GAN along with L_1 -distance is used to super
282 resolve CT scan imagery [28].

283 **A.2 Accelerated Gradients**

284 The idea of accelerated training has long been studied. An elegant line of research focuses on
285 variance reduction that aims to address stochastic and finite sum problems by averaging the stochastic
286 noise [29, 30]. Among momentum based acceleration, much theoretical progress has been made to
287 accelerate any smooth convex optimization [31, 32]. Further, many efforts have been made towards
288 changing the step size across iterations based on estimated gradient norm [33, 34, 35].

289 **B More Theoretical Analysis**

290 **B.1 Main Results: Sub-Optimality Gap**

291 Here, we analyze the continuous time gradient flow in both approaches. In this analysis, we define
292 each iterate, $\theta(t)$ at a continuous time, t . The optimal set of parameters is denoted by θ^* . The sub-
293 optimality gap of generator and discriminator are defined by $\kappa(t) = \kappa(\theta(t)) := l(\theta(t)) - l(\theta^*)$ and
294 $\pi(t) = \pi(\theta(t)) := g(\theta^*) - g(\theta(t))$, respectively. In adversarial setting, $l(\cdot)$ is a convex downward
295 and $g(\cdot)$ is a convex upward function. For clarity, we first analyze the gradient flow in sole supervision
296 using common theoretic tools and then extend this analysis to adversarial regularization.

297 **Theorem 4.** *In purely supervised learning, the sub-optimality gap at the average over all iterates in*
298 *a trajectory of T time steps is upper bounded by*

$$\mathcal{O}\left(\frac{\|\theta(0) - \theta^*\|^2}{2T}\right).$$

299 *Proof.* Refer to Appendix D.8.

300 **Theorem 5.** *In supervised learning with adversarial regularization, the sub-optimality gap at the*
301 *average over all iterates in a trajectory of T time steps is upper bounded by*

$$\mathcal{O}\left(\frac{\|\theta(0) - \theta^*\|^2}{2T} - \pi\left(\frac{1}{T}\int_0^T \theta(t)dt\right)\right).$$

302 *Proof.* Refer to Appendix D.9.

303 According to **Theorem 4** and **5**, the distance to optimal solution decreases rapidly in augmented
304 objective when compared with purely supervised objective. Since sub-optimality gap is a non-negative
305 quantity and $\pi\left(\frac{1}{T}\int_0^T \theta(t)dt\right) \geq 0$, adversarial regularization has a *tighter* sub-optimality gap. The
306 tightness is controlled by the sub-optimality gap of adversary, $\pi(\cdot)$ at the average over all iterates in
307 the same trajectory. Also, these theorems do not require all iterates to be within the tiny landscape
308 of optimal empirical risk. The genericness of these theorems provides further evidence of better
309 empirical risk in adversarial regularization.

310 **B.2 Main Results: Provable Convergence**

311 In this section, we analyze the convergence guarantee of the minimax adversarial training under
312 strongly-convex-strongly-concave and smooth nonconvex-nonconcave criteria. In this regard, we
313 assume finite α -moment of estimated stochastic gradients as the unbounded variance has a profound
314 impact on optimization process [36, 37]. At each iteration $k = 1, \dots, T$, we denote unbiased

315 stochastic gradient by $\mathbf{g}_k = \mathbf{g}(\theta_k) := \nabla l(\theta_k, \xi) - \nabla g(\theta_k, \xi)$, where ξ represents stochasticity.
 316 Here, we analyze rates for global clipping. Similar analyses can also be made for coordinate-wise
 317 clipping [36].

318 **Assumption 7.** (Existence of α -moment) Suppose we have access to gradients at each iteration.
 319 There exist positive real numbers $\alpha \in (1, 2]$ and $G > 0$, such that $\mathbb{E}[\|\mathbf{g}(\theta)\|^\alpha] \leq G^\alpha$ for all θ .

320 **Theorem 6.** (Strongly-convex-strongly-concave convergence) Suppose **Assumption 7** holds. Let
 321 $l(\theta_k) \triangleq l(\theta_k) - g(\theta_k)$ is a μ -strongly convex function. Let $\{\theta_k\}$ be the sequence of iterates obtained
 322 using global clipping on SGD with momentum $\beta = 0$. Define the output to be k -weighted combination
 323 of iterates: $\bar{\theta} = \frac{\sum_{k=1}^T k\theta_{k-1}}{\sum_{k=1}^T k}$. If adaptive clipping $\tau_k = Gk^{\frac{1}{\alpha}}\mu^{\frac{1}{\alpha}}$ and step size $\eta_k = \frac{5}{2\mu(k+1)}$, then
 324 the output iterate $\bar{\theta}$ satisfies

$$\mathbb{E}[l(\bar{\theta})] - l(\theta^*) \leq \mathcal{O}\left(G^2(\mu(T+1))^{\frac{2-2\alpha}{\alpha}} - (g(\theta^*) - \mathbb{E}[g(\bar{\theta})])\right).$$

325

326 *Proof.* Refer to Appendix D.10.

327 Observe that when we eliminate adversarial regularization and set $\alpha = 2$, we recover exactly the
 328 SGD rate, i.e., $\mathcal{O}\left(\frac{G^2}{\mu T}\right)$ [37]. Thus, adversarial regularization does converge under strongly-convex-
 329 strongly-concave criterion. However, it is determined by the convergence of the inner maximization
 330 loop in minimax optimization.

331 **Theorem 7.** (Nonconvex-nonconcave convergence) Suppose **Assumption 3.1** and **3.2** hold. Let
 332 $l(\theta_k) \triangleq l(\theta_k) - g(\theta_k)$ is a possible L -smooth function and $\{\theta_k\}$ be the sequence of iterates obtained
 333 using global clipping on SGD with momentum, $\beta = 0$. Given constant clipping $\tau_k = G(\eta_k L)^{\frac{1}{\alpha}}$ and
 334 constant step size $\eta_k = \left(\frac{R_0^\alpha L^{2-2\alpha}}{G^2 T^\alpha}\right)^{\frac{1}{3\alpha-2}}$, where $R_0 = l(\theta_0) - l(\theta^*)$, the sequence $\{\theta_k\}$ satisfies

$$\frac{1}{T} \sum_{k=1}^T \mathbb{E}[\|\nabla l(\theta_{k-1})\|^2] \leq \mathcal{O}\left(G^{\frac{2\alpha}{3\alpha-2}} \left(\frac{R_0 L}{T}\right)^{\frac{2\alpha-2}{3\alpha-2}} - \frac{1}{T} \sum_{k=1}^T \mathbb{E}[\|\nabla g(\theta_{k-1})\|^2]\right).$$

335

336 *Proof.* Refer to Appendix D.11.

337 By setting $\alpha = 2$ and discarding adversarial acceleration, we obtain the standard SGD rate, $\mathcal{O}\left(\frac{G}{\sqrt{T}}\right)$.
 338 It is important to heed the fact adversarial regularization converges under nonconvex-nonconcave
 339 criterion as well. To this end, we have established that augmented objective is *guaranteed* to converge
 340 under strongly-convex-strongly-concave and nonconvex-nonconcave criteria provided the assumptions
 341 are satisfied. These convergence guarantees provide additional insights to our understanding of
 342 adversarial regularization in practice. It is necessary to highlight the fact that such analysis is far
 343 from being conclusive. While this paper studies minimax optimization under nonconvex-smooth
 344 settings, it will be interesting to derive convergence guarantees under nonconvex-nonsmooth settings.

