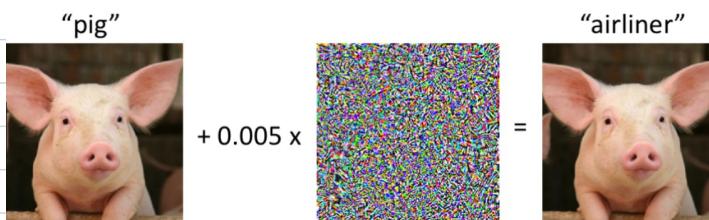


# Adversarial Machine Learning

Adversarial Example:



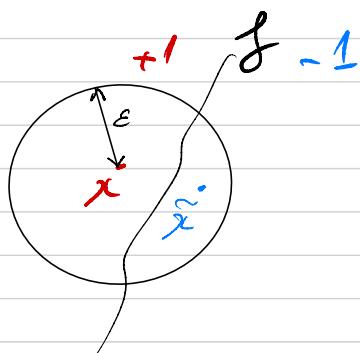
- classification problem  $(X, y) \sim D; y \in \{-1, 1\}$

- Standard Risk:

$$Z(f) = \mathbb{E}_D [\mathbb{1}_{f(x) \neq y}] = P_D [f(x) \neq y]$$

- Adversarial Risk:

$$R_\epsilon(f) = \mathbb{E}_D \left[ \max_{\tilde{x} \in \mathcal{X}, \|x - \tilde{x}\| \leq \epsilon} \mathbb{1}_{f(\tilde{x}) \neq y} \right]$$



Questions?

- How should we define the adversary power?  
 $\ell_2, \ell_\infty, \ell_1, \ell_0, \dots$

= Threat model (what the adversary can do)

- If  $Z(f) \leq \delta$ , How large can  $R_\epsilon(f)$  be?

- How can I design classifier so that it's robust?

- Given a non-robust classifier, can we make it robust?

- How can we find adversarial example.

- Which access we have on the NN to attack it.

- Why NN are non-robust?

## Adversarial examples:

### White Box Attack

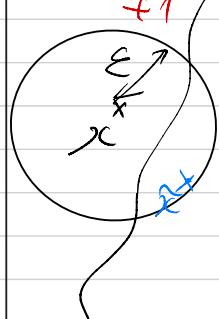
Task: Given  $\mathcal{X}$ ,  $\mathcal{Y}$ ,  $f$  and  $x \in \mathcal{X}$ , a metric, and  $\epsilon$ -Budget  
 → find an adversarial input  $\tilde{x}$  such that  $\|x - \tilde{x}\| \leq \epsilon$   
 $\hookrightarrow f(\tilde{x}) \neq y$

• Separable Delta,  $\exists h: \mathcal{X} \rightarrow \{-1, 1\}$  s.t  $h(x) = y$

• Given  $x$ , if  $f(x) \neq y \Rightarrow x = \tilde{x}$

let's assume that  $f(x) = y$

• 1) Ideally:  $B_\epsilon(x) \cap \{f(x) = -1\} \neq \emptyset$



• 2) Use gradient descent (gradient w.r.t input)

• we assume, we have  $g \in [0, 1]$

$$(g(x) = P(y=1|x)) \quad \begin{cases} f(x) = +1 & y \frac{1}{2} < g(x) < 1 \\ -1 & \text{o.w.} \end{cases}$$

• Wrong direction to follow will be:

• Smooth loss for classif.  $\ell(y, g(x))$

$$\frac{d}{dx} \ell(y, g(x)) = \underbrace{\ell'(y, g(x)) \cdot y}_{< 0} \cdot \nabla_x g(x) \quad (\text{loss is } \downarrow)$$

(Max. loss) → dir. to follow  $\propto -y \nabla_x g(x)$

$$\bullet h(x) \nabla_x g(x) \in \mathbb{R}^d$$

$$\Rightarrow \tilde{x} = x - \epsilon h(x) \frac{\nabla_x g(x)}{\|\nabla_x g(x)\|_2} \quad \text{under L2 const}$$

optimal local update.

Taylor:

$$b(\tilde{x}) \approx b(x) + \nabla_x b(x)(\tilde{x} - x)$$

$$\rightarrow \max_{\|x - \tilde{x}\| \leq \epsilon} b(\tilde{x}) \Leftrightarrow \max_{\|\tilde{x} - x\| \leq \epsilon} \nabla_x b(x)^T (\tilde{x} - x)$$

$$\Leftrightarrow \max_{\|v\| \leq \epsilon} \nabla_x b(x)^T v$$

$$\Rightarrow v^* = \frac{\nabla_x b(x)}{\|\nabla_x b(x)\|} \epsilon$$

Apply this to  $b(x) = -g(x)$

Rq

Sometimes we do multiple steps of this procedure.

