

## Announcements

- (1) HW#1 graded on gradescope
- (2) Reminder: project meetings,  
and HW#2

# Adversarial Robustness

Deep nets achieve human-level accuracy on many image classification tasks

**But are they robust?**

The seminal work of Szegedy et al. introduced the notion of adversarial examples

"For most images  $x$ , can find a perturbation  $\Delta$  st. the deep net misclassifies  $x + \Delta$ "

The perturbation is often imperceptible to humans, e.g.  $\|\Delta\|_{\infty} \leq 0.05$

Even worse, attacks can be

① targeted - can choose what it is misclassified as, appropriate choice of  $\Delta$

② black box - some attacks don't need to know the parameters of deep net they're attacking

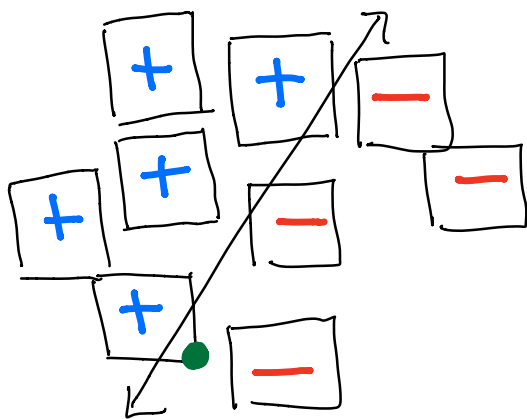
Main Question: Can we robustify deep learning?

Madry et al. introduced a popular defense based on adversarial training:

$$\min_{\theta} \mathbb{E}_{(x,y) \sim D} \left[ \max_{\|\Delta\|_p \leq \delta} L(x+\Delta, y; \theta) \right]$$

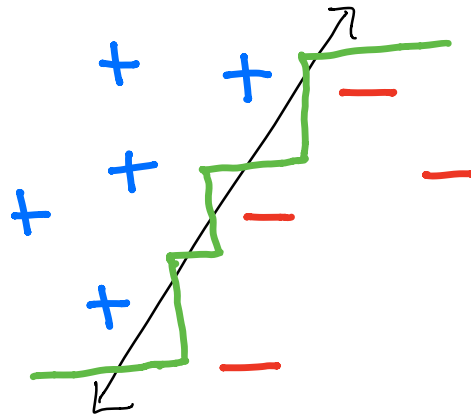
Essentially this is robust optimization applied to model fitting

Some intuition: Let's start with linearly separable data



Around each point, draw a  $\delta$ -ball in  $\mathbb{R}^d$  wherever the ball crosses the decision boundary, we get an adversarial example

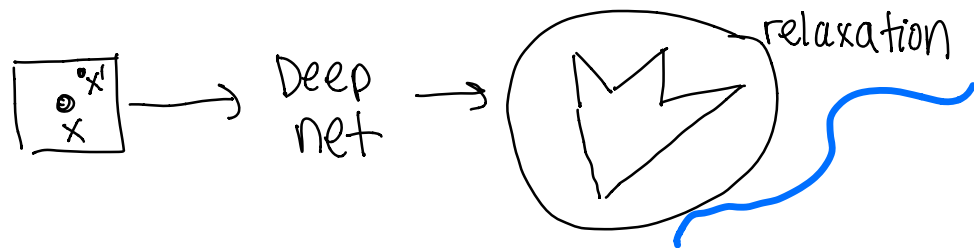
We can still find a robust decision boundary, but it becomes more complex



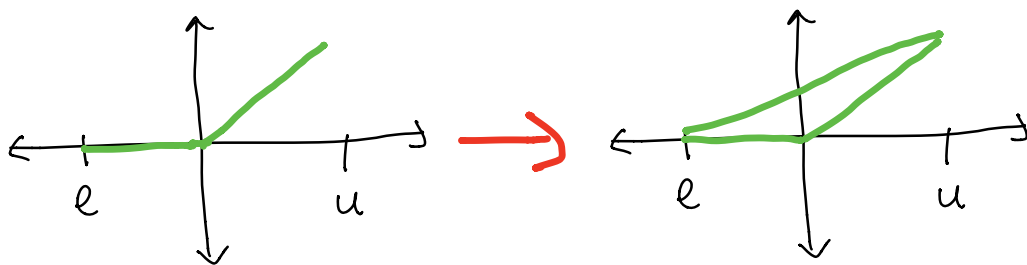
Nevertheless we use SGD to train

Adversarial training leads to significant robustness improvements, but generally at the cost of accuracy

