

Must Learn AI Security

aka.ms/MustLearnAISecurity

ROD TRENT
SENIOR PROGRAM MANAGER
MICROSOFT

This post is part of an ongoing series to educate about new and known security vulnerabilities against AI.

The full series index (including code, queries, and detections) is located here:

<https://aka.ms/MustLearnAISecurity>

The book version (pdf) of this series is located here: Coming...

The book will be updated when each new part in this series is released.

This book is updated every time a new part of this series is posted. The most current edition of this book will always be located at: https://github.com/rod-trent/OpenAISecurity/tree/main/Must_Learn/Book_Version

Book release ver. 0.27, August 24, 2023 8:00am EST

Contents

Must Learn AI Security Series: Introduction	7
Giving AI Security the Must Learn Treatment.....	7
Azure Open AI Weekly Community Copilot.....	9
Must Learn AI Security Part 1: Prompt Injection Attacks Against AI	10
Chapter 1	10
What is a prompt?	10
What is a prompt injection attack?	10
Why it matters	11
How it might happen.....	11
Real-world Example	12
How to monitor	13
What to capture	14
How to mitigate	14
EXTRA: Content Filtering	16
Must Learn AI Security Part 2: Data Poisoning Attacks Against AI	17
Chapter 2	17
What is a Data Poisoning attack?	17
How it works	17
Types of Data Poisoning attacks.....	18
Why it matters	18
Why it might happen	19
Real-world Examples.....	19
How to mitigate	19
How to monitor	20
What to capture	21
Must Learn AI Security Part 3: Adversarial Attacks Against AI	22
Chapter 3	22
What is an Adversarial attack?	22
How it works	22
Types of Adversarial attacks.....	23
Why it might happen	24
Real-world Example	25

How to Mitigate	25
How to Monitor and what to capture.....	27
Must Learn AI Security Part 4: Trojan Attacks Against AI	29
Chapter 4	29
What is a Trojan attack against AI?	29
How it works	29
Types of Trojan attacks	30
Why it matters	31
How it might happen.....	32
Real-world Example	33
How to Mitigate	33
How to monitor	34
What to capture	35
Must Learn AI Security Part 5: Evasion Attacks Against AI.....	38
Chapter 5	38
What is an Evasion attack against AI?	38
How it works	38
Types of Evasion attacks	39
Why it matters	40
Why it might happen	41
Real-world Example	42
How to Mitigate	42
How to monitor	43
What to capture	44
Must Learn AI Security Part 6: Model Inversion Attacks Against AI	46
Chapter 6	46
What is a Model Inversion attack against AI?	46
How it works	46
Types of Model Inversion attacks	47
Why it matters	48
Why it might happen	49
Real-world Example	50
How to Mitigate	51

How to monitor/What to capture	51
Must Learn AI Security Part 7: Membership Inference Attacks Against AI	53
Chapter 7	53
What is a Membership Inference attack against AI?	53
How it works	53
Types of Membership Inference attacks	54
Why it matters	55
Why it might happen	57
Real-world Example	58
How to Mitigate	59
How to monitor/What to capture	60

Must Learn AI Security Series: Introduction

Giving AI Security the Must Learn Treatment

Let's start with a joke...

Q: How can you tell when there is a Microsoft person in the room?

A: The conversation turns to AI within 5 minutes.

To be honest, this is actually more a truism than a joke, but hopefully you get the gist.

AI is everywhere right now.

In a rapidly evolving digital landscape, the power and potential of Artificial Intelligence (AI) have ignited what seems like - bubble or not - a technological revolution. The latest "hotness" of Generative AI has catapulted AI into the forefront of every technology conversation and as a Microsoft person, that joke absolutely applies to me. As a technology person, I was drawn to ChatGPT early on but as a security person I was immediately worried about the security of the shiny, new thing.

This blog series begins with the haunting realization that AI, like any other technology, is susceptible to exploitation and abuse. The very intelligence that empowers AI to make autonomous decisions can also be manipulated to execute malevolent actions. Within the virtual realm, a breach of AI security can have profound real-world consequences, endangering privacy, economy, and even human lives.

Throughout the chapters that follow, I'll delve into the multifaceted dimensions of AI security. I'll explore the challenges posed by

adversarial attacks and attempt to provide prescriptive guidance on how to monitor, capture, and mitigate each type of AI harm.

Additionally, I'll examine the need for transparent and accountable AI systems that respect user privacy and uphold ethical standards. As AI becomes an integral part of our daily lives, it is imperative that we confront the ethical implications of its proliferation.

You'll see very quickly that applying standard best practices for security - overlaying existing templates for good security - will work in most cases. You'll also come to realize that most AI security is focused on first, writing secure code and then, ensuring data sources are protected.

We're all in this together. I'm by no means an expert in AI security, but I'm working toward that because it's the next important thing. We can all learn from one another. As I embark on this quest, I'll extend an invitation to all here to provide feedback. You can provide feedback through the chat system here on this site, but also through the GitHub repository that will house the queries, detections, and other collateral for this series.

GitHub repo for Must Learn AI Security: <https://aka.ms/MustLearnAISecurity>

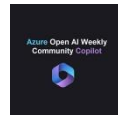
So, welcome to "[Must Learn AI Security](#)," a continuing and evolving comprehensive exploration of where AI and security intersect. As with the original [Must Learn KQL series](#), the content here will be made available as a series, following a logical design to enable you to get up to speed on the topics and concepts of a secure AI. I have no clue when it will complete - or if it will. Just like Must Learn KQL it will end when it ends. As the chapters grow, you can expect a downloadable PDF book version you can read with your favorite eReader (browser, Kindle, etc.).

Obviously, I work for Microsoft and many of my references and examples will be around Microsoft Security platform products like Microsoft Sentinel, Defender, and Azure OpenAI. But I'll try my best to keep the references to a minimum as this is an important topic for everyone - not just Microsoft customers.

I look forward to working with you all.

Lastly, here's some recommended resources to get you started and keep you informed:

1. **Azure OpenAI community group on LinkedIn:** <https://www.linkedin.com/groups/14241561/>
2. **Book:** [Not with a Bug, But with a Sticker: Attacks on Machine Learning Systems and What To Do About Them](#)
3. **Weekly Azure OpenAI newsletter:**



[Azure Open AI Weekly Community Copilot](#)

[Welcome to the Microsoft Azure Open AI community newsletter!](#)

Must Learn AI Security Part 1: Prompt Injection Attacks Against AI

Chapter 1

To understand the Prompt Injection Attack for AI, it helps to first understand what a *Prompt* is.

What is a prompt?

When we interact with AI language models, such as ChatGPT, Google Bard, and others, we provide a *prompt* in the form of a question, sentence, or short paragraph. The *prompt* is what is fed to the AI. It is our desired information that the model should analyze and then produce a result in the form of a task or response. It acts like a conversation starter or cue that helps create the desired output. Prompts let us control the conversation and direct it in a certain way.

What is a prompt injection attack?

A prompt injection attack refers to the act of *maliciously* manipulating the input prompts given to an AI system to trick, subvert, or exploit its behavior. The goal of such an attack could vary depending on the context, but some potential objectives could include:

1. **Bias Injection:** Injecting biased or harmful prompts to influence the AI's outputs in a way that promotes misinformation, hate speech, or discriminatory content.
2. **Data Poisoning:** Introducing tainted or misleading prompts during the AI training process to compromise the model's performance and cause it to produce erroneous results.

3. **Evasion:** Crafting prompts specifically designed to evade the AI's security or detection mechanisms, enabling malicious activities to go unnoticed.
4. **Model Exploitation:** Manipulating the prompts to cause the AI model to perform actions it was not designed for, such as revealing sensitive information or performing unauthorized tasks.
5. **Adversarial Attacks:** Crafting adversarial prompts that exploit vulnerabilities in the AI model, causing it to make incorrect or unintended decisions.

Why it matters

A prompt is crucial in shaping the output generated by the language model. It provides the initial context, specific instructions, or the desired format for the response. The quality and precision of the prompt can influence the relevance and accuracy of the model's output.

For example, if you ask (your *prompt*), "What's the best cure for poison ivy?", the model, as you should expect, is designed to concentrate on health-related information. The response should offer solutions based on the data sources that was used to train the model. It should probably provide common methods of a cure and a warning that they might not work for everyone. And should end with advising to consult a doctor. However, if someone has tampered with the language model by adding harmful data, users could receive incorrect or unsafe information.

How it might happen

A great, current example of how this might happen is related in a recent security issue reported by Wired Magazine. In the article, [A New Attack Impacts Major AI Chatbots—and No One Knows How to](#)

[Stop It](#), it talks about someone using a string of nonsense characters to trick ChatGPT into responding in a way it normally wouldn't.

Reading the article, a user could take the supplied nonsense string (copy) and tack it onto (paste) their own prompt and cause ChatGPT to respond differently or issue a response that would normally be disallowed by policy.

In one sense, I guess, you could say the author of the article is a threat actor using a prompt injection attack. We're just left to determine if it was *malicious* or not.

