

# ECE595 / STAT598: Machine Learning I

## Lecture 33.2: Adversarial Attack - Basic Terminologies

Spring 2020

Stanley Chan

School of Electrical and Computer Engineering  
Purdue University



# Outline

- **Lecture 33 Overview**
- Lecture 34 Min-distance attack
- Lecture 35 Max-loss attack and regularized attack

## Today's Lecture

- What are adversarial attacks?
  - The surprising findings by Szegedy (2013) and Goodfellow (2014)
  - Examples of attacks
  - Physical attacks
- **Basic terminologies**
  - **Defining attack**
  - **Multi-class problem**
  - **Three forms of attack**
  - **Objective function and constraint sets**

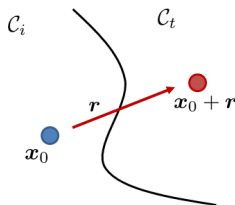
## Definition: Additive Adversarial Attack

### Definition (**Additive** Adversarial Attack)

Let  $\mathbf{x}_0 \in \mathbb{R}^d$  be a data point belong to class  $\mathcal{C}_i$ . Define a target class  $\mathcal{C}_t$ . An **additive** adversarial attack is an addition of a perturbation  $\mathbf{r} \in \mathbb{R}^d$  such that the perturbed data

$$\mathbf{x} = \mathbf{x}_0 + \mathbf{r}$$

is misclassified as  $\mathcal{C}_t$ .



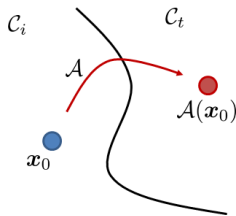
## Definition: General Adversarial Attack

### Definition (Adversarial Attack)

Let  $\mathbf{x}_0 \in \mathbb{R}^d$  be a data point belong to class  $\mathcal{C}_i$ . Define a target class  $\mathcal{C}_t$ . An **adversarial attack** is a mapping  $\mathcal{A} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  such that the perturbed data

$$\mathbf{x} = \mathcal{A}(\mathbf{x}_0)$$

is misclassified as  $\mathcal{C}_t$ .



# Example: Geometric Attack

## Fast Geometrically-Perturbed Adversarial Faces (WACV 2019)

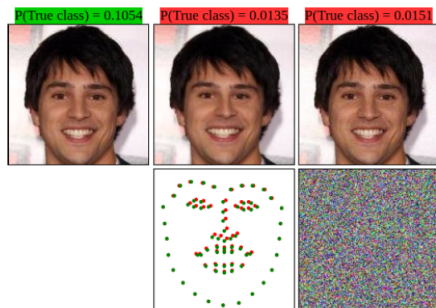
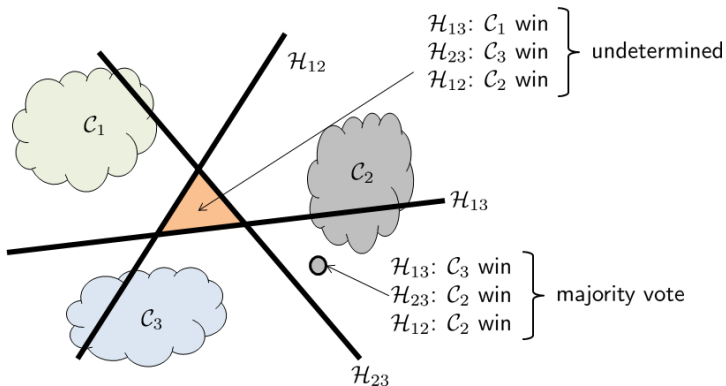


Figure 1. Comparison of the proposed attack to an intensity-based attack. First column: the ground truth image, which is correctly classified. Second column: the spatially transformed adversarial image wrongly classified and the corresponding adversarial landmark locations computed by our method. Third column: the adversarial image wrongly classified and the corresponding perturbation generated by the fast gradient sign method [7]. The proposed method leads to natural adversarial faces which are clean from additive noise.

