



THEORY OF MACHINE LEARNING

LECTURE 25

ROBUSTNESS (PART 2)



ANNOUNCEMENTS

- **Homework 4 - due last on Tuesday Apr 26**
 - Discussion → Issues in analyzing convergence of a non-convex opt problem.
- **Project presentations:** starting next week!
 - Please sign up! → On Tuesday, saw how to structure your pres.
 - Optional - submit presentation pdf on canvas for smoother transitions

RECAP: LEARNING IN THE PRESENCE OF ADVERSARIES

- Training time versus test time
- Training time: adversary corrupts small fraction of inputs
- Test time: adversary evaluates model on inputs with "imperceptible error" added (can be viewed as input distribution vs test)
- Former has multiple models - Benign noise, Huber's corruption model, data poisoning
- Field of robust statistics



ALGORITHMS AT HIGH LEVEL

- Use entire data, but limit *influence* of outliers */ not noisy points.*

- Median instead of mean (low dimensions)

- Truncated gradients

(clustering, mean estimation)

- "Inlier pursuit": key idea is that inliers "reinforce" one another

- RANSAC algorithm

- More sophisticated "filtering" algorithms

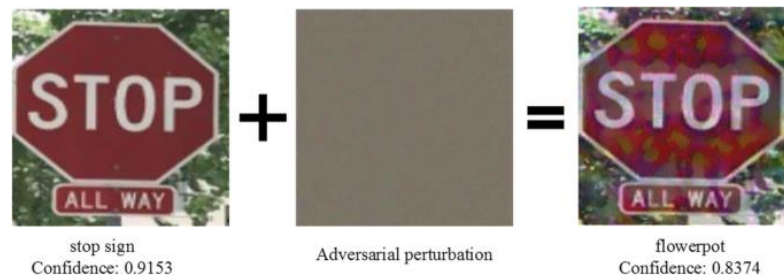
Promising idea/hypothesis: if you had only inliers a "simpler" model is possible.

- Main problem of study - robust mean estimation

- Can be used as a subroutine in other algorithms (use robust mean estimation for gradients!)

ROBUSTNESS OF TRAINED MODELS

- “Intriguing” property of deep learning models - models that generalize well are surprisingly brittle! [Szegedy et al. 2013]

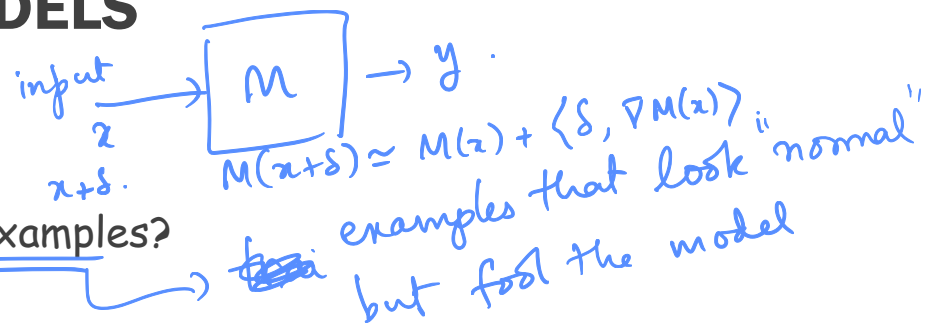


deep learning
models have tons of
params

- Obvious consequences
- Why possible? (statistical explanation)

out of dist.
input? so
generalization guarantees
don't mean anything.

ROBUSTNESS OF TRAINED MODELS



- How can we "generate" such adversarial examples?

- Now standard approaches

- Fast Gradient Sign Method (FGSM) - way to use the gradients to carefully choose direction of movement
- Projected Gradient Descent (PGD) - iterative procedure to "maximally" affect loss function

(implemented in standard DL libraries).

assume you have access to model..