Real-world Example

A real-world example of a prompt injection attack against AI is a study conducted by researchers from the University of Maryland and Stanford University in 2021. They explored the vulnerabilities of OpenAI's GPT-3 language model to prompt injection attacks, also known as "Trojan attacks" in the context of NLP models.

In their experiment, they demonstrated that an attacker could exploit the vulnerabilities of GPT-3 by manipulating the input prompt in a way that the AI model would produce malicious or harmful content as output. For instance, if GPT-3 is used as a code generation tool, an attacker could craft the input prompt in such a way that the generated code includes a hidden backdoor, allowing unauthorized access to a system or application.

This example shows that AI-powered language models like GPT-3 can be susceptible to prompt injection attacks, where an attacker manipulates the input prompt to make the AI system generate malicious or undesirable content. To mitigate such risks, AI developers and users need to be aware of the potential vulnerabilities and implement appropriate security measures, such

as input validation, prompt filtering, and monitoring the generated content for malicious activities.

How to monitor

Continuously monitoring and logging application activities is necessary to detect and respond to potential security incidents quickly. Monitoring should produce a based model of accurate prompts and any outliers should be identified and resolved through ongoing mitigation.

Monitoring can be accomplished through a data aggregator that analyzes for outliers. A good example is a modern SIEM, like [Microsoft Sentinel](#), which enables organizations to collect and analyze data and then create custom detections from alerts to notify security teams when prompts are outside norms or organization policies.

For the growing library of queries, detections, and more for Microsoft Sentinel see: [OpenAISecurity/Security/Sentinel at main · rod-trent/OpenAISecurity \(github.com\)](#)

One big note here. You need to identify if your AI provider allows monitoring of prompts. As its early days, most currently don't. They do capture the prompts - some for a shorter, some for longer retention periods - they just don't expose it to customers for logging purposes. The idea is that prompts are user-specific and private and instead it's monitoring what the result or response is that matters most. Personally, I don't agree. That should be left to the organization. But there are content filtering mechanisms available to help curb what user can and cannot enter as prompts.

There are also other mechanisms that can be used, such as filtering usage data through a proxy (CASB) or ensuring that your organization develops its own interface to the AI provider and use

the API instead of direct prompts so that you can better control what users can do.

What to capture

Once you've identified the data available in the log stream, you can start to focus on the specific pieces of artifact (evidence) that will be useful in capturing potential attackers and creating detections.

Here's a few things to consider capturing:

1. IP Addresses (internal and external)
2. Human and non-human accounts
3. Geographical data - this is important to match up to known threats (nation state or otherwise)
4. Success AND failures

Consider creating a watchlist of known entities (users, IPs) that should be able to access your AI and one for approved geographical locations. Using an editable watchlist enables you to quickly adjust your detections should the threat landscape change.

Microsoft Sentinel users, see: [Monitor Azure Open AI Deployments with Microsoft Sentinel](#)

How to mitigate

Mitigation for this type of attack is generally considered precautionary steps to avoid it in the first place. To mitigate a prompt injection attack, both developers and users should take appropriate precautions. Here are some steps to follow:

1. **Input validation:** Implement strict validation checks on user inputs to filter out malicious content, ensuring only valid and safe prompts are passed to the AI model.

2. **Sanitization:** Sanitize user inputs to remove or neutralize potentially harmful elements before processing them in the AI system.
3. **Rate limiting:** Apply rate limiting on user requests to prevent excessive or rapid attempts at injecting malicious prompts, making it harder for attackers to exploit the system.
4. **Monitoring and logging:** Monitor and log user inputs and AI responses to identify suspicious patterns, enabling early detection of potential prompt injection attacks.
5. **Regular updates and patches:** Keep your AI models and related software up to date, applying security patches and updates to minimize vulnerabilities.
6. **User education:** Educate users about the risks of prompt injection attacks, encouraging them to be cautious when providing input to AI systems and to report any suspicious behavior.
7. **Secure AI model training:** Ensure your AI models are trained on high-quality, diverse, and reliable data sources to reduce the chances of the model producing harmful outputs.
8. **Phish Resistant MFA:** For organizations developing their own AI apps, make sure to use proper identity mechanisms.
9. **Trusted devices/applications:** Ensure only trusted devices and applications are granted access.
10. **Data Loss Protection (DLP):** Protect sensitive corporate data from leaving the company due to user negligence, mishandling of data, or malicious intent.

By implementing these measures, you can reduce the risk of prompt injection attacks and enhance the overall security of your AI systems. To defend against prompt injection attacks, developers and researchers need to employ robust security measures, conduct thorough testing and validation, and implement mechanisms to detect and mitigate potential risks associated with manipulated prompts.

EXTRA: Content Filtering

One other thing to consider is developing a strong content filtering policy. As the data flows into the modern SIEM, outliers are identified, and detections and alerts are created, a part of mitigation that can help is to develop a better content filtering strategy. Azure OpenAI, for example, provides a stock feature for quickly adjusting content filtering. To create a fully customizable version, customers need to request full access to their own filtering.

See:

[Azure OpenAI Service content filtering](#)

[Preventing abuse and harmful content generation](#)

Create content filtering configuration



Content filtering configurations are created within a Resource and can be associated with Deployments.

[Learn more about configurability here.](#)

The default content filtering configuration is set to filter at the medium severity threshold for all four content harms categories for both, prompts and completions. That means that content that is detected at severity level medium or high is filtered, while content detected at severity level low is not filtered by the content filters.

Create custom configuration name

CustomContentFilter148

Set severity levels

		User prompts (Input)				Model completions (Output)		
		Low	Medium	High		Low	Medium	High
Hate ⓘ	<input checked="" type="checkbox"/> On	✓	⊖	⊖		<input checked="" type="checkbox"/> On	✓	⊖
Sexual ⓘ	<input checked="" type="checkbox"/> On	✓	⊖	⊖		<input checked="" type="checkbox"/> On	✓	⊖
Self-harm ⓘ	<input checked="" type="checkbox"/> On	✓	⊖	⊖		<input checked="" type="checkbox"/> On	✓	⊖
Violence ⓘ	<input checked="" type="checkbox"/> On	✓	⊖	⊖		<input checked="" type="checkbox"/> On	✓	⊖

[Learn more about content filters here](#)

Save

Cancel

Must Learn AI Security Part 2: Data Poisoning Attacks Against AI

Chapter 2

What is a Data Poisoning attack?

A Data Poisoning attack is a type of malicious activity aimed at machine learning models. A successful attack results in incorrect or misleading information being fed into the training data. The objective of this attack is to skew the model's learning process, causing it to make incorrect predictions or classifications.

As you can imagine from the description, data protection is key to protecting against this method of attack. While external forces are definitely a threat, more than often this type of attack is the result of an internal threat enacted by someone with either proper or hacked credentials.

How it works

As noted just prior, access to the data source used for training is the primary element of this attack and generally follows these steps:

1. **Model Targeting:** The attacker first identifies a target model that they wish to manipulate.
2. **Injecting Poisoned Data:** The attacker then injects poisoned data into the training set. This data is carefully crafted to look normal but contains misleading features or labels that are intended to mislead the learning algorithm.
3. **Training on Poisoned Data:** The targeted model is trained or retrained using the contaminated training data. The model

learns from both the authentic and poisoned data, which can subtly or substantially alter its behavior.

4. **Exploiting the Compromised Model:** Once the model has been trained on the poisoned data, it may behave in ways that benefit the attacker. For example, it might systematically misclassify certain types of inputs, or it could leak sensitive information.

Types of Data Poisoning attacks

Data poisoning against AI is an ongoing and evolving area of security. While both the methods used to conduct these attacks and the techniques to defend against them continue to evolve, it's still essential knowledge. Currently, the following types of attacks have been identified and categorized.

1. **Targeted Attacks:** These attacks are aimed at specific misclassification or a particular wrong behavior of the model. The attacker may want the model to misclassify images of a certain type or favor one class over another.
2. **Random Attacks:** These attacks aren't targeted at any particular misbehavior. Instead, they aim to reduce the overall performance of the model by injecting random noise or incorrect labels into the training data.

Why it matters

Data poisoning attacks can have serious consequences, such as:

1. **Loss of Integrity:** The model may lose its reliability and start making incorrect predictions or decisions.
2. **Loss of Confidentiality:** In some cases, attackers may use data poisoning to infer sensitive information about the training data or the individuals involved in the training process.

3. **Reputation Damage:** If a poisoned model is widely used, it may lead to the erosion of trust in both the system and the organization responsible for it.

Why it might happen

Other than providing information for nefarious and dangerous purposes, this type of attack is generally considered more frequently for political purposes through the delivery of “fake” information to alter or steer election results. But imagine if an attacker recategorized “not safe for work” images so that they were viewable to get someone fired.

Real-world Examples

One example of a data poisoning attack against AI is manipulating the training data of the model to corrupt its learning process. This can be done by intentionally inserting incorrect, misleading, or manipulated data into the model's training dataset to skew its behavior and outputs. An example of this would be to add incorrect labels to images in a facial recognition dataset to manipulate the system into purposely misidentifying faces.

