

ADVERSARIALLY TRAINED NEURAL POLICIES IN THE FOURIER DOMAIN

Ezgi Korkmaz

Electrical Engineering and Computer Science School,
KTH Royal Institute of Technology,
Stockholm, Sweden
ezgikorkmazk@gmail.com

ABSTRACT

Reinforcement learning policies based on deep neural networks are vulnerable to imperceptible adversarial perturbations to their inputs, in much the same way as neural network image classifiers. Recent work has proposed several methods for adversarial training for deep reinforcement learning agents to improve robustness to adversarial perturbations. In this paper, we study the effects of adversarial training on the neural policy learned by the agent. In particular, we compare the Fourier spectrum of minimal perturbations computed for both adversarially trained and vanilla trained neural policies. Via experiments in the OpenAI Atari environments we show that minimal perturbations computed for adversarially trained policies are more focused on lower frequencies in the Fourier domain, indicating a higher sensitivity of these policies to low frequency perturbations. We believe our results can be an initial step towards understanding the relationship between adversarial training and different notions of robustness for neural policies.

1 INTRODUCTION

Deep neural networks (DNNs) have been initially employed in deep reinforcement learning by Mnih et al. (2015) to approximate the state-action value function for large action size or state size MDPs. With this initial success deep reinforcement learning became an emerging subfield with many applications such as robotics Kalashnikov et al. (2018), financial trading Noonan (2017) and medical Daochang & Jiang (2018); Huan-Hsin et al. (2017).

While the successes of DNNs grew, a line of research focused on their reliability and robustness. Initially, Szegedy et al. (2014) demonstrated that it is possible to fool image classifiers by adding visually imperceptible perturbations to neural network inputs. Follow up work by Goodfellow et al. (2015) showed that these perturbations demonstrate that neural networks are learning approximately linear functions. Several studies focused on overcoming this susceptibility towards specifically crafted visually imperceptible perturbations, and proposed training neural networks against these perturbations Madry et al. (2018). More recent study showed that much of the work focusing on adversarial training suffers from something called the obfuscated gradient problem Athalye et al. (2018). Further, while there is a significant amount of study focusing on adversarial training, several other works suggest that adversarial perturbations may be inevitable Dohmatob (2019); Mahlouljifar et al. (2019); Gourdeau et al. (2019).

For these reasons in this paper we focus on adversarially trained neural policies and make the following contributions:

- We investigate the frequency domain of the minimal perturbations produced by the Carlini & Wagner (2017) formulation for both adversarially trained models and vanilla trained models.
- We conduct multiple experiments in the OpenAI Atari Baselines.
- We show that the perturbations produced from adversarially trained models are suppressed in high frequencies and more concentrated in lower frequencies in the Fourier domain compared vanilla trained neural policies.

2 BACKGROUND

2.1 ADVERSARIAL EXAMPLES

Manipulating the output of neural networks by introducing imperceptible perturbations has been introduced by Szegedy et al. (2014) based on a box constrained optimization method. While this proposed method is computationally expensive, Goodfellow et al. (2015) proposed a faster and simpler method based gradients in a nearby ϵ -ball,

$$x_{\text{adv}} = x + \epsilon \cdot \frac{\nabla_x J(x, y)}{\|\nabla_x J(x, y)\|_p}, \quad (1)$$

where x represents the input, y represents the labels, and $J(x, y)$ represents the cost function used to train the network. Kurakin et al. (2016) propose to search iteratively inside this ϵ -ball with the fast gradient sign method (FGSM) proposed by Goodfellow et al. (2015).

$$x_{\text{adv}}^0 = x, \quad (2)$$

$$x_{\text{adv}}^{N+1} = \text{clip}_{\epsilon}(x_{\text{adv}}^N + \alpha \text{sign}(\nabla_x J(x_{\text{adv}}^N, y))) \quad (3)$$

