# SECURING THE DIGITAL FORTRESS: ADVERSARIAL MACHINE LEARNING CHALLENGES AND COUNTERMEASURES IN CYBERSECURITY

## Atharv Raotole[1], Ayush Deshmukh[2]

[1]*Department of Information Technology, Sardar Patel Institute of Technology, Mumbai, India*
[2]*Department of Computer Engineering, Dwarkadas J. Sanghvi College of Engg, Mumbai, India*

---------------------------------------------------------------------***---------------------------------------------------------------------

**Abstract** - *Machine learning models are increasingly integral to a wide array of applications, from autonomous vehicles to online security and voice recognition systems. However, the pervasive use of machine learning also exposes these systems to adversarial threats. In this paper we discuss the multifaceted domain of Adversarial Machine Learning in Cybersecurity, with a dual focus. Firstly, it explores how machine learning models can be manipulated, presenting an in-depth analysis of the diverse range of adversarial attacks, from evasion attacks to data poisoning. Second, it endeavors to develop and propose effective methods and strategies for defending against these adversarial attacks, thereby bolstering the resilience of machine learning systems in the context of cybersecurity. The impacts of successful adversarial attacks are experimented, emphasizing the far-reaching consequences on security, integrity, and trust in AI systems.*

***Key Words***: **Adversarial Machine Learning, Cybersecurity, Adversarial Attacks, Robustness, Defense Mechanisms, Privacy-Preserving Machine Learning.**

## 1. INTRODUCTION

In an era marked by the relentless march of technology, and with the ever-growing reliance on machine learning models across various domains, cybersecurity faces an ever-evolving and formidable adversary: adversarial machine learning. As artificial intelligence and machine learning systems seamlessly integrate into critical applications like fraud detection, autonomous vehicles, and malware identification, the vulnerability of these systems to adversarial attacks has become a substantial concern.

Adversarial machine learning involves the deliberate manipulation of machine learning models by malicious actors who aim to undermine their functionality and accuracy. These adversaries adeptly exploit the inherent vulnerabilities of AI and machine learning systems to craft inputs that deceive, mislead, or compromise the performance of these models. In the context of cybersecurity, adversarial attacks pose a severe threat, as the very systems designed to protect against threats can themselves become the targets of exploitation.

The ramifications of successful adversarial attacks are profound, casting shadows on the security, integrity, and trustworthiness of artificial intelligence systems. As our reliance on these systems continues to grow, the imperative of safeguarding them against adversarial threats becomes increasingly apparent. This paper aims to shine a light on the critical importance of addressing this issue, and as you delve into its pages, you will find an exploration of the multifaceted domain of adversarial machine learning in cybersecurity, from understanding the diversity of adversarial attacks to proposing effective methods and strategies for bolstering the resilience of machine learning systems in the context of security.

## 2. LITERATURE REVIEW

In the paper by Anthi et al. [1], adversarial machine learning techniques in the context of power system security are explored, with a focus on intrusion detection systems (IDS) and the introduction of the Jacobian-based Saliency Map Attack (JSMA). The study demonstrates the success of JSMA attacks in evading detection but is limited in its applicability to the power system domain and lacks comprehensive real-world evaluations of proposed defense mechanisms.

Alotaibi and Rassam [2] present a survey on adversarial machine learning attacks against intrusion detection systems (IDS), providing insights into various strategies and defense mechanisms. However, the paper lacks a comprehensive evaluation of the effectiveness of these defense strategies.

Rosenberg et al. [3] introduce adversarial attack methods and their computational costs, but the practical implications of these attacks and evaluated defense methods are not extensively discussed.

Tygar [4] introduces the field of adversarial machine learning, discussing vulnerabilities in machine learning algorithms and potential countermeasures, but empirical evidence supporting the effectiveness of these countermeasures is lacking.

In the systematic review by Martins et al. [5], the impact of adversarial attacks on intrusion and malware detection is explored, emphasizing the dependence on

factors such as data quality, classifier complexity, and adversary sophistication. However, the practicality of tested attacks and adversarial defenses in intrusion scenarios is not extensively addressed.

