

SRN Aleksander Madry

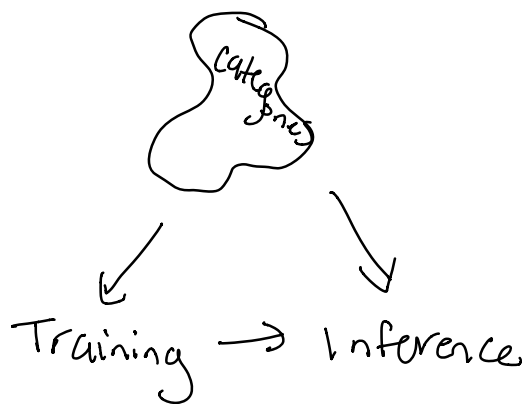
Machine Learning: A Robustness Perspective

Is it truly ready for real-world deployment?

- can be manipulated to exploit host
- ex: tesla car drove straight into overturned car

Robust ML: A Challenge

- ImageNet: An ML Home Run → over 1 mil images so that computer can predict image
 - ↓ error over the years
- Reflect:
 - A limitation of the (Supervised) ML Framework



measure of performance:
fraction of mistakes
during testing

But: In reality, the
distribution we use
ML are NOT the
ones we train it on

What can go wrong?

ML Predictions are (mostly) accurate
but brittle

ex: ^{w/ confidence of}
pig (91% recognition) but w/ a little
noise in image, the computer ~~but~~
thinks it's an airliner

ex: rotation + translation suffices to
fool state-of-the-art vision models
- data augmentation ~~does~~ NOT
seem to help here

Why is the brittleness of ML a problem?

- Security (recognize wrong criminals
w/ sunglasses or noise)
- Safety (car, miss something dangerous
w/ self-driving cars)
- ML Alignment (need to understand
the "failure modes" of ML)

ML pipeline (via adversarial lens)

data collection \rightarrow training \rightarrow inference \rightarrow deployment

Is ML inherently not reliable?

No! But need to rethink it?

Why are our models brittle?

$\rightarrow d \Rightarrow \infty$ (wired dimensions of objects)

\rightarrow only optimizing for avg case but not worst case

Why are adv. perturbations bad?

dog + meaningless perturbation = cat

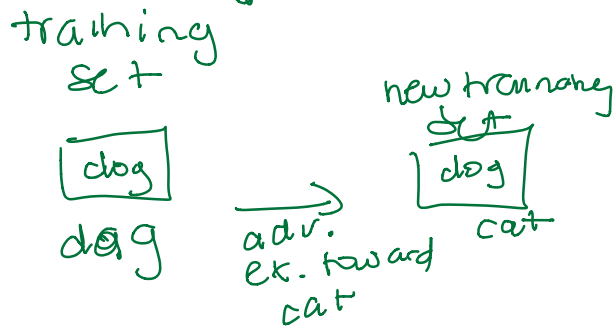
(we can tell it is still a dog, but ML thinks it's a cat)

- ML has no concept of dog

\swarrow
image is meaningless

\searrow
classes are meaningless

Are adv. perturbations just meaningless?



- 1) Make adversarial example toward the other class
- 2) Relabel the image as the target class
- 3) Train w/ new dataset but test on the original test set

So we train on a "totally mislabeled" dataset but expect performance on a "correct" dataset

What will happen?

- we get a nontrivial accuracy on the original classification task

What's going on?

The Robust Features Model

Robust Features:

correlated w/
label even w/
adversary

Non-robust features:

Correlated w/ label on
avg but can be flipped
w/in

All robust models are misleading but...