



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

*Distributed
Computing*



Rethinking Adversarial Examples

Master's Thesis

Yahya Jabary

yjabary@ethz.ch

Computer Engineering and Networks Laboratory
ETH Zürich

Supervisors:

Prof. Dr. Roger Wattenhofer

Prof. Dr. Shahram Dustdar

December 3, 2024

Acknowledgements

...

Abstract

...

Keywords: Robustness, Alignment, Security in Machine Learning

Contents

| | |
|---|------------|
| Acknowledgements | i |
| Abstract | ii |
| 1 Introduction | 1 |
| 1.1 Definition | 1 |
| 1.2 Mental Models | 3 |
| 1.3 Counterintuitive Properties | 4 |
| 2 Methodology | 5 |
| 3 Results | 6 |
| 4 Stuff | 7 |
| 4.1 First Section Title | 8 |
| 4.1.1 First Subsection Title | 8 |
| Bibliography | 9 |
| A First Appendix Chapter Title | A-1 |

Introduction

We have two goals in writing this document. One: fulfilling the requirements for a master’s degree by presenting and extending our original research [1] in thesis form. Two: offering a fresh and cohesive perspective on the rapidly evolving and, in our view, really exciting field of adversarial machine learning. To our knowledge, this is the first attempt to introduce this topic as a gateway for a broader audience, and we hope it will be valuable to those interested. No prior knowledge is assumed, but a basic understanding of linear algebra will be helpful.

1.1 Definition

Adversarial Examples are closely related to the concept of perturbation methods.

The origin of perturbations can be traced back to the early days of computational geometry by Seidel et al. in 1998 [2]. Perturbation techniques in computational geometry address a fundamental challenge: handling “degeneracies” in geometric algorithms. These are special cases that occur when geometric primitives align in ways that break the general position assumptions the algorithms rely on.

Example: Perturbation scheme for a linear classifier.

Consider a simple case of determining whether a point lies above or below a line [3]. While this classification appears straightforward, numerical issues arise when the point lies exactly on the line. Such degeneracies can cascade into algorithm failures or inconsistent results. The elegant solution is to imagine slightly moving (perturbing) the geometric objects to eliminate these special cases. Formally, we can express symbolic perturbation as $p_\varepsilon(x) = x + \varepsilon \cdot \delta(x)$ where x is the original input, ε is an infinitesimally small positive number the exact value of which is unimportant, and $\delta(x)$ is the perturbation function to break degeneracies.

A perturbation scheme should be (1) consistent, meaning that the same input always produces the same perturbed output (2) infinitesimal, such that perturbations are small enough not to affect non-degenerate cases and (3) effective, in breaking all possible degeneracies.

One powerful perturbation approach is Simulation of Simplicity (SoS) [4, 5, 6, 7, 8, 9]. SoS systematically perturbs input coordinates using powers of a symbolic infinitesimal. For a point $p_i = (x_i, y_i)$, the perturbed coordinates become:

$$(\tilde{x}_i, \tilde{y}_i) = (x_i + \varepsilon^{2i}, y_i + \varepsilon^{2i+1}) = p_i + \varepsilon^{2i} \cdot (1, \varepsilon)$$

This scheme ensures that no two perturbed points share any coordinate, effectively eliminating collinearity and other degeneracies.

The beauty of perturbation methods lies in their ability to handle degeneracies without explicitly detecting them, making geometric algorithms both simpler and more robust.

Adversarial examples follow the same principles as perturbation methods, but with the opposite goal of misleading machine learning models. This concept was first introduced by Szegedy et al. in 2014 [10]. Intuitively adversarial examples can be understood as seeking the closest point in the input space that lies on the “wrong side” of a decision boundary.

Example: Fast Gradient Sign Method (FGSM)

FGSM is one of the earliest and most widely recognized adversarial attack techniques, introduced by Goodfellow et al. [11]. The perturbation is controlled by a parameter $\varepsilon > 0$, which determines the magnitude of the change based on the direction of change for each pixel or feature in the input x . The model’s loss function is denoted by J , θ represents the model’s parameters, and y is the true target label.

FGSM works by calculating the gradient of the loss function with respect to the input, $\nabla_x J(\theta, x, y)$, and then adjusting the input in the direction of this gradient. The sign of the gradient, $\text{sign}(\nabla_x J(\theta, x, y))$, is used to ensure that the perturbation is small, while the ℓ_∞ -norm constraint ensures that the change to the input remains imperceptible to human observers [12].