Patil et al. [6] investigate adversarial attacks on AI-based malware detection models and introduce an effective adversarial training defense mechanism. While focusing on a specific type of adversarial attack, the paper lacks a comprehensive discussion on real-world implications.

Siva Kumar et al. [10] identify gaps in securing ML systems compared to traditional software security and outline a research agenda for enhancing security practices in adversarial ML.

Ahsan et al. [9] examine the implementation of machine learning techniques in cybersecurity and discuss their effectiveness in countering existing threats. However, the paper primarily discusses theoretical aspects and lacks empirical validation of machine learning techniques in real-world cyber security settings.

Sarker [12] explores adversarial machine learning in cybersecurity and defense mechanisms, emphasizing the importance of working with real-world data. Yet, the empirical validation of machine learning techniques in real-world settings and further exploration of adversarial attacks on machine learning-based cybersecurity systems are needed.

Kurakin, Goodfellow, and Bengio [11] address adversarial attacks on machine learning models and introduce adversarial training at scale, with a particular focus on ImageNet. Their contributions include recommendations for scaling adversarial training to large models and datasets, insights into enhanced model robustness against single-step attacks, and an exploration of the reduced transferability of multi-step attacks. They find that single-step attacks are notably effective in black-box scenarios. Additionally, the study resolves a "label leaking" effect, ultimately improving adversarial model performance. This research offers valuable insights into the challenges and considerations of large-scale adversarial training.

## 3. ADVERSARIAL ATTACKS ON MACHINE LEARNING MODELS

The landscape of machine learning security is characterized by an incessant arms race, where machine learning models must confront an array of adversarial attacks designed to exploit their vulnerabilities. This section delves into the intricacies of adversarial attacks, elucidating various types, providing real-world incidents for context, and exploring the motivations behind adversarial actions.
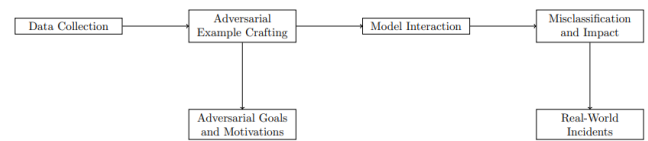


**Fig -1:.** Adversarial Attack Workflow

A. *Types of Adversarial Attacks*

Adversarial attacks can manifest in several forms, each with distinct characteristics and objectives.

- Evasion Attacks:

  Evasion attacks, also known as adversarial examples, in- volve manipulating input data to mislead machine learning models. Adversaries introduce imperceptible perturbations into the input, causing the model to make incorrect predictions. These perturbations, while seemingly inconsequential to hu- man observers, can be highly effective in causing model misclassification.

- Poisoning Attacks:

  In poisoning attacks, adversaries attempt to compromise the integrity of the training data by introducing malicious data points during the training phase. These poisoned data points subtly skew the model's learned parameters, rendering it vulnerable to specific attacks during deployment.

- Data-Driven Attacks:

  Data-driven attacks involve manipulating the data distribution upon which the model is trained. Attackers can induce a shift in the data distribution, making the model less robust to real-world variations. This is particularly concerning in applications like anomaly detection and fraud prevention.

- Model Inversion Attacks:

  Model inversion attacks focus on reversing the prediction process to gain insights into the training data or other sensitive information. By submitting queries to the model and analyzing the model's responses, adversaries attempt to reconstruct details about the data the model was trained on, potentially revealing sensitive information.

- Membership Inference Attacks:

  Membership inference attacks determine whether a specific data point was part of the training dataset used to build a ma- chine learning model.

Attackers exploit the model's responses to query data points, aiming to identify if a particular data instance was in the training set or not. This has implications for data privacy and confidentiality

- Model Stealing Attacks:

  Model stealing attacks target the theft of a machine learning model itself. Attackers query the model and use its responses to reconstruct a clone of the target model. This can have severe implications for intellectual property and proprietary model protections.