Madry et al. (2018) explained adversarial training in terms of the theory of robust optimization, and referred to this class of iterative methods used to produce adversarial perturbations as projected gradient descent (PGD). Carlini & Wagner (2017) formulated this problem in a more targeted way, and proposed a method based on distance minimization for a given label in image classification. For deep reinforcement learning this formulation is based on distance minimization for a given a target action which is not equal to the best action decided by the trained policy,

$$\begin{aligned} \min_{s_{\text{adv}} \in D_{\epsilon, p}(s)} \quad & \|s_{\text{adv}} - s\|_p \\ \text{subject to} \quad & a^*(s) \neq a^*(s_{\text{adv}}), \end{aligned}$$

Note that $a^*(s)$ denotes the action taken in state s , and $a^*(s_{\text{adv}})$ denotes the action taken in state s_{adv} . Athalye et al. (2018) showed that the Carlini & Wagner (2017) adversarial formulation can break several proposed defenses. For this reason, in this paper we will focus on perturbations produced by the Carlini & Wagner (2017) formulation.

2.2 PERTURBATION FORMULATIONS AND ADVERSARIAL TRAINING

The first studies on adversarial examples for deep reinforcement learning were Huang et al. (2017) and Kos & Song (2017). Both papers focused on using FGSM perturbations applied to the state observations in order to degrade the performance of trained neural policies. In the other direction, Mandlekar et al. (2017) introduced a form of adversarial training for deep neural policies by modifying the inputs at training time with perturbations based on the gradient of the cost function. Since reinforcement learning involves interaction with the environment there has also been effort to incorporate this interaction into adversarial training. Both Pinto et al. (2017) and Gleave et al. (2020) use game-theoretic models of the interaction between the deep reinforcement learning agent and then design training algorithms based on playing against the adversary in a zero-sum game in order to improve robustness. Recent work by Zhang et al. (2020) involves modifying the original MDP to create what they term a state-adversarial MDP. The authors then design theoretically-motivated adversarial training algorithms for deep neural policies based on training within the state-adversarial MDP.

3 EXPERIMENTAL SETUP

In our experiments we use OpenAI Brockman et al. (2016) Atari baselines Bellemare et al. (2013). Our models are trained with DDQN Wang et al. (2016) and SA-DDQN Zhang et al. (2020). We test trained policies averaged over 10 episodes. Note that SA-DDQN is certified against ℓ_{∞} -norm

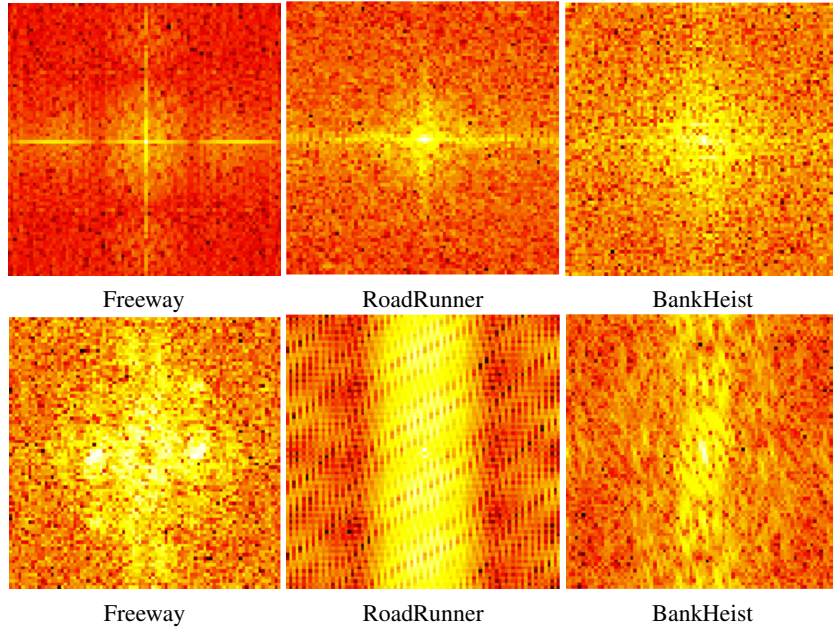


Figure 1: Fourier spectrum of the perturbations computed via Carlini & Wagner (2017) for state-of-the-art adversarially trained models and vanilla trained models. First Row: Adversarially trained. Second Row: Vanilla trained.

bounded perturbations at $1/255$. Therefore, we also bound the perturbation by this threshold and find the perturbations with ℓ_∞ -norm lower than this value.