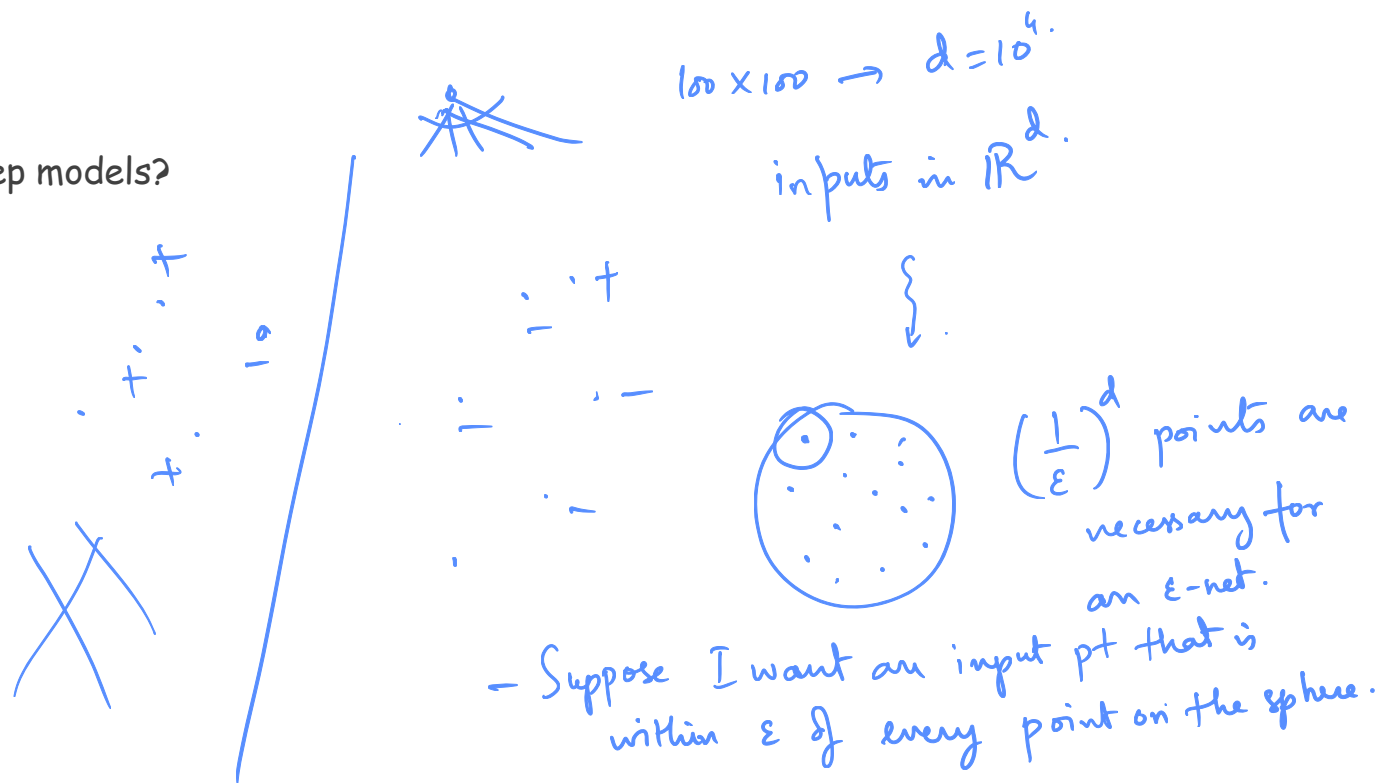
→ Same model arch ; ImageNet.

"generalizable" adversarial examples.

WHY DO ADVERSARIAL EXAMPLES EXIST?

- I.e., why only for deep models?

- Margin theory



"PROVABLY" DEFENDING?

- Can we show that no corruption of small magnitude can hurt the classifier?

- "Adversarial training":

- Instead of minimizing empirical risk, minimize empirical "robust risk"

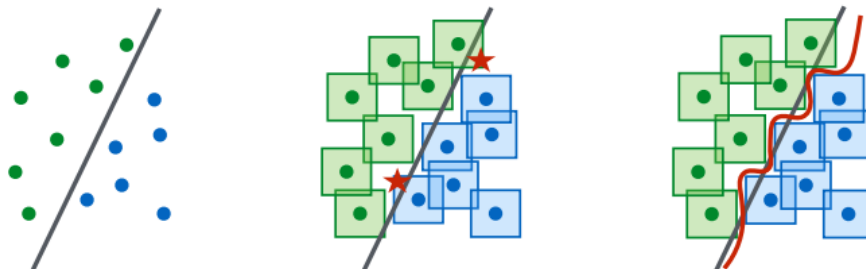
minimize $\mathbb{E}_{x \sim \mathcal{D}_{\text{inputs}}}$ $\mathbb{E}_{y \sim \text{label}(x)}$ $l(x, y; w)$ parameters.

min $\mathbb{E}_{x \sim \mathcal{D}_{\text{inputs}}} \max_{\delta, \|\delta\| \leq \epsilon} l(x + \delta, \text{label}(x); w)$. [Wald].

Saddle point opts.

- How to solve this optimization problem? (uses theorem of Danskin - gives a way of solving min-max opt problems) \rightarrow [Madry, ... '2018].

- Do we need "richer models"?



Trade off between robustness & "regular" accuracy.

Tradeoff between robustness & accuracy.

$x_1, x_2, \dots, x_N.$

clean training
data.

model M \rightarrow clean accuracy $\rightarrow 80\%.$ \rightarrow noisy accuracy $\rightarrow 5\%.$

model M' trained ~~on~~ using min max training \rightarrow noisy accuracy $\rightarrow 60\%.$
clean accuracy $\rightarrow 65\%.$

"PROVABLY" DEFENDING – CONNECTION TO PRIVACY

Differential privacy:

dataset $\underline{X} \rightarrow$ 

Clever formulation

by [work. Nissim, Naor]
Smith.



Suppose my "trusted party" ensures
that \forall query that analyst gives,

$\text{Answer}(x) \ \& \ \text{Answer}(x')$, where
you flipped
are indistinguishable (you bit).

$$\begin{array}{c} \text{---} x \quad \text{---} \frac{1}{0} \text{---} \\ \text{---} \quad \frac{0}{0} \text{---} \end{array}$$

$$\text{Answer}(x) \simeq \text{Answer}(x + \delta)$$



MORE NUGGETS

- Do we need more data for obtaining robust models?
- State of the art accuracies
- Expressibility vs trainability vs robustness