Kolter and Wang gave certified defenses against adversarial examples via convex relaxations:



Their idea was to convexify the ReLU on a restricted domain:



to get an outer bound on

we'll study randomized smoothing, an approach for handling  $l_2$ -perturbations

LeCuyer et al were the first to give rigorous guarantees

The main ideas were

(1) Show that <sup>any</sup> differentially private classifier with a margin is provably robust

(2) Add a layer of noise to make the classifier differentially private

Here we mean differential privacy w.r.t. the pixel values

First we need a generalization of our earlier definition of privacy

def: Let  $p$  be a metric on databases. We say a randomized algorithm  $A$  is  $(\epsilon, \delta)$ -differentially private w.r.t.  $p$  if  $\forall$  databases  $m, m'$  with  $p(m, m') \leq 1$

$$\mathbb{P}(A(m) \in S) \leq e^\epsilon p(A(m') \in S) + \delta$$

$\forall S \subseteq \mathcal{O} \leftarrow \text{output space}$

Remember, for us databases  $\equiv$  images

Now consider a deep net  $f: \mathbb{R}^d \rightarrow [k]$   
set of labels  $\uparrow$

Most deep nets compute a score function

$$\varphi: \mathbb{R}^d \rightarrow [0, 1]^k$$

where the coordinates of  $\varphi(x)$  are interpreted as probabilities assigned to each label, and  $\sum_{i=1}^k \varphi(x)_i = 1$

Then  $f(x) = \operatorname{argmax}_{i \in [k]} \varphi(x)_i$

i.e. it outputs the "most probable" label

The main proposition is:



Proposition: Suppose  $\mathcal{A}^k$  is  $(\epsilon, \delta)$ -DP <sup>randomized scoring function</sup>  
w.r.t.  $\mathcal{P} \equiv \ell_p$ . Then for any input  $x$ , if  
for some  $i \in [k]$  we have

$$\underbrace{\mathbb{E}[\mathcal{A}_i(x)]}_{\phi(x)_i} > e^{2\epsilon} \max_{j \neq i} \mathbb{E}[\mathcal{A}_j(x)] + (1 + e^\epsilon) \delta$$

then  $f(x) = \operatorname{argmax}_{i \in [k]} \mathbb{E}[\mathcal{A}_i(x)]$  is robust  
to attacks  $\Delta$  with  $\|\Delta\|_p \leq 1$  on  $x$

Note: The classification rule is deterministic,  
because it's based on expectation of  $\mathcal{A}$

First we will prove a helper lemma that  
says:

"If  $\mathcal{A}$  is differentially private,  
then so is its expectation"

In particular:

Lemma 1: Suppose  $A$  is a randomized function with outputs in  $[0, b]$  that satisfies  $(\epsilon, \delta)$ -DP w.r.t.  $\mathcal{P} \equiv \mathcal{L}_p$ . Then

$$\mathbb{E}[A(x)] \leq e^\epsilon \mathbb{E}[A(x+\Delta)] + b\delta$$

$\forall \Delta$  with  $\|\Delta\|_p \leq 1$

Proof: We can write

$$\mathbb{E}[A(x)] \leq \int_0^b \mathbb{P}[A(x) > t] dt$$

Now invoking the definition of DP

$$\leq e^\epsilon \left( \int_0^b \mathbb{P}[A(x+\Delta) > t] dt \right) + \int_0^b \delta dt$$

$$= e^\epsilon \mathbb{E}[A(x+\Delta)] + b\delta$$



Now we are ready to prove the main proposition

Proof: Consider any  $\Delta$  with  $\|\Delta\|_p \leq 1$  and let  $x' = x + \Delta$ .