B. *Real World Incidents:*

The impact of adversarial attacks is not confined to hypothetical scenarios but has been vividly illustrated through various real-world incidents, demonstrating the tangible risks and vulnerabilities associated with these attacks. Some notable examples include:

1. Stop Sign Manipulation in Autonomous Vehicles:

   In 2017, researchers showcased a vulnerability in object recognition systems employed in autonomous vehicles. They were able to manipulate stop signs by strategically placing stickers on them, causing the signs to be misclassified as yield signs. This misclassification could lead to catastrophic consequences if exploited by malicious actors, as it might cause autonomous vehicles to disregard stop signs.

2. Deep Learning-Based CAPTCHA Solvers for Fraudulent Account Creation:

   Cybercriminals have leveraged deep learning-based CAPTCHA solvers to automate the creation of fraudulent accounts on websites and online platforms. These solvers employ machine learning models to circumvent CAPTCHA challenges, thus facilitating the mass creation of accounts for various nefarious purposes, including spamming, identity theft, and fraudulent activities.

3. Adversarial Image Attacks in Social Media:

   Adversarial attacks have been employed on social media platforms to manipulate images and videos. For example, attackers have used deep learning techniques to generate deepfake videos, convincingly altering the appearance and speech of individuals, often with malicious intent, such as spreading disinformation or defamation.

4. Speech Recognition Manipulation for Voice Assistants:

   Voice-activated virtual assistants, such as Amazon Alexa and Google Assistant, have been susceptible to adversarial attacks that manipulate voice commands. Attackers have successfully devised audio attacks that subtly modify spoken commands to trick these systems into executing unintended actions, posing security and privacy risks.

C. *Motivations and Goals of Adversaries:*

To effectively mitigate adversarial threats, comprehend- ing the intricate motivations and goals of adversaries is of paramount importance. The adversaries targeting machine learning models encompass a wide spectrum of intentions, each driving their actions with distinct objectives:

1. Financial Gain and Fraudulent Activities:

   Many adversaries are primarily motivated by financial gain. They exploit machine learning vulnerabilities to subvert systems such as fraud detection, enabling activities like credit card fraud, insurance fraud, or manipulating stock trading algorithms for illicit profits. By undermining these safeguards, adversaries aim to amass wealth through fraudulent means.

2. Privacy Invasion and Data Exploitation:

   For some, the lure of sensitive data is irresistible. Adversaries aim to compromise the privacy of individuals and organizations. They exploit vulnerabilities in recommendation systems, user profiles, and databases to breach user data. The compromised data can be used for identity theft, blackmail, or sold on the dark web, raising profound concerns about data privacy and security.

3. Strategic Espionage and National Security:

   In a more ominous realm, nation-state actors target machine learning models for strategic or intelligence purposes. They seek to compromise models used in military applications, critical infrastructure, or governmental systems. These at- tacks have significant implications for national security, as adversaries may aim to disrupt critical operations, manipulate intelligence data, or gain a strategic advantage.

4. Malicious Manipulation and Misinformation:

Some adversaries are motivated by a desire to manipulate machine learning models for nefarious purposes. For instance, they may seek to create deep fake content, spreading disinformation, and undermining trust in media sources. In addition, they may attempt to manipulate autonomous systems, such as self-driving cars, to engage in criminal activities or acts of sabotage.

In understanding the motivations and goals behind adversarial attacks on machine learning models, we lay the foundation for the development of effective defense strategies. As we progress through this paper, we will delve deeper into these motivations and propose corresponding defense mechanisms. This comprehensive analysis contributes to our overarching goal of enhancing the resilience of machine learning systems against a diverse range of adversarial challenges

# 4. METHODS FOR ADVERSARIAL ATTACK GENERATION

This section delves into the intricate techniques employed by adversaries to manipulate machine learning models, providing a comprehensive exploration of adversarial attack gen- eration methods. It encompasses not only the 'how' but the 'why' of attack strategies, offering insight into the tools and algorithms at the forefront of adversarial machine learning.

