

ADVERSARIAL ROBUSTNESS

Recall: ERM = minimize θ $\mathbb{E}_{x,y \sim \text{training}}$ $l(f_{\theta}(x); y)$

Goal: maximize acc. ~~on~~ on test data.

Usually in standard ML, train, test $\overset{\text{iid}}{\sim}$

↓
not always true

① DA / DG

ROBUSTNESS

↙
"Distributions"

→
"Adversarial"
imperceptible

$x + \delta = x'$

input "some" noise adversarial example

DEFINITION: x' is an adv. version of x
(given model $f_{\theta}(\cdot)$)

if

① $f_{\theta}(x) \neq f_{\theta}(x')$

② $f_{\text{human}}(x) = f_{\text{human}}(x')$

$f_{\text{det}}(x) = f_{\text{det}}(x')$

pred. not same
human pred. identical
generally

③ for a human, the diff. betⁿ x and x'
 $d_{\text{human}}(x, x') = 0$

↓
 assumption

for some distance measure d
 $d(x, x') \leq \epsilon$

Q: ① what noise δ should we add?

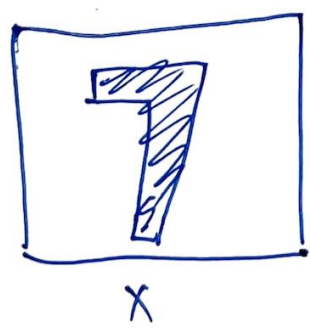
Q: ② can we quantify "imperceptible"
 what is $d(\cdot)$ (What is the THREAT MODEL)

e.g.

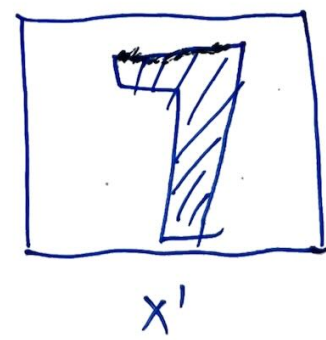
L_2 norm	$\ x' - x\ _2 \leq \epsilon$
L_1 norm	$\ x' - x\ _1 \leq \epsilon$
L_∞ norm	$\ x' - x\ _\infty \leq \epsilon$

aside: L_∞ norm
 $\|a\|_\infty = \max_i |a_i|$

L_p norm threat models



rotate it
 by 0.5°



OPTIMAL
 TRANSPORT
 COST

$\|x' - x\|_p \leq \epsilon$

ATTACK : creating adv. examples
(when applied to the model,
accuracy drops)

DEFENSE : preventing attacks
improving robustness

ATTACK problem formulation

Goal : given input x , model $f_{\theta}(\cdot)$
find δ
 $x' = x + \delta$

st. ① $f_{\theta}(x') \neq f_{\theta}(x) = y \rightarrow g.t.$

② $\|x - x'\|_p \leq \epsilon$

TYPES OF ATTACK BASED ON OUTCOME

- ① targeted \rightarrow make the model predict label y_{target}
- ② untargeted \rightarrow make the model predict some wrong label

BASED ON WHEN THE MODEL IS ATTACKED

- ① inference / test time attacks \rightarrow TODAY
- ② training time attacks "poisoning"

WHAT DOES THE ATTACKER HAVE ACCESS TO?

- ① Black Box attack : only access to I/O
- ② White Box attack : I/O, model weights, architecture
...

ATTACK GOAL : make classifier fail

$$\text{maximize}_{\delta} \ell(f_{\theta}(x), y)$$

$$\text{s.t. } \|\delta\|_p \leq \epsilon$$

$$\delta \leftarrow \delta + \eta \nabla \ell$$

- we need to optimize
- we know G.D.
- but constraint?

METHOD Fast Gradient Sign Method (FGSM)

①

- simple idea
- one-step search
- look @ direction of gradient
- go in that direction
↳ but don't violate the constraint

$$\delta_{\text{FGSM}} = \max_{\|\delta\|_{\infty} < \epsilon} \langle \delta, \nabla_x \ell(f_{\theta}(x), y) \rangle$$

$$\delta_{\text{FGSM}} = \epsilon \cdot \text{sign}(\nabla_x \ell(f_{\theta}(x), y))$$

METHOD

(2)

Projected G-P

-iterative version FGSM

$$x^0 = x \quad (\text{input})$$

$$\hat{x}^{t+1} = x^t + \eta g^t \quad \longrightarrow \text{from FGSM}$$

problem! what about the constraint

"do a projection"

$$x^{t+1} = \Pi(\hat{x}^{t+1})$$

↓
projection

for L_∞ : $x^{t+1} = \text{clip}(\hat{x}^{t+1}, x-\epsilon, x+\epsilon)$

for L_2 : normalize

ALTERNATIVE FORMULATION

i.e. make attack imperceptible

equiv. $\left\{ \begin{array}{l} \text{minimize } \|\delta\|_p \\ \text{st. } f_\theta(x+\delta) \neq y \end{array} \right.$

$$\text{minimize } \|\delta\|_p - \lambda \cdot \ell(f_\theta(x+\delta), y)$$

ADVERSARIAL DEFENSE

- remember adv. attack is at test-time
- ~~model is already trained~~
- need diff. training method

recall: ERM $\min_{\theta} \mathbb{E}_{x,y \sim \text{train}} \ell(f_{\theta}(x), y)$

Adversarial Training

$$\min_{\theta} \mathbb{E}_{x,y \sim \text{training}} \left[\max_{x' \substack{x' = x + \delta \\ \text{s.t. } \|x - x'\| \leq \epsilon}} \ell(f_{\theta}(x'), y) \right]$$

- instead of minimizing loss on "clean" training data

- ① find worst perturbations (adv. examples)
- ② minimize the worst-case loss

Inner Max: attack problem (solved using PGD, FGSM, ...)

Outer Min: G-D.

"min-max" problems

Alternating Optimization