Another example is the manipulation of images to deceive image classification models. An early example of this is Tay, Microsoft's Twitter chatbot released in 2016. Twitter intended for Tay to be a friendly bot that Twitter users could interact with. However, within 24 hours of its release, Tay was transformed into a racist and sexist bot due to data poisoning attacks.

How to mitigate

Defending against data poisoning attacks can be complex, but some general strategies include:

1. **Monitoring Data Access:** Using a monitoring mechanism, record user logins and access. Use a Watchlist of trusted users to monitor against.
2. **Monitoring Data Application Activity:** Using the same monitoring mechanism, set a baseline for normal activity (time, schedule) and alert on outliers.
3. **Data Validation and Cleaning:** Regularly reviewing and cleaning the training data to detect and remove any anomalies or inconsistencies.
4. **Robust Learning Algorithms:** Designing algorithms that can detect and mitigate the effects of anomalous data.
5. **Monitoring Model Behavior:** Continuously monitoring the model's behavior and performance can help detect unexpected changes that might indicate a poisoning attack.

How to monitor

Continuously monitoring and logging data access and data application activities are necessary to detect and respond to potential security incidents quickly. Monitoring should produce a based model of accurate prompts and any outliers should be identified and resolved through ongoing mitigation.

Monitoring can be accomplished through a data aggregator that analyzes for outliers. A good example is a modern SIEM, like [Microsoft Sentinel](#), which enables organizations to collect and analyze data and then create custom detections from alerts to notify security teams when prompts are outside norms or organization policies.

For the growing library of queries, detections, and more for Microsoft Sentinel see: [OpenAISecurity/Security/Sentinel at main · rod-trent/OpenAISecurity \(github.com\)](https://github.com/rod-trent/OpenAISecurity)

What to capture

It should be noted that “hallucinations” can sometimes be mistaken for Data Poisoning or [Prompt Injection](#) attacks. This is why monitoring for activity and outliers is so important to identify an actual attack versus a misconfiguration.

For more on hallucinations, see: [Using Azure AI Studio to Reduce Hallucinations](#)

Once you’ve identified the data available in the log stream, you can start to focus on the specific pieces of artifact (evidence) that will be useful in capturing potential attackers and creating detections.

Here’s a few things to consider capturing:

1. IP Addresses (internal and external)
2. Logins: anomalous activity, time elements
3. Potentially compromised accounts (general access, data application access)
4. Human and non-human accounts
5. Geographical data - this is important to match up to known threats (nation state or otherwise)
6. Data modeling success AND failures

Microsoft Sentinel users, see: [Monitor Azure Open AI Deployments with Microsoft Sentinel](#)

Must Learn AI Security Part 3: Adversarial Attacks Against AI

Chapter 3

What is an Adversarial attack?

Adversarial attacks against AI are like throwing a wrench in the gears of a well-oiled machine. These attacks involve crafting sneaky input data to confuse AI systems, making them produce incorrect or misleading results. It's like someone giving you a fake treasure map and watching you dig holes all day. These attacks can expose vulnerabilities in AI systems and, if not addressed, can have some serious consequences, like a bad hair day for AI. So, it's crucial to develop robust AI models that can withstand these pesky adversarial attacks.

How it works

Adversarial attacks against AI are like a sneaky game of trick-or-treat. The attacks aim to fool an AI model by making small, crafty changes to the input data and generally happens in the following order:

1. First, the attacker identifies an AI model's weakness.
2. Next, they create an adversarial example, which is the input data with some subtle alterations. To the human eye, the changes are barely noticeable, but they're just enough to send the AI model into a tizzy.
3. The altered input data is then fed into the AI model which tries to make sense of it. But because of the adversarial example, the model ends up making incorrect predictions or classifications.

4. The attacker then sits back and watches the chaos unfold, like a mischievous kid who's just tied everyone's shoelaces together.

Keep in mind that not all AI models are defenseless against these attacks. Like much of AI security right now, ways to protect AI systems from adversarial attacks are being developed, like training them with adversarial examples or building more robust models.

Types of Adversarial attacks

There's a whole variety of adversarial attacks against AI, just like there are many ways to ruin a perfectly good pie. Here are a few common types:

1. **[Fast Gradient Sign Method \(FGSM\)](#)**: This one generates adversarial examples by adding small, malicious changes to the input data that confuse the AI model.
2. **[Projected Gradient Descent \(PGD\)](#)**: This attack iteratively adjusts the input data to maximize the AI model's error, making it confused.
3. **[Carlini & Wagner \(C&W\) Attack](#)**: This sneaky attack is like slipping onions into a fruit salad. It optimizes the input data to minimize the difference between the original and adversarial examples while still fooling the AI model.
4. **[DeepFool](#)**: This attack is like a game of hide-and-seek with a twist. It finds the smallest possible perturbation to the input data, making it almost invisible to the AI model while still causing it to make incorrect predictions.
5. **[One-Pixel Attack](#)**: This one alters just one pixel of an image to confuse AI models in image classification tasks, showing that even the tiniest change can trip up these fancy AI systems.

AI models need to be designed and tested to withstand these adversarial attacks.

Why it might happen

People have different reasons for launching adversarial attacks against AI. It's like asking why someone would put salt in a sugar bowl. Some reasons include:

1. **Exploiting vulnerabilities:** Just like some folks get a kick out of finding loopholes, attackers might want to expose weaknesses in an AI system and use them to their advantage.
2. **Sabotage:** Some attackers might want to undermine a competitor's AI system or cause reputational damage.
3. **Security research:** Not all adversarial attacks are malicious. Some researchers use these attacks to study AI systems' vulnerabilities and develop more robust and secure models. It's like testing the locks on your doors to make sure no one can break in.
4. **Bypassing security systems:** Some attackers might use adversarial attacks to fool AI-powered security systems, like facial recognition or spam filters. It's like wearing a disguise to sneak past the bouncer at a nightclub.
5. **Stealing sensitive data:** By attacking AI models, some folks might be trying to access confidential information or intellectual property.

While there are benefits for the attackers, these actions can have serious consequences for others. That's why it's essential to develop AI models that can stand their ground against these sneaky attacks.

Real-world Example

A real-world example of an adversarial attack against AI is a research experiment conducted by a team of researchers at Google Brain, OpenAI, and Stanford University in 2017. They demonstrated that by slightly modifying an image, they could deceive an AI-based image recognition system into misclassifying it.

In this particular experiment, they used a technique called "fast gradient sign method" (FGSM) to create adversarial examples. They took an image of a panda and added a small amount of carefully calculated noise, which is imperceptible to humans. This noise caused the AI image recognition system to misclassify the panda as a gibbon with a high confidence level, even though the altered image still appeared to be a panda to humans.

This example highlights the vulnerability of AI systems, particularly deep neural networks, to adversarial attacks. By making subtle changes to the input data, attackers can manipulate the AI system's output, potentially leading to incorrect decisions or unintended actions. Adversarial attacks can pose significant risks in various applications, including autonomous vehicles, security systems, and medical diagnostics, among others.

How to Mitigate

Mitigating adversarial attacks against AI systems typically involves a combination of approaches, as no single method can guarantee complete protection. Some potential methods to mitigate adversarial attacks include:

1. **Data Augmentation:** Enhance the training dataset by adding adversarial examples generated using various attack methods,

which can help the AI system to learn and recognize these perturbations and improve its robustness against such attacks.

2. **Adversarial Training:** Train the AI model using a combination of clean and adversarial examples, allowing the model to learn from both types of data and improve its resilience against adversarial attacks.
3. **Gradient Masking:** Regularize the model during training by adding noise or applying other transformations to the gradient, making it harder for an attacker to compute the gradient and generate adversarial examples.
4. **Defensive Distillation:** Train a second model that learns to mimic the output probabilities of the original model, effectively smoothing the decision boundaries and making it more difficult for an attacker to find adversarial examples.
5. **Randomization:** Introduce randomization during the inference stage, such as by applying random transformations to input data or randomly selecting subsets of the model for evaluation. This can make it more challenging for an attacker to generate adversarial examples that consistently fool the AI system.
6. **Detection Methods:** Employ techniques to detect adversarial examples at runtime, such as by comparing the input's features to known clean and adversarial examples or monitoring the model's behavior during inference.
7. **Ensemble Methods:** Use multiple AI models or an ensemble of models to make predictions. This can reduce the likelihood of a single adversarial example fooling all models simultaneously.
8. **Robust Model Architectures:** Design AI models with in-built robustness to adversarial attacks, such as by incorporating attention mechanisms, dropout layers, or other architectural components that can help the model withstand adversarial perturbations.

9. **Regularization Techniques:** Apply regularization techniques like L1 or L2 regularization during training to reduce model complexity and improve generalization, which can make the model less susceptible to adversarial attacks.
10. **Ongoing Research and Collaboration:** Stay up to date with the latest research in adversarial robustness and collaborate with other researchers and practitioners to develop and share effective mitigation techniques.