4 NEURAL POLICY PERTURBATIONS IN FOURIER DOMAIN

In this paper we conduct an investigation on the frequency domain of the perturbations computed from vanilla trained agents and adversarially trained agents. In particular we conduct the following experiment: For a given state compute a minimum length perturbation which causes the agent to change its optimal learned action. If the minimal perturbation has norm smaller than a given threshold, compute the Fourier transform of the perturbation and record this data. By comparing the results of this experiment for adversarially trained versus vanilla trained agents, we can understand the affects of adversarial training on the directions to which the neural policy is sensitive. We now describe the details of the experimental setup.

For the adversarially trained agents, we focus on the state-of-the-art adversarial training algorithm SA-DQN proposed by Zhang et al. (2020). In this study the authors model the interaction between the neural policy and the introduced perturbations as a state-adversarial modified Markov Decision Process (MDP). The authors claim that the agents trained in the SA-MDP with the proposed algorithm SA-DQN are more robust to adversarial perturbations and natural noise introduced to the agent’s perception system. Furthermore, the authors demonstrate the robustness of SA-DQN against perturbations produced by the PGD attack proposed by Madry et al. (2018).

To compute the minimum length perturbations for our experiments we use the Carlini & Wagner (2017) formulation, which searches for a perturbation of minimum length that causes the agent to change its optimal action. For each state we use the Carlini & Wagner (2017) method to produce a perturbation $\eta = s_{\text{adv}} - s$. Note that the certified bound for Zhang et al. (2020) is $\frac{1}{255}$; therefore, we ensure the perturbation produced has norm $\|\eta\|_\infty < \frac{1}{255}$. We then compute the Fourier transform of η and add it to our dataset. Note that it is possible to break the certified defense via Carlini & Wagner (2017) perturbations. The “certified defense” proposed by Zhang et al. (2020) only holds for a fraction of states.

In Figure 1 we visualize the Fourier transform of a minimal perturbation for both vanilla trained and adversarially trained agents in RoadRunner, BankHeist and Freeway. The center of each image corresponds to the Fourier basis function where both spatial frequencies are zero, and the magnitude of the spatial frequencies increases as one moves outward from the center. From these visualizations it is clear that the perturbations for the adversarially trained models have their Fourier transform concentrated at lower frequencies than those of the vanilla trained models. To make this claim formal, for each number f we compute the total energy $\mathcal{E}(f)$ of the Fourier transform for basis functions whose maximum spatial frequency is equal to f . In Figure 2 we plot the average of $\mathcal{E}(f)$ over the minimal perturbations computed in our experiments. We find that the minimal perturbations computed for adversarially trained neural policies are indeed shifted towards lower frequencies when compared to those for vanilla trained neural policies. This shift in the frequency domain of the computed perturbations implies that adversarially trained neural policies may be more sensitivity towards low frequency perturbations.

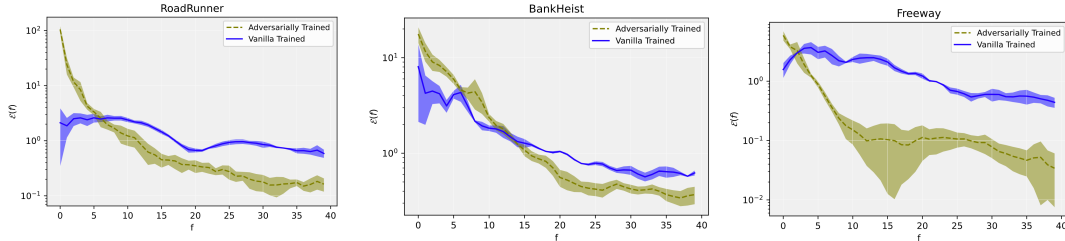


Figure 2: Perturbations computed via Carlini & Wagner (2017) for adversarially trained models and vanilla trained models in Fourier domain.