345 B.3 Main Results: Generalization Error

346 Motivated by the role of over-parametrization in generalization [38, 39, 40], we study the general-
 347 ization behavior of adversarial regularization. We use Rademacher complexity to get a bound on
 348 generalization error. Since it depends on hypothesis class, we use a set of restricted parameters of
 349 trained networks to get a tighter bound on generalization. The restricted set of parameters is defined
 350 as

$$\mathcal{W} = \{(V, U) \mid V \in \mathbb{R}^{d_y \times h}, U \in \mathbb{R}^{h \times d_x}, \|v_i\| \leq \alpha_i, \|u_i - u_i^0\| \leq \beta_i\},$$

351 where $i = 1, 2, \dots, h$. Here, $v_i \in \mathbb{R}^{d_y}$ and $u_i \in \mathbb{R}^{d_x}$ denote vector representation of each neuron in
 352 the top layer and hidden layer, respectively. Thus, the restricted hypothesis class becomes

$$\mathcal{F}_{\mathcal{W}} = \{V[Ux]_+ \mid (V, U) \in \mathcal{W}\},$$

353 where $[\cdot]_+$ represents ReLU activation. For any hypothesis class \mathcal{F} , let $l \circ \mathcal{F}$ denote the composition
 354 of loss function and hypothesis class. The following generalization bound holds for any $f \in \mathcal{F}_{\mathcal{W}}$
 355 over m training samples with probability $1 - \delta$.

$$\mathbb{E}_{(x,y) \sim \mathcal{D}} [l \circ f] \leq \frac{1}{m} \sum_{i=1}^m l(f(x_i); y) + 2\mathcal{R}_{\mathcal{S}}(l \circ \mathcal{F}_{\mathcal{W}}) + 3\sqrt{\frac{\ln(2/\delta)}{2m}},$$

356 where $\mathcal{R}_{\mathcal{S}}(\mathcal{H})$ is the Rademacher complexity of a hypothesis class \mathcal{H} with respect to training set \mathcal{S} .

$$\mathcal{R}_{\mathcal{S}}(\mathcal{H}) = \frac{1}{m} \mathbb{E}_{\xi_i \in \{\pm 1\}^m} \left[\sup_{f \in \mathcal{H}} \sum_{i=1}^m \xi_i f(x_i) \right].$$

357 **Relative Generalization Error:** We define relative generalization error as

$$e_{gen,r} = \left(\mathbb{E}_{(x,y) \sim \mathcal{D}} [l \circ f] - \frac{1}{m} \sum_{i=1}^m l(f(x_i); y) \right) \times N^*.$$

358 To be consistent with [40] while studying generalization, we assume $l(f(\theta; x); y)$ be a locally K -
 359 Lipschitz function, i.e. given $y \in Y$, $\|\nabla l(f(\theta; x); y)\| \leq K$, $\forall \theta$. Using K -Lipschitz property of
 360 loss function l in **Lemma 9** of [40], one can easily prove that the Rademacher complexity of $l \circ \mathcal{F}_{\mathcal{W}}$
 361 is bounded as

$$\begin{aligned} & \mathcal{R}_{\mathcal{S}}(l \circ \mathcal{F}_{\mathcal{W}}) \\ & \leq \frac{2K\sqrt{d_y}}{m} \sum_{j=1}^h \alpha_j \left(\beta_j \|X\|_F + \|u_j^0 X\|_2 \right) \\ & \leq \frac{2K\sqrt{d_y}}{\sqrt{m}} \|\alpha\|_2 \left(\|\beta\|_2 \sqrt{\frac{1}{m} \sum_{i=1}^m \|x_i\|_2^2} + \sqrt{\frac{1}{m} \sum_{i=1}^m \|U^0 x_i\|_2^2} \right). \end{aligned}$$

362 Adapted to current setting, the generalization error becomes

$$\mathcal{O} \left(\|U^0\|_2 \|V\|_F + \|U - U^0\|_F \|V\|_F + \sqrt{h} \right).$$

363 Next, we empirically verify the required assumptions and corresponding theoretical results.

364 C Experiments

365 This section contains empirical results to support the theoretical findings on sole supervision and
 366 adversarial regularization.

367 C.1 Training Details

368 The majority of the experiments are conducted on two layer neural networks with ReLU activation
 369 function. For completeness however, we experiment with practical neural network architectures.
 370 We do not use weight decay, dropout or normalization in these networks. We experiment on both
 371 MNIST and CIFAR10 datasets. We use SGD with momentum 0.9, batch size 64 and fixed learning
 372 rate of 0.01 for MNIST and CIFAR10. We use mean square error of 0.001 for MNIST and 0.02 for
 373 CIFAR10 as our convergence criteria. We train on both datasets for a maximum of 1000 epochs, or
 374 until convergence. In these settings, we train 13 architectures on both datasets in which the number
 375 of hidden units (h) range from 2^3 to 2^{15} . We use PyTorch to design all these experiments. All
 376 parameters are initialized with uniform distribution.

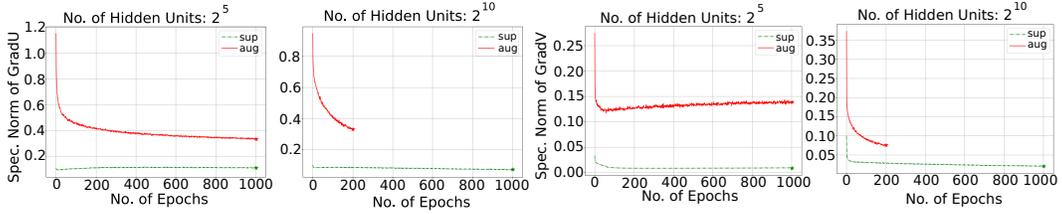


Figure 1: Comparison of gradient updates between supervised (sup) and augmented (aug) objective in the *hidden layer* (left) and *top layer* (right) on MNIST.

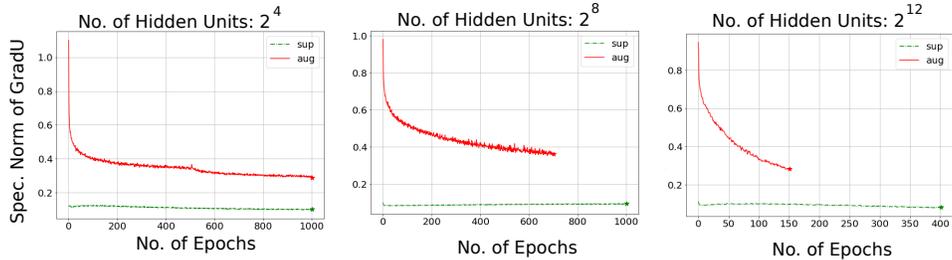


Figure 2: Comparison of gradient updates between supervised and augmented objective as observed in the *hidden layer* on MNIST.

377 C.2 Experimental Results

378 C.2.1 Results on MNIST

379 Figure 1 provides empirical evidence of the vanishing gradient issue and how adversarial regular-
 380 ization helps circumvent this. In all the experimented architectures, the spectral norm of gradients
 381 estimated by purely supervised objective is smaller than adversarial learning. This is consistent with
 382 the theoretical analysis in Section 3. The main reason for such non-vanishing gradient is the feedback
 383 signal from discriminator. Further, we observe that the rate of convergence is at least as good as sole
 384 supervision, as marked by \star in Figure 1.

385 Figure 5 offers experimental support to better empirical risk in adversarial setting. We observe the
 386 significance of near optimal region, i.e. ϵ with 32 hidden units in Figure 5. Since the expressive
 387 power of such networks is very small in both approaches, evidently neither meets the convergence
 388 criteria. However, as the capacity increases the supervised cost, which is common in both approaches,
 389 guides them to a tiny landscape around optimal risk and thereby, it satisfies the condition of **Theorem**
 390 **1**. Under this circumstance, the optimal empirical risk attained by augmented objective can be
 391 provably better than sole supervision as predicted by our theory. Figure 5 and 4 supports this theory as
 392 augmented objective consistently achieves better performance either by risk or by rate of convergence
 393 for networks with sufficient expressive power.

394 As shown in Figure 2 and 3, the estimated gradient in SGD+momentum vanishes within the tiny
 395 landscape of optimal empirical risk. Further, the adversarial regularization accelerates gradient
 396 updates and attains minimal empirical risk compared to sole supervision. It is evident from Figure 4
 397 where we observe this particular phenomenon across a wide variety of architectures.

398 Furthermore, we compare the optimal empirical risk and iteration complexity with different number
 399 of hidden units in Figure 6. One can infer from Figure 6 (a) (left) that the value of ϵ in **Theorem 1**
 400 is approximately equal to 0.005^2 . The number of epochs required to attain optimum in adversarial
 401 learning is always less than or equal to supervised learning, which validates our theorems.

²The value of ϵ is more relevant to the present body of analysis as it performs the inverse mapping in practical scenarios. Moreover, it is not hard to estimate δ where discriminator acts as the mapping function.

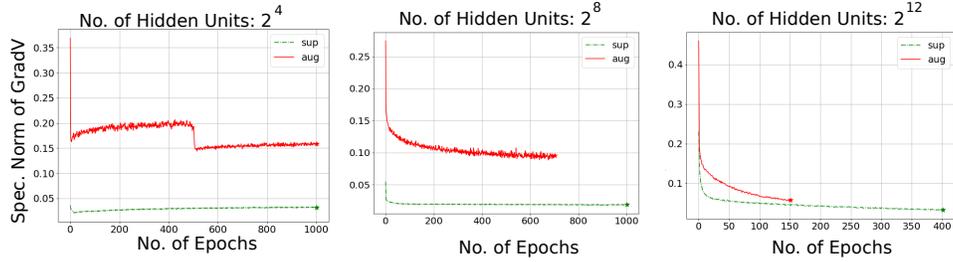


Figure 3: Comparison of gradient updates between supervised and augmented objective as observed in the *top layer* on MNIST.

