



Original Article

Adversarial Machine Learning: Exploring Security Vulnerabilities in AI-Driven Systems

Neha Iyer

Blockchain Specialist, Accenture, Singapore

Abstract - Adversarial Machine Learning (AML) has emerged as a critical area of research in the field of artificial intelligence (AI) and cybersecurity. As AI-driven systems become increasingly integrated into various sectors, including finance, healthcare, and autonomous vehicles, the security of these systems is of paramount importance. AML explores the vulnerabilities of machine learning (ML) models to adversarial attacks, where malicious actors manipulate input data to deceive the models. This paper provides a comprehensive overview of AML, including its theoretical foundations, types of attacks, defense mechanisms, and real-world implications. We also discuss the challenges and future directions in this rapidly evolving field, emphasizing the need for robust and secure AI systems.

Keywords - Adversarial Machine Learning, Evasion Attacks, Poisoning Attacks, Model Robustness, Adversarial Training, Defensive Distillation, Federated Learning, Explainable AI (XAI), Deep Neural Networks (DNNs), Cybersecurity in AI

1. Introduction

1.1 Background and Motivation

1.1.1. The Growing Influence of AI and Machine Learning

The proliferation of machine learning (ML) and artificial intelligence (AI) has fundamentally transformed industries ranging from finance and healthcare to autonomous systems and cybersecurity. These technologies enable automation, enhance decision-making, and improve efficiency by learning patterns from vast amounts of data. AI-powered applications such as fraud detection systems, medical diagnosis tools, and self-driving cars are now becoming indispensable in daily life. However, as AI models continue to evolve and integrate deeper into critical infrastructure, they also become high-value targets for cyber threats.

1.1.2. Emerging Security Threats in AI Systems

Despite the remarkable progress in ML, these systems remain vulnerable to a range of adversarial attacks. Unlike traditional software vulnerabilities, where security loopholes stem from coding flaws, AI security risks often arise from the inherent design and training of ML models. Attackers exploit weaknesses in learning algorithms by manipulating input data, deceiving models into making incorrect predictions, or even degrading their performance over time. This growing field of threats has given rise to the domain of **Adversarial Machine Learning (AML)**, which aims to understand, analyze, and defend against these attacks.

1.1.3. Types of Adversarial Threats

Adversarial attacks on AI systems can take multiple forms, each posing distinct risks:

- **Evasion Attacks:** Attackers subtly modify input data at inference time to trick AI models into making wrong predictions (e.g., fooling an image classifier into misidentifying objects).
- **Poisoning Attacks:** Attackers manipulate training data to degrade the model's performance, embedding backdoors or biases that can be exploited in the future.
- **Model Extraction and Inversion Attacks:** Attackers attempt to reverse-engineer AI models, stealing their functionality or reconstructing sensitive training data.

1.1.4. Real-World Implications of Adversarial Attacks

The consequences of adversarial attacks are not just theoretical but have real-world implications:

- **Autonomous Vehicles:** Attackers can manipulate road signs using minor alterations, causing self-driving cars to misinterpret traffic signals.
- **Healthcare and Medical Diagnosis:** Adversarial perturbations in medical images can lead to incorrect diagnoses, jeopardizing patient safety.

- **Cybersecurity and Fraud Detection:** Attackers can generate adversarial examples to evade security systems, enabling fraudulent transactions or bypassing authentication mechanisms.

1.1.5. The Need for Robust and Secure AI Models

As AI-driven systems become increasingly embedded in critical applications, ensuring their security and reliability is of paramount importance. Developing robust AI models that can withstand adversarial attacks is an urgent challenge for researchers, engineers, and policymakers. Various defense mechanisms, including adversarial training, input preprocessing, and model robustness techniques, are actively being explored to mitigate these threats.

