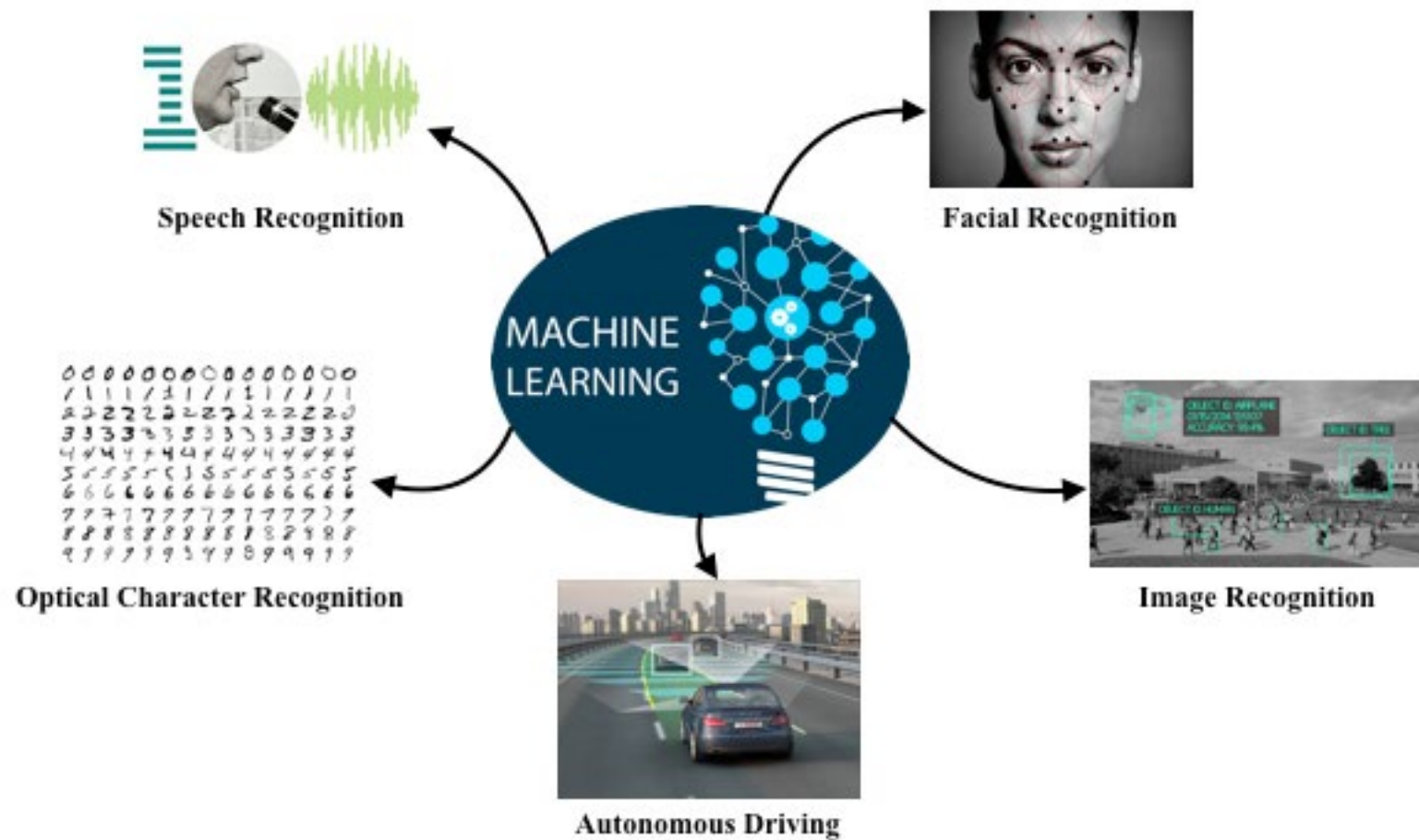


Adversarial Machine Learning

- ✱ Attack algorithms

What is AML?