<https://arxiv.org/pdf/1809.08999.pdf>

# The Multi-Class Problem

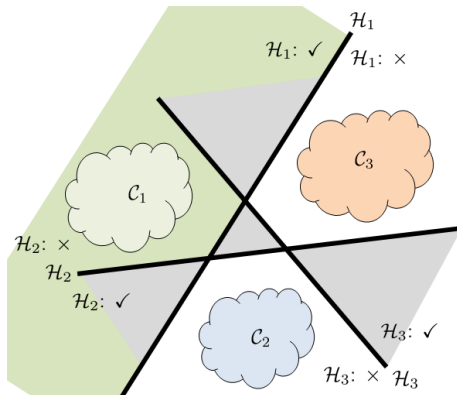
## Approach 1: One-on-One



- Class  $i$  VS Class  $j$
- Give me a point, check which class has more votes
- There is an undetermined region

# The Multi-Class Problem

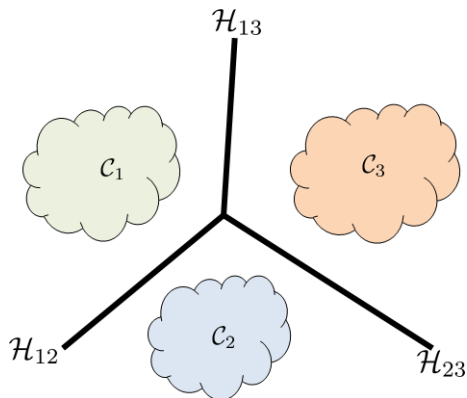
## Approach 2: One-on-All



- Class  $i$  VS not Class  $i$
- Give me a point, check which class has no conflict
- There are undetermined regions

# The Multi-Class Problem

## Approach 3: Linear Machine



- Every point in the space gets assigned a class.
- You give me  $\mathbf{x}$ , I compute  $g_1(\mathbf{x}), g_2(\mathbf{x}), \dots, g_K(\mathbf{x})$ .
- If  $g_i(\mathbf{x}) \geq g_j(\mathbf{x})$  for all  $j \neq i$ , then  $\mathbf{x}$  belongs to class  $i$ .



## Correct Classification

- We are mostly interested the linear machine problem.
- Let us try to simplify the notation. The statement:  
If  $g_i(\mathbf{x}) \geq g_j(\mathbf{x})$  for all  $j \neq i$ , then  $\mathbf{x}$  belongs to class  $i$ .  
is equivalent to (asking everyone to be less than 0)

$$g_1(\mathbf{x}) - g_i(\mathbf{x}) \leq 0$$

⋮

$$g_k(\mathbf{x}) - g_i(\mathbf{x}) \leq 0,$$

- and is also equivalent to (asking the worst guy to be less than 0)

$$\max_{j \neq i} \{g_j(\mathbf{x})\} - g_i(\mathbf{x}) \leq 0$$

- Therefore, if I want to launch an **adversarial attack**, I want to move you to class  $t$ :

$$\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0.$$

# Our Approach

Here is what we are going to do

- First, we will preview the three **equivalent** forms of attack:
  - Minimum Distance Attack: Minimize the perturbation magnitude while accomplishing the attack objective
  - Maximum Loss Attack: Maximize the training loss while ensuring perturbation is controlled
  - Regularization-based Attack: Use regularization to control the amount of perturbation
- Then, we will try to understand the **geometry** of the attacks.
- We will look at the **linear classifier** case to gain insights.

# Minimum Distance Attack

## Definition (Minimum Distance Attack)

The **minimum distance attack** finds a perturbed data  $\mathbf{x}$  by solving the optimization

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \|\mathbf{x} - \mathbf{x}_0\| \\ & \text{subject to} && \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0, \end{aligned} \tag{1}$$

where  $\|\cdot\|$  can be any norm specified by the user.