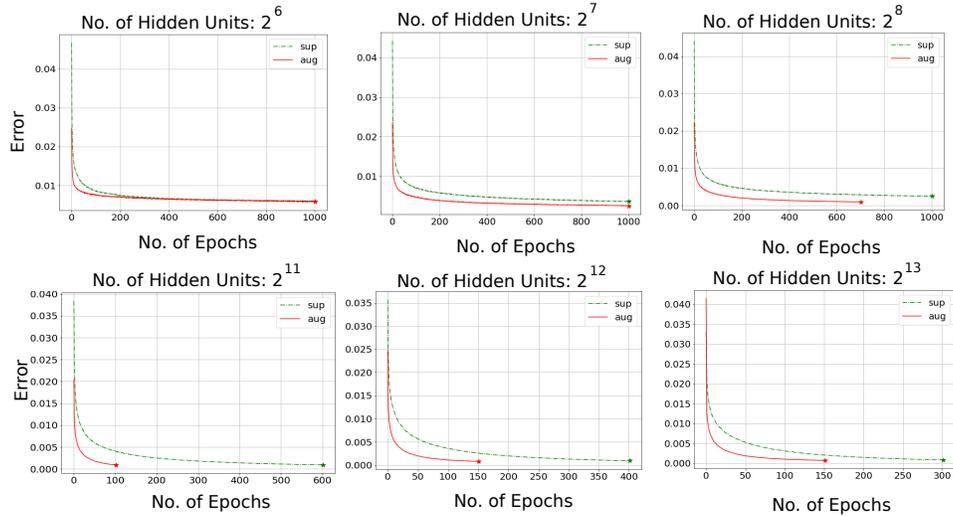


Figure 4: Comparison of optimal empirical risk on MNIST.

402 C.2.2 Results on CIFAR10

403 These theorems are also verified on CIFAR10 dataset. As shown in Figure 6 (right), supervised
 404 learning with adversarial regularization performs better than sole supervision both in terms of optimal
 405 empirical risk and rate of convergence. Here, we find ϵ to be approximately equal to 0.06.

406 Similar to our analysis on MNIST, we also observe vanishing gradient issue on CIFAR10 which is
 407 shown in Figure 7 and 8. Figure 9 illustrates how model capacity correlates with empirical risk and
 408 thereby, satisfies the condition of **Theorem 1**. Across a wide variety of architectures, we verify that
 409 supervised learning with adversarial regularization can be better than sole supervision both in terms
 410 of optimal empirical risk and iteration complexity as predicted by our theory. As shown in Figure 9,
 411 though both methods start with almost same initial empirical risk, augmented objective traverses
 412 through a shorter path and attains minimal risk upon convergence.

413 C.2.3 Results on Various Network Configurations

414 To study the impact of these findings on more realistic scenarios, we experiment on various network
 415 configurations. As shown in Figure 10 and 11, the issue of vanishing gradient is persistent across these
 416 experimented configurations. Furthermore, the discussion on adversarial acceleration is also supported
 417 by Figure 12. In addition, Table 1 shows that the proposed hypothesis: *adversarial regularization*
 418 *achieves tighter ϵ -stationary point at an optimal rate* holds under practical circumstances. More
 419 specifically, we observe accelerated gradient updates not only in two layer ReLU networks, but
 420 also in deep MLP with exponential linear activations, convolution layers, skip connections, dense
 421 connections, L_1 regularized networks, and L_2 regularized networks. Thus, augmented objective owes
 422 its performance benefits to adversarial learning at a fundamental level.

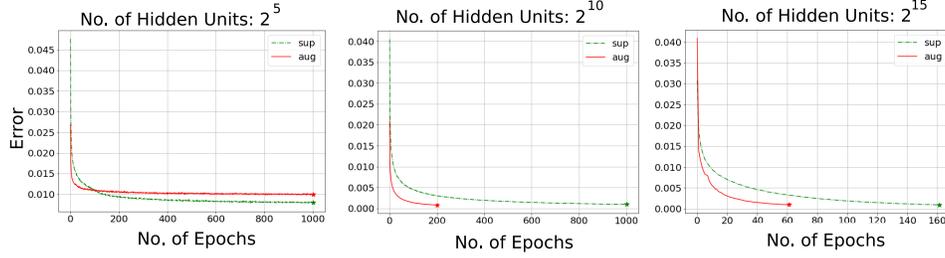


Figure 5: Comparison of optimal empirical risk on MNIST.

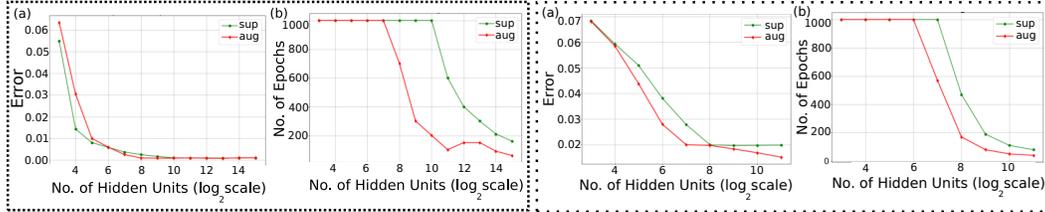


Figure 6: Comparison on MNIST (left) and CIFAR10 (right). (a) Optimal empirical risk. (b) Iteration Complexity. Adversarial regularization attains tighter ϵ -stationary point at an optimal rate.

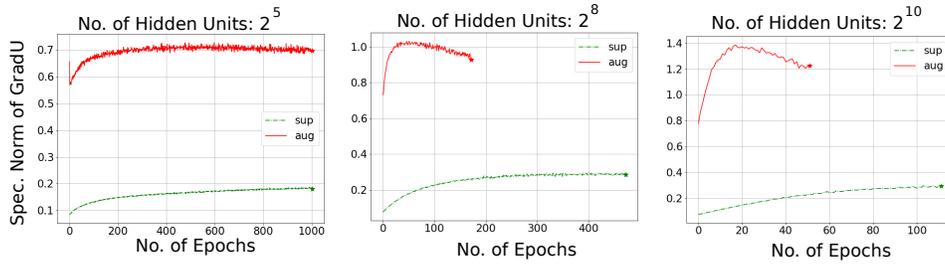


Figure 7: Comparison of gradient updates between supervised and augmented objective as observed in the *hidden layer* on CIFAR10.

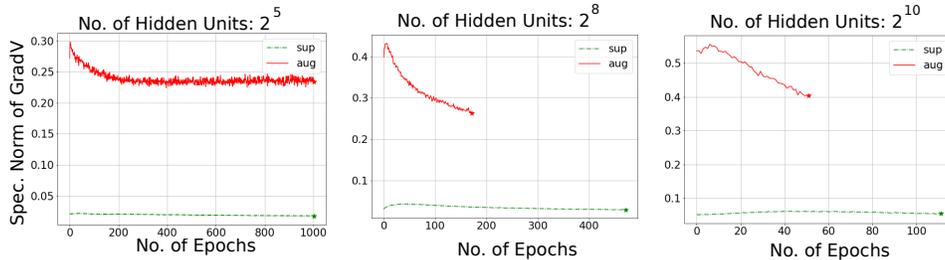


Figure 8: Comparison of gradient updates between supervised and augmented objective as observed in the *top layer* on CIFAR10.