Adversarial Examples

※ Attack algorithms



Classified as panda



Small adversarial noise



Classified as gibbon

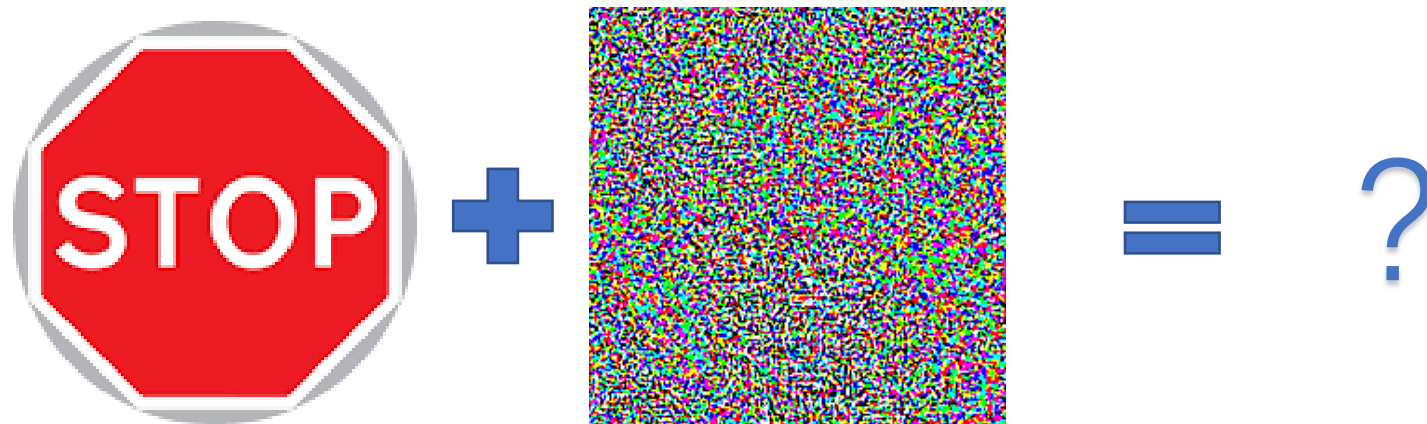


?

Who cares panda?

Adversarial Examples

※ Attack algorithms



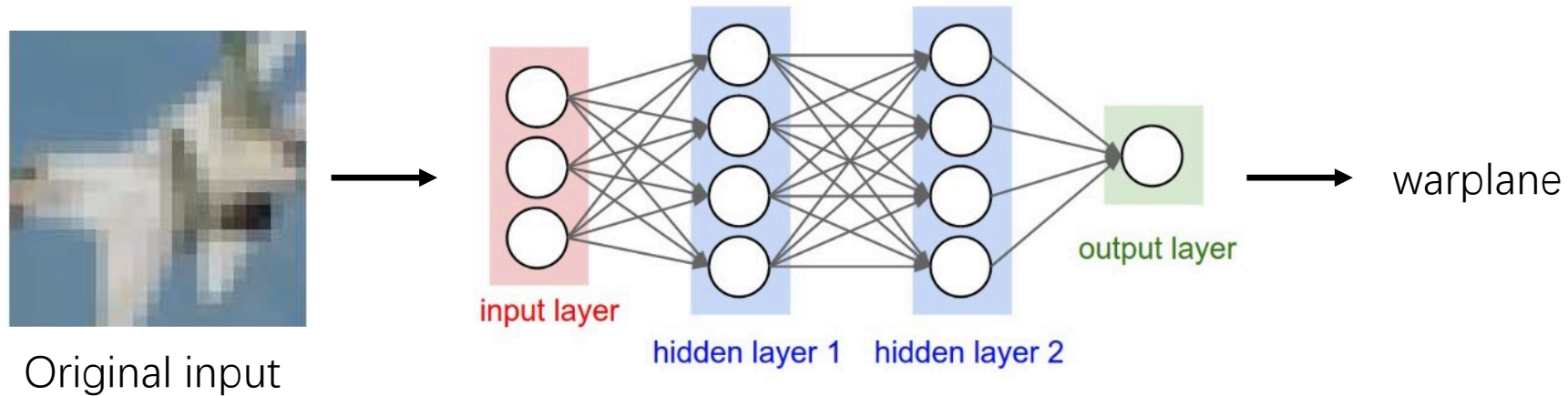
Small adversarial noise

Outline

- Attack
 - Formulation
 - Distance metrics
- Attack algorithms
 - L-BFGS
 - Fast Gradient Sign
 - AdvGAN
 - One pixel attack