2. Theoretical Foundations of Adversarial Machine Learning

2.1 Machine Learning Basics

Machine learning (ML) is a branch of artificial intelligence (AI) that enables systems to learn patterns from data and make predictions or decisions without being explicitly programmed. It is widely used in applications such as image recognition, natural language processing, and fraud detection. ML models can be broadly classified into three main types:

2.1.1. Supervised Learning

In supervised learning, models are trained on labeled datasets, where each input is paired with the correct output. The goal is to learn a mapping function that generalizes well to unseen data. Examples include:

- Classification: Predicting discrete labels (e.g., spam detection, disease diagnosis).
- Regression: Predicting continuous values (e.g., stock price forecasting, temperature prediction).

2.1.2. Unsupervised Learning

Unsupervised learning involves training models on unlabeled data to uncover hidden structures or patterns. These methods are commonly used for:

- Clustering: Grouping similar data points (e.g., customer segmentation, anomaly detection).
- Dimensionality Reduction: Reducing data complexity while preserving essential features (e.g., PCA, autoencoders).

2.1.3. Reinforcement Learning

Reinforcement learning (RL) is a learning paradigm where an agent interacts with an environment, making sequential decisions to maximize cumulative rewards. It is widely applied in robotics, game playing (e.g., AlphaGo), and autonomous systems. The learning process is guided by:

- State: The current situation of the agent.
- Action: The possible moves or decisions the agent can take.
- Reward: A numerical signal indicating the success or failure of an action.

2.2 Adversarial Attacks

Adversarial attacks exploit vulnerabilities in ML models by introducing small perturbations in input data, leading to incorrect predictions. These attacks pose significant security risks, especially in critical applications like healthcare, finance, and autonomous driving. The two primary types of adversarial attacks are: In the Model Training phase, the figure illustrates two major types of attacks: poisoning attacks and model extraction attacks. Poisoning attacks involve injecting maliciously crafted data into the training set, leading to incorrect model learning. This is represented in the image by poisoned data being fed into the feature extraction pipeline, influencing the training algorithm and resulting in a biased or manipulated model. Meanwhile, model extraction attacks involve adversaries stealing the model itself by querying it repeatedly to approximate its decision boundaries. Additionally, model inversion attacks allow attackers to infer sensitive information about the training data by exploiting vulnerabilities in the learned model.

The Model Prediction phase focuses on adversarial attacks, where attackers craft perturbed input samples designed to mislead the trained model. The image depicts how an attacker manipulates the input data in real time, causing the model to

generate incorrect predictions. This is particularly dangerous in safety-critical applications, where even minor perturbations can have severe consequences, such as misleading an autonomous vehicle into misinterpreting road signs.

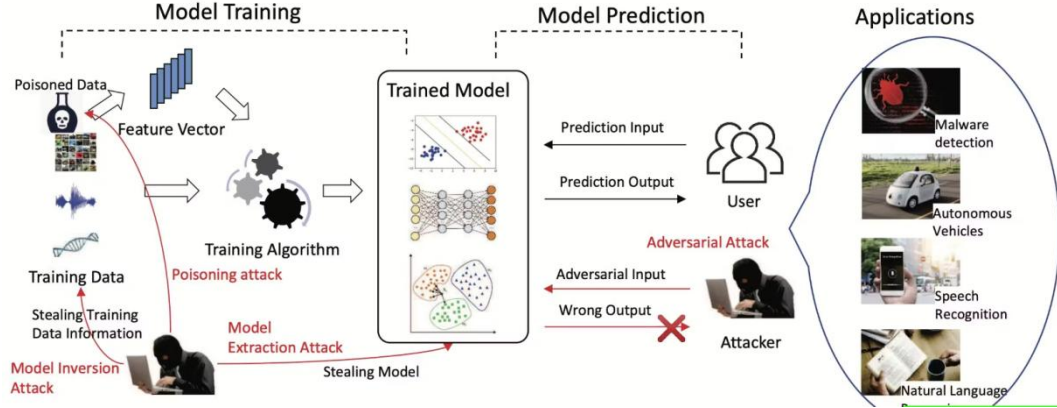


Fig 1: Model Training, Model Prediction, and Applications

The applications section highlights real-world scenarios where adversarial attacks pose significant risks. These include malware detection, where attackers can bypass security mechanisms by slightly modifying malware signatures, and autonomous vehicles, where adversarial perturbations can cause misinterpretations of road conditions. Other vulnerable areas include speech recognition and natural language processing, where attackers can manipulate input signals to deceive AI-driven assistants and chatbots.