423 **C.3 Results on Generalization Error**

424 The generalization trend in sole supervision is shown in Figure 13(a) and 13(c). As per equation (B.3),
 425 the combined measure of Frobenius norm of top layer, i.e. $\|V\|_F$ and distance from initialization of
 426 hidden layer, i.e. $\|U - U^0\|_F$ explains the generalization gap on MNIST and CIFAR10. We verify
 427 this measure in our experimental setting and study whether it can explain generalization in adversarial
 428 learning. Note that adversarial learning and sole supervision share exactly same mapping function (f),
 429 learning algorithm (SGD+momentum) and empirical data distribution (S). Thus the generalization
 430 bound, which is derived for a purely supervised objective, is expected to explain the generalization

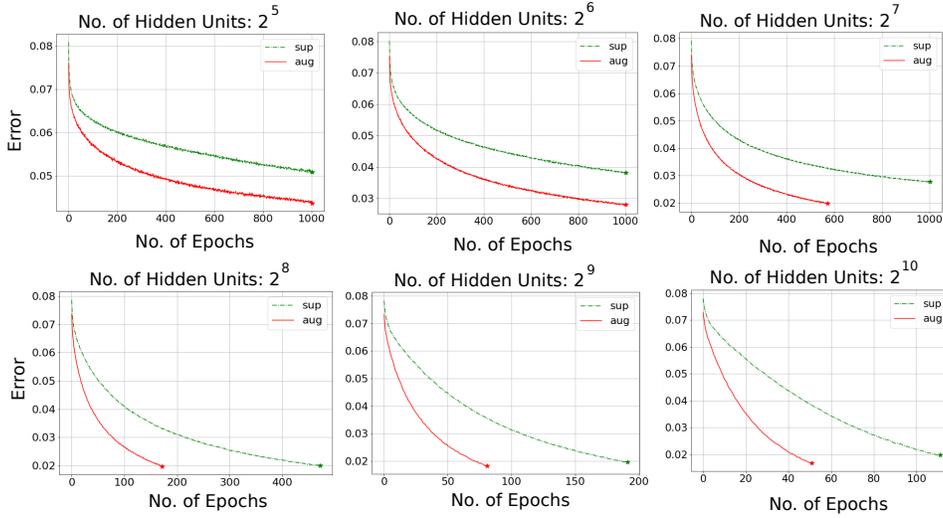


Figure 9: Comparison of optimal empirical risk on CIFAR10.

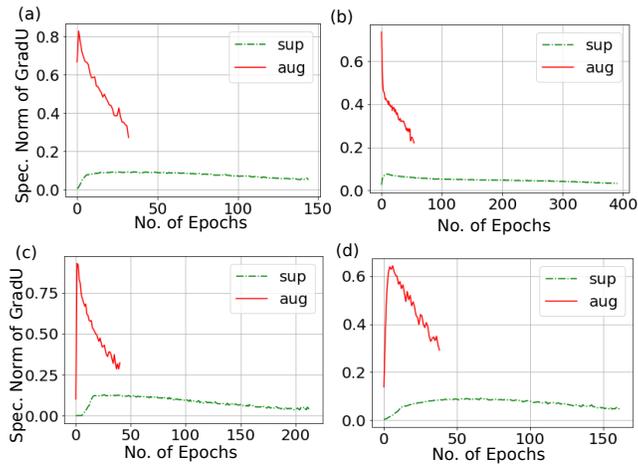


Figure 10: Comparison of gradient updates between supervised and augmented objective as observed in the *first layer* on MNIST. (a) Multi-Layer Perceptron. (b) Exponential Activation. (c) Residual Network. (d) Dense Network.

431 error in adversarial learning with expert regularization. However, as shown in Figure 13(b) and 13(d),
 432 this bound does not fully explain the generalization error observed in adversarial learning.

433 In Figure 14, we observe that the relative generalization error of adversarial regularization can be
 434 better than sole supervision. This is feasible for a network with sufficient expressive power to achieve
 435 near optimal convergence.

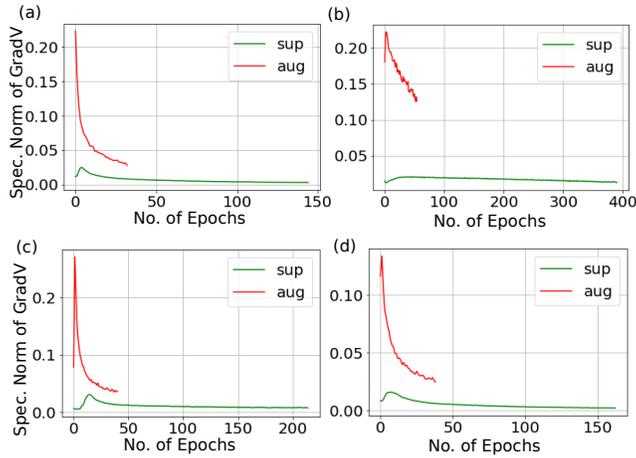


Figure 11: Comparison of gradient updates between supervised and augmented objective as observed in the *last layer* on MNIST. (a) Multi-Layer Perceptron. (b) Exponential Activation. (c) Residual Network. (d) Dense Network.

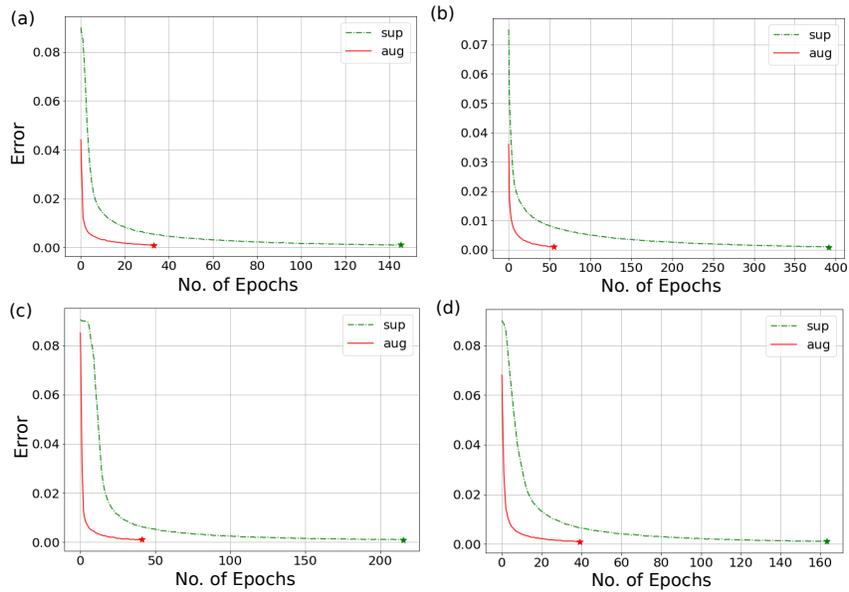


Figure 12: Comparison of optimal empirical risk on MNIST. (a) Multi-Layer Perceptron. (b) Exponential Activation. (c) Residual Network. (d) Dense Network.

Table 1: Hypothesis Testing on Various Network Configurations

Architecture	No. Layer	Activation	No. ResBlock	No. DenseBlock	No. Epoch Sup	No. Epoch Aug	Hypothesis
MLP-Deep	6	ELU	2	0	391	55	✓
CNN-ResNet	6	ReLU	2	0	215	41	✓
CNN-DenseNet	6	ReLU	2	1	163	39	✓
CNN-DenseNet-L1	6	ReLU	2	1	1000	39	✓
CNN-DenseNet-L2	6	ReLU	2	1	155	39	✓
CNN-ResNet-AvgPool	6	ReLU	2	0	109	29	✓

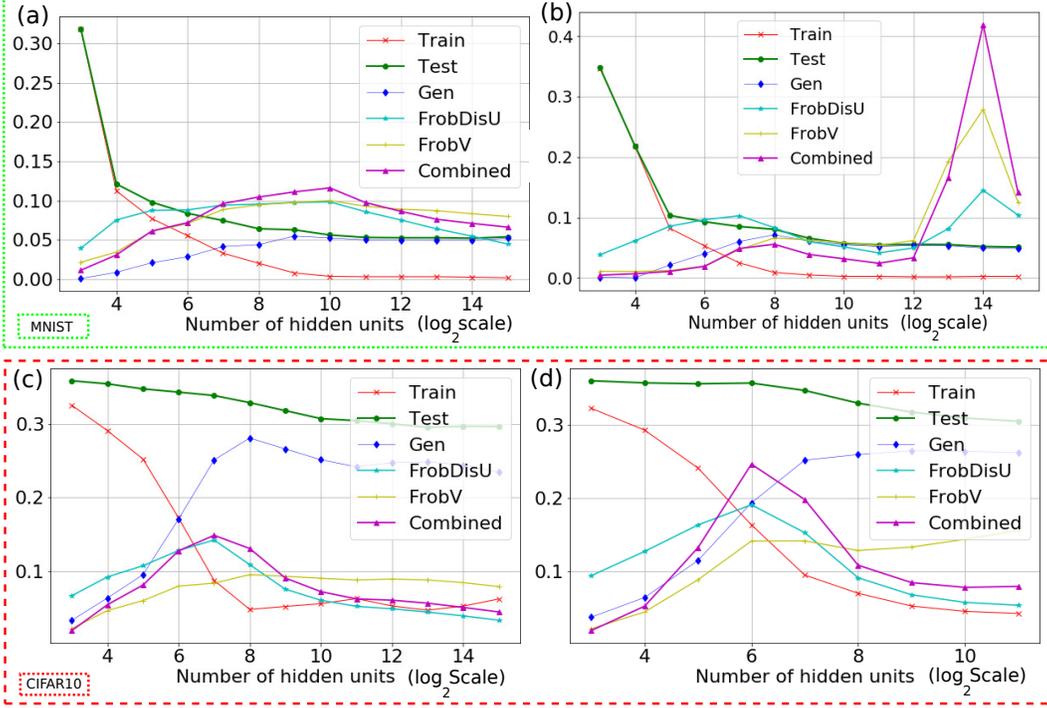


Figure 13: Generalization error on MNIST and CIFAR10.