$\Rightarrow$  this is called **white box attack**!

## Black Box Attack:

- $x - \boxed{g} \rightarrow g(x) \Rightarrow$  we can approx  $\nabla_x g(x)$  using finite difference.

$$\Rightarrow g(x + \delta e_i) - g(x) \sim \delta \nabla_i g(x) \rightarrow \nabla_x g(x)$$

- If we only have  $h+1, -1$ ; Build  $S = \{x_i, y_i\}$  ( $y_i = f(x_i)$ )  
 $\rightarrow$  use  $S$  to train a different N.N

$\rightarrow$  Attack this NN with White Box Attack  
 $\rightarrow$  So very often we can transfer it to the unknown network

### Attack of a physical object



(i.e. Stop sign)

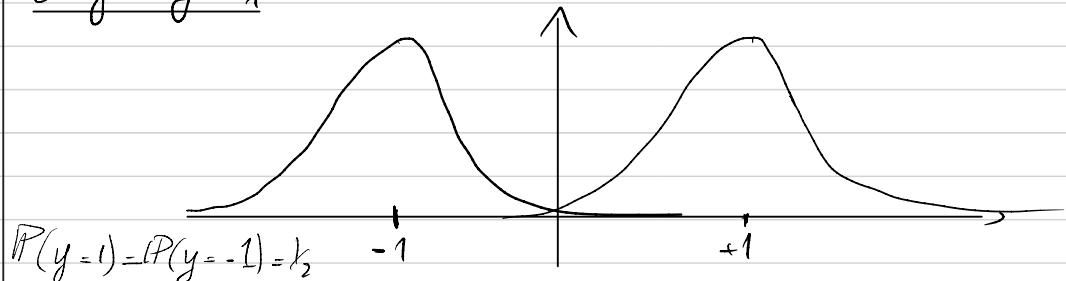
Figure 2: To stop or not to stop.

## • Robust vs Non-Robust features:

$$\tilde{x} \xrightarrow{\tau} x \in \mathbb{R}^D \quad x_1 = y + z_1 \\ x_2 = y \sqrt{\frac{\log d}{d-1}} + z_2 \quad z_i \sim \mathcal{N}(0,1) \\ \vdots \\ x_d = y \sqrt{\frac{\log d}{d-1}} + z_d \quad ("SNR")$$

$x_1$  has a very high signal  
 $(x_2, \dots, x_d)$  has a very low signal (but we have a lot of them)

• Only using  $x_1$ :



$$P(y=1) = P(y=-1) = \frac{1}{2}$$

$$\Rightarrow P(y|x_1) = \frac{P(x_1|y)P(y)}{P(x_1)} \quad . \text{How often will you be wrong?}$$

$$\Rightarrow \max_{y \in \{-1, 1\}} P(x_1|y)$$

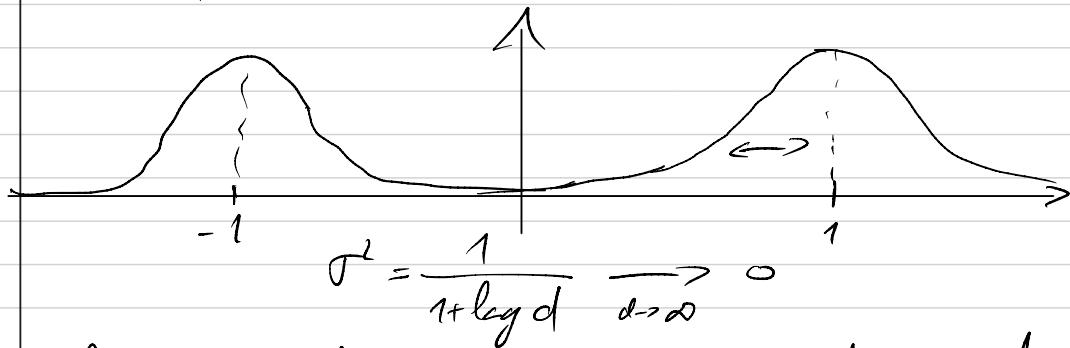
Prob of making a mistake  $\sim 0,16$

- Using all components:

$$\begin{aligned}
 \underset{y \in \{-1, 1\}}{\operatorname{argmax}} P(y | x) &= \underset{y}{\operatorname{argmax}} \frac{P(x|y)P(y)}{P(x)} \\
 &= \underset{y}{\operatorname{argmax}} \prod_{i=1}^d P(x_i | y) \\
 &= \underset{y}{\operatorname{argmax}} \prod_{i=1}^d \log \left( \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x_i - y_i \alpha_i)^2} \right) \\
 &= \underset{y}{\operatorname{argmin}} \sum_{i=1}^d (x_i - y_i \alpha_i)^2 \\
 \alpha_1 = 1, \alpha_i &= \sqrt{\frac{\log d}{d-1}}, i > 1 \\
 &= \underset{y}{\operatorname{argmax}} \sum_{i=1}^d x_i \alpha_i \\
 &= \underset{y}{\operatorname{argmax}} y \left( \sum_{i=1}^d \alpha_i^2 \right) + y \sum_{i=1}^d x_i \alpha_i \\
 &= \underset{y}{\operatorname{argmax}} y (1 + \log d) + y \sum_{i=1}^d x_i \alpha_i
 \end{aligned}$$

recalling  $\rightarrow \operatorname{argmax}_y y \hat{y} + \hat{y}^2$  with  $\hat{y} \sim N(0, \frac{1}{1+\log d})$

So it's equivalent;



$d$ : If we train well we can expect to have a standard error close to 0 when  $d$  grows.

## adversarial Risk =

- assume we allow an adversary to move every component by  $\epsilon$ . When  $\epsilon = 2 \sqrt{\frac{\log d}{d-1}}$

- . what can the adversary do?
  - First, learn the classifier.
  - then set  $\tilde{x}_i = x_i - y \cdot 2 \sqrt{\frac{\log d}{d-1}}$

$$\tilde{x}_i = -y \sqrt{\frac{\log d}{d-1}} + z_i$$

- . The optimal classifier will predict  $y$  whereas the true label was  $y$ .  
 $\rightarrow$  Adversarial Risk = 1

→ We could have only considered  $x_1$  and we still get a risk of 0,16

d There is a trade-off between having a small:

- adv Risk
- Standard Risk

## Summary

Standard Risk:  $Z(f) = \mathbb{E}_{x,y \sim P} [\mathbb{1}_{f(x) \neq y}]$ .

Adversarial example:  $\tilde{x}$  s.t.  $\|x - \tilde{x}\| \leq \epsilon$   
 $f(x) \neq y$

Adversarial Risk:  $R_\epsilon(f) = \mathbb{E}_{x,y \sim P} [\max_{\substack{x, \|x - \tilde{x}\| \leq \epsilon \\ f(\tilde{x}) \neq y}} \mathbb{1}_{f(\tilde{x}) \neq y}]$

→ It is possible to have  $Z(f) \lll 1$  and  $R_\epsilon(f) = 1$

## Why are ML algos Vulnerable to attack?

. Robust vs Non-Robust features

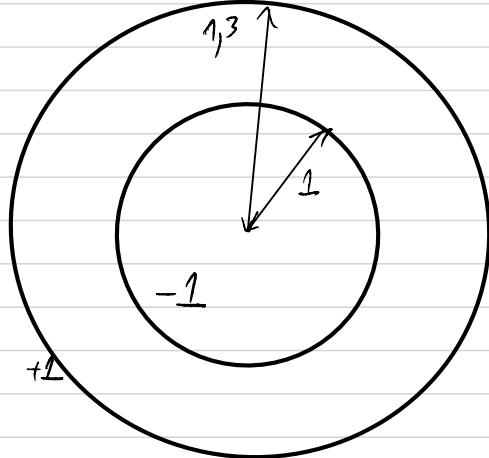
$$d-1 \left\{ \begin{array}{l} x_1 = y + z_1 \\ x_2 = dy + z_2 \\ \vdots \\ x_{d-1} = dy + z_d \end{array} \right. \quad \alpha = \sqrt{\frac{\log d}{d-1}}, \quad z_i \sim N(0, 1)$$

conclusion:

- . If we use all features, we can make  $Z(f) \rightarrow 0$ , But if we consider  $\infty$  perturbation  $\epsilon = 2\alpha$  then  $R_\epsilon(f) \rightarrow 1$
- . If we use only the first feature of  $x$ ,  $Z(f) = R_\epsilon(f) = 0,16$

## Adversarial Sphere:

$$D = 500, X \in \mathbb{R}^D, y \in \{-1\}$$



$X$  is uniform on the sphere

⚠  $D$  is very large

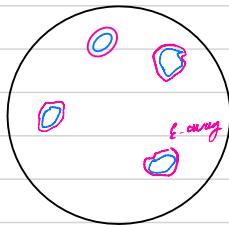
Experiment:  $|S_{\text{Train}}| = 50 \text{ Million}$ , NN 2 hidden layers, ReLU act.  
500 hidden nodes