2.2.1. Evasion Attacks

Evasion attacks occur at inference time, where an adversary modifies an input instance to fool the model without altering its perceived meaning. Examples include:

- Image Perturbation: Slightly modifying pixels in an image to cause misclassification.
- Malware Evasion: Crafting malicious software samples that bypass detection systems.

2.2.2. Poisoning Attacks

Poisoning attacks target the training phase, where adversaries inject manipulated data into the dataset to degrade model performance or insert backdoors. Common strategies include:

- Label Manipulation: Altering the labels of training samples to mislead the model.
- Feature Corruption: Introducing noise into critical features to obscure meaningful patterns.

2.3 Mathematical Formulation

Adversarial attacks aim to find a small perturbation δ that, when added to an input sample x , causes the model f to misclassify it. The attack can be formulated as an optimization problem:

$$\min_{\delta} L(f(x + \delta), y)$$

where:

x is the original input.

δ is the adversarial perturbation.

f is the ML model.

y is the true label of the input.

L is the loss function (e.g., cross-entropy loss).

The attacker aims to find the minimal perturbation δ that maximizes model misclassification while remaining imperceptible to humans.

2.4 Example: Fast Gradient Sign Method (FGSM)

One of the simplest and most effective adversarial attack techniques is the Fast Gradient Sign Method (FGSM). This method exploits the gradient of the loss function with respect to the input to craft adversarial examples. The FGSM attack is formulated as:

$$\delta = \epsilon \cdot \text{sign}(\nabla_x L(f(x), y))$$

where:

ϵ is the perturbation magnitude, controlling the attack's strength.

$(\nabla_x L(f(x), y))$ is the gradient of the loss function with respect to the input.

$\text{sign}(\cdot)$ extracts the sign of the gradient, ensuring directional modification.

The adversarial example is then generated as:

$$x_{adv} = x + \delta x$$

This method effectively forces the model towards incorrect predictions with minimal perturbations, making it one of the most widely studied attack strategies.

3. Types of Adversarial Attacks

3.1 Evasion Attacks

Evasion attacks occur at the test phase of a machine learning model, where an adversary carefully modifies input samples to deceive the model into making incorrect predictions. Unlike poisoning attacks, which target the training process, evasion attacks do not require direct access to the training dataset. Instead, attackers generate adversarial examples—perturbed inputs that appear unchanged to humans but cause significant errors in the model's predictions. These attacks are particularly concerning in critical AI applications such as facial recognition, malware detection, and autonomous driving, where incorrect decisions can have severe consequences.

3.1.1 White-Box Attacks

In white-box attacks, the attacker has complete knowledge of the target model, including its architecture, parameters, and training data distribution. This access allows them to exploit the model's vulnerabilities systematically. One of the most widely used white-box attack methods is the Fast Gradient Sign Method (FGSM), which calculates the gradient of the loss function with respect to the input and perturbs the input in the direction that maximally increases the loss. White-box attacks are particularly dangerous because they enable precise adversarial example generation, making defense mechanisms more challenging to implement.

3.1.2 Black-Box Attacks

Black-box attacks, in contrast, assume that the attacker has no prior knowledge of the model's internal structure. Instead, the attacker must rely on querying the model and analyzing its outputs to infer decision boundaries. A well-known black-box attack is the Transferability Attack, where the attacker first trains a substitute model that approximates the target model's behavior. Once the substitute model is trained, adversarial examples generated against it are transferred to the target model, often with high success rates. This property of adversarial transferability poses a significant threat, as attackers do not need direct access to proprietary or protected AI systems to launch an effective attack.

3.2 Poisoning Attacks

Poisoning attacks target the model during its training phase by introducing manipulated data into the dataset. The goal is to degrade the model's performance or implant hidden behaviors that only activate under specific conditions. Since machine learning models heavily depend on the quality and integrity of training data, poisoning attacks can have long-lasting and severe consequences. They are especially dangerous in scenarios where models are retrained periodically using crowdsourced or publicly available data, such as spam detection and recommendation systems.