D. Understanding the Attack Landscape:

To effectively counter adversarial threats, it is essential to comprehend the landscape of techniques attackers employ. Adversarial attack generation involves a multifaceted approach, encompassing evasion, poisoning, and data-driven attacks. It is imperative to understand how these techniques are harnessed to exploit machine learning models and their vulnerabilities

E. Crafting Adversarial Examples:

A central element of adversarial attack generation is the crafting of adversarial examples. Adversaries manipulate in- put data to introduce perturbations that lead to model mis- classification. The creation of these adversarial examples is underpinned by an intimate knowledge of the target model, its architecture, and the nature of the data it processes.

F. Common Attack Algorithms:

To orchestrate these adversarial attacks, adversaries employ a variety of algorithms.

1. Fast Gradient Sign Method (FGSM):

FGSM is a fundamental one-step attack method. It calculates the gradient of the loss function with respect to the input data and then perturbs the input data in the direction of the gradient. This perturbation is typically scaled by a small constant ($\epsilon$) to ensure it remains within a certain distortion threshold.

Let's denote the input data as x, the loss function as J, the model as f and the perturbation as D. The FGSM attack can be described as:

$$D = \epsilon.sign(\nabla_x J(f(x), y)) \qquad (1)$$

Where:

- $\epsilon$ is a small positive constant, determining the magnitude of perturbation.

- $\nabla_x$ denotes the gradient with respect to x.

- J(f(x),y) represents the loss associated with the model's prediction f(x) and true label y.

2. Projected Gradient Descent (PGD):

PGD is an iterative attack that extends FGSM over multiple iterations. In each iteration, small perturbations are applied to the input, and the result is projected back onto a defined $\epsilon$-neighborhood. This process continues for a set number of iterations, increasing the robustness of the attack against defenses. The PGD attack can be shown as follows

$$x^{(t+1)} = Clip_{(x,\epsilon)}(x^t + \alpha.sign(\nabla_x J(f(x), y))) \qquad (2)$$

Where:

- $x^t$ represents the input data at iteration t.

- $\alpha$ is a step size, determining the size of the perturbations in each iteration.

- $Clip_{(x,\epsilon)}$ enforces the $\epsilon$-neighborhood constraint, ensur- ing the perturbed example remains within the desired distortion bounds

3. Carlini and Wagner (C&W) Attack:

The C&W attack is known for its versatility and effective- ness. It formulates adversarial examples as an optimization problem, aiming to minimize perturbation while ensuring misclassification. The C&W attack involves solving a complex optimization problem, which may vary depending on the specifics of the attack. The optimization problem often includes the objective of

minimizing perturbation while ensuring that the model's output for the adversarial example x' satisfies the desired misclassification criteria. It's typically solved using iterative optimization techniques such as the Adam optimizer.

4. DeepFool:

DeepFool is an attack algorithm designed to minimize the L2 norm of perturbations required to push an input data point across the decision boundary of a machine learning model. It aims to find the smallest perturbation necessary for misclassification.

The DeepFool algorithm aims to minimize the following objective function:

$$\arg\min \frac{1}{2}||r||_2^2$$

Subject to constraint:

$$f(x + r) = f(x)$$

Where:

- r represents preturbation.

- f(x+r) is the model's prediction for the perturbed input x+r.

## 5. IMPACTS ON ADVERSARIAL ATTACKS

A. Consequences of Successful Adversarial Attacks: Adversarial attacks on machine learning models yield a myriad of detrimental consequences, including but not limited to:

Misclassification and Model Compromise: Successful adversarial attacks lead to model misclassification, resulting in incorrect predictions. Adversaries can exploit these misclassifications to bypass security mechanisms, such as malware detection or intrusion detection systems, rendering them vulnerable to exploitation.

