

# Robustness

- Model Misspecification
- Measurement error
- Dirty data
- Adversary

Simple example:

$$X_1, \dots, X_n \sim N(\mu, 1)$$

Goal: Est  $\mu$  from  $X_1, \dots, X_n$



$\hat{\mu} = \frac{1}{n} \sum X_i$ . If  $n$  is large enough,  $\hat{\mu} \approx \mu$ .

Adversary: Can add pts to dataset (few)

If  $\bullet$  is very large ( $> 100n + \mu$ ) then  $|\hat{\mu} - \mu|$  will be large ( $> 100$ )

Defenses:

- Prune dataset
- Median, instead of mean

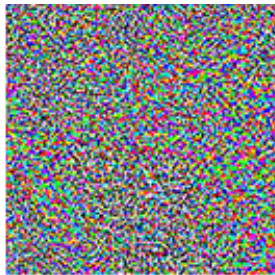
Today: Robustness to adversarial Examples



"panda"

57.7% confidence

+  $\epsilon$



=



"gibbon"

99.3% confidence

↑ Carefully crafted

Setting: Model trained on some training data.

At test time, adversary can modify each point "a bit"

Adversary goal: Reduce test accuracy as much as possible

$x'$  is An adversarial example for  $x$  on model  $f_\theta$  if

1. (Informal)  $x \approx x' \Leftrightarrow d(x, x')$  is small  $\Leftrightarrow x$  and  $x'$  have same label according to human
2.  $f_\theta(x) \neq f_\theta(x')$

Proxy for "human perception" in 1:

- Use dist instead

Common:  $l_p$  dist ( $\|x - x'\|_p = \left( \sum (x_i - x'_i)^p \right)^{1/p}$ )

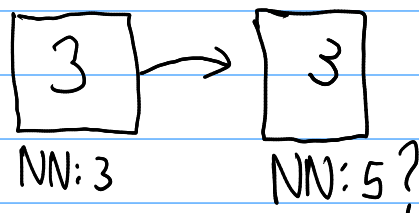
Adversary may output  $x'$  in  $\{y : \|x - y\|_p \leq \epsilon\}$

Common:  $p = 0, 2, \infty$   $\leftarrow$  Today: Can change each pixel by  $\leq \epsilon$ .   
  $\uparrow$  small

$\uparrow$  Can change  $\epsilon$  pts arbitrarily

Other dists

- Wasserstein
- Translations, Rotations, Resizing



## Attacker:

Given trained model  $f_\theta$ , test example  $x$ .

Construct  $x'$ :

1.  $\|x - x'\|_\infty \leq \epsilon$
2.  $f_\theta(x) \neq f_\theta(x')$

Notes:

- a) White-box
- b) Untargeted vs Targeted attacks  $\leftarrow (f_\theta(x') = c \neq f_\theta(x) \text{ for target } c)$

Optimize NN:

$$\min_{\theta} \frac{1}{n} \sum l(x_i, y_i, \theta)$$

Data fixed, optimize  $\theta$

$$\theta \leftarrow \theta - \eta \frac{1}{n} \sum \nabla_{\theta} l(x_i, y_i, \theta)$$

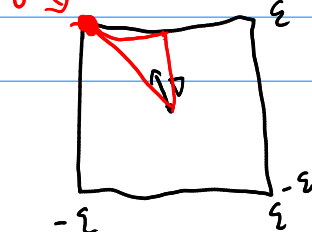
Generate Adversarial example

$$s' = \arg \max_{s} l(x + s, y, \theta) \text{ s.t. } \|s\|_\infty \leq \epsilon. \Rightarrow x' = x + s'$$

Gradient-Based opt

- Simple: Fast Gradient Sign Method (FGSM)

$\nabla_s l(x + s, y, \theta) \Rightarrow$  Take <sup>biggest</sup> 1 step in  $\nabla$  direction, s.t. constraint



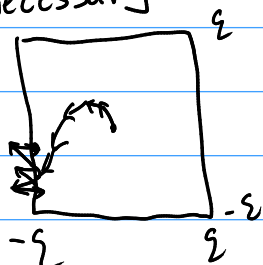
$$s^* = \epsilon \cdot \text{sign}(\nabla_s l(x + s, y, \theta)) \in \{\pm 1\}^d$$

Better: Projected Gradient Descent (PGD)

$$\delta^{t+1} = \text{Proj} \left( \delta^t + \eta \text{Sign} \left( \nabla_{\delta} l(x + \delta^t, y, \theta) \right) \right)$$

Projects into  $[-\epsilon, \epsilon]^d$  if necessary

Same as before, but  $\eta$  is hyper param



$$\text{Proj}([3\epsilon, -2\epsilon, 0.5\epsilon])$$

$$= [\epsilon, -\epsilon, 0.5\epsilon]$$

Untargeted:  $\max_{\delta} l(x + \delta, y, \theta)$  s.t.  $\|\delta\|_{\infty} \leq \epsilon$

Targeted to  $c$ :  $\max_{\delta} l(x + \delta, y, \theta) - l(x + \delta, c, \theta)$  s.t.  $\|\delta\|_{\infty} \leq \epsilon$

## Defenses

### Adversarial Training

- Usual goal:  $\min_{\theta} \mathbb{E}_{(x,y) \sim p} [l(x,y,\theta)]$

- Robust setting:  $\min_{\theta} \mathbb{E}_{(x,y) \sim p} \left[ \max_{\delta: \|\delta\|_{\infty} \leq \epsilon} l(x + \delta, y, \theta) \right]$

$$\hookrightarrow \min_{\theta} \frac{1}{n} \sum \max_{\delta_i: \|\delta_i\|_{\infty} \leq \epsilon} l(x_i + \delta_i, y_i, \theta)$$

How to solve

1. Draw minibatch  $B$

2. For each  $(x_i, y_i)$  in  $B$ , compute  $\delta_i^* = \underset{\delta_i: \|\delta_i\|_{\infty} \leq \epsilon}{\text{argmax}} l(x_i + \delta_i, y_i, \theta)$

3.  $\theta \leftarrow \theta - \eta \frac{1}{|B|} \sum_{i \in B} \nabla_{\theta} l(x_i + \delta_i^*, y_i, \theta)$

(use methods above)

4. Repeat

- Athalye, Carlini, Wagner '18

- ICLR '18 accepted 9 papers on adv defs.

Defense	Dataset	Distance	Accuracy
Buckman et al. (2018)	CIFAR	0.031 ( $l_\infty$ )	0%*
Ma et al. (2018)	CIFAR	0.031 ( $l_\infty$ )	5%
Guo et al. (2018)	ImageNet	0.005 ( $l_2$ )	0%*
Dhillon et al. (2018)	CIFAR	0.031 ( $l_\infty$ )	0%
Xie et al. (2018)	ImageNet	0.031 ( $l_\infty$ )	0%*
Song et al. (2018)	CIFAR	0.031 ( $l_\infty$ )	9%*
Samangouei et al. (2018)	MNIST	0.005 ( $l_2$ )	55%**
Madry et al. (2018)	CIFAR	0.031 ( $l_\infty$ )	47%
Na et al. (2018)	CIFAR	0.015 ( $l_\infty$ )	15%

## Backdoor attacks