3.2.1 Data Poisoning

Data poisoning involves injecting misleading or malicious samples into the training dataset to bias the learning process. Attackers may introduce mislabeled data points, causing the model to generalize incorrectly. For instance, in a spam detection

system, an attacker might add emails containing spam content but labeled as "not spam," tricking the model into classifying actual spam emails incorrectly. This type of attack can significantly reduce the reliability of AI-driven systems, leading to security vulnerabilities in applications like fraud detection, medical diagnosis, and financial forecasting.

3.2.2 Backdoor Attacks

Backdoor attacks introduce hidden triggers into the training data, which, when present, cause the model to behave maliciously while otherwise functioning normally. An attacker may insert a specific pattern—such as a checkerboard overlay—into images during training. The model learns to associate this pattern with a particular label, allowing the attacker to manipulate predictions at will. For example, a backdoored facial recognition system might correctly identify users under normal conditions but misclassify an attacker's face when they wear glasses with a particular pattern. These attacks are particularly challenging to detect because the model performs well on standard validation datasets, making them an insidious security risk.

3.3 Example: Carlini and Wagner Attack

One of the most advanced and effective adversarial attack methods is the Carlini and Wagner (C&W) attack, which uses an optimization-based approach to craft adversarial examples with minimal perturbations. Unlike simpler attacks such as FGSM, which apply fixed-magnitude perturbations, the C&W attack formulates adversarial generation as an optimization problem. It seeks to find the smallest possible modification to the input that causes the model to misclassify it while maintaining human perceptual similarity. The Carlini and Wagner (C&W) attack is a more sophisticated method for generating adversarial examples. It uses an optimization-based approach to find the smallest perturbation that causes the model to misclassify the input. The C&W attack is defined as:

$$\min_{\delta} L(\delta) + c \cdot L(f(x + \delta), y)$$

- $\|\delta\|$ is the norm of the perturbation.
- (c) is a hyperparameter that balances the trade-off between the perturbation size and the loss.

4. Defense Mechanisms

As adversarial attacks on machine learning models continue to evolve, researchers and practitioners have developed various defense mechanisms to mitigate their impact. These defenses aim to enhance the robustness of AI models by either modifying their training process, preprocessing inputs to remove adversarial perturbations, or designing architectures that are inherently resistant to adversarial manipulations. Despite these efforts, no single defense mechanism is foolproof, as adversaries continuously develop new techniques to bypass them. Nonetheless, combining multiple defense strategies can significantly improve the resilience of AI systems.

4.1 Adversarial Training

Adversarial training is one of the most widely used techniques for improving model robustness against adversarial attacks. This method involves augmenting the training dataset with adversarial examples, forcing the model to learn representations that are more resilient to perturbations. By exposing the model to adversarially perturbed inputs during training, it learns to generalize better and become less susceptible to small, targeted modifications. For instance, adversarial examples generated using the Fast Gradient Sign Method (FGSM) or the Carlini and Wagner (C&W) attack can be included in the training process. The model is then trained to correctly classify both the original and adversarially modified inputs. While adversarial training improves robustness, it comes at the cost of increased computational overhead and may not always generalize well against unseen attack strategies.

4.2 Input Transformation

Another effective defense strategy involves preprocessing the input data to remove adversarial perturbations before feeding them into the model. Input transformation techniques modify or filter the input in ways that diminish the effectiveness of adversarial modifications. Some commonly used transformation methods include:

- **Noise Injection:** Adding small amounts of random noise to the input data can disrupt the adversarial perturbations while maintaining the essential features needed for classification. However, excessive noise can degrade model performance on legitimate inputs.
- **Smoothing Filters:** Applying Gaussian smoothing or median filtering to input images can blur adversarial perturbations, reducing their impact. These techniques work particularly well against pixel-based perturbations but may struggle against more sophisticated attacks.