The process for generating an adversarial example with FGSM can be expressed as:

$$x' = x + \varepsilon \cdot \text{sign}(\nabla_x J(\theta, x, y))$$

In the targeted version of FGSM, the perturbation is designed to minimize the loss with respect to a chosen target class, making the model predict the wrong class deliberately. In the untargeted version, the perturbation is designed to increase the loss for the correct class, causing the model to misclassify the input. In both cases, the perturbation is added to the original input, x , to create the adversarial example.

1.2 Mental Models

In 2013 Szegedy et al. [13] discovered that deep neural networks are vulnerable to adversarial examples – inputs with imperceptible perturbations that cause misclassification, revealing blind spots where the models’ decision boundaries are brittle, despite appearing to generalize well on normal inputs. Since then we have made a lot of progress in both finding attack vectors [11] [14] [15] and adversarially training models to be more robust [16] [14] [17] to these attacks.

However, despite the progress, we still don’t fully understand:

- What they are, why they exist and how they work.
- How to defend against them without sacrificing accuracy (robustness vs accuracy trade-off).
- Why they transfer between different models, datasets, architectures, training procedures and even to the human visual system [18].

There have been many attempts at explaining adversarial examples, each with limitations and assumptions – some complementary, some contradictory ¹.

For example, the *dimpled manifold hypothesis* [20] suggests that the decision boundary of deep neural networks is close to the data manifold, making it easy to find adversarial examples, while the *non-robust features hypothesis* [19] suggests that models exploit non-robust features that are imperceptible to humans, leading to a vulnerability against small perturbations

The *dimpled manifold hypothesis* is a particularly controversial one, as it was criticized by Yannik Kilcher [21] in 2021, who also provided a counterexample in less than 100 lines of code [22]. Despite the lack of generalizability, a master’s student from the University of Vienna, Lukas Karner, successfully verified and replicated all experiments detailed in the dimples paper [23] in 2023. Lukas allowed me to use his results in my work. He mentioned that there is currently

¹For a comprehensive overview of the hypotheses, see the Addendum of Ilyas et al. [19]

no paper or thesis based on his work and that he would be happy to see his results being used in a meaningful way.

This leaves us with the possibility that the experiments carried out are correct in themselves, but that the chain of reasoning is inconclusive and therefore doesn't generalize. Investigating this further would require more rigor by formalizing falsifiable hypotheses based on the paper and conducting experiments to test them.

Another aspect of the *dimpled manifold hypothesis* worth exploring is the idea of projecting the perturbations on the data manifold before applying them to the input space. Reducing the dimensionality of the perturbations or compressing them makes them visible to the human eye and interpretable as demonstrated by Karner [23].

1.3 Counterintuitive Properties

Methodology

Results

CHAPTER 4

Stuff



Figure 4.1: This is an example graphic.

“The stuff is what the stuff is, brother. Okay. We don’t ask questions about the weights. We just wake up, we go to work, we use the weights, we go back home. Okay. If we change the weights, the predictions would be different and less good, probably... depending on the weather... so we don’t ask about the weights.”

— James Mickens, *USENIX Security 18* [24]

4.1 First Section Title

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua.

4.1.1 First Subsection Title

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua.

Theorem 4.1 (First Theorem). *This is our first theorem.*

Proof. And this is the proof of the first theorem with a complicated formula and a reference to Theorem 4.1. Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua.

$$\frac{d}{dx} \arctan(\sin(x^2)) = -2 \cdot \frac{\cos(x^2)x}{-2 + (\cos(x^2))^2} \quad (4.1)$$

□

And here we cite some external documents [25, 26]. An example of an included graphic can be found in Figure 4.1. Note that in L^AT_EX, “quotes” do not use the usual double quote characters.