$\rightarrow f$   
 $|S_{\text{Test}}| = 20 \text{ Million} \rightarrow \text{No err on the test set !!}$   
 $L(f)$  is very small

But the adv. Risk  $R_E \sim 1$  for  $\epsilon = \frac{1}{\sqrt{D}}$

### Effect of high dimension:

$X$  is uniform



$E$ : set of  $x \in \mathcal{X}, \|x\| = 1$   
and  $f(x) = 1$  "easy"

- Standard Risk:

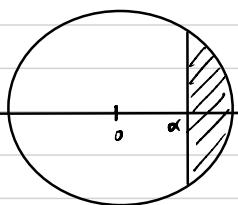
$$L(f) \approx P(\|x\|=1, f(x)=1) = \frac{\text{Area of } E}{\text{Area of the Sphere}}$$

- Adversarial Risk:

$$R_E(f) \sim \frac{\text{Area in } E}{\text{Area of the Sphere}}$$

Goal: Show that  $L(f) \ll 1$  but  $R_E(f) \sim 1$

• E: Spherical cap:



$$E = \{x, \|x\|=1, x_1 > \alpha\}; \alpha \in [0, 1]$$

Probability of area P,  $0 < P < \frac{1}{2}$

$$\text{i)} \alpha = 0 \quad P = \frac{1}{2}$$

$$\text{ii)} \alpha = 1 \quad P = 0$$

$\rightarrow \alpha = fct(p)$ ? intuition: proportional.

- P is very small, D is very large

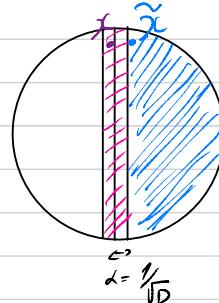
$$\alpha = \frac{c \ln(p)}{\sqrt{D}}$$

$\rightarrow$  whatever p, in large dim, the cap is almost half the sphere.  
 $\Rightarrow$  all the mass you have on the surface of the sphere is on the equator!!

P small, D large

$$E = 2\alpha = O\left(\frac{1}{\sqrt{D}}\right)$$

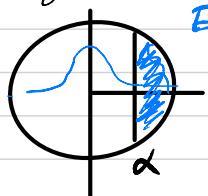
$$\|x - \tilde{x}\| \leq \epsilon$$



$$E \alpha = \frac{1}{\sqrt{D}}$$

$$Z(f) = p \text{ (small)}, R_E(f) = 1-p \sim 1$$

Proof:



simply:  $x \sim \text{Unif}(\text{sph})$

$$X \sim N(0, \frac{1}{D} I), \sigma^2 = \frac{1}{D}$$

X is spherical,  $\|x\| \sim 1$  with high prob.

$$P = P(x_1 > \alpha) = \frac{1}{2\pi r^2} \int_{\alpha}^{\infty} e^{-\frac{y^2}{2r^2}} dy,$$

$$= Q(\alpha \sqrt{D})$$

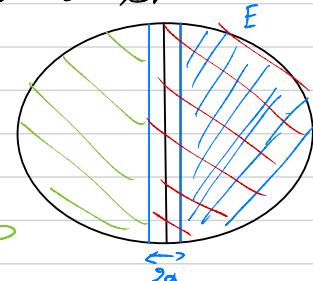
$$\text{where } Q(x) = \int_x^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} dy$$

$$\Rightarrow \alpha = \frac{Q^{-1}(P)}{\sqrt{D}}$$

$\Rightarrow$  all the mass is on the equator

$D$  large  $P$  small  $\alpha$  small

Let  $x$  uniform on the sphere  
with proba  $1-p$ ,  
I want to be on  $E$ .



What is the proba of obtaining a point  $E$ -away of the cap?

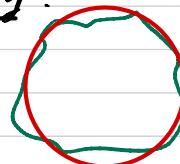
$$E = 2a \sim \frac{1}{\sqrt{D}}$$

(Standard nib =  $p$ )

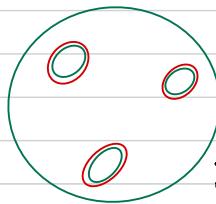
The adversarial risk is  $R(f) = 1 - p \sim 1$

### - Isoperimetric inequality :

In our case:



- The fixed volume body with the smallest surface area is the sphere



given a shape, we can add points  $E$  away

→ the spherical cap is the shape which minimizes the total area with a fixed initial area.