How to Monitor and what to capture

Monitoring adversarial attacks against AI systems involves detecting and analyzing unusual or malicious activities that target the AI models. This can be achieved through a combination of techniques including the following:

1. **Input Monitoring:** Analyze input data for anomalies, unexpected patterns, or changes in distribution that might indicate an adversarial attack. This can be done using statistical methods, machine learning algorithms, or deep learning techniques to detect and flag suspicious inputs.
2. **Model Behavior Monitoring:** Track the AI model's behavior, such as its confidence in predictions or output probabilities, to identify anomalies that could suggest an adversarial attack. An unusually high or low confidence level or a sudden change in the model's behavior may be indicative of an attack.
3. **Performance Metrics Tracking:** Continuously monitor the AI system's performance metrics, such as accuracy, precision, recall, and F1 score, to identify any sudden or unexpected drops in performance that could be the result of an adversarial attack.
4. **Log Analysis:** Collect and analyze logs from the AI system and its surrounding infrastructure to identify unusual activities, patterns, or access attempts that might suggest an attack.

5. **Intrusion Detection Systems (IDS):** Implement intrusion detection systems that monitor network traffic, system activities, or application-level events to detect and report potential adversarial attacks.
6. **Runtime Verification:** Employ runtime verification techniques to ensure that the AI model's behavior adheres to a predefined set of properties or specifications, which can help detect deviations caused by adversarial attacks.
7. **Periodic Model Evaluation:** Regularly evaluate the AI model using test datasets and validation sets to assess its performance and robustness against known and unknown adversarial examples.
8. **Audit Trails:** Maintain detailed audit trails of all activities, including data access, model updates, and system configurations, to support the investigation and analysis of potential adversarial attacks.
9. **Incident Response Plan:** Develop a comprehensive incident response plan to address potential adversarial attacks, including steps to detect, analyze, contain, eradicate, and recover from an attack.
10. **Collaboration and Information Sharing:** Collaborate with other organizations, researchers, and practitioners to share information about adversarial attacks, detection techniques, and best practices for monitoring and mitigating such attacks. This can help improve the overall security posture of AI systems across the community.

Must Learn AI Security Part 4: Trojan Attacks Against AI

Chapter 4

What is a Trojan attack against AI?

Much like any type of Trojan attack in the security realm, a Trojan attack against AI is a type of cyber-attack where a malicious actor disguises a piece of malware as a legitimate software program or data file. Once the Trojan is installed on an AI system, it can give the attacker unauthorized access to the system, steal sensitive data, or cause other types of damage. In the case of AI, Trojan attacks can be particularly damaging because they can manipulate the algorithms that make decisions based on data, leading to incorrect or even dangerous outcomes.

How it works

The general steps taken in a Trojan attack against AI can vary, but here are some common steps that attackers may take:

1. **Reconnaissance:** The attacker does research on the target AI system to identify vulnerabilities and weaknesses.
2. **Delivery:** The attacker delivers a Trojan to the AI system, often through email phishing, social engineering or through infected software.
3. **Installation:** The Trojan is installed on the AI system, allowing the attacker access to the system.
4. **Command and Control:** The attacker establishes a command-and-control infrastructure to remotely control the Trojan and carry out malicious actions.
5. **Exploitation:** The attacker exploits the Trojan to carry out malicious actions, which can include stealing sensitive data,

manipulating algorithms to produce incorrect results, or causing other types of damage.

6. **Cover-up:** The attacker may attempt to cover up their tracks to avoid detection and continue their malicious activities.

These steps take a similar approach to the adversary tactics and techniques of the **MITRE ATT&CK Matrix for Enterprise**. If not already, you should become very familiar with these threat models and methodologies.

ATT&CK Matrix for Enterprise

layout: side ▾ show sub-techniques hide sub-techniques

Reconnaissance 10 techniques	Resource Development 8 techniques	Initial Access 9 techniques	Execution 14 techniques	Persistence 19 techniques	Privilege Escalation 13 techniques	Defense Evasion 42 techniques	Credential Access 17 techniques	Discovery 31 techniques	Lateral Movement 9 techniques	Collection 17 techniques	Command and Control 16 techniques
Active Scanning (2)	Acquire Access	Drive-by Compromise	Cloud Administration Command	Account Manipulation (2)	Abuse Elevation Control Mechanism (4)	Abuse Elevation Control Mechanism (4)	Adversary-in-the-Middle (2)	Account Discovery (4)	Exploitation of Remote Services	Adversary-in-the-Middle (2)	Application Layer Protocol (4)
Gather Victim Host Information (4)	Acquire Infrastructure (3)	Exploit Public-Facing Application	Command and Scripting Interpreter (9)	BITS Jobs	Access Token Manipulation (2)	Access Token Manipulation (2)	Brute Force (4)	Application Window Discovery	Internal Spearphishing	Archive Collected Data (3)	Communication Through Removable Media
Gather Victim Identity Information (2)	Compromise Accounts (3)	External Remote Services	Container Administration Command	Boot or Logon Autostart Execution (14)	Boot or Logon Autostart Execution (14)	Build Image on Host	Credentials from Password Stores (3)	Browser Information Discovery	Lateral Tool Transfer	Audio Capture	Data Encoding (2)
Gather Victim Network Information (6)	Compromise Infrastructure (7)	Hardware Additions	Deploy Container	Boot or Logon Initialization Scripts (5)	Boot or Logon Initialization Scripts (5)	Debugger Evasion	Exploitation for Credential Access	Cloud Infrastructure Discovery	Remote Service Session Hijacking (2)	Automated Collection	Data Obfuscation (3)
Gather Victim Org Information (4)	Develop Capabilities (4)	Phishing (3)	Exploitation for Client Execution	Browser Extensions	Create or Modify System Process (4)	Deobfuscate/Decode Files or Information	Forced Authentication	Cloud Service Dashboard	Remote Services (7)	Clipboard Data	Dynamic Resolution (3)
Phishing for Information (3)	Establish Accounts (3)	Replication Through Removable Media	Inter-Process Communication (2)	Compromise Client Software Binary	Domain Policy Modification (2)	Deploy Container	Forge Web Credentials (2)	Cloud Service Discovery	Replication Through Removable Media	Data from Cloud Storage	Encrypted Channel (2)
Search Closed Sources (2)	Obtain Capabilities (6)	Supply Chain Compromise (2)	Native API	Create Account (3)	Event Triggered Execution (14)	Direct Volume Access	Input Capture (4)	Cloud Storage Object Discovery	Container and Resource Discovery	Data from Configuration Repository (2)	Fallback Channels
Search Open Technical Databases (5)	Stage Capabilities (6)	Trusted Relationship	Scheduled Task/Job (5)	Create or Modify System Process (4)	Exploitation for Privilege Escalation	Execution Guardrails (1)	Modify Authentication Process (4)	Container and Resource Discovery	Debugger Evasion	Data from Information Repositories (2)	Ingress Tool Transfer
Search Open Websites/Domains (3)		Valid Accounts (4)	Serverless Execution	Event Triggered Execution (14)		Exploitation for Defense Evasion	Multi-Factor Authentication Interception	Device Driver Discovery	File and Directory Discovery	Data from Local System	Non-Application
Search Victim-Owned Websites			Shared Modules			File and Directory Permissions Modification		Domain Trust Discovery			
								File and Directory Discovery			

MITRE ATT&CK Matrix

See: [*ATT&CK Matrix for Enterprise*](#)

Types of Trojan attacks

There are different types of Trojan attacks against AI. Here are a few examples:

1. **Data Poisoning:** In this type of attack, the attacker injects incorrect or malicious data into an AI system, which can manipulate the system's decision-making process.

2. **Model Stealing:** In this type of attack, the attacker steals the AI model used by a company or organization, which can allow the attacker to replicate the model and use it for malicious purposes.
3. **Backdoor Access:** In this type of attack, the attacker gains unauthorized access to an AI system by exploiting a vulnerability or creating a backdoor.
4. **Adversarial Attacks:** In this type of attack, the attacker creates adversarial inputs that can cause an AI system to produce incorrect or unexpected outputs.
5. **Malware Injection:** In this type of attack, the attacker injects malware into an AI system through a Trojan, which can allow the attacker to control the system and carry out malicious activities.

It's important to be aware of these different types of Trojan attacks against AI and take appropriate measures to prevent them.

Why it matters

The negative results of a Trojan attack against AI can be severe and can vary depending on the type and severity of the attack. Here are some possible negative results:

1. **Data Theft:** Attackers can use Trojan attacks to steal sensitive data from an AI system, such as personal information, financial data, or intellectual property.
2. **Manipulation of Algorithms:** Attackers can use Trojan attacks to manipulate the algorithms used by an AI system, which can result in incorrect or biased decisions.
3. **System Disruption:** Trojan attacks can disrupt the functioning of an AI system, which can cause it to malfunction or stop working altogether.

4. **Financial Loss:** Trojan attacks can result in financial loss for companies or organizations, either through theft of funds or loss of revenue due to system disruption.
5. **Reputation Damage:** If a company or organization is the victim of a Trojan attack, it can damage their reputation and erode trust with customers and partners.