Erosion of User Privacy: Adversarial attacks on recommendation systems and user profiling can erode user privacy. By manipulating the recommendations or exploiting vulnerabilities in these systems, adversaries can gain unauthorized access to sensitive user information, undermining confidentiality and trust.

Financial Losses: Financial institutions, stock markets, and online payment systems rely on machine learning for fraud detection and algorithmic trading. Successful adversarial attacks can lead to substantial financial losses, as fraudsters exploit vulnerabilities in these systems for personal gain.

B. Financial and Reputational Costs:

Successful adversarial attacks exact both financial and reputational costs on organizations and individuals. These include:

Monetary Losses: Organizations experience direct financial losses when fraud and security breaches occur as a result of adversarial attacks. Remediation efforts, legal fees, and compensation to affected parties can be substantial.

Reputational Damage: Adversarial attacks tarnish an organization's reputation. News of security breaches and system failures undermines trust in products and services, affecting customer loyalty and investor confidence.

Compliance and Regulatory Penalties: In some cases, regulatory bodies may impose penalties and fines on organizations that fail to protect against adversarial attacks. Non-compliance with data protection and cybersecurity regulations can have far-reaching legal and financial consequences.

By elucidating the comprehensive spectrum of consequences, along with real-world case studies and the financial and reputational costs, this section underscores the urgency of addressing adversarial threats. It emphasizes the critical need for robust defense mechanisms and countermeasures to mitigate these profound impacts.

## 6. MECHANISMS AND STRATEGIES

This section provides a comprehensive examination of the multifaceted world of defense mechanisms and strategies against adversarial attacks in machine learning. It underscores the critical importance of model robustness and presents a spectrum of countermeasures to safeguard machine learning systems against adversarial manipulation.
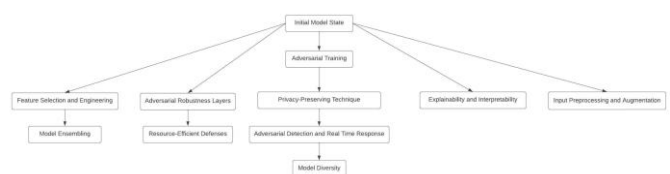


**Fig. 2:** Landscape of defense mechanism

A. The Imperative of Robustness:

Robustness in machine learning serves as the linchpin in defense against adversarial attacks. It constitutes the model's ability to maintain its performance and integrity even in the face of adversarial input. Robust models possess inherent resistance to manipulation and can

reliably generate predictions even when subjected to maliciously crafted data.

1. Achieving Robustness: Achieving robustness necessitates a multi-pronged approach, embracing several key strategies:

   a. Adversarial Training:

   Adversarial training stands at the forefront of fortifying models against adversarial attacks. This technique involves augmenting the training dataset with carefully crafted adversarial examples. By training on adversarial inputs, the model acquires the ability to recognize and defend against such manipulations during inference. The adversarial examples act as an additional training signal, making the model more resilient.

   b. Feature Selection and Engineering:

   Prudent feature selection and engineering are pivotal com- ponents of model robustness. Feature selection aims to retain only the most relevant and robust input features while discarding those that are more susceptible to adversarial perturbations. It reduces the model's attack surface by focusing on data elements that are less prone to manipulation.

   c. Model Ensembling:

   Model ensembling is an effective defense strategy that involves combining the outputs of multiple machine learning models or predictions. Ensembles can include models with varying architectures or training data. The diversity of these ensemble members makes it significantly more challenging for adversaries to craft adversarial examples that fool the entire ensemble. It enhances the overall robustness and accuracy of the system.

B. Countermeasures against Adversarial Attacks:

Countermeasures against adversarial attacks encompass a diverse array of strategies, each designed to mitigate vulnerabilities and enhance the resilience of machine learning models.

1. Adversarial Training:

   Adversarial training, as a central defense mechanism, in- volves training models on adversarial examples. During the training process, adversarial perturbations are introduced to the input data. By learning to adapt to these perturbations, the model becomes more robust against adversarial attacks during inference.