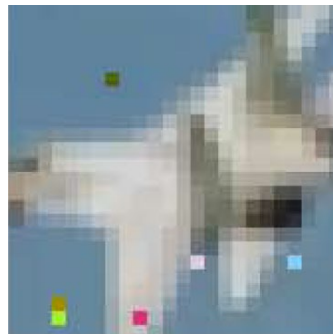
Attack

※ Attack algorithms




Attack: find a new input (similar to original input) but classified as another class t (untargeted or targeted)


Attacker knows the classifier




How to find adversarial examples

※ Attack algorithms

minimize $\mathcal{D}(x, x + \delta)$  distance between x and $x + \delta$

such that $C(x + \delta) = t$  $x + \delta$ is classified as target class t

$x + \delta \in [0, 1]^n$  each element of $x + \delta$ in $[0, 1]$ (for a valid image)

Distance Metrics

※ Attack algorithms

Two images: x and x'

- L_0 : measures the number of coordinates such that $x_i \neq x'_i$
 - corresponds to the number of pixels that have been changed in an image
- L_2 : Euclidean distance
- L_∞ : $\max(|x_1 - x'_1|, \dots, |x_n - x'_n|)$
 - measures maximum change to any of the elements



L-BFGS

※ Attack algorithms

$$\begin{array}{ll}\text{minimize} & \|x - x'\|_2^2 \\ \text{such that} & C(x') = l \\ & x' \in [0, 1]^n\end{array}$$



$$\begin{array}{ll}\text{minimize} & c \cdot \|x - x'\|_2^2 + \text{loss}_{F,l}(x') \\ \text{such that} & x' \in [0, 1]^n\end{array}$$

Initial formulation

$$\begin{array}{ll}\text{minimize} & \mathcal{D}(x, x + \delta) \\ \text{such that} & C(x + \delta) = t \\ & x + \delta \in [0, 1]^n\end{array}$$