While input transformations provide a lightweight and easily implementable defense, adversaries can often develop countermeasures that bypass them by generating perturbations that survive the transformations.

4.3 Model Robustness

Improving the robustness of the model itself is another approach to mitigating adversarial attacks. Several strategies can enhance the inherent resilience of machine learning models:

- **Using deeper architectures:** More complex models with additional layers and non-linearities can make it harder for adversaries to find effective perturbations. However, this comes with increased computational cost.
- **Ensemble Learning:** Combining multiple models to form an ensemble can increase robustness by making it more difficult for a single adversarial attack to deceive all models simultaneously. Ensemble-based defenses aggregate predictions from multiple independent classifiers, reducing the likelihood of misclassification.
- **Regularization Techniques:** Implementing strategies like weight decay, dropout, and batch normalization can improve generalization and make the model less sensitive to small perturbations in the input space.

Although these techniques strengthen model resilience, no model architecture is completely immune to adversarial attacks, emphasizing the need for ongoing research and hybrid defense approaches.

4.4 Example: Defensive Distillation

One of the more innovative defense mechanisms against adversarial attacks is Defensive Distillation. This technique involves training a neural network in two stages to make its decision boundaries smoother, thereby reducing its vulnerability to adversarial perturbations. In the first stage, a teacher model is trained using the original dataset. Instead of training a student model directly on the hard labels (e.g., categorical class labels), the student model is trained on the soft probabilities output by the teacher model. These probability distributions provide a more gradual learning signal, helping the student model learn more generalized and robust decision boundaries.

Mathematically, if the teacher model outputs a probability distribution $p(y|x)$ over classes instead of a hard label, the student model is trained using the same dataset but with the smoothed probability distribution. This reduces the model's sensitivity to small perturbations, making it harder for adversaries to craft effective adversarial examples.

Despite its effectiveness, defensive distillation is not a perfect solution, as later research has shown that strong attacks can still bypass it. Nevertheless, it remains a valuable component in the broader landscape of adversarial defenses, particularly when combined with other strategies such as adversarial training and input transformations.

5. Real-World Implications

The growing reliance on machine learning and AI across industries means that adversarial machine learning (AML) attacks can have significant real-world consequences. From cybersecurity threats to safety-critical systems, the vulnerability of AI models to adversarial manipulation raises concerns about their reliability, security, and ethical implications. As AI systems become more deeply integrated into critical applications, adversarial attacks pose risks that extend beyond theoretical research into practical, high-stakes domains. One of the most concerning implications of adversarial attacks is in autonomous vehicles and transportation systems. Self-driving cars rely on AI models for object detection, lane following, and decision-making. However, researchers have demonstrated that small, strategically placed perturbations on road signs or lane markings can mislead these models, causing potentially catastrophic accidents. An attacker could, for example, alter a stop sign with minor pixel modifications, leading a vehicle to misinterpret it as a speed limit sign, resulting in dangerous outcomes.

Another major area of concern is cybersecurity and fraud detection. Many AI-powered cybersecurity systems use machine learning models to detect malware, phishing attempts, and fraudulent activities. Adversarial attacks can fool these detection systems, allowing malicious actors to bypass security measures. In the financial sector, fraud detection algorithms that rely on transaction patterns can be manipulated using adversarial techniques to evade detection, leading to significant financial losses for institutions and individuals alike. Beyond security risks, adversarial attacks also pose a serious challenge to healthcare and medical diagnostics. AI-driven medical imaging systems, such as those used for detecting cancer or analyzing radiology scans, are vulnerable to adversarial perturbations. If an attacker were to manipulate medical images, a benign tumor could be classified as malignant or vice versa, leading to incorrect diagnoses and potentially life-threatening consequences. Ensuring the robustness of AI in healthcare is critical, as these systems assist doctors in making high-stakes decisions that impact patient outcomes.

These real-world implications highlight the urgent need for stronger defenses against adversarial attacks. As AI continues to expand into critical industries, researchers, policymakers, and technology developers must work together to enhance the security and reliability of AI systems. The future of AI security depends on developing more resilient models, integrating real-time attack detection mechanisms, and enforcing regulatory standards that ensure AI technologies are safe and trustworthy.