436 D Technical Proofs

437 D.1 Proof of Lemma 1

438 Using Jensen's inequality,

$$\begin{aligned}
 \|\nabla_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y)]\|^2 &\leq \mathbb{E}_{(x,y) \sim \mathcal{P}} \left[\|\nabla_{\theta} l(f(\theta; x); y)\|^2 \right] \\
 &\leq \mathbb{E}_{(x,y) \sim \mathcal{P}} \left[\|\nabla_p l(p; y) \nabla_{\theta} f(\theta; x)\|^2 \right], \text{ where } p = f(\theta; x) \\
 &\leq \mathbb{E}_{(x,y) \sim \mathcal{P}} \left[\underbrace{\|\nabla_p l(p; y)\|^2 \|\nabla_{\theta} f(\theta; x)\|^2}_{\text{Cauchy-Schwarz inequality}} \right] \\
 &\leq L^2 \mathbb{E}_{(x,y) \sim \mathcal{P}} \left[\|\nabla_p l(p; y)\|^2 \right]
 \end{aligned}$$

439 Let $p = f(\theta; x)$ and $q = f(\theta^*; y)$. Using β -smoothness and L -Lipschitz property, we get

$$\|\nabla_p l(p; y)\| - \|\nabla_q l(q; y)\| \leq \|\nabla_p l(p; y) - \nabla_q l(q; y)\| \leq \beta \|p - q\| \leq \beta L \|\theta - \theta^*\|.$$

440 Since $\|\theta - \theta^*\| \leq \epsilon$,

$$\|\nabla_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y)]\|^2 \leq L^2 \mathbb{E}_{(x,y) \sim \mathcal{P}} \left[(\|\nabla_q l(q; y)\| + L\beta\epsilon)^2 \right].$$

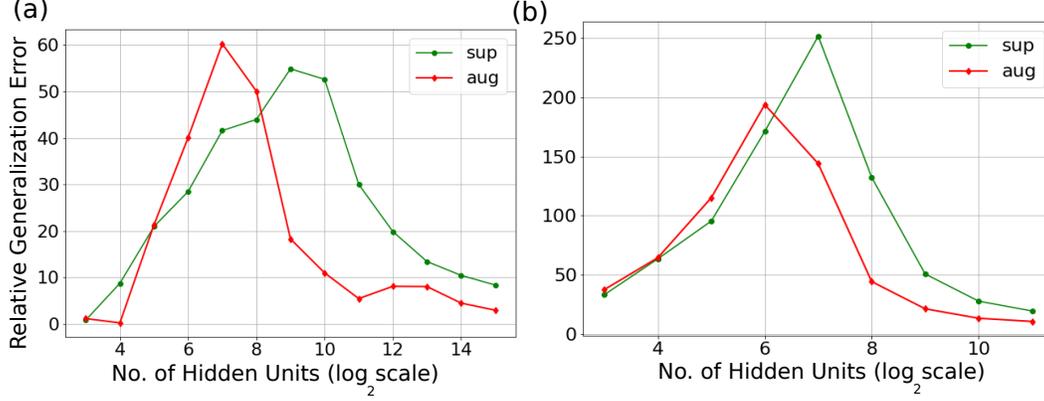


Figure 14: Relative generalization. (a) MNIST. (b) CIFAR10.

441 Upon substituting optimality condition, i.e. $\|\nabla_{\theta} l(q; y)\| = 0$, the above expression simplifies to

$$\|\nabla_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y)]\| \leq L^2 \beta \epsilon.$$

442 This completes the proof of the theorem. □

443 D.2 Proof of Lemma 2

444 Using similar arguments from Appendix D.1,

$$\begin{aligned} \|\nabla_{\theta} \mathbb{E}_{x \sim \mathcal{P}_X} [g(\psi; f(\theta; x))]\|^2 &\leq \mathbb{E}_{x \sim \mathcal{P}_X} \left[\|\nabla_{\theta} g(\psi; f(\theta; x))\|^2 \right] \\ &\leq \mathbb{E}_{x \sim \mathcal{P}_X} \left[\|\nabla_p g(\psi; p)\|^2 \|\nabla_{\theta} f(\theta; x)\|^2 \right], \text{ where } p = f(\theta; x) \\ &\leq L^2 \mathbb{E}_{x \sim \mathcal{P}_X} \left[\|\nabla_p g(\psi; p)\|^2 \right] \\ &\leq L^2 \mathbb{E}_{x \sim \mathcal{P}_X} \left[(\|\nabla_p g(\psi; p)\| + \delta)^2 \right] \\ &\leq L^2 \delta^2 \end{aligned}$$

445 Taking square root, $\|\nabla_{\theta} \mathbb{E}_{x \sim \mathcal{P}_X} [g(\psi; f(\theta; x))]\| \leq L\delta$, which finishes the proof. □

446 D.3 Proof of Theorem 1

447 By applying triangle inequality after simplification,

$$\begin{aligned} \|\nabla_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y) - g(\psi; f(\theta; x))]\| &\leq \|\nabla_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [l(f(\theta; x); y)]\| + \|\nabla_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [g(\psi; f(\theta; x))]\| \\ &\leq L^2 \beta \epsilon + L\delta \text{ (Lemma 1 and Lemma 2),} \end{aligned}$$

448 which completes the statement of the theorem. □

449 D.4 Proof of Theorem 2

450 We parameterize the path between θ_k and θ_{k+1} as following:

$$\gamma(t) = t\theta_{k+1} + (1-t)\theta_k \forall t \in [0, 1]. \quad (4)$$

451 By fixed step gradient descent, the iterate $\theta_{k+1} = \theta_k - h_k \nabla l(\theta_k)$. Using Taylor's expansion,

$$\begin{aligned} l(\theta_{k+1}) &= l(\theta_k) + \nabla l(\theta_k) (\theta_{k+1} - \theta_k) + \frac{1}{2} (\theta_{k+1} - \theta_k)^T \nabla^2 l(\theta_k) (\theta_{k+1} - \theta_k) \\ &= l(\theta_k) - h_k \|\nabla l(\theta_k)\|^2 + \frac{1}{2} (\theta_{k+1} - \theta_k)^T \nabla^2 l(\theta_k) (\theta_{k+1} - \theta_k), \text{ } (\because \theta_{k+1} - \theta_k = -h_k \nabla l(\theta_k)). \end{aligned}$$

452 Using Cauchy-Schwarz inequality and integrating over parameterized curve $\gamma(t)$,

$$\begin{aligned} l(\theta_{k+1}) &\leq l(\theta_k) - h_k \|\nabla l(\theta_k)\|^2 + \frac{1}{2} \|(\theta_{k+1} - \theta_k)\| \|\nabla^2 l(\theta_k)(\theta_{k+1} - \theta_k)\| \\ &\leq l(\theta_k) - h_k \|\nabla l(\theta_k)\|^2 + \frac{1}{2} \|(\theta_{k+1} - \theta_k)\|^2 \int_0^1 \|\nabla^2 l(\gamma(t))\| dt. \end{aligned}$$

453 We know by **Assumption 5**

$$\|\nabla^2 l(\theta)\| \leq L_0 + L_1 \|\nabla l(\theta)\|.$$

454 Then using descent rule and arguments of **Theorem 1**, we obtain the following inequality:

$$\begin{aligned} l(\theta_{k+1}) &\leq l(\theta_k) - h_k \|\nabla l(\theta_k)\|^2 + \frac{h_k^2 \|\nabla l(\theta_k)\|^2}{2} \int_0^1 (L_0 + L_1 \|\nabla l(\gamma(t))\|) dt \\ &\leq l(\theta_k) - h_k \|\nabla l(\theta_k)\|^2 + \frac{h_k^2 \|\nabla l(\theta_k)\|^2}{2} \int_0^1 (L_0 + L_1 L^2 \beta \epsilon) dt \\ &\leq l(\theta_k) - h_k \|\nabla l(\theta_k)\|^2 + \frac{h_k^2 \|\nabla l(\theta_k)\|^2}{2} (L_0 + L_1 L^2 \beta \epsilon). \end{aligned}$$