5 CONCLUSION

In this paper we focused on perturbations computed from state-of-the-art adversarially trained neural policies. We conducted several experiments and compared adversarially trained models to vanilla trained models. We investigated the frequency domain of the perturbations computed from state-of-the-art adversarially trained neural policies and vanilla trained neural policies. We found that the perturbations computed from adversarially trained models are more concentrated in lower frequencies compared to the vanilla trained neural policies. We believe this initial work outlines the vulnerabilities of adversarially trained neural policies and can be initial step towards building robust and reliable deep reinforcement learning agents.

REFERENCES

- Anish Athalye, Nicholas Carlini, and David A. Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In Jennifer G. Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 274–283. PMLR, 2018.
- Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael. Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research.*, pp. 253–279, 2013.
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv:1606.01540*, 2016.
- Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 39–57, 2017.
- Liu Daochang and Tingting. Jiang. Deep reinforcement learning for surgical gesture segmentation and classification. In *International conference on medical image computing and computer-assisted intervention.*, pp. 247–255. Springer, Cham, 2018.

- Elvis Dohmatob. Generalized no free lunch theorem for adversarial robustness. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 1646–1654. PMLR, 09–15 Jun 2019.
- Adam Gleave, Michael Dennis, Cody Wild, Kant Neel, Sergey Levine, and Stuart Russell. Adversarial policies: Attacking deep reinforcement learning. *International Conference on Learning Representations ICLR*, 2020.
- Ian Goodfellow, Jonathan Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *International Conference on Learning Representations*, 2015.
- Pascale Gourdeau, Varun Kanade, Marta Kwiatkowska, and James Worrell. On the hardness of robust classification. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 7444–7453, 2019.
- Sandy Huang, Nicholas Papernot, Yan Goodfellow, Ian an Duan, and Pieter Abbeel. Adversarial attacks on neural network policies. *Workshop Track of the 5th International Conference on Learning Representations*, 2017.
- Tseng Huan-Hsin, Sunan Cui, Yi Luo, Jen-Tzung Chien, Randall K. Ten Haken, and Issam El. Naqa. Deep reinforcement learning for automated radiation adaptation in lung cancer. *Medical physics* 44, pp. 6690–6705, 2017.
- Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, and Sergey. Levine. Qt-opt: Scalable deep reinforcement learning for vision-based robotic manipulation. *arXiv preprint arXiv:1806.10293*, 2018.
- Jernej Kos and Dawn Song. Delving into adversarial attacks on deep policies. *International Conference on Learning Representations*, 2017.
- Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial examples in the physical world. *arXiv preprint arXiv:1607.02533*, 2016.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
- Saeed Mahloujifar, Xiao Zhang, Mohammad Mahmoody, and David Evans. Empirically measuring concentration: Fundamental limits on intrinsic robustness. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 5210–5221, 2019.
- Ajay Mandlekar, Yuke Zhu, Animesh Garg, Li Fei-Fei, and Silvio Savarese. Adversarially robust policy learning: Active construction of physically-plausible perturbations. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 3932–3939, 2017.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, arc G Bellemare, Alex Graves, Martin Riedmiller, Andreas Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Antonoglou, Helen King, Dharmashan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518: 529–533, 2015.
- Laura Noonan. Jpmorgan develops robot to execute trades. *Financial Times*, pp. 1928–1937, July 2017.

- Lerrel Pinto, James Davidson, Rahul Sukthankar, and Abhinav Gupta. Robust adversarial reinforcement learning. *International Conference on Learning Representations ICLR*, 2017.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dimutru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *In Proceedings of the International Conference on Learning Representations (ICLR)*, 2014.
- Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Van Hasselt, Marc Lanctot, and Nando. De Freitas. Dueling network architectures for deep reinforcement learning. *International Conference on Machine Learning ICML*, pp. 1995–2003, 2016.
- Huan Zhang, Hongge Chen, Chaowei Xiao, Bo Li, Mingyan Liu, Duane S. Boning, and Cho-Jui Hsieh. Robust deep reinforcement learning against adversarial perturbations on state observations. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.