Note that these two are not equivalent optimization problems

Fast Gradient Sign

※ Attack algorithms

$$x' = x - \epsilon \cdot \text{sign}(\nabla \text{loss}_{F,t}(x))$$

ϵ is chosen to be sufficiently small so as to be undetectable

fast rather than optimal

Fast Gradient Sign

✧ Attack algorithms

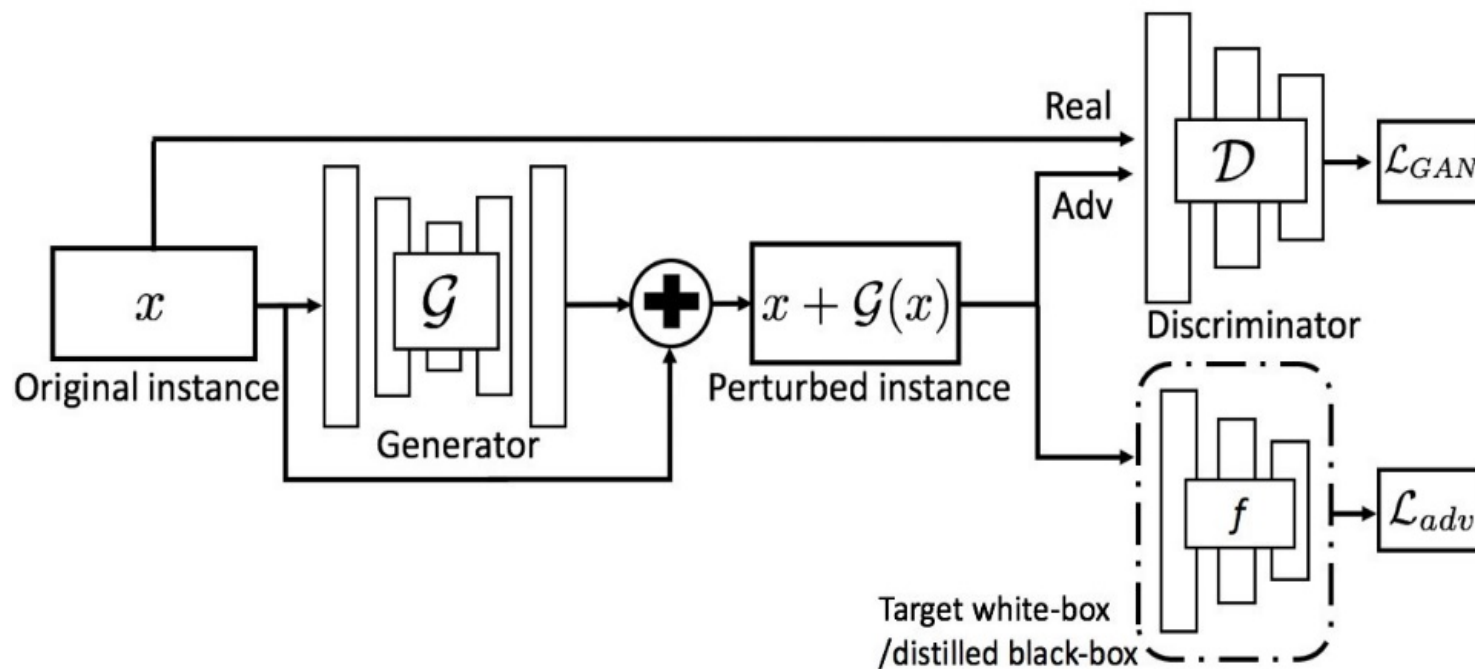


MNIST

Adversarial Image	Perturbation
Pred: 4	eps: 38
Pred: 7	eps: 60
Pred: 8	eps: 42
Pred: 8	eps: 12
Pred: 9	eps: 17

AdvGAN

※ Attack algorithms



$$\mathcal{L}_{GAN} = \mathbb{E}_x \log \mathcal{D}(x) + \mathbb{E}_x \log(1 - \mathcal{D}(x + \mathcal{G}(x))).$$

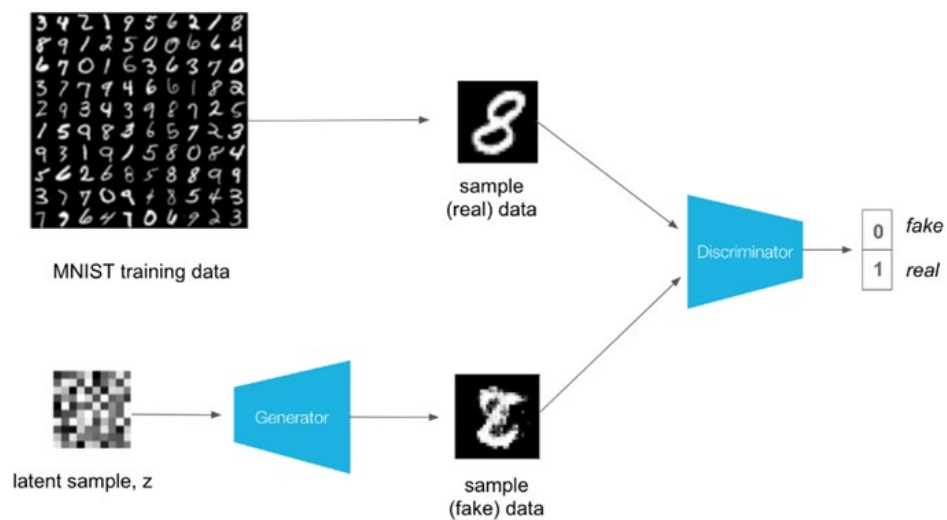
$$\mathcal{L}_{adv}^f = \mathbb{E}_x \ell_f(x + \mathcal{G}(x), t),$$

$$\mathcal{L}_{hinge} = \mathbb{E}_x \max(0, \|\mathcal{G}(x)\|_2 - c),$$

$$\mathcal{L} = \mathcal{L}_{adv}^f + \alpha \mathcal{L}_{GAN} + \beta \mathcal{L}_{hinge},$$