2. Feature Selection and Engineering:

   Feature selection and engineering strive to reduce the vulnerability of machine learning models to adversarial manipulation by selecting and engineering input features carefully. This process can be guided by domain knowledge to identify and retain only the most informative and robust features, minimizing the model's exposure to adversarial perturbations.

3. Model Ensembling:

   Model ensembling is a potent defense mechanism that leverages diversity. It involves combining the predictions of multiple models, each with distinct strengths and weaknesses. Adversaries find it more challenging to craft adversarial examples that can deceive the entire ensemble. Model ensembling enhances the overall robustness and reliability of machine learning systems.

4. Adversarial Robustness Layers:

   Adversarial robustness layers are specialized components integrated into the neural network architecture with the primary objective of enhancing the model's resilience against adversarial attacks. These layers act as a proactive defense mechanism by detecting and mitigating adversarial inputs before they can propagate through the network and influence the model's predictions.

   The key functions of adversarial robustness layers include:

   - Preprocessing and Normalization: Adversarial robustness layers often start with preprocessing and normalization steps to prepare the input data. These steps aim to enhance the data's consistency and reduce its vulnerability to adversarial perturbations.

   - Anomaly Detection: Anomaly detection techniques, such as autoencoders or other unsupervised learning methods, are employed to scrutinize the input data. These layers identify deviations or irregularities that align with adversarial patterns. Anomalies are identified by comparing the input's representations at different layers within the network.

   - Feature Extraction: Specialized feature extraction layers may be utilized to extract salient and robust features from the input data. These features are specifically chosen for their ability to resist adversarial manipulation.

   - Threshold-Based Filters: After preprocessing and anomaly detection, threshold-based filters are

applied. These filters as- sess the extent to which the input data deviates from expected norms. Inputs exceeding predefined thresholds for abnormality are either discarded or flagged for further scrutiny.

- Feedback Mechanisms: Adversarial robustness layers can incorporate feedback mechanisms that enable the neural net- work to adapt and improve its ability to detect and mitigate adversarial attacks. When an adversarial input is detected, the network may adjust its parameters and decision boundaries to enhance future detection capabilities.

- The specific algorithms and techniques used within these layers can be tailored to the unique requirements of the machine learning model and the nature of the data. Research in this area continues to evolve, with a focus on developing more effective and adaptable adversarial robustness layers to enhance model security and reliability.

5. Input Preprocessing and Augmentation:

Input preprocessing and augmentation techniques introduce noise or other perturbations to the input data. By distorting the input in a controlled manner, these techniques make it significantly more challenging for adversaries to craft effective adversarial examples. Input preprocessing and augmentation contribute to enhanced model robustness. These defense mechanisms and strategies are pivotal in ensuring the reliability and trustworthiness of AI systems in the face of evolving adversarial challenges.

## 7. EVALUATING MODEL ROBUSTNESS

Evaluating the robustness of machine learning models against adversarial attacks is a fundamental aspect of securing AI systems. This systematic evaluation process incorporates critical components, beginning with the selection of adversarial benchmark datasets. These datasets are crucial for assessing model resilience, as they encompass a diverse range of adversarial examples crafted with various attack techniques, ensuring a comprehensive assessment. To comprehensively evaluate model robustness, the establishment of an adversarial attack taxonomy is essential, categorizing attacks into evasion and data poisoning categories, enabling systematic assessment of vulnerabilities. The selection of appropriate evaluation metrics, such as accuracy, robust accuracy, and area under the receiver operating characteristic curve (AUC-ROC), plays a pivotal role in gauging a model's performance under adversarial conditions. To rigorously test a model's resilience against

adversarial attacks, several well-defined evaluation methodologies are employed. White-box testing, where attackers possess complete knowledge of the model, allows for a thorough examination of vulnerabilities. In contrast, black-box testing simulates real-world attack scenarios with minimal insights into the model's internal workings. Transferability testing explores whether adversarial examples crafted for one model can successfully fool other models, shedding light on the generalizability of adversarial attacks and their potential impact on broader machine learning ecosystems. This comprehensive evaluation framework is essential to fortify machine learning models against adversarial threats and ensure the security of AI systems.