455 Let us choose $h_k = \frac{1}{L_0 + L_1 L^2 \beta \epsilon}$. Now,

$$\begin{aligned} l(\theta_{k+1}) &\leq l(\theta_k) - \frac{h_k \|\nabla l(\theta_k)\|^2}{2} \\ &\leq l(\theta_k) - \frac{\|\nabla l(\theta_k)\|^2}{2(L_0 + L_1 \lambda M)}. \end{aligned}$$

456 Assume that it takes T iterations to reach ϵ -stationary point, i.e., $\epsilon \leq \|\nabla l(\theta_k)\|$ for $k \leq T$. By a
457 telescopic sum over k ,

$$\begin{aligned} \sum_{k=0}^{T-1} l(\theta_{k+1}) - l(\theta_k) &\leq \frac{-T\epsilon^2}{2(L_0 + L_1 \lambda M)} \\ \implies T &\leq \frac{2(l(\theta_0) - l^*)(L_0 + L_1 L^2 \beta \epsilon)}{\epsilon^2}. \end{aligned}$$

458 Therefore, we get

$$\sup_{\theta_0 \in \{\mathbb{R}^h \times d_x, \mathbb{R}^{d_y \times h}\}, l \in \mathcal{L}} \mathcal{T}_\epsilon(A_h[l, \theta_0], l) = \mathcal{O}\left(\frac{(l(\theta_0) - l^*)(L_0 + L_1 L^2 \beta \epsilon)}{\epsilon^2}\right)$$

459 which finishes the proof. □

460 D.5 Proof of Corollary 1

461 Using the arguments made in the proof of **Theorem 2** and first-order Taylor's expansion, we get

$$\begin{aligned} l(\theta_{k+1}) &= l(\theta_k) - h_k \|\nabla l(\theta_k)\|^2 \\ &\leq l(\theta_k) - h_k \epsilon^2. \end{aligned}$$

462 By telescopic sum,

$$\begin{aligned} \sum_{k=0}^{T-1} l(\theta_{k+1}) - l(\theta_k) &\leq -T h_k \epsilon^2 \\ \implies T &\leq \frac{(l(\theta_0) - l^*)}{h_k \epsilon^2}. \end{aligned}$$

463 So,

$$\sup_{\theta_0 \in \{\mathbb{R}^h \times d_x, \mathbb{R}^{d_y \times h}\}, l \in \mathcal{L}} \mathcal{T}_\epsilon(A_h[l, \theta_0], l) = \mathcal{O}\left(\frac{(l(\theta_0) - l^*)}{h \epsilon^2}\right)$$

464 which finishes the proof. □

465 **D.6 Proof of Theorem 3**

466 Recall that the target function $l(\theta)$ remains identical in both settings except for additional cost of
467 discriminator over generator in augmented objective. In this setting, the parameters are updated as

$$\theta_{k+1} = \theta_k - h_k \nabla (l(\theta_k) - g(\psi; f(\theta_k; x))). \quad (5)$$

468 Using Taylor's expansion, the triangle and Cauchy-Schwarz inequality as in Appendix D.4, we obtain

$$l(\theta_{k+1}) \leq l(\theta_k) - h_k \|\nabla l(\theta_k)\|^2 - h_k \|\nabla l(\theta_k)\| \|\nabla g(\psi; f(\theta_k; x))\| + \frac{h_k^2 \|\nabla (l(\theta_k) - g(\psi; f(\theta_k; x)))\|^2}{2} \int_0^1 \|\nabla^2 l(\gamma(t))\| dt.$$

469 By **Assumption 5** and **6**,

$$l(\theta_{k+1}) \leq l(\theta_k) - h_k \|\nabla l(\theta_k)\|^2 - h_k \|\nabla l(\theta_k)\| \zeta + \frac{h_k^2 \|\nabla l(\theta_k) - \nabla g(\psi; f(\theta_k; x))\|^2}{2} \int_0^1 (L_0 + L_1 \|\nabla l(\gamma(t))\|) dt.$$

470 Upon simplification using arguments of Appendix D.4 and applying Minkowski's inequality,

$$l(\theta_{k+1}) \leq l(\theta_k) - h_k \|\nabla l(\theta_k)\|^2 - h_k \|\nabla l(\theta_k)\| \zeta + \frac{h_k^2 (\|\nabla l(\theta_k)\|^2 + \|\nabla g(\psi; f(\theta_k; x))\|^2)}{2} (L_0 + L_1 \lambda M).$$

471 Using $h_k = \frac{1}{L_0 + L_1 L^2 \beta \epsilon}$, we get

$$\begin{aligned} l(\theta_{k+1}) &\leq l(\theta_k) - \frac{h_k \|\nabla l(\theta_k)\|^2}{2} - h_k \|\nabla l(\theta_k)\| \zeta + \frac{h_k \|\nabla g(\psi; f(\theta_k; x))\|^2}{2} \\ &\leq l(\theta_k) - \frac{h_k \|\nabla l(\theta_k)\|^2}{2} - h_k \|\nabla l(\theta_k)\| \zeta + \frac{h_k L^2 \delta^2}{2}, \text{ (from Lemma 2).} \end{aligned}$$

472 Assuming T iterations to reach ϵ -stationary point, i.e., $\epsilon \leq \|\nabla l(\theta_k)\|$ for $k \leq T$. By a telescopic
473 sum over k ,

$$\begin{aligned} \sum_{k=0}^{T-1} l(\theta_{k+1}) - l(\theta_k) &\leq \frac{-T(\epsilon^2 + 2\epsilon\zeta - L^2\delta^2)}{2(L_0 + L_1 L^2 \beta \epsilon)} \\ \implies T &\leq \frac{2(l(\theta_0) - l^*)(L_0 + L_1 L^2 \beta \epsilon)}{\epsilon^2 + 2\epsilon\zeta - L^2\delta^2}. \end{aligned}$$

474 Therefore, we obtain

$$\sup_{\theta_0 \in \{\mathbb{R}^{h \times d_x}, \mathbb{R}^{d_y \times h}\}, l \in \mathcal{L}} \mathcal{T}_\epsilon(A_h[l, \theta_0], l) = \mathcal{O}\left(\frac{(l(\theta_0) - l^*)(L_0 + L_1 \lambda M)}{\epsilon^2 + 2\epsilon\zeta - \delta^2 M^2}\right)$$

475 which finishes the proof. \square

476 **D.7 Proof of Corollary 2**

477 Using the arguments made in the proof of **Theorem 3** and first-order Taylor's approximation, we get

$$\begin{aligned} l(\theta_{k+1}) &= l(\theta_k) - h_k \|\nabla l(\theta_k)\|^2 - h_k \|\nabla l(\theta_k)\| \|\nabla g(\psi; f(\theta_k; x))\| \\ &\leq l(\theta_k) - h_k \epsilon^2 - h_k \epsilon \zeta. \end{aligned}$$

478 By telescopic sum,

$$\begin{aligned} \sum_{k=0}^{T-1} l(\theta_{k+1}) - l(\theta_k) &\leq -T h_k \epsilon^2 - T h_k \epsilon \zeta \\ \implies T &\leq \frac{(l(\theta_0) - l^*)}{h_k \epsilon^2 + h_k \epsilon \zeta}. \end{aligned}$$

479 Therefore,

$$\sup_{\theta_0 \in \{\mathbb{R}^{h \times d_x}, \mathbb{R}^{d_y \times h}\}, l \in \mathcal{L}} \mathcal{T}_\epsilon(A_h[l, \theta_0], l) = \mathcal{O}\left(\frac{(l(\theta_0) - l^*)}{h \epsilon^2 + h \epsilon \zeta}\right)$$

480 which finishes the proof. \square

481 **D.8 Proof of Theorem 4**

482 In sole supervision, the parameters are updated by $\frac{d\theta(t)}{dt} = -\nabla l(\theta(t))$. We define distance to optimal
 483 solution as $r^2(t) = \frac{1}{2} \|\theta(t) - \theta^*\|^2$. Now differentiating both sides, we get

$$\begin{aligned} \frac{dr^2(t)}{dt} &= \left\langle \frac{d\theta(t)}{dt}, \theta(t) - \theta^* \right\rangle \\ &= \langle -\nabla l(\theta(t)), \theta(t) - \theta^* \rangle. \end{aligned}$$