6. Challenges and Future Directions

Adversarial Machine Learning (AML) presents several challenges that hinder the development of fully secure and robust AI models. One of the major challenges is scalability, as defending against adversarial attacks in large-scale AI systems requires significant computational resources. Deploying real-time defenses in high-frequency applications such as autonomous vehicles, financial fraud detection, and real-time surveillance is a daunting task due to the high processing power needed to detect and counter adversarial inputs. Additionally, the cost of implementing robust security measures can be prohibitive for many organizations, limiting widespread adoption.

Another fundamental challenge is the dynamic nature of adversarial attacks. Attackers continuously evolve their attack strategies, rendering existing defense mechanisms obsolete over time. As security measures improve, adversaries devise more sophisticated attack techniques that exploit new vulnerabilities. This constant cat-and-mouse game between attackers and defenders necessitates continuous research and updates to defensive strategies, making AML a highly dynamic and evolving field.

The complexity of modern AI models further exacerbates the difficulty of defending against adversarial attacks. Deep neural networks, which power many AI applications, operate in high-dimensional input spaces, making them susceptible to even small perturbations. The lack of interpretability in these models makes it challenging to understand why they are vulnerable to specific attacks and how to effectively mitigate these weaknesses. This has led to a growing demand for explainable AI (XAI) techniques that can provide insights into how adversarial attacks exploit model weaknesses.

6.2 Future Directions

To address these challenges, researchers are exploring several promising directions. One such approach is the development of hybrid defense mechanisms that combine multiple techniques to enhance model robustness. For example, adversarial training can be complemented with input transformation techniques and anomaly detection systems to create multi-layered defenses that are more resilient against a wide range of adversarial attacks.

Another important area of research is explainability in adversarial defense. By developing methods that provide deeper insights into why models are vulnerable, researchers can design more targeted defenses. Explainability techniques such as saliency maps and feature attribution methods can help identify weak points in models and suggest ways to reinforce them against adversarial manipulations. This research direction aims to bridge the gap between AML security and interpretability, ensuring that AI models are both robust and transparent.

Furthermore, collaborative defense strategies are gaining traction in the AML community. Sharing adversarial attack datasets, defensive techniques, and security findings across organizations and research institutions can accelerate progress in

securing AI systems. Open-source initiatives and collaborative platforms can facilitate the rapid development of improved AML defenses, creating a more unified approach to tackling adversarial threats.

6.3 Example: Federated Learning

A promising technique for mitigating adversarial attacks, especially data poisoning attacks, is federated learning. Unlike traditional centralized training methods, federated learning distributes the training process across multiple devices or servers without sharing raw data. This decentralized approach enhances privacy and security by preventing attackers from directly tampering with the centralized training dataset. Federated learning is particularly useful in sensitive applications such as healthcare and finance, where data security is paramount. However, challenges remain in ensuring that federated models are not vulnerable to adversarial attacks introduced at the local device level.

7. Conclusion

Adversarial Machine Learning (AML) is an essential and rapidly growing field that addresses the security vulnerabilities of AI-driven systems. As AI technologies become more prevalent across various industries, ensuring their robustness against adversarial attacks is critical. This paper has explored the theoretical foundations of AML, different types of adversarial attacks, defense mechanisms, and the real-world implications of these security threats.

Despite significant progress in adversarial defenses, challenges such as scalability, evolving attack strategies, and model complexity continue to pose serious threats to AI security. However, emerging research directions, including hybrid defense mechanisms, explainability techniques, and collaborative security efforts, offer promising solutions to strengthen AI systems. Federated learning, in particular, presents a novel approach to mitigating adversarial threats while preserving data privacy.

As AI continues to transform industries such as healthcare, finance, transportation, and cybersecurity, securing these systems from adversarial attacks must remain a top priority. Ongoing research and collaboration among academia, industry, and policymakers will be crucial in developing resilient AI models that can withstand adversarial threats. By fostering a proactive approach to AI security, we can build trustworthy and robust AI systems that benefit society while minimizing risks associated with adversarial manipulations.

