# Adversarial Attacks on Different DL Models

Monday, 2 September 2024     9:57 PM

**Understanding Adversarial Examples**

Adversarial examples are carefully crafted inputs designed to deceive deep neural networks into making incorrect predictions. These inputs are typically derived from normal, unaltered data by introducing subtle perturbations (noise/change) — tiny modifications that are often imperceptible to the human eye. Despite these changes being nearly invisible to humans, they can drastically alter the output of a neural network, causing it to misclassify the input or produce erroneous results. The key characteristic of adversarial examples is that they exploit the weaknesses in the model's decision boundaries, leading to confident yet incorrect predictions. This vulnerability poses a significant challenge to the reliability of deep learning models, particularly in critical applications.
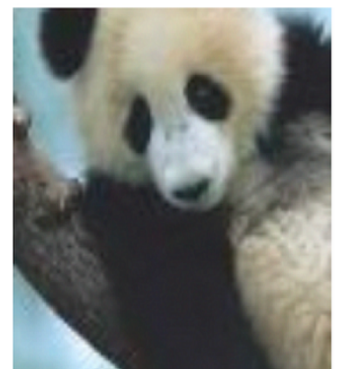


$$x$$
"panda"

$$+\,.007\times$$

$$\text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$
"nematode"

$$=$$

$$\boldsymbol{x} + \epsilon\,\text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$
"gibbon"

57.7% confidence               8.2% confidence               99.3 % confidence

Types of Adversarial Attacks: An Overview

The research paper categorizes adversarial attacks into several types based on different criteria, including the attacker's knowledge, the nature of the target model, and the goals of the attack. Here's an overview of the main types:

1. **White-box Attacks**: In white-box attacks, the attacker has complete knowledge of the target neural network, including its architecture, parameters, and training data. This allows the attacker to precisely calculate the gradients and create highly effective adversarial examples that can easily fool the model. Common methods in this category include the Fast Gradient Sign Method (FGSM) and the Jacobian-based Saliency Map Attack (JSMA).

2. **Black-box Attacks**: In black-box attacks, the attacker has no direct access to the internal workings of the target model. Instead, they can only observe the model's outputs in response to certain inputs and create adversarial examples. The attacker might use a substitute model to generate adversarial examples and then transfer these to the target model, exploiting the transferability property of adversarial examples. Methods like Zeroth Order Optimization (ZOO) are designed for this scenario.

3. **Targeted Attacks**: Targeted attacks are designed to mislead the neural network into classifying an input as a specific and give incorrect label. For example, in an image recognition task, a targeted attack might modify an

4. **Non-targeted Attacks**: Non-targeted attacks are opposite to targeted

image of a dog so that the network classifies it as a cat.

4.  **Non-targeted Attacks**: Non-targeted attacks are opposite to targeted attacks which aim to cause any misclassification, without targeting a specific incorrect label. The goal here is simply to degrade the accuracy of the model by pushing inputs outside of their correct classifications. This type of attack is generally easier to perform, as there is more flexibility in the resulting misclassification.

Methods for Generating Adversarial Examples

The paper outlines several methods for generating adversarial examples, each with unique approaches and objectives.

Here are some of the most prominent methods:

**Fast Gradient Sign Method (FGSM)**:

FGSM is one of the earliest and most well-known techniques for generating adversarial examples. This method works by calculating the gradient of the loss function with respect to the input data and then applying a small perturbation in the direction of the gradient. The perturbation is designed to maximize the model's error, effectively pushing the input towards a decision boundary that results in a misclassification. The adversarial example is generated using the formula:

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

## Projected Gradient Descent (PGD)

# Projected Gradient Descent (PGD)

Core of machine learning optimization lies the fundamental concept of gradient descent. This iterative algorithm fine-tunes model parameters to minimize a given loss function. Mathematically, the update rule is expressed as $\Theta_{t+1} = \Theta_t - \alpha \cdot \nabla J(\Theta_t)$, where $\Theta_t$ represents the parameters at iteration $t$, $\alpha$ is the learning rate, and $\nabla J(\Theta_t)$ is the gradient of the loss function.

Projected Gradient Descent (PGD) builds upon this foundation, introducing thoughtful constraints to enhance its effectiveness in crafting adversarial examples. In the context of adversarial attacks, the objective is to perturb input data to deliberately mislead the model.

PGD incorporates a perturbation budget ($\epsilon$) and a step size ($\alpha$) to control the amount and direction of perturbation. The update rule for PGD is defined as $x'_{t+1} = \Pi(x_t + \alpha \cdot sign(\nabla_x J(\Theta, x_t, y)))$, where, $x_t$ is the input at iteration $t$, $\alpha$ is the step size, $\nabla_x J(\Theta, x_t, y)$ is the gradient of the loss with respect to the input, and $\Pi$ is the projection operator ensuring perturbed input stays within predefined bounds. Understanding this working mechanism is pivotal. It not only deepens our grasp of how models learn but also sheds light on the nuanced process of adversarial example generation through the lens of PGD. To see a practical implementation of PGD and adversarial

attacks, refer to our example notebook [here](#).

**What is PGD Attack?**

Projected Gradient Descent is an iterative method used in adversarial machine learning to create adversarial examples. These are inputs to machine learning models that an attacker has intentionally designed to cause the model to make a mistake. PGD is known for its effectiveness and sophistication compared to simpler methods like the Fast Gradient Sign Method (FGSM).

**How Does PGD Work?**

PGD operates by applying small but iteratively adjusted perturbations to the input data, aimed at maximizing the model's prediction error. It's a more refined approach where the attacker can control the size of these perturbations and the number of iterations, allowing a careful balance between making the adversarial examples effective and keeping them undetectable.

**Key Features of PGD:**

- Iterative Nature: Allows thorough exploration of the model's loss surface.

- Fine-grained Control: Adjustments in step size and iterations for balance.

- Evasion of Defenses: Capability to bypass defences designed for simpler attacks.

- Flexibility: Supports both non-targeted and targeted attacks.

**Importance in Adversarial Machine Learning**

PGD serves as a benchmark in the security of machine learning models. Its ability to craft more sophisticated adversarial examples makes it a standard tool for testing model robustness. Understanding and defending against PGD attacks is crucial for developing models that are secure in real-world applications.

**Defending Against PGD Attacks**

Defence against PGD is challenging but essential. Strategies include adversarial training, input transformations, and robust optimization techniques. Developing models resilient to PGD attacks ensures better security against a wide range of adversarial threats.

**Conclusion**

PGD attacks represent a significant challenge in the field of adversarial machine learning,

**Conclusion**

offering both a tool for attackers and a benchmark for defenders. By understanding and preparing for PGD attacks, we can advance towards creating more secure and reliable machine learning systems.