

# PART 6

2022-04-29

# L15 Adversarial Attack

$$x \Rightarrow [] \Rightarrow \dots \Rightarrow [] \Rightarrow y \quad \left. \begin{array}{c} \\ \\ \end{array} \right\} p(y|x)$$

$\approx p(y|z)$

Classification:

$$c(x) = \arg \max_y p(y|x)$$

Benign example:  $(x, y) \quad c(x) = y$

Targeted attack:  $x \rightarrow x'$ ,  $D(x, x') \leq \delta$   
 $c(x') = t \neq y \quad \rightarrow L_1, L_2, L_\infty$

Untargeted attack:  $c(x') \neq y$

## General Principle

$$D_0 L(x, y, \theta)$$

Training:  $\min_{\theta} E_{x,y} [-\log p(y|x, \theta)]$

choose  $\theta$  to max prob of true class  $y$

After training:  $\theta$  fixed  $(x, y)$

Untargeted Attack

$$\max_{x'} -\log p(y|x', \theta) \quad \left. \begin{array}{l} \uparrow \text{prob of } y \\ \downarrow \text{prob of } y \end{array} \right.$$

$$D(x', x) \leq \delta$$

Targeted attack

$$\min_{x'} -\log p(t|x', \theta) \quad \left. \begin{array}{l} \uparrow \text{prob of } t \\ \downarrow \text{prob of } t \end{array} \right.$$

$$D(x', x)$$

$$D_x L(x, y, \theta) \quad D_x L(x, t, \theta)$$

## Fast Gradient Sign Method (FGSM)

Target :  $y \rightarrow t$        $\min_{x'} L(\bar{x}', t)$   
 $D(x, x') \leq \delta$

$$\bar{x}' = \bar{x} - \varepsilon \nabla_x L(x, t)$$

(0.1, 100, -200, --)

problem: Magnitude of perturbations NOT controlled

$$\bar{x}' = \bar{x} - \varepsilon \operatorname{Sign} \nabla_x L(x, t)$$

What  $t$  to choose?      (1, 1, -1, --)

$x \Rightarrow$	cat	dog ..	Tower	<u>Least likely class</u>
	0.7	0.1	0.000001	

## Basic Iterative Method (BIM)

$$x'_0 = x$$

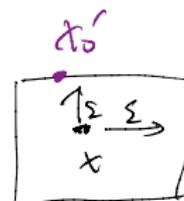
$$x'^{i+1} = \text{Clip}_{\bar{x}, \varepsilon} [x'_i - \varepsilon \text{sign}(\nabla_x L(x'_i, t))]$$

If  $x'_0$  random point

on the  $L_\infty$ -ball of  $\bar{x}$ ,

projected Gradient Descent

(PGD)



2022-05-04

First Attack algo: Box-constrained L-BFGS (2014)

$$\min_r \|r\|_2^2$$

$$\text{s.t. } C(x') = t \quad x' = x + r \quad x' \in [0, 1]^d$$

$\Rightarrow$  Hard to solve. Reformulated as follows

$$\min_r C\|r\|_2^2 + L(x+r, t)$$

$$L: -\log p(\cdot)$$

$$\text{s.t. } x+r \in [0, 1]^d$$

strong attack : CW or Margin-Based (2017)

$$x \Rightarrow \dots \Rightarrow \begin{bmatrix} z_1 \\ \vdots \\ z_t \\ \vdots \\ z_c \end{bmatrix} \text{ logits}$$

ideas

①  $\max z_t, \min \max_{i \neq t} z_i : \delta = z_t - \max_{i \neq t} z_i$

$\max \delta$

②  $\delta > k$ , turn attention to  $\min \|x - x'\|_2^2$

CW objective

$$\min_{x'} c \|x - x'\|_2^2 + \max \left\{ \max_{i \neq t} z_i - z_t, -k \right\}$$

$\underbrace{-\sigma}_{\leq -k} \quad \underbrace{\}_{> -k} \quad \textcircled{2} \quad \textcircled{1}$

## Linearity Hypothesis (Goodfellow et al 2015)

Linear classifier :  $\underline{\bar{w}^T \bar{x}} > 0$   
 $\bar{r} = -\varepsilon \bar{w}$

$$\bar{x}' = \bar{x} + \bar{r}$$

$$\bar{w}^T \bar{x}' = \underline{\bar{w}^T \bar{x}} - \varepsilon \bar{w}^T \bar{w}$$

Small

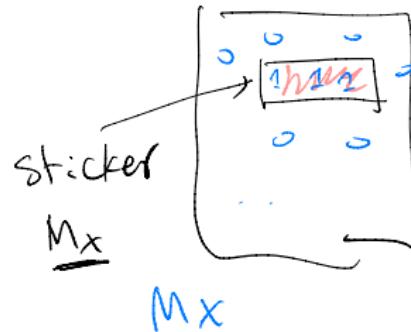
$$\underbrace{w_1^2 + w_2^2 + \dots + w_d^2}_{large}$$

$d = 65,000$

can easily cross the decision boundary

works even if  $\bar{r} = -\varepsilon \text{Sign}(\bar{w})$

## sticker attack



$$\min_{\tilde{r}} \text{clip}(\tilde{r} \cdot \text{mult}^L L(X + r \cdot M_x, t))$$

# Adversarial Training / Data Augmentation

Data  $\Rightarrow$  Model 1



A-E.

mixup

drop

Adversarial



Continuity

Data + A.E.  $\Rightarrow$  Model 2

~~Model 2~~



More robust against  
attack

L16 : XAI

2022-05-06

## Sensitivity vs Attribution

$$Z_c = 0.1 x_1 + 10 x_2 + 2$$

Sensitivity:

$$\frac{x_1}{0.1} \quad \frac{x_2}{10}$$

$$\frac{\partial Z_c}{\partial x_i}$$

Attribution:

(100, 0.1)	10	1	$x_i \frac{\partial Z_c}{x_i}$
(1, 1)	0.1	10	
Input			
$x_1 \quad x_2$			

LIME

$$z = [1, 0, \dots, 1]$$

$$g(z) = \bar{w}^T z$$

(x)

$$\min_w L(w) = \sum_{z \in \text{Samples around } x} t_z (f(x_z) - g(z))^2$$



z: Samples  
around  
 $x_z$

w

$$[0.1 \ 0.3 \ 0.5 \ \dots] \quad [1, \dots, 1]$$

$$[1, 1, 1, 1, 1, 1, 1]$$

$$[1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 0]$$

$$[0 \ 1 \ 0 \ 0 \ 1 \ 1 \ 0]$$

More important  
than



