Then applying Lemma 1 coordinate wise, we have:

$$\mathbb{E}[A_i(x)] \leq e^\epsilon \mathbb{E}[A_i(x')] + \delta \quad (1)$$

$$\mathbb{E}[A_i(x')] \leq e^\epsilon \mathbb{E}[A_i(x)] + \delta \quad (2)$$

Now we will use (1) + (2) to give a lower bound on the score of label  $i$  on input  $x'$ , and an upper bound for every other label:

$$\mathbb{E}[A_i(x')] \stackrel{(1)}{\geq} \frac{\mathbb{E}[A_i(x)] - \delta}{e^\epsilon}$$

Using the hypothesis of the prop.,  
we have:

$$> \frac{e^{2\epsilon} \max_{j \neq i} \mathbb{E}[A_j(x)] + (1+e^\epsilon)\delta - \delta}{e^\epsilon}$$

$$= e^\epsilon \max_{j \neq i} \mathbb{E}[A_j(x)] + \delta$$

$$\stackrel{(2)}{\geq} \max_{j \neq i} \mathbb{E}[A_j(x')]$$

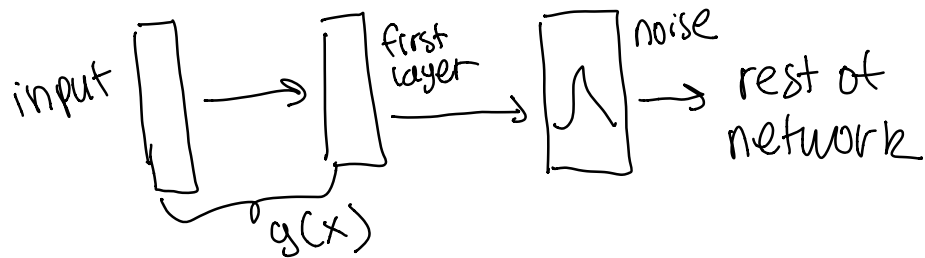
Thus we have

$$\mathbb{E}[A_i(x')] > \max_{j \neq i} \mathbb{E}[A_j(x')]$$

which means we still output label  $i$ .  $\square$

But what do we do if the randomized scoring function isn't DP?

Let's add a noise layer, e.g.



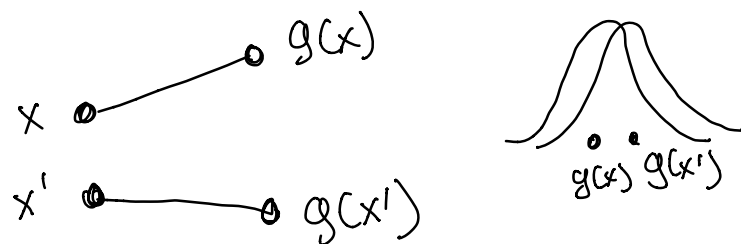
You could add it at the input layer, but Lucuyer et al. find that this is better

Main Question: How much noise should we add?

It's driven by the sensitivity, e.g.

$$N_{p,q} = \max_{x \neq x'} \frac{\|g(x) - g(x')\|_q}{\|x - x'\|_p}$$

Intuitively, the larger the sensitivity, the more noise you have to add



Note: If the weights in the first layer are  $W$ , and there is no nonlinearity, we have

$$N_{p,q} = \|W\|_{p,q} \triangleq \sup_{\|x\|_p \leq 1} \|Wx\|_q$$

Some  $p, q$ -norms are hard to compute, and they use relaxations instead

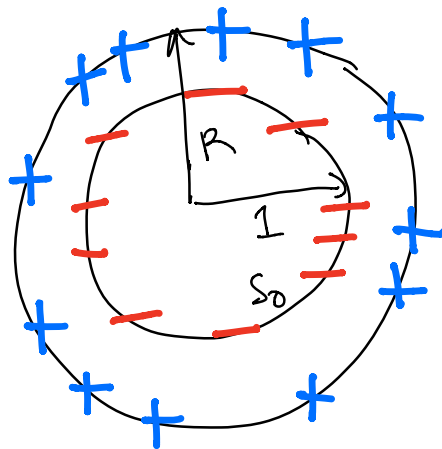
Can also train the network in a manner that makes it less sensitive

Epilogue: Yang et al. showed how optimal randomized smoothing schemes come from Wolff crystals

Interlude

Main Question: Are there situations where adversarial examples are inevitable?