484 Using convexity and integrating over all iterates in a trajectory of T time steps,

$$\begin{aligned} \frac{1}{T} \int_0^T \frac{dr^2(t)}{dt} dt &\leq \frac{1}{T} \int_0^T -\kappa(t) dt \\ \implies \frac{1}{T} (r^2(T) - r^2(0)) &\leq -\frac{1}{T} \int_0^T \kappa(t) dt \\ \implies \frac{1}{T} \int_0^T \kappa(\theta(t)) dt &\leq \frac{r^2(0)}{T}. \end{aligned}$$

485 By Jensen's inequality,

$$\kappa \left(\frac{1}{T} \int_0^T \theta(t) dt \right) \leq \frac{1}{T} \int_0^T \kappa(\theta(t)) dt.$$

486 Therefore, $\kappa \left(\frac{1}{T} \int_0^T \theta(t) dt \right) = \mathcal{O} \left(\frac{\|\theta(0) - \theta^*\|^2}{2T} \right)$ which finishes the proof. \square

487 **D.9 Proof of Theorem 5**

488 In supervised learning with adversarial regularization, the parameters are updated by $\frac{d\theta(t)}{dt} =$
 489 $-\nabla l(\theta(t)) + \nabla g(\theta(t))$. Using arguments of Appendix D.8, we obtain

$$\frac{dr^2(t)}{dt} = \langle -\nabla l(\theta(t)), \theta(t) - \theta^* \rangle + \langle \nabla g(\theta(t)), \theta(t) - \theta^* \rangle.$$

490 Since $l(\cdot)$ is a convex downward and $g(\cdot)$ is a convex upward function, we get

$$\begin{aligned} \frac{1}{T} \int_0^T \frac{dr^2(t)}{dt} dt &\leq -\frac{1}{T} \int_0^T \kappa(t) dt - \frac{1}{T} \int_0^T \pi(t) dt \\ \implies \frac{1}{T} (r^2(T) - r^2(0)) &\leq -\frac{1}{T} \int_0^T \kappa(t) dt - \frac{1}{T} \int_0^T \pi(t) dt \\ \implies \frac{1}{T} \int_0^T \kappa(\theta(t)) dt &\leq \frac{r^2(0)}{T} - \frac{1}{T} \int_0^T \pi(\theta(t)) dt. \end{aligned}$$

491 Now, using Jensen's inequality on both $\kappa(\cdot)$ and $\pi(\cdot)$

$$\kappa \left(\frac{1}{T} \int_0^T \theta(t) dt \right) = \mathcal{O} \left(\frac{\|\theta(0) - \theta^*\|^2}{2T} - \pi \left(\frac{1}{T} \int_0^T \theta(t) dt \right) \right)$$

492 which finishes the proof. \square

493 **D.10 Proof of Theorem 6**

494 For simplicity, let us denote the bias $b_k = \mathbb{E}[\hat{\mathbf{g}}_k] - \nabla l(\theta_k)$.

$$\begin{aligned} \|\theta_k - \theta^*\|^2 &= \|\theta_{k-1} - \eta_k \hat{\mathbf{g}}_{k-1} - \theta^*\|^2 \\ &= \|\theta_{k-1} - \theta^*\|^2 - 2\eta_k \langle \theta_{k-1} - \theta^*, \hat{\mathbf{g}}_{k-1} \rangle + \eta_k^2 \|\hat{\mathbf{g}}_{k-1}\|^2 \\ &= \|\theta_{k-1} - \theta^*\|^2 - 2\eta_k \langle \theta_{k-1} - \theta^*, \nabla l(\theta_{k-1}) \rangle - 2\eta_k \langle \theta_{k-1} - \theta^*, b_{k-1} \rangle + \eta_k^2 \|\hat{\mathbf{g}}_{k-1}\|^2 \\ &\leq \|\theta_{k-1} - \theta^*\|^2 - 2\eta_k \langle \theta_{k-1} - \theta^*, \nabla l(\theta_{k-1}) \rangle + \underbrace{2\eta_k \|\theta_{k-1} - \theta^*\| \|b_{k-1}\|}_{\text{By Cauchy-Schwarz inequality}} + \eta_k^2 \|\hat{\mathbf{g}}_{k-1}\|^2 \\ &\leq \|\theta_{k-1} - \theta^*\|^2 - 2\eta_k \langle \theta_{k-1} - \theta^*, \nabla l(\theta_{k-1}) \rangle + \underbrace{\eta_k \left(\|\theta_{k-1} - \theta^*\|^2 + \|b_{k-1}\|^2 \right)}_{\text{By AM-GM inequality}} + \eta_k^2 \|\hat{\mathbf{g}}_{k-1}\|^2 \end{aligned}$$

495 By μ -strong convexity, it is required that there exist positive constants μ such that for all (x, y) ,
 496 $\mathfrak{l}(y) \geq \mathfrak{l}(x) + \langle y - x, \nabla \mathfrak{l}(x) \rangle + \frac{\mu}{2} \|y - x\|^2$. Using strong-convexity at θ_{k-1} and θ^* , we get

$$\begin{aligned} \|\theta_k - \theta^*\|^2 &\leq \|\theta_{k-1} - \theta^*\|^2 - 2\eta_k (\mathfrak{l}(\theta_{k-1}) - \mathfrak{l}(\theta^*)) - \eta_k \mu \|\theta_{k-1} - \theta^*\|^2 + \eta_k \left(\|\theta_{k-1} - \theta^*\|^2 + \|b_{k-1}\|^2 \right) + \eta_k^2 \|\hat{\mathfrak{g}}_{k-1}\|^2 \\ &\leq \|\theta_{k-1} - \theta^*\|^2 (1 - \eta_k \mu + \eta_k) - 2\eta_k (\mathfrak{l}(\theta_{k-1}) - \mathfrak{l}(\theta^*)) + \eta_k \|b_{k-1}\|^2 + \eta_k^2 \|\hat{\mathfrak{g}}_{k-1}\|^2. \end{aligned}$$

497 **Lemma 3.** Suppose **Assumption 7** holds for any $\mathfrak{g}(\theta)$ and $\alpha \in (1, 2]$. With global clipping parameter
 498 $\tau \geq 0$, the variance and bias of the estimator $\hat{\mathfrak{g}}$ are upper bounded as:

$$\mathbb{E} \left[\|\hat{\mathfrak{g}}(\theta)\|^2 \right] \leq G^\alpha \tau^{2-\alpha} \text{ and } \|\mathbb{E} [\hat{\mathfrak{g}}(\theta)] - \nabla \mathfrak{l}(\theta) + \nabla g(\theta)\|^2 \leq G^{2\alpha} \tau^{2-2\alpha}.$$

499

500 One can easily prove this using **Lemma 2** of [36]. Upon rearranging, taking expectation of both
 501 sides, and using **Lemma 3**,

$$\mathbb{E} [\mathfrak{l}(\theta_{k-1})] - \mathfrak{l}(\theta^*) \leq \mathbb{E} \left[\left(\frac{\eta_k^{-1} - \mu + 1}{2} \right) \|\theta_{k-1} - \theta^*\|^2 - \frac{\eta_k^{-1}}{2} \|\theta_k - \theta^*\|^2 \right] + \frac{1}{2} G^{2\alpha} \tau^{2-2\alpha} + \frac{\eta_k}{2} G^\alpha \tau^{2-\alpha}.$$

502 Let us choose $\frac{\eta_k^{-1} - \mu + 1}{2} = k - 1$ and $\frac{\eta_k^{-1}}{2} = k + 1$. After simplification, $\eta_k = \frac{5}{2\mu(k+1)}$. Now,
 503 substitute $\tau_k = G k^{\frac{1}{\alpha}} \mu^{\frac{1}{\alpha}}$, $\eta_k = \frac{5}{2\mu(k+1)}$ and multiply k both sides. Thus,

$$k \mathbb{E} [\mathfrak{l}(\theta_{k-1})] - k \mathfrak{l}(\theta^*) \leq \mathbb{E} \left[k(k-1) \|\theta_{k-1} - \theta^*\|^2 - k(k+1) \|\theta_k - \theta^*\|^2 \right] + \frac{G^2 k^{\frac{2-\alpha}{\alpha}} \mu^{\frac{2-2\alpha}{\alpha}}}{2} \left[\frac{5}{2} \left(\frac{k}{k+1} \right) + 1 \right].$$

504 Since $\frac{k}{k+1} < 1$ for $k = 1, \dots, T$, we get

$$k \mathbb{E} [\mathfrak{l}(\theta_{k-1})] - k \mathfrak{l}(\theta^*) \leq \mathbb{E} \left[k(k-1) \|\theta_{k-1} - \theta^*\|^2 - k(k+1) \|\theta_k - \theta^*\|^2 \right] + \frac{7G^2 k^{\frac{2-\alpha}{\alpha}} \mu^{\frac{2-2\alpha}{\alpha}}}{4}.$$