## 8. FUTURE DIRECTIONS AND CHALLENGES

Persistent challenges include achieving robustness in multimodal learning, developing defense mechanisms with broad applicability, improving model interpretability, and creating resource-efficient defenses. Future research directions include the development of comprehensive adversarial resilience frameworks, extending robustness to federated learning, and exploring the intersection of adversarial robustness and privacy preservation. Additionally, there is a need for real-time adversarial detection systems, especially in critical applications like autonomous vehicles and healthcare. These challenges and opportunities underscore the growing significance of secure, resilient, and privacy-conscious AI systems in an increasingly AI-driven era.

## 9. CONCLUSIONS

As the field of machine learning and artificial intelligence continues to evolve, so too will adversarial attacks. By gaining a deep understanding of the adversarial landscape, we are better equipped to anticipate and counteract the evolving strategies of adversaries. By comprehensively analyzing the multifaceted aspects of adversarial threats in the machine learning landscape, this research paper contributes to the development of effective defenses and strategies to enhance the resilience of machine learning systems.

## REFERENCES

[1] E. Anthi, L. Williams, M. Rhode, P. Burnap, and A. Wedgbury, "Adversarial attacks on machine learning cybersecurity defences in industrial control systems, Journal of Information Security and Applications, vol. 58, p. 102717, 2021. [Online].

[2] A. Alotaibi and M. A. Rassam, "Adversarial machine learning attacks against intrusion detection systems: A survey on strategies and defense," Future Internet, vol. 15, no. 2, 2023. [Online].

[3] I. Rosenberg, A. Shabtai, Y. Elovici, and L. Rokach, "Adversarial machine learning attacks and defense methods in the cyber security domain," ACM Comput. Surv., vol. 54, no. 5, may 2021. [Online].

[4] J. Tygar, "Adversarial machine learning," Internet Computing, IEEE, vol. 15, pp. 4 – 6, 11 2011.

[5] N. Martins, J. Cruz, T. Cruz, and P. Henriques Abreu, "Adversarial ma- chine learning applied to intrusion and malware scenarios: A systematic review," IEEE Access, vol. PP, pp. 1–1, 02 2020.

[6] S. Patil, V. Varadarajan, D. Walimbe, S. Gulechha, S. Shenoy, A. Raina, and K. Kotecha, "Improving the robustness of ai-based malware detection using adversarial machine learning," Algorithms, vol. 14, no. 10, 2021. [Online].

[7]                                    G. Apruzzese, M. Andreolini, L. Ferretti, M. Marchet ti,  and M. Colajanni, "Modeling realistic adversarial attacks against network intrusion detection systems," Digital Threats, vol. 3, no. 3, feb 2022. [Online].

[8] S. Soni and B. Bhushan, "Use of machine learning algorithms for designing efficient cyber security solutions," in 2019 2nd International Conference on Intelligent Computing, Instrumentation and Control Technologies (ICICICT), vol. 1, 2019, pp. 1496–1501.

[9] M. Ahsan, K. E. Nygard, R. Gomes, M. M. Chowdhury, N. Rifat, and J. F. Connolly, "Cybersecurity threats and their mitigation approaches using machine learningmdash;a review," Journal of Cybersecurity and Privacy, vol. 2, no. 3, pp. 527–555, 2022. [Online].

[10] R. S. Siva Kumar, M. Nyström, J. Lambert, A. Marshall, M. Goertzel, A. Comissoneru, M. Swann, and S. Xia, "Adversarial machine learning- industry perspectives," in 2020 IEEE Security and Privacy Workshops (SPW), 2020, pp. 69–75.

[11] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial machine learn- ing at scale," 2017.

[12] I. H. Sarker, "Machine Learning for Intelligent Data Analysis and Automation in Cybersecurity: Current and Future Prospects," sep 19 2022.