These negative results can have long-lasting consequences for companies or organizations that fall victim to Trojan attacks against AI, which is why it's important to take preventative measures to secure these systems.

How it might happen

A Trojan attack against AI can happen in several ways, but here are some common methods:

1. **Social Engineering:** Attackers may use social engineering tactics to trick users into downloading and installing Trojan malware, often through phishing emails or other types of social engineering attacks.
2. **Software Vulnerabilities:** Attackers may exploit vulnerabilities in software or operating systems used by an AI system to gain access and install Trojan malware.
3. **Third-Party Software:** Attackers may target third-party software components or libraries used by an AI system, which can contain vulnerabilities that can be exploited to install Trojan malware.
4. **Malicious Websites:** Attackers can use malicious websites to exploit vulnerabilities in a user's browser or operating system, which can allow them to install Trojan malware on the AI system.
5. **Physical Access:** Attackers may gain physical access to an AI system and install Trojan malware directly onto the system.

Once the Trojan malware is installed on the AI system, the attacker can use it to remotely control the system, steal data, or manipulate algorithms to produce incorrect or biased results.

Real-world Example

A real-world example of a Trojan attack against AI occurred in 2019, when researchers from the University of California, Berkeley, published a paper detailing how they inserted backdoor trojans into deep learning models. They demonstrated that an attacker could train the model to recognize a specific, seemingly innocuous trigger, like a small watermark, patch, or a specific color pattern. When the AI model encounters this trigger in the input data, it will produce a specific, predefined incorrect output, which can be controlled by the attacker.

In their experiment, the researchers inserted a backdoor into a facial recognition system. They trained the model to recognize a specific pattern of glasses on a person's face. When the AI system encountered a face with these glasses, it would incorrectly classify the person as a specific individual, regardless of their actual identity. This could be exploited to bypass security systems, falsely incriminate someone, or cause other unintended consequences.

This example highlights the risk of trojan attacks in AI systems, where an attacker can manipulate the training process or insert malicious code into the model itself, causing the system to behave in unintended and potentially harmful ways when exposed to specific triggers.

How to Mitigate

There are several ways to mitigate Trojan attacks against AI, including:

1. **Use of Antivirus and Firewall Software:** Antivirus and firewall software can help prevent Trojan malware from being installed on an AI system and can detect and block malicious activity.
2. **Regular Software Updates:** Regular software updates can help fix vulnerabilities in the software or operating system used by the AI system, making it more difficult for attackers to exploit these vulnerabilities.
3. **Strong Access Controls:** Implementing strong access controls, such as limiting user access to only what is necessary and requiring multi-factor authentication, can help prevent unauthorized access to the AI system.
4. **Employee Education:** Educating employees on how to recognize and prevent social engineering attacks, such as phishing emails, can help prevent Trojan malware from being installed on the AI system.
5. **Adversarial Training:** Adversarial training involves training an AI system to recognize and defend against adversarial attacks, such as adversarial inputs or data poisoning.
6. **Regular System Audits:** Regular system audits can help identify vulnerabilities and weaknesses in the AI system, allowing them to be addressed before they can be exploited by attackers.

By implementing these mitigation strategies, companies and organizations can better protect their AI systems from Trojan attacks and other types of cyber threats.

How to monitor

To monitor against Trojan attacks against AI, here are some steps you can take:

1. **Implement Real-Time Monitoring:** Implementing real-time monitoring of AI systems can help detect and alert security teams to any unusual activity or attempts to access the system.
2. **Implement Intrusion Detection and Prevention:** Intrusion detection and prevention systems can help detect and prevent unauthorized access to AI systems, including Trojan attacks.
3. **Use Machine Learning:** Machine learning can be used to detect anomalies in the behavior of an AI system and flag any suspicious activity that could be indicative of a Trojan attack.
4. **Conduct Regular Penetration Testing:** Regular penetration testing can help identify vulnerabilities in an AI system, allowing them to be addressed before they can be exploited by attackers.
5. **Monitor Network Traffic:** Monitoring network traffic can help detect any attempts to exfiltrate data from an AI system or any suspicious activity that could be indicative of a Trojan attack.
6. **Implement User Behavior Analytics:** User behavior analytics can help detect any unusual or suspicious behavior by users of an AI system, which could be indicative of a Trojan attack.

By implementing these monitoring strategies, companies and organizations can better protect their AI systems from Trojan attacks and other types of cyber threats. It's important to continually evaluate and update monitoring strategies to ensure that they are effective and up to date with the latest threats.

What to capture

To identify when a Trojan attack against AI is happening, you should capture the following types of data:

1. **Network Traffic:** Monitoring network traffic can help detect any unusual traffic patterns that could be indicative of a

Trojan attack. This includes capturing data on the volume and frequency of data transfers, the source and destination IP addresses, and the type of data being transferred.

2. **System Logs:** System logs can provide valuable information on user activity, system performance, and security events. Capturing data on user logins, system activity, and system errors can help detect any unusual or suspicious activity that could be indicative of a Trojan attack.
3. **User Behavior Analytics:** Capturing data on user behavior, such as the types of files accessed, the frequency of access, and the times of day when access occurs, can help detect any unusual or suspicious behavior that could be indicative of a Trojan attack.
4. **AI Model Performance Metrics:** Capturing data on the performance of an AI model, such as accuracy, precision, and recall, can help detect any unusual or unexpected changes in the performance of the model that could be indicative of a Trojan attack.
5. **Security Alerts:** Capturing data on security alerts generated by intrusion detection and prevention systems, firewalls, and antivirus software can help detect any attempted or successful Trojan attacks.

By capturing and analyzing this data, companies and organizations can better detect and respond to Trojan attacks against AI, helping to mitigate their impact and reduce the risk of data theft, system disruption, and other negative consequences.

To prevent Trojan attacks against AI, it's important to maintain strong cybersecurity practices, including regular software updates, strong passwords, and employee education about phishing and social engineering tactics. To prevent Trojan attacks, it's important to keep software up-to-date, use strong passwords, and be cautious when downloading files from unknown sources. This can include

using secure coding practices, regularly updating software and systems, and implementing strong access controls and monitoring.

Must Learn AI Security Part 5: Evasion Attacks Against AI

Chapter 5

What is an Evasion attack against AI?

An Evasion attack against AI involves attempting to bypass or deceive an AI system's defenses or detection mechanisms in order to manipulate or exploit the system. This can be achieved through techniques such as altering or obscuring input data, using adversarial examples, or employing other tactics that make it difficult for the AI system to accurately classify or make decisions based on the input. Evasion attacks can be a serious security concern in applications such as cybersecurity, fraud detection, and autonomous vehicles.

How it works

An evasion attack against AI typically involves the following steps:

1. **Adversary identifies the target AI system and its vulnerabilities:** The attacker first identifies the target AI system and analyzes its vulnerabilities. They may also try to gather information about the system's algorithms and the types of data it uses to make decisions.
2. **Adversary generates adversarial examples:** The attacker generates adversarial examples, which are inputs that have been specifically crafted to deceive the AI system. These examples are designed to look similar to legitimate inputs, but with subtle modifications that cause the AI system to misclassify or produce incorrect outputs.
3. **Adversary submits adversarial examples to the AI system:** The attacker then submits the adversarial examples to

the AI system, either directly or by embedding them in legitimate data.

4. **AI system produces incorrect output:** When the AI system processes the adversarial examples, it produces incorrect outputs. This can have serious consequences, depending on the application of the AI system. For example, in the case of a fraud detection system, an evasion attack could allow a fraudster to bypass the system and carry out fraudulent activities undetected.
5. **Adversary refines the attack:** If the initial attack is unsuccessful, the attacker may refine their approach by using more sophisticated techniques or by testing the AI system's responses to different types of inputs.

Evasion attacks against AI can be difficult to detect and defend against, as attackers can use a variety of techniques to evade detection and deceive the system.

Types of Evasion attacks

Here are some common types of evasion attacks against AI:

1. **Adversarial examples:** Adversarial examples are inputs that have been specifically crafted to deceive an AI system. These examples are designed to look similar to legitimate inputs, but with subtle modifications that cause the AI system to misclassify or produce incorrect outputs.
2. **Input perturbation:** Input perturbation involves adding noise or random perturbations to the input data in order to bypass the AI system's detection mechanisms.
3. **Feature manipulation:** Feature manipulation involves modifying the input data in a way that changes the features or attributes that the AI system uses to make decisions. This can

be done in a way that is difficult to detect or that causes the AI system to misclassify the input.

4. **Model inversion:** Model inversion involves using the output of an AI system to reverse-engineer the model and learn sensitive information about the data that was used to train the model.
5. **Data poisoning:** Data poisoning involves injecting malicious data into the training data used to train an AI system. This can cause the AI system to learn incorrect or biased models, which can be exploited by attackers.

These are just a few examples of the many types of evasion attacks that can be carried out against AI systems. As AI technology continues to advance, it is likely that attackers will develop new and more sophisticated techniques for evading detection and exploiting vulnerabilities in AI systems.