References

- [1] Goodfellow, I. J., Shlens, J., & Szegedy, C. (2014). Explaining and Harnessing Adversarial Examples. arXiv preprint arXiv:1412.6572.
- [2] Carlini, N., & Wagner, D. (2017). Towards Evaluating the Robustness of Neural Networks. In 2017 IEEE Symposium on Security and Privacy (SP) (pp. 39-57).
- [3] Papernot, N., McDaniel, P., Jha, S., Fredrikson, M., Celik, Z. B., & Swami, A. (2016). The Limitations of Deep Learning in Adversarial Settings. In 2016 IEEE European Symposium on Security and Privacy (EuroS&P) (pp. 372-387).
- [4] Liu, Y., Chen, X., Liu, C., & Song, D. (2017). Delving into Transferable Adversarial Examples and Black-box Attacks. In 31st AAAI Conference on Artificial Intelligence (AAAI-17).
- [5] Papernot, N., McDaniel, P., Wu, X., Jha, S., & Swami, A. (2016). Distillation as a Defense to Adversarial Perturbations Against Deep Neural Networks. In 2016 IEEE Symposium on Security and Privacy (SP) (pp. 582-597).
- [6] TechTarget. Adversarial machine learning: Definition & overview. Retrieved from <https://www.techtarget.com/searchenterpriseai/definition/adversarial-machine-learning>
- [7] Title of the article. Journal of Engineering Research and Reports, Volume(Issue), Pages. Retrieved from <https://journaljerr.com/index.php/JERR/article/view/1413>
- [8] Palo Alto Networks. What are adversarial attacks on AI & machine learning? Retrieved from <https://www.paloaltonetworks.com/cyberpedia/what-are-adversarial-attacks-on-AI-Machine-Learning>
- [9] Lakera AI. Adversarial machine learning. Retrieved from <https://www.lakera.ai/blog/adversarial-machine-learning>
- [10] Coursera. Adversarial machine learning: Understanding and defending against attacks. Retrieved from <https://www.coursera.org/articles/adversarial-machine-learning>

- [11] Viso AI. Adversarial machine learning: An overview. Retrieved from <https://viso.ai/deep-learning/adversarial-machine-learning/>
- [12] Title of the paper. IEEE Transactions on Neural Networks and Learning Systems, Volume(Issue), Pages. Retrieved from <https://ieeexplore.ieee.org/iel7/9739/5451756/09887796.pdf>
- [13] DataCamp. Adversarial machine learning: A data science perspective. Retrieved from <https://www.datacamp.com/blog/adversarial-machine-learning>

Algorithms

Algorithm 1: Fast Gradient Sign Method (FGSM)

```
def fgsm_attack(model, x, y, epsilon):
    # Compute the gradient of the loss function with respect to
    # the input
    x_grad = compute_gradient(model, x, y)

    # Generate the adversarial example
    x_adv = x + epsilon * torch.sign(x_grad)

    return x_adv
```

Algorithm 3: Adversarial Training

```
def adversarial_training(model, train_data, epsilon,
num_epochs):
    for epoch in range(num_epochs):
        for x, y in train_data:
            # Generate adversarial examples
            x_adv = fgsm_attack(model, x, y, epsilon)

            # Compute the loss on the adversarial examples
            loss = model_loss(model(x_adv), y)

            # Backpropagate and update the model
            loss.backward()
            optimizer.step()
            optimizer.zero_grad()
```

Algorithm 2: Carlini and Wagner Attack

```
def cw_attack(model, x, y, c, learning_rate, max_iterations):
    # Initialize the perturbation
    delta = torch.zeros_like(x, requires_grad=True)

    for _ in range(max_iterations):
        # Compute the loss
        loss = torch.norm(delta, p=2) + c * model_loss(model(x + delta), y)

        # Compute the gradient
        loss.backward()

        # Update the perturbation
        delta.data -= learning_rate * delta.grad.data
        delta.grad.data.zero_()

    return x + delta
```