Gilmer et al studied "adversarial spheres", where data is uniform on concentric spheres



They considered the following measure of adversarial robustness:

Let  $E = \{x \mid \text{model makes a mistake on } x\}$

$$\text{Let } d(E) = \frac{\int_x [\text{dist}(x, E)]}{x}$$

i.e. average Euclidean distance to the nearest error

They showed:

Theorem For any model that achieves accuracy  $\frac{1}{2} \leq p < 1$  on  $S_0$ , then

$$d(E) \leq O\left(\frac{\Phi^{-1}(p)}{\sqrt{d}}\right)$$

where  $d$  is the dimension,  $\Phi^{-1}$  is the inverse normal c.d.f.



The main idea is to use concentration of measure, e.g.

Theorem: Let  $E \subseteq S_0$  be measurable, with  $\mu(E) \geq \frac{1}{2}$ . Let  $E_\delta$  be the  $\delta$ -neighborhood, i. e.

$$E_\delta = \left\{ x \in S_0 \mid \exists y \in E \text{ with } \|x - y\|_2 \leq \delta \right\}$$

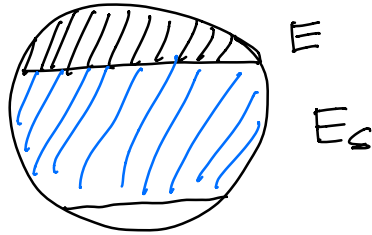
Then we have

$$\mu(E_\delta) \geq 1 - 2e^{-\frac{d\delta^2}{2}}$$

In particular, if  $\delta = \Theta\left(\frac{1}{\sqrt{d}}\right)$ , we cover almost all of the sphere

Note: There are versions that start from the assumption  $\mu(E) = p$ , as in our application

The key intuition is that, subject to  $\mu(E) = p$ , the set that maximizes  $d(E)$  is a spherical cap



And by concentration of measure, most of the sphere is within distance  $O(\frac{1}{\sqrt{d}})$

Does this explain why being robust to adversarial examples is difficult in practice?

Likely not!

Theorem [informal]: There are accurate classifiers for natural images that are robust to adversarial examples

Proof: Humans



Bubeck, Price and Razenshteyn showed:

Theorem [informal] If there exist robust classifiers, there is an exp. time algorithm to find one using just polynomially many samples

But efficiently finding one can be hard!

In particular, they give an instance where:

- ① robust learning is information-theoretically possible
- ② there is a computationally efficient non-robust learner
- ③ No SQ algorithm can find a robust classifier

Note: what makes this problem diff. than the usual ones in agnostic learning is that the notion of robustness has to do with perturbations of the input, not label corruptions

Follow-up work of Bubeck, Lee, Price and Razenshteyn showed cryptographic hardness, rather than just  $\Omega$  lower bounds

## Taking a Step Back

Main Question: what makes a classifier robust in real world applications?

Adversarial robustness started off with  $\ell_\infty$  perturbations, but could it be neither necessary nor sufficient?

Many reasons for wanting to be robust come from shifts in the distr.

## Data Poisoning Attacks

Lastly, we consider another threat model:

"If an adversary can corrupt an  $\epsilon$ -fraction of our training data, can we learn a good model?"

These questions were first posed in robust statistics

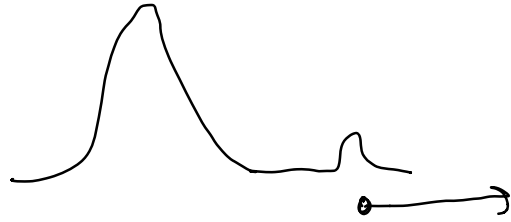
A canonical problem of robustly learning a Gaussian:

① we take  $N$  samples from  $\mathcal{N}(\mu, \sigma^2)$

② An adversary is allowed to corrupt an  $\epsilon$ -fraction arbitrarily

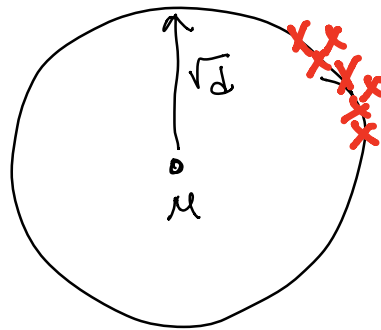


Intuitively, the reason this works is



As you send the  $\epsilon$ -bump to  $\infty$ , the empirical mean / variance diverge but the median and MAD don't