Why it matters

Evasion attacks against AI can have various negative results, including:

1. **Compromised security:** Evasion attacks can compromise the security of AI systems, making it easier for attackers to carry out malicious activities such as data theft, fraud, and cyberattacks.
2. **Inaccurate decisions:** Evasion attacks can cause AI systems to make inaccurate decisions, which can have serious consequences in applications such as healthcare, finance, and autonomous vehicles.
3. **Bias and discrimination:** Evasion attacks can be used to introduce bias and discrimination into AI systems, which can have negative impacts on individuals or groups that are unfairly targeted or excluded.

4. **Reduced trust in AI:** Evasion attacks can reduce public trust in AI technology by highlighting its vulnerabilities and limitations.
5. **Higher costs and reduced efficiency:** Evasion attacks can increase the costs of developing and deploying AI systems by requiring additional resources to detect and defend against attacks. They can also reduce the efficiency of AI systems by introducing errors and false positives that require additional human intervention to correct.

Evasion attacks against AI pose a serious threat to the security, accuracy, and fairness of AI systems, and it is important to develop effective defenses and detection mechanisms to mitigate these risks.

Why it might happen

An attacker might use an evasion attack against AI for various purposes, such as:

1. **Malicious intent:** An attacker might use an evasion attack to carry out a malicious activity such as stealing sensitive data, bypassing security systems, or disrupting critical infrastructure.
2. **Financial gain:** An attacker might use an evasion attack to gain financial advantage by manipulating AI models used in trading or investments.
3. **Privacy violation:** An attacker might use an evasion attack to violate an individual's privacy by manipulating AI models used for personal identification or profiling.
4. **Competitive advantage:** An attacker might use an evasion attack to gain a competitive advantage by manipulating AI models used in business operations such as pricing, product recommendations, or demand forecasting.

5. **Research purposes:** An attacker might use an evasion attack to conduct research on the vulnerabilities and limitations of AI systems.

The motivations behind an evasion attack against AI can vary depending on the attacker's goals and objectives. However, regardless of the attacker's intent, evasion attacks can have serious negative consequences for the security, accuracy, and fairness of AI systems.

Real-world Example

A real-world example of an evasion attack against AI is the case of stickers being used to trick an AI-powered self-driving car. In 2018, researchers at Tencent Keen Security Lab demonstrated an evasion attack against Tesla's Autopilot system. They placed small, innocuous-looking stickers on the road surface in a specific pattern. These stickers confused the AI system, causing it to identify a non-existent lane and subsequently follow an incorrect path.

This example highlights how small, targeted changes in the input data (in this case, the stickers on the road) can manipulate the behavior of an AI system, potentially leading to dangerous consequences. Such attacks exploit the vulnerabilities in AI models and can be used to deceive the system into making incorrect decisions or taking unintended actions.

How to Mitigate

There are several ways an organization can mitigate an evasion attack against AI. Here are some examples:

1. **Robust defenses:** Organizations can deploy robust defenses such as intrusion detection and prevention systems, firewalls, and antivirus software to detect and prevent attacks.
2. **Regular vulnerability assessments:** Organizations can perform regular vulnerability assessments to identify vulnerabilities and weaknesses in their AI systems.
3. **Data integrity checks:** Organizations can implement data integrity checks to ensure that the data used to train and test AI models is accurate and free from manipulation.
4. **Adversarial training:** Organizations can use adversarial training techniques to train AI models to recognize and defend against adversarial attacks.
5. **Defense in depth:** Organizations can use a defense-in-depth approach, which involves layering multiple defenses to provide redundancy and increase the difficulty of evading detection.
6. **Human oversight:** Organizations can incorporate human oversight into their AI systems to provide an additional layer of defense against adversarial attacks.
7. **Regular updates and patches:** Organizations should keep their AI systems up to date with the latest security patches and updates to mitigate known vulnerabilities.

Mitigating an evasion attack against AI requires a proactive approach that involves a combination of technical solutions, process improvements, and human oversight. By implementing these measures, organizations can reduce the risk of successful evasion attacks and increase the security, accuracy, and fairness of their AI systems.

How to monitor

Organizations can monitor against evasion attacks against AI in several ways. Here are some examples:

1. **Anomaly detection:** Organizations can use anomaly detection techniques to identify deviations from normal behavior in AI systems. This can help detect abnormal inputs or outputs that may indicate an evasion attack.
2. **Model monitoring:** Organizations can monitor the performance of AI models to detect changes in their behavior that may indicate an evasion attack.
3. **Data lineage tracking:** Organizations can track the lineage of data used in AI models to detect any changes or manipulations to the data that could result in an evasion attack.
4. **Adversarial testing:** Organizations can conduct adversarial testing to identify vulnerabilities and weaknesses in their AI systems.
5. **Network monitoring:** Organizations can monitor network traffic to detect any suspicious activity that may indicate an evasion attack.
6. **Human review:** Organizations can incorporate human review into their AI systems to provide an additional layer of defense against evasion attacks.
7. **Continuous evaluation:** Organizations can continuously evaluate the performance of their AI systems to ensure that they are functioning as intended and to detect any anomalies or deviations from normal behavior.

Monitoring against evasion attacks requires a proactive approach that involves a combination of technical solutions and human oversight.

What to capture

During monitoring to detect evasion attacks against AI, several things should be captured. Here are some examples:

1. **Input data:** The input data used to train and test the AI models should be captured and monitored to detect any changes or manipulations that may indicate an evasion attack.
2. **Output data:** The output data produced by the AI models should be captured and monitored to detect any anomalies or deviations from normal behavior that may indicate an evasion attack.
3. **Model behavior:** The behavior of the AI models should be captured and monitored to detect any changes or deviations from normal behavior that may indicate an evasion attack.
4. **Network traffic:** Network traffic should be captured and monitored to detect any suspicious activity that may indicate an evasion attack.
5. **System logs:** System logs should be captured and monitored to detect any unusual or abnormal activity that may indicate an evasion attack.
6. **Adversarial testing results:** The results of adversarial testing should be captured and monitored to identify vulnerabilities and weaknesses in the AI systems that may be exploited by attackers.

Capturing and monitoring these types of data can help organizations detect and respond to evasion attacks against AI in a timely manner, reducing the risk of successful attacks and increasing the security, accuracy, and fairness of their AI systems.

Must Learn AI Security Part 6: Model Inversion Attacks Against AI

Chapter 6

What is a Model Inversion attack against AI?

A Model Inversion attack against AI refers to the process where an attacker attempts to reconstruct the original data used for training a machine learning model by only having access to the model's output. This type of attack poses a significant risk to the privacy of the data used for training the model.

How it works

A Model Inversion attack against AI works by exploiting the information leakage from the machine learning model's outputs to reconstruct or approximate the original training data. An attacker uses the model's predictions and confidence scores to iteratively refine their input to generate a close approximation of the original data.

Here's a step-by-step explanation of how a Model Inversion attack works:

1. **Access to the model:** The attacker needs access to the AI model, which can be through a public API, a stolen copy of the model, or any other means of interacting with the model's predictions.
2. **Identifying target:** The attacker selects a target individual or data point whose information they want to reconstruct from the model.

3. **Generating initial input:** The attacker starts with an initial input that could be random or based on some prior knowledge of the target domain.
4. **Analyzing model outputs:** The attacker inputs the generated data into the model and collects the model's predictions and confidence scores.
5. **Refining input:** Using the information from the model's outputs, the attacker iteratively refines the input data to maximize the model's confidence in the target label or class. This process involves optimization techniques like gradient descent or genetic algorithms.
6. **Convergence:** The attacker repeats steps 4 and 5 until the input converges to a close approximation of the target data point, or the confidence scores reach a certain threshold.
7. **Reconstruction:** The attacker now has a data point that closely resembles the original training data point, effectively compromising the privacy of the target individual or data point.

It's worth noting that Model Inversion attacks are more likely to be successful in cases where the model is overfitted, as it may have memorized specific training examples.

Types of Model Inversion attacks

There are two primary types of Model Inversion attacks against AI: the black-box attack and the white-box attack. Both types aim to reconstruct or approximate the original training data, but they differ in the level of access the attacker has to the AI model.

1. **Black-box Model Inversion attack:** In this type of attack, the attacker only has access to the model's input-output pairs, meaning they can input data and receive the corresponding predictions. However, they have no knowledge of the model's

architecture, parameters, or the training data. The attacker generates inputs, analyzes the model's outputs, and iteratively refines the inputs based on the obtained information. This process continues until the attacker is able to approximate the original training data.

2. **White-box Model Inversion attack:** In a white-box attack, the attacker has more extensive access to the AI model, including its architecture, parameters (such as weights and biases), and possibly partial knowledge of the training data. This additional information allows the attacker to exploit the inner workings of the model more effectively and reconstruct the original training data with higher accuracy. In this scenario, the attacker can use gradient-based optimization techniques to maximize the model's confidence in the target label or class, ultimately converging to a close approximation of the target data point.

Both black-box and white-box Model Inversion attacks pose significant threats to the privacy of the data used in training AI models.

Why it matters

A Model Inversion attack against AI can lead to several negative consequences, mainly related to the breach of data privacy and potential misuse of sensitive information.