By assuming the spherical cap, we have a lower bound on the adversarial risk.

## Adversarial training:

$$\begin{aligned} \hat{\theta} &= \min \mathcal{L}(f_{\theta}), \text{ But we don't want small } \mathcal{L} \\ &\Rightarrow \hat{\theta} = \min R_E(f_{\theta}) \end{aligned}$$

Pb: • We don't know D  $\Rightarrow$  we approx. the true risk by a sample avg.

- The Pass function is not smooth  $\Rightarrow$  approx. by a smooth loss

$$\text{i.e } g(x) = P(y=1|x); f(x) = \begin{cases} +1 & \text{if } g(x) \in S_1 \\ -1 & \text{if } g(x) \in S_2 \end{cases}$$

→ you can consider the loss  $l(y, x, g) = \frac{1}{2} - y(g(x) - \frac{1}{2})$

$$\left\{ \begin{array}{l} \text{if } y = 1 \rightarrow l = 1 - g(x) \\ \text{if } y = -1 \rightarrow l = g(x) \end{array} \right. \quad \text{prob. of wrong pred.}$$

Goal:

$$\min_{\Theta \in \mathbb{R}^d} \frac{1}{N} \sum_{i=1}^N \max_{\|\tilde{x}_i - x_i\| \leq \epsilon} \ell(y_i, \tilde{x}_i, \Theta)$$

How:

At  $\theta_t$ : 1) for every  $i$ , find  $\tilde{x}_i = \arg \max_{x_i} \ell(y_i, x_i, \theta)$   
 $\|x_i - \tilde{x}_i\| \leq \epsilon$

2) Comput. grad:

$$\nabla_{\theta} l(y_i, \tilde{x}_i, \theta)$$

and do a grad. step

(min)

$\Rightarrow$  This is called **adversarial Training**, and obtains

state of the art result for robust classifier.

# Randomized Smoothing:

- How to make a given classifier robust?

Let  $f: \mathbb{R}^d \rightarrow \mathcal{Y}$  a given classifier.

$\Rightarrow$  then consider  $g$  s.t.  $g(x)$  will be the class the most likely returned by  $f$  when evaluated at  $x + z$  with  $z \sim N(0, \sigma^2 I)$

•  $g$  is the smooth classifier

$\Rightarrow$  consider  $x$  and  $y$ , we assume  $g(x) = +1$

For  $\tilde{x}$  s.t.  $\|x - \tilde{x}\| \leq \epsilon$ , we see that  $g(\tilde{x}) = +1$

**Robust Training**: an equivalent method is to consider a ball of radius  $\sqrt{D}\sigma$  (with a uniform dist.)

- Idea: consider the binary case, and original classifier  $f$ .

Take a point  $X$ , and let's say that label  $y=1$  is most likely for  $X + Z$ .

Now, if we take  $\tilde{X} + Z$ , where  $\|X - \tilde{X}\|_2 \leq \epsilon$  for some small  $\epsilon$ .

It is intuitive that in average we are still most likely to return  $y=1$  rather than  $y=-1$   
 $\rightarrow$  i.e. we will return the same label for points not too far away.

Example:  $D=1$ ,



- let  $p = \mathbb{E} [\frac{1}{\mathbb{I}_{f(x+z)=1}}]$  and assume  $\frac{1}{2} < p \leq 1$

$\Rightarrow$  min. of the points in the neighborhood of  $x$  are given  $y=1$   
now  $\tilde{p} = \mathbb{E} [\frac{1}{\mathbb{I}_{f(\tilde{x}+z)=1}}]$  when  $\|x - \tilde{x}\| \leq \epsilon$  rather than  $y=-1$

assume we move  $x$  to the right by  $\epsilon$ , how small can  $\tilde{p}$  be?

$\rightarrow$  in the worst case all the points labeled  $y=-1$  would be on the right of  $x$  in the tail of the gaussian. More precisely, the worst case happens if all the points to the left of  $x + \sigma Q^{-1}(1-p)$  are labeled  $y=1$  and all the points to the right  $y=-1$ .

$\Rightarrow$  how far can we move  $\tilde{x}$  from  $x$  to the right so still the majority is  $y=1$ ?

$\Rightarrow$  we see that we can move it at most by  $\sigma Q^{-1}(1-p)$

$\Rightarrow$  the longer the  $p$ , (more biased) the original average was, the more we are adversarial robust.

- Why not just take a very large  $\sigma$ ?

The averaging (smoothing) will in general increase the Standard Risk.  
(Trade off)