- I want to make you to class  $\mathcal{C}_t$ .
- So the constraint needs to be satisfied.
- But I also want to minimize the attack strength. This gives the objective.

# Maximum Loss Attack

## Definition (Maximum Loss Attack)

The **maximum loss attack** finds a perturbed data  $\mathbf{x}$  by solving the optimization

$$\begin{aligned} & \underset{\mathbf{x}}{\text{maximize}} && g_t(\mathbf{x}) - \max_{j \neq t} \{g_j(\mathbf{x})\} \\ & \text{subject to} && \|\mathbf{x} - \mathbf{x}_0\| \leq \eta, \end{aligned} \tag{2}$$

where  $\|\cdot\|$  can be any norm specified by the user, and  $\eta > 0$  denotes the attack strength.

- I want to bound my attack  $\|\mathbf{x} - \mathbf{x}_0\| \leq \eta$
- I want to make  $g_t(\mathbf{x})$  as big as possible
- So I want to maximize  $g_t(\mathbf{x}) - \max_{j \neq t} \{g_j(\mathbf{x})\}$
- This is equivalent to

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \\ & \text{subject to} && \|\mathbf{x} - \mathbf{x}_0\| \leq \eta, \end{aligned}$$

# Regularization-based Attack

## Definition (Regularization-based Attack)

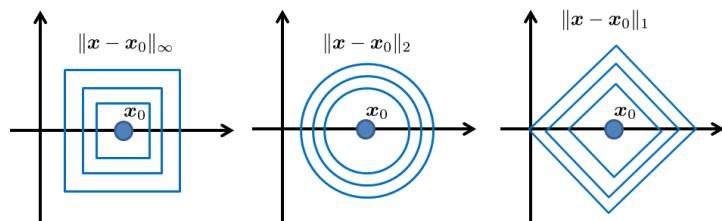
The **regularization-based attack** finds a perturbed data  $\mathbf{x}$  by solving the optimization

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\| + \lambda (\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x})) \quad (3)$$

where  $\|\cdot\|$  can be any norm specified by the user, and  $\lambda > 0$  is a regularization parameter.

- Combine the two parts via regularization
- By adjusting  $(\epsilon, \eta, \lambda)$ , all three will give the same optimal value.

## Understanding the Geometry: Objective Function



- $l_0$ -norm:  $\varphi(x) = \|x - x_0\|_0$ , which gives the most sparse solution. Useful when we want to limit the number of attack pixels.
- $l_1$ -norm:  $\varphi(x) = \|x - x_0\|_1$ , which is a convex surrogate of the  $l_0$ -norm.
- $l_\infty$ -norm:  $\varphi(x) = \|x - x_0\|_\infty$ , which minimizes the maximum element of the perturbation.

## Understanding the Geometry: Constraint

- The constraint set is

$$\Omega = \{\mathbf{x} \mid \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0\}$$

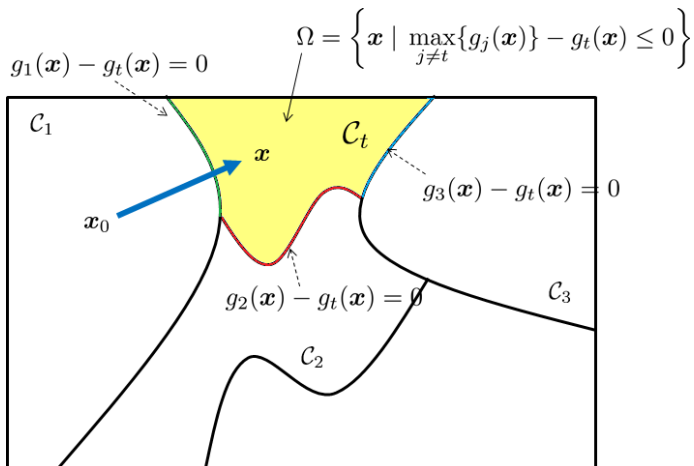
- We can write  $\Omega$  as

$$\Omega = \left\{ \mathbf{x} \mid \begin{array}{l} g_1(\mathbf{x}) - g_t(\mathbf{x}) \leq 0 \\ g_2(\mathbf{x}) - g_t(\mathbf{x}) \leq 0 \\ \vdots \\ g_k(\mathbf{x}) - g_t(\mathbf{x}) \leq 0 \end{array} \right\}$$