505 Taking telescopic sum over $k = 1, \dots, T$, we obtain

$$\sum_{k=1}^T k \mathbb{E} [\mathfrak{l}(\theta_{k-1})] - \mathfrak{l}(\theta^*) \sum_{k=1}^T k \leq \mathbb{E} \left[-T(T+1) \|\theta_T - \theta^*\|^2 \right] + \frac{7G^2 \mu^{\frac{2-2\alpha}{\alpha}}}{4} \sum_{k=1}^T k^{\frac{2-\alpha}{\alpha}}.$$

506 Using $\sum_{k=1}^T k^{\frac{2-\alpha}{\alpha}} \leq \int_0^{T+1} k^{\frac{2-\alpha}{\alpha}} dk \leq (T+1)^{\frac{2}{\alpha}}$,

$$\sum_{k=1}^T k \mathbb{E} [\mathfrak{l}(\theta_{k-1})] - \mathfrak{l}(\theta^*) \frac{T(T+1)}{2} \leq \frac{7G^2 \mu^{\frac{2-2\alpha}{\alpha}}}{4} (T+1)^{\frac{2}{\alpha}}.$$

507 Now, dividing both sides by $\frac{T(T+1)}{2}$ and using $T^{-1} \leq 2(T+1)^{-1}$ for $T \geq 1$,

$$\frac{\sum_{k=1}^T k \mathbb{E} [\mathfrak{l}(\theta_{k-1})]}{\sum_{k=1}^T k} - \mathfrak{l}(\theta^*) \leq 7G^2 \mu^{\frac{2-2\alpha}{\alpha}} (T+1)^{\frac{2-2\alpha}{\alpha}}.$$

508 By Jensen's inequality,

$$\mathbb{E} \left[\mathfrak{l} \left(\frac{\sum_{k=1}^T k \theta_{k-1}}{\sum_{k=1}^T k} \right) \right] - \mathfrak{l}(\theta^*) \leq \mathcal{O} \left(G^2 (\mu(T+1))^{\frac{2-2\alpha}{\alpha}} \right)$$

509 Substituting $\mathfrak{l}(\theta) = l(\theta) - g(\theta)$, we get

$$\mathbb{E} [l(\bar{\theta})] - l(\theta^*) \leq \mathcal{O} \left(G^2 (\mu(T+1))^{\frac{2-2\alpha}{\alpha}} - (g(\theta^*) - \mathbb{E} [g(\bar{\theta})]) \right),$$

510 which finishes the proof. \square

511 **D.11 Proof of Theorem 7**

512 The notations of l and b_k follow from Appendix D.10. Using L -smooth property of l , we get

$$\begin{aligned}
l(\theta_k) &\leq l(\theta_{k-1}) + \langle \nabla l(\theta_{k-1}), \theta_k - \theta_{k-1} \rangle + \frac{L}{2} \|\theta_k - \theta_{k-1}\|^2 \\
&\leq l(\theta_{k-1}) + \langle \nabla l(\theta_{k-1}), -\eta_k \hat{\mathbf{g}}_{k-1} \rangle + \frac{\eta_k^2 L}{2} \|\hat{\mathbf{g}}_{k-1}\|^2 \\
&\leq l(\theta_{k-1}) - \eta_k \|\nabla l(\theta_{k-1})\|^2 - \eta_k \langle \nabla l(\theta_{k-1}), b_{k-1} \rangle + \frac{\eta_k^2 L}{2} \|\hat{\mathbf{g}}_{k-1}\|^2 \\
&\leq l(\theta_{k-1}) - \eta_k \|\nabla l(\theta_{k-1})\|^2 + \underbrace{\eta_k \|\nabla l(\theta_{k-1})\| \|b_{k-1}\|}_{\text{By Cauchy-Schwarz inequality}} + \frac{\eta_k^2 L}{2} \|\hat{\mathbf{g}}_{k-1}\|^2 \\
&\leq l(\theta_{k-1}) - \eta_k \|\nabla l(\theta_{k-1})\|^2 + \underbrace{\frac{\eta_k}{2} \left(\|\nabla l(\theta_{k-1})\|^2 + \|b_{k-1}\|^2 \right)}_{\text{By AM-GM inequality}} + \frac{\eta_k^2 L}{2} \|\hat{\mathbf{g}}_{k-1}\|^2
\end{aligned}$$

513 Taking expectation of both sides,

$$\mathbb{E} [l(\theta_k) - l(\theta_{k-1})] \leq \mathbb{E} \left[\frac{-\eta_k}{2} \|\nabla l(\theta_{k-1})\|^2 \right] + \frac{\eta_k}{2} G^{2\alpha} \tau^{2-2\alpha} + \frac{\eta_k^2 L}{2} G^\alpha \tau^{2-\alpha}.$$

514 Upon rearranging and taking telescopic sum over $k = 1, \dots, T$, we obtain

$$\frac{1}{T} \sum_{k=1}^T \mathbb{E} \left[\|\nabla l(\theta_{k-1})\|^2 \right] \leq \frac{2\eta_k^{-1}}{2} (l(\theta_0) - l(\theta^*)) + G^{2\alpha} \tau^{2-2\alpha} + \eta_k L G^\alpha \tau^{2-\alpha}.$$

515 By choosing $\tau = G (\eta_k L)^{-\frac{1}{\alpha}}$,

$$\frac{1}{T} \sum_{k=1}^T \mathbb{E} \left[\|\nabla l(\theta_{k-1})\|^2 \right] \leq \frac{2\eta_k^{-1} R_0}{T} + 2G^2 (\eta_k L)^{\frac{2\alpha-2}{\alpha}}.$$

516 Let us choose $\eta_k = \left(\frac{R_0^\alpha L^{2-2\alpha}}{G^2 T^\alpha} \right)^{\frac{1}{3\alpha-2}}$. Thus,

$$\frac{1}{T} \sum_{k=1}^T \mathbb{E} \left[\|\nabla l(\theta_{k-1})\|^2 \right] \leq 4G^{\frac{2\alpha}{3\alpha-2}} \left(\frac{R_0 L}{T} \right)^{\frac{2\alpha-2}{3\alpha-2}}$$

517 Now, substituting $l(\theta) = l(\theta) - g(\theta)$, we get

$$\frac{1}{T} \sum_{k=1}^T \mathbb{E} \left[\|\nabla l(\theta_{k-1})\|^2 + \|\nabla g(\theta_{k-1})\|^2 - 2\langle \nabla l(\theta_{k-1}), \nabla g(\theta_{k-1}) \rangle \right] \leq 4G^{\frac{2\alpha}{3\alpha-2}} \left(\frac{R_0 L}{T} \right)^{\frac{2\alpha-2}{3\alpha-2}}.$$

518 Since the gradients received from $l(\theta)$ and $g(\theta)$ are negatively correlated at any instant during the optimization process, the above expression simplifies to

$$\frac{1}{T} \sum_{k=1}^T \mathbb{E} \left[\|\nabla l(\theta_{k-1})\|^2 + \|\nabla g(\theta_{k-1})\|^2 + 2\|\nabla l(\theta_{k-1})\| \|\nabla g(\theta_{k-1})\| \right] \leq 4G^{\frac{2\alpha}{3\alpha-2}} \left(\frac{R_0 L}{T} \right)^{\frac{2\alpha-2}{3\alpha-2}}.$$

520 Therefore,

$$\frac{1}{T} \sum_{k=1}^T \mathbb{E} \left[\|\nabla l(\theta_{k-1})\|^2 \right] + \frac{1}{T} \sum_{k=1}^T \mathbb{E} \left[\|\nabla g(\theta_{k-1})\|^2 \right] \leq 4G^{\frac{2\alpha}{3\alpha-2}} \left(\frac{R_0 L}{T} \right)^{\frac{2\alpha-2}{3\alpha-2}}.$$

521 Upon simplification,

$$\frac{1}{T} \sum_{k=1}^T \mathbb{E} \left[\|\nabla l(\theta_{k-1})\|^2 \right] \leq \mathcal{O} \left(G^{\frac{2\alpha}{3\alpha-2}} \left(\frac{R_0 L}{T} \right)^{\frac{2\alpha-2}{3\alpha-2}} - \frac{1}{T} \sum_{k=1}^T \mathbb{E} \left[\|\nabla g(\theta_{k-1})\|^2 \right] \right)$$

522 which finishes the proof. \square