AdvGAN

※ Attack algorithms



Untargeted

Pred: 9	Pred: 3	Pred: 8	Pred: 8	Pred: 4	Pred: 3	Pred: 8	Pred: 3	Pred: 3	Pred: 8

Targeted

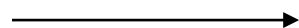
Target: 0	Target: 1	Target: 2	Target: 3	Target: 4	Target: 5	Target: 6	Target: 7	Target: 8	Target: 9
Pred: 0	Pred: 1	Pred: 2	Pred: 3	Pred: 4	Pred: 5	Pred: 6	Pred: 7	Pred: 8	Pred: 9
Pred: 0	Pred: 1	Pred: 2	Pred: 3	Pred: 4	Pred: 5	Pred: 6	Pred: 7	Pred: 8	Pred: 9
Pred: 0	Pred: 1	Pred: 2	Pred: 3	Pred: 4	Pred: 5	Pred: 6	Pred: 7	Pred: 8	Pred: 9

One pixel attack

※ Attack algorithms

$$\begin{array}{ll}\underset{e(\mathbf{x})^*}{\text{maximize}} & f_{adv}(\mathbf{x} + e(\mathbf{x})) \\ \text{subject to} & \|e(\mathbf{x})\| \leq L\end{array}$$

modify a part of all dimensions



$$\begin{array}{ll}\underset{e(\mathbf{x})^*}{\text{maximize}} & f_{adv}(\mathbf{x} + e(\mathbf{x})) \\ \text{subject to} & \|e(\mathbf{x})\|_0 \leq d,\end{array}$$

modify d dimensions

One pixel attack

Method	Success rate	Confidence	Number of pixels	Network
Our method	35.20%	60.08%	1 (0.098%)	NiN
Our method	31.40%	53.58%	1 (0.098%)	VGG
LSA[15]	97.89%	72%	33 (3.24%)	NiN
LSA[15]	97.98%	77%	30 (2.99%)	VGG
FGSM[11]	93.67%	93%	1024 (100%)	NiN
FGSM[11]	90.93%	90%	1024 (100%)	VGG

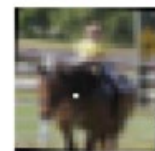
TABLE IX

COMPARISON OF NON-TARGETED ATTACK EFFECTIVENESS BETWEEN THE PROPOSED METHOD AND TWO PREVIOUS WORKS. THIS SUGGESTS THAT ONE PIXEL IS ENOUGH TO CREATE ADVERSARIAL IMAGES FROM MOST OF THE NATURAL IMAGES.

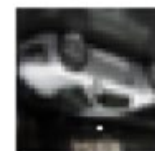
AllConv



SHIP
CAR(99.7%)



HORSE
DOG(70.7%)



CAR
AIRPLANE(82.4%)



DEER
AIRPLANE(49.8%)



HORSE
DOG(88.0%)

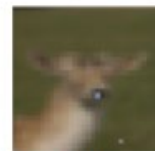
NiN



HORSE
FROG(99.9%)



DOG
CAT(75.5%)



DEER
DOG(86.4%)

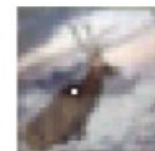


BIRD
FROG(88.8%)

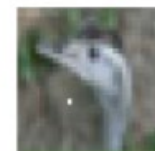


SHIP
AIRPLANE(62.7%)

VGG



DEER
AIRPLANE(85.3%)



BIRD
FROG(86.5%)



CAT
BIRD(66.2%)



SHIP
AIRPLANE(88.2%)



CAT
DOG(78.2%)