Bibliography

- [1] Y. Jabary, A. Plesner, T. Kuzhagaliyev, and R. Wattenhofer, “Seeing through the mask: Rethinking adversarial examples for captchas,” *arXiv preprint arXiv:2409.05558*, 2024.
- [2] R. Seidel, “The nature and meaning of perturbations in geometric computing,” *Discrete & Computational Geometry*, vol. 19, pp. 1–17, 1998.
- [3] M. De Berg, *Computational geometry: algorithms and applications*. Springer Science & Business Media, 2000.
- [4] W. R. Franklin and S. V. G. de Magalhães, “Implementing simulation of simplicity for geometric degeneracies,” *arXiv preprint arXiv:2212.08226*, 2022.
- [5] Edelsbrunner, Letscher, and Zomorodian, “Topological persistence and simplification,” *Discrete & computational geometry*, vol. 28, pp. 511–533, 2002.
- [6] H. Edelsbrunner and D. Guoy, “Sink-insertion for mesh improvement,” in *Proceedings of the seventeenth annual symposium on Computational geometry*, 2001, pp. 115–123.
- [7] H. Edelsbrunner and E. P. Mücke, “Simulation of simplicity: a technique to cope with degenerate cases in geometric algorithms,” *ACM Transactions on Graphics (tog)*, vol. 9, no. 1, pp. 66–104, 1990.
- [8] B. Lévy, “Robustness and efficiency of geometric programs: The predicate construction kit (pck),” *Computer-Aided Design*, vol. 72, pp. 3–12, 2016.
- [9] P. Schorn, “An axiomatic approach to robust geometric programs,” *Journal of symbolic computation*, vol. 16, no. 2, pp. 155–165, 1993.
- [10] C. Szegedy, “Intriguing properties of neural networks,” *arXiv preprint arXiv:1312.6199*, 2013.
- [11] I. J. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and harnessing adversarial examples,” *arXiv preprint arXiv:1412.6572*, 2014.
- [12] J. Zhang and C. Li, “Adversarial examples: Opportunities and challenges,” *IEEE transactions on neural networks and learning systems*, vol. 31, no. 7, pp. 2578–2593, 2019.
- [13] E. D. Cubuk, B. Zoph, S. S. Schoenholz, and Q. V. Le, “Intriguing properties of adversarial examples,” *arXiv preprint arXiv:1711.02846*, 2017.

- [14] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, “Towards deep learning models resistant to adversarial attacks,” *arXiv preprint arXiv:1706.06083*, 2017.
- [15] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, and A. Swami, “The limitations of deep learning in adversarial settings,” in *2016 IEEE European symposium on security and privacy (EuroS&P)*. IEEE, 2016, pp. 372–387.
- [16] A. Shafahi, M. Najibi, M. A. Ghiasi, Z. Xu, J. Dickerson, C. Studer, L. S. Davis, G. Taylor, and T. Goldstein, “Adversarial training for free!” *Advances in neural information processing systems*, vol. 32, 2019.
- [17] N. Papernot, P. McDaniel, X. Wu, S. Jha, and A. Swami, “Distillation as a defense to adversarial perturbations against deep neural networks,” in *2016 IEEE symposium on security and privacy (SP)*. IEEE, 2016, pp. 582–597.
- [18] G. Elsayed, S. Shankar, B. Cheung, N. Papernot, A. Kurakin, I. Goodfellow, and J. Sohl-Dickstein, “Adversarial examples that fool both computer vision and time-limited humans,” *Advances in neural information processing systems*, vol. 31, 2018.
- [19] A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran, and A. Madry, “Adversarial examples are not bugs, they are features,” *Advances in neural information processing systems*, vol. 32, 2019.
- [20] A. Shamir, O. Melamed, and O. BenShmuel, “The dimpled manifold model of adversarial examples in machine learning,” *arXiv preprint arXiv:2106.10151*, 2021.
- [21] Y. Kilcher, “The Dimpled Manifold Model of Adversarial Examples in Machine Learning (Research Paper Explained),” https://www.youtube.com/watch?v=k_hUdZJNzkU/, 2021, [Online; accessed 17-July-2024].
- [22] Y. Kilcher, “dimple test,” <https://gist.github.com/yk/de8d987c4eb6a39b6d9c08f0744b1f64/>, 2021, [Online; accessed 17-July-2024].
- [23] L. Karner, “The Dimpled Manifold Revisited,” <https://github.com/LukasKarner/dimpled-manifolds/>, 2023, [Online; accessed 17-July-2024].
- [24] J. Mickens, “Q: Why do keynote speakers keep suggesting that improving security is possible? a: Because keynote speakers make bad life decisions and are poor role models,” in *27th USENIX Security Symposium (USENIX Security 18)*. Baltimore, MD: USENIX Association, Aug. 2018. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity18/presentation/mickens>

- [25] A. One and A. Two, “A theoretical work on computer science,” in *30th Symposium on Comparative Irrelevance, Somewhere, Some Country*, Jun. 1999.
- [26] A. One and A. Two, “A theoretical work on computer science,” in *30th Symposium on Comparative Irrelevance, Somewhere, Some Country*, Jun. 1999.

First Appendix Chapter Title