Some of these negative effects include:

1. **Privacy violation:** The primary concern of a Model Inversion attack is the potential exposure of sensitive information contained in the original training data. This can lead to a violation of individuals' privacy rights and cause harm to those whose data has been reconstructed.

2. **Identity theft:** In cases where the AI model involves personally identifiable information (PII), such as facial recognition or biometric data, a successful Model Inversion attack can lead to identity theft. Attackers may use the reconstructed data to impersonate individuals or gain unauthorized access to personal accounts and services.
3. **Loss of trust:** Model Inversion attacks can undermine trust in AI systems, as individuals and organizations become concerned about the security and privacy risks associated with using AI models trained on their data.
4. **Legal and regulatory issues:** Companies that experience a Model Inversion attack may face legal and regulatory consequences if they fail to protect users' data privacy according to established laws and regulations, such as the General Data Protection Regulation (GDPR) in the European Union.
5. **Misuse of sensitive information:** If the reconstructed data contains sensitive information, such as medical records or financial data, attackers could use this information for malicious purposes, including extortion, fraud, or targeted advertising.
6. **Damage to reputation:** An organization that suffers a Model Inversion attack may experience damage to its reputation, as users and clients may view the organization as having inadequate security and privacy measures in place.

Implementing robust security measures and following best practices can help protect AI systems from potential attacks.

Why it might happen

There are several reasons why someone might perform a Model Inversion attack against AI, which can range from monetary gain to competitive advantage or even just curiosity.

Some of these motivations include:

1. **Financial gain:** Attackers might attempt a Model Inversion attack to acquire sensitive information such as credit card details, social security numbers, or other financial data, which they could use for fraudulent activities or sell on the dark web.
2. **Identity theft:** In cases where the AI model involves personally identifiable information (PII), such as facial recognition or biometric data, attackers could use the reconstructed data to impersonate individuals, gain unauthorized access to personal accounts, or commit various forms of identity theft.
3. **Competitive advantage:** Competitors might perform a Model Inversion attack to gain insights into a company's proprietary data or trade secrets, which could give them a competitive advantage in the market.
4. **Corporate espionage:** Attackers might use Model Inversion attacks to gather sensitive information about a company's business strategies, product plans, or customer data, which could be used for corporate espionage or market manipulation.
5. **Curiosity or intellectual challenge:** Some attackers might be driven by curiosity or the intellectual challenge of successfully performing a Model Inversion attack, rather than having a specific malicious intent.
6. **Exposing vulnerabilities:** In some cases, security researchers or ethical hackers might perform a Model Inversion attack to demonstrate the potential vulnerabilities in an AI system and encourage the development of more secure and privacy-preserving AI models.

Real-world Example

Consider a facial recognition AI model that has been trained using a large dataset of individuals' images, along with their corresponding names. The AI model can recognize a person's face and output the

person's name when given an input image.

An attacker, with no access to the original dataset, wants to uncover the image of a specific person, say, Meredith. The attacker starts by inputting random images into the model and analyzing the output probabilities of the model recognizing Meredith. By iteratively refining the input images and optimizing them based on the model's output probabilities, the attacker can eventually generate an image that closely resembles Meredith's face.

In this example, the attacker has successfully performed a Model Inversion attack, compromising the privacy of Meredith's image without having direct access to the original dataset. This highlights the importance of privacy-preserving techniques like differential privacy and secure multi-party computation in the development of AI models.

How to Mitigate

To mitigate the risk of Model Inversion attacks, privacy-preserving techniques like differential privacy, federated learning, and secure multi-party computation can be employed in the development and deployment of AI models.

How to monitor/What to capture

Detecting a Model Inversion attack against AI can be challenging, as the attacker often has limited access to the model and may not leave obvious traces. However, there are several indicators and activities you can monitor to detect potential Model Inversion attacks:

1. **Unusual query patterns:** Monitor the usage patterns of your AI model to identify any unusual or suspicious behavior. For

instance, a high number of queries from a single source or an unexpected increase in queries during specific time periods could indicate a potential attack.

2. **Atypical inputs:** Keep an eye on the inputs provided to the model. If you notice a series of inputs that seem unusual, random, or unrelated to the typical use case, it could be an attempt to perform a Model Inversion attack.
3. **High-confidence incorrect predictions:** If your model starts generating high-confidence predictions that are incorrect or don't align with the expected output, it could be a sign that someone is trying to reverse-engineer the model by refining inputs based on the model's output probabilities.
4. **Access logs and user behavior:** Regularly review access logs and user behavior to identify any unauthorized or suspicious access to the AI model or attempts to exfiltrate model parameters.
5. **Rate-limiting violations:** Implement rate-limiting on your AI model's API to prevent excessive queries in a short period. Monitor for violations of these rate limits, as they could indicate an attacker trying to gather information for a Model Inversion attack.
6. **Multiple failed login attempts:** Track failed login attempts and unauthorized access attempts to your AI system, as attackers may try to gain access to the model itself or related resources to facilitate a Model Inversion attack.

Model Inversion attacks against AI involve exploiting information leakage from machine learning models to reconstruct or approximate the original training data. These attacks can compromise data privacy and lead to various negative consequences, such as identity theft, loss of trust, legal issues, and misuse of sensitive information.

Must Learn AI Security Part 7: Membership Inference Attacks Against AI

Chapter 7

What is a Membership Inference attack against AI?

A Membership Inference Attack against AI refers to a type of privacy breach where an attacker tries to determine if a specific data point was part of the training dataset used to build a machine learning model. In this attack, the adversary queries the AI model and analyzes the output, such as the model's confidence in its predictions, to infer whether the data point was included in the training data or not.

How it works

A Membership Inference Attack against AI typically happens in the following steps:

1. **Data collection:** The attacker gathers data samples that they believe could be part of the target AI model's training dataset. They may also collect additional data samples that are unlikely to be part of the training data.
2. **Model access:** The attacker needs to have query access to the AI model, either through an API or by interacting with a service that uses the model. They do not need direct access to the model's parameters or the actual training dataset.
3. **Creating a shadow model:** The attacker trains a "shadow model" using their collected data, attempting to replicate the target AI model's behavior. They may create multiple shadow models with different subsets of data to improve their chances of success.

4. **Analyzing model outputs:** The attacker queries the target AI model and their shadow models with their collected data samples. They analyze the model outputs, such as prediction confidence scores or class probabilities, to identify patterns that may indicate membership in the training dataset.
5. **Inference:** Based on the analysis of the model outputs, the attacker makes an educated guess about whether a specific data point was part of the training dataset or not. If their inference is accurate, they have successfully executed a membership inference attack.

It is important to note that the success of a Membership Inference Attack depends on various factors, such as the target model's architecture, the quality of the attacker's shadow models, and the availability of sufficient data samples for analysis.

Types of Membership Inference attacks

There are several types of Membership Inference Attacks against AI, which can be broadly categorized into two classes: *passive* attacks and *active* attacks.

1. **Passive attacks:** In passive attacks, the attacker only relies on the available information and their observations of the AI model's behavior. They do not try to manipulate the model or its training process. Passive attacks can be further divided into:
 - a. **Black-box attacks:** The attacker has no knowledge of the model's architecture, parameters, or training data, and only has access to the model's input-output behavior through an API or a service. They use this limited information to create shadow models and infer membership.

- b. **White-box attacks:** The attacker has more information about the target AI model, such as its architecture and parameters. This additional information can help the attacker create better shadow models and improve the accuracy of their membership inference.
- 2. **Active attacks:** In active attacks, the attacker tries to manipulate the AI model's training process or its behavior to gain insights into the training data. Some examples of active attacks include:
 - a. **Data poisoning:** The attacker injects carefully crafted data samples into the model's training data, aiming to influence the model's behavior and make it easier to infer membership.
 - b. **Model inversion:** The attacker exploits the model's parameters or architecture to recreate or approximate the training data, which can then be used to perform a membership inference attack.

Each type of Membership Inference Attack has its own challenges and success rates, depending on factors such as the target model's architecture, the quality of the attacker's shadow models, and the availability of data samples for analysis.

Why it matters

Membership Inference Attacks against AI can lead to several negative consequences, including:

1. **Privacy violations:** If an attacker can successfully infer that a specific data point was part of an AI model's training dataset, they may be able to reveal sensitive information about individuals or organizations. This could include personal information such as health records, financial data, or social

media activity, potentially leading to identity theft, discrimination, or other privacy breaches.

2. **Data leakage:** A successful attack can expose proprietary or confidential information that a company or organization intended to keep secret. This could compromise trade secrets, intellectual property, or other valuable data, leading to financial losses or reputational damage.
3. **Regulatory and legal risks:** Privacy breaches resulting from Membership Inference Attacks could lead to non-compliance with data protection regulations, such as the General Data Protection Regulation (GDPR) or the California Consumer Privacy Act (CCPA). Non-compliance could result in fines, legal action, and reputational damage for the affected organization.
4. **Erosion of trust:** Users and stakeholders may lose trust in AI systems and the organizations that develop and deploy them if they believe that their privacy is not being adequately protected. This loss of trust could hinder the adoption of AI technologies and limit their potential benefits.
5. **Effects on data sharing:** Concerns about the potential for Membership Inference Attacks may discourage individuals and organizations from contributing data to AI projects, limiting the availability of high-quality training data and hindering AI research and development.