- Remark: If you want to replace max by  $i^*$ , then  $i^*$  is a function of  $\mathbf{x}$ :

$$\Omega = \{\mathbf{x} \mid g_{i^*(\mathbf{x})}(\mathbf{x}) - g_t(\mathbf{x}) \leq 0\}.$$

# Understanding the Geometry: Constraint





# Linear Classifier

- Let us take a closer look at the linear case.
- Each discriminant function takes the form

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i,0}.$$

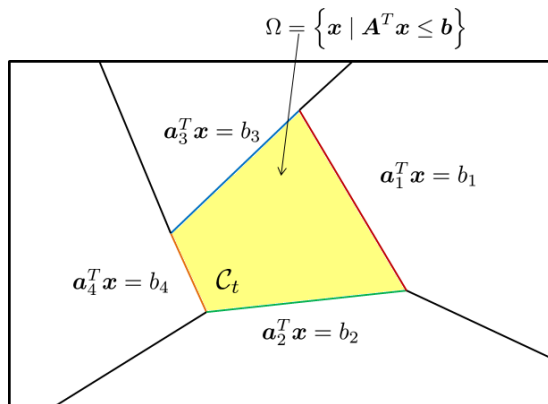
- The decision boundary between the  $i$ -th class and the  $t$ -th class is therefore

$$g(\mathbf{x}) = (\mathbf{w}_i - \mathbf{w}_t)^T \mathbf{x} + w_{i,0} - w_{t,0} = 0.$$

- The constraint set  $\Omega$  is

$$\begin{bmatrix} \mathbf{w}_1^T - \mathbf{w}_t^T \\ \vdots \\ \mathbf{w}_{t-1}^T - \mathbf{w}_t^T \\ \mathbf{w}_{t+1}^T - \mathbf{w}_t^T \\ \vdots \\ \mathbf{w}_k^T - \mathbf{w}_t^T \end{bmatrix} \mathbf{x} + \begin{bmatrix} w_{1,0} - w_{t,0} \\ \vdots \\ w_{t-1,0} - w_{t,0} \\ w_{t+1,0} - w_{t,0} \\ \vdots \\ w_{k,0} - w_{t,0} \end{bmatrix} \leq \mathbf{0} \Leftrightarrow \mathbf{A}^T \mathbf{x} \leq \mathbf{b}$$

# Linear Classifier



- You can show  $\Omega = \{ \mathbf{A}^T \mathbf{x} \leq \mathbf{b} \}$  is convex.
- But the complement  $\Omega^c = \{ \mathbf{A}^T \mathbf{x} > \mathbf{b} \}$  is not convex.
- So targeted attack is easier to analyze than untargeted attack.

## Attack: The Simplest Example

The optimization is:

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \|\mathbf{x} - \mathbf{x}_0\| \\ & \text{subject to} && \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0, \end{aligned}$$

- Suppose we use  $\ell_2$ -norm, and consider **linear** classifiers, then
- the attack is given by

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 \quad \text{subject to} \quad \mathbf{A}^T \mathbf{x} \leq \mathbf{b},$$

- This is a **quadratic programming** problem.
- We will discuss how to solve this problem analytically.

# Summary

- Adversarial attack is a universal phenomenon for **any** classifier.
- Attacking deep networks are popular because people think that they are unbeatable.
- There is really nothing too magical behind adversarial attack.
- All attacks are based on one of the three forms of attacks.
- Deep networks are trickier, as we will see, because the internal model information is not easy to extract.
- We will learn the basic principles of attacks, and try to gain insights from linear models.