In high-dimensions, things are more tricky



Just an  $\epsilon$ -fraction of corruptions can move the model parameters by  $\epsilon\sqrt{d}$



There are approaches based on considering all 1-d projections, e.g.

def: The Tukey depth of point  $x$  among a sample set  $S$  is

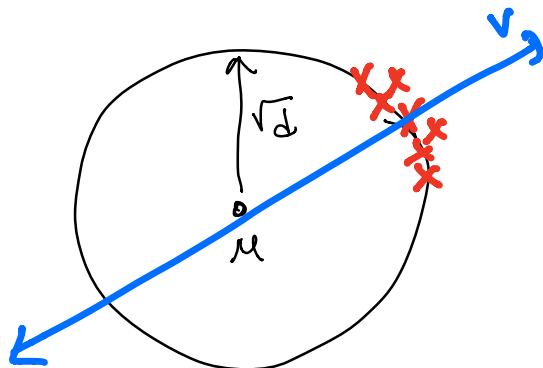
$$\min_{1\text{-d proj. } \pi_v} \min_{\substack{y \in S \\ \pi_v y < \pi_v x}} (|\{y \mid \pi_v y < \pi_v x\}|, \dots)$$

i.e. the worst case over all 1-d proj's. of the fewest points to left or right

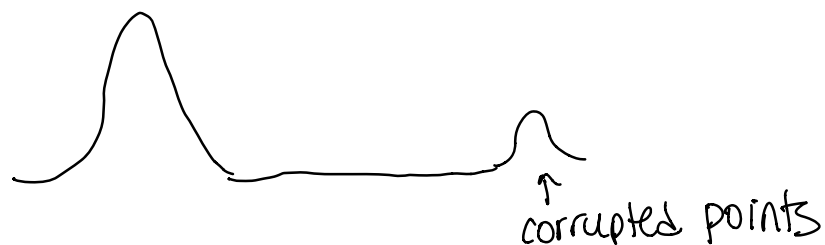
def: The Tukey median of a sample set  $S$  is the point  $x$  (not nec. in  $S$ ) that maximizes Tukey depth

Now it is easy to see that the Tukey median would be robust if we could compute it

For example, consider



Projecting onto  $v$ , we'd get:



Un fortunately:

Theorem [Johnson, Preparata] The Tukey median is NP-hard to compute in high-dimensions

Nevertheless we can get around this

Diakonikolas, Kamath, Kane, Li, Moitra and Stewart showed

Theorem: Given  $\epsilon$ -corrupted samples from a high-dimensional Gaussian  $\mathcal{N}(\mu, \Sigma)$  there is a poly time algorithm with poly. sample complexity that outputs  $\hat{\mu}, \hat{\Sigma}$  with

$$d_{TV}(\mathcal{N}(\mu, \Sigma), \mathcal{N}(\hat{\mu}, \hat{\Sigma})) \leq O(\epsilon \log^{3/2} 1/\epsilon)$$

Lai, Rao and Vempala gave related results independently

Tran, Li, Madry showed how to use these results to defend against data poisoning attacks in deep nets:

- repeat →
- ① Train a deep network
  - ② For each label  $i \in [k]$ , robustly compute the covariance of  $\{x \mid f(x) = i, x \in S\}$  and filter out unlikely  $x$ 's using  $\hat{\Sigma}$

They called these patterns in poisoned data spectral signatures

Intuition: Even when data is not Gaussian, looking for non-Gaussian projections can reveal interesting structure

Applications to Interpretability

Koh and Liang applied another

formalism from robust statistics,  
called the influence curve

"the change in an estimator from  
changing a single point"

They studied how a model's prediction  
is influenced by up/down weighting  
a sample in the training set

## Applications

- ① When you get a point wrong, can you blame some of your training data?
- ② Are you making a prediction because of meaningful similarities?

E.g. is your prediction most influenced by accidents (i.e. similar lighting / background) or the actual content