To minimize these negative consequences, it is essential for AI developers and organizations to implement privacy-preserving techniques, such as differential privacy, federated learning, and secure multi-party computation, and to follow best practices for data protection and model development.

Why it might happen

When an attacker successfully performs a Membership Inference Attack against AI, they can gain valuable information and insights, such as:

1. **Membership status:** The primary goal of the attack is to determine whether a specific data point was part of the AI model's training dataset. Knowing this information may be valuable in itself, especially if the data is sensitive or confidential.
2. **Privacy-sensitive information:** If the attacker can infer membership, they may be able to expose sensitive information about individuals or organizations associated with the data points. This could include personal details, health records, financial data, or other private information that could be exploited for malicious purposes, such as identity theft or targeted attacks.
3. **Proprietary or confidential data:** Successful attacks can reveal proprietary or confidential information that a company or organization intended to keep secret. Attackers could use this information for corporate espionage, intellectual property theft, or to gain a competitive advantage.
4. **Insight into AI model's behavior:** By analyzing the AI model's responses during the attack, the attacker may gain insights into the model's behavior, weaknesses, and potential biases. This information could be used to launch further attacks or exploit vulnerabilities in the AI system.
5. **[Evasion](#) and [Adversarial](#) attacks:** Information obtained from a successful Membership Inference Attack can potentially be used to craft adversarial examples or devise evasion strategies that target the AI model's specific weaknesses, making it more difficult for the model to detect or classify the attacker's malicious inputs.

Overall, a successful Membership Inference Attack can provide the attacker with valuable information and insights that they can exploit for various malicious purposes or gain a strategic advantage.

Real-world Example

While there haven't been many publicly reported real-world cases of successful Membership Inference Attacks, researchers have demonstrated the feasibility of such attacks in various experimental settings. One notable example is the study conducted by Shokri et al. in 2017, titled "[Membership Inference Attacks Against Machine Learning Models](#)."

In this study, the researchers demonstrated how an attacker could perform Membership Inference Attacks against machine learning models trained on real-world datasets, including the CIFAR-100 image classification dataset and the Adult Income dataset from the UCI Machine Learning Repository. The researchers used black-box attacks, meaning they had no knowledge of the target models' architecture or parameters and only had access to their input-output behavior.

The attack involved creating shadow models to mimic the target models and analyzing the prediction confidence scores to infer membership. The researchers found that their attacks were successful in determining whether a data point was part of the training dataset with significantly higher accuracy than random guessing. This study illustrated the potential risks associated with Membership Inference Attacks and the importance of adopting privacy-preserving techniques to protect sensitive data used in AI systems.

While the study serves as a theoretical example, it highlights the

potential real-world risks that AI systems might face if they do not implement adequate privacy protections.

How to Mitigate

To mitigate Membership Inference Attacks against AI, developers and organizations can employ several techniques and best practices to protect sensitive data and enhance model privacy:

1. **Differential privacy:** Implementing differential privacy adds controlled noise to the model's outputs or during the training process, making it difficult for attackers to infer membership based on the model's responses. This technique can help protect the privacy of individual data points without significantly compromising the model's accuracy.
2. **Federated learning:** In federated learning, the AI model is trained on decentralized data sources without requiring the data to be centralized. This approach reduces the risk of membership inference attacks, as the attacker will have limited access to the distributed data and the model's global parameters.
3. **Model generalization:** Improve the generalization of the AI model by using techniques such as early stopping, regularization, and dropout during training. A model with better generalization is less likely to overfit to the training data and leak information about individual data points.
4. **Limit model access:** Restrict the number of queries or rate at which users can access the AI model. This can make it more difficult for an attacker to gather enough information to perform a successful attack.
5. **Monitoring and auditing:** Regularly monitor and audit the AI model's behavior to detect any anomalies or signs of potential attacks. This can help identify and respond to threats proactively.

6. **Data anonymization:** Remove or anonymize personally identifiable information (PII) from the training dataset to reduce the risk of privacy breaches and limit the potential impact of a successful attack.
7. **Secure multi-party computation:** Use secure multi-party computation techniques to protect the privacy of data during the training process. This approach allows multiple parties to collaboratively train an AI model without revealing their individual data.
8. **Train multiple models:** Instead of using a single model, consider training multiple models on different subsets of data. This can make it more challenging for an attacker to perform a successful attack, as they would need to attack multiple models to gain the desired information.

By implementing these techniques and best practices, developers and organizations can significantly reduce the risk of Membership Inference Attacks against AI systems and better protect sensitive data and user privacy.

How to monitor/What to capture

To identify an active Membership Inference Attack, you should monitor and audit various aspects of the AI model's behavior, user access, and system performance. Here are some key indicators to watch for:

1. **Unusual query patterns:** Keep an eye on the rate, volume, and type of queries made to the AI model. An unusually high number of queries or a sudden spike in requests may indicate an attacker is probing the model.
2. **High-confidence predictions on unusual inputs:** If the model produces high-confidence predictions on atypical inputs or synthetic data, it may suggest an attacker is testing the

model's behavior to gather information for a Membership Inference Attack.

3. **Repeated queries with slight variations:** Monitor for repeated queries with slight variations in input data, which could indicate an attacker is trying to understand the model's decision boundaries or confidence scores.
4. **Unusual user access patterns:** Track user access logs to identify any unusual patterns, such as unauthorized access, multiple failed login attempts, or access from suspicious IP addresses.
5. **Anomalies in model performance:** Keep track of the AI model's performance metrics, such as accuracy, precision, and recall, to identify any unexpected fluctuations or anomalies that could be linked to an attack.
6. **Unusual data access patterns:** Monitor data access logs to detect any irregularities in data access patterns, such as unauthorized access to training data or attempts to inject malicious data into the training set.
7. **System resource usage:** Monitor system resources, such as CPU, memory, and network usage, to identify any unusual spikes or patterns that may indicate an ongoing attack.
8. **[Model inversion](#) or [data poisoning](#) attempts:** Look for signs of model inversion or data poisoning, where an attacker tries to manipulate the AI model's training process or exploit its parameters to recreate or approximate the training data.

By monitoring these indicators and setting up alerts for suspicious activity, you can proactively detect and respond to potential Membership Inference Attacks, helping to protect your AI system and the sensitive data it relies on.

LEARNING ABOUT SECURITY AND AI IS LIKE HAVING YOUR CAKE AND EATING IT TOO – A FABULOUS COMBO THAT EVERYONE NEEDS IN THEIR LIVES. HERE ARE A FEW REASONS WHY THIS TOPIC IS SO IMPORTANT:

CYBERSECURITY THREATS ARE THE NEW BLACK: WITH TECHNOLOGY ADVANCING FASTER THAN YOU CAN SAY "HOLD MY LATTE," CYBERCRIMINALS ARE FINDING NEW, INNOVATIVE WAYS TO WREAK HAVOC. GETTING YOUR HANDS DIRTY IN THE WORLD OF SECURITY AND AI HELPS YOU STAY ONE STEP AHEAD OF THOSE PESKY HACKERS AND PROTECT YOUR PRECIOUS DATA.

PRIVACY, PLEASE: UNLESS YOU WANT YOUR PERSONAL INFORMATION TO BECOME THE TALK OF THE DIGITAL TOWN, YOU BETTER KNOW A THING OR TWO ABOUT SECURING YOUR AI SYSTEMS. KNOWLEDGE IS POWER, AND POWER KEEPS YOUR PRIVATE LIFE AS PRIVATE AS IT SHOULD BE.

THE RISE OF THE AI EMPIRE: AI IS PRACTICALLY EVERYWHERE THESE DAYS, FROM YOUR PHONE'S VOICE ASSISTANT TO SELF-DRIVING CARS. BY UNDERSTANDING HOW SECURITY AND AI INTERSECT, YOU'LL BE BETTER EQUIPPED TO NAVIGATE THIS BRAVE NEW WORLD AND MAKE INFORMED DECISIONS, BOTH PERSONALLY AND PROFESSIONALLY.

ETHICAL CONUNDRUMS GALORE: AS AI CONTINUES TO STRUT ITS STUFF, ETHICAL QUESTIONS ARE POPPING UP LIKE NOBODY'S BUSINESS. LEARNING ABOUT SECURITY AND AI HELPS YOU WRAP YOUR HEAD AROUND THESE COMPLEX ISSUES AND PARTICIPATE IN MEANINGFUL CONVERSATIONS THAT SHAPE THE FUTURE OF TECHNOLOGY.

EMBRACE THE FABULOUS WORLD OF SECURITY AND AI, AND PIVOT CONFIDENTLY INTO A FUTURE WHERE YOU'RE IN CONTROL, WELL-INFORMED, AND READY TO KICK SOME CYBER-BUTT!