Adversarial Examples

Neil Gong

Today's lecture

- What is adversarial example
- Why do we care
- How to find adversarial example

Adversarial Examples





Normal example: digit 0

Adversarial example: predicted to be 9

Adversarial Examples

- Classifier C
- Normal example x
 - Image, text, audio, graph, software
- Perturb x to x'
 - Preserving semantics
- C misclassifies x'
 - *Targeted*: *C*(*x*')=*t*, an attacker-chosen target label
 - Untargeted: $C(x') \neq C(x)$

Why do we care?



Stop sign to speed limit

Malware -> benign software

Spam -> non-spam

Privacy protection

Guiding design of ML

Attacker's Background Knowledge

Access to prediction API Learning algorithm Black box Hyperparameter, e.g., neural network architecture Weakest attacker Training data

Model parameters White box Strongest attacker

How to Find Adversarial Examples - Image Domain

- Perturb x to x'
 - Preserving semantics
 - Human perceives x' and x as the same
 - d(x,x') is small

 L_0, L_2, L_∞ norm of the noise x'-x

Minimize d(x,x')Subject to (1) C(x') = t or $C(x') \neq C(x)$ (2) x' is still an image

Solving the optimization problem

minimize $\|\delta\|_p + c \cdot f(x+\delta)$ such that $x + \delta \in [0, 1]^n$ Box constraints

Loss function

$$f_{1}(x') = -\log_{F,t}(x') + 1$$

$$f_{2}(x') = (\max_{i \neq t} (F(x')_{i}) - F(x')_{t})^{+}$$

$$f_{3}(x') = \operatorname{softplus}(\max_{i \neq t} (F(x')_{i}) - F(x')_{t}) - \log(2)$$

$$f_{4}(x') = (0.5 - F(x')_{t})^{+}$$

$$f_{5}(x') = -\log(2F(x')_{t} - 2)$$

$$f_{6}(x') = (\max_{i \neq t} (Z(x')_{i}) - Z(x')_{t})^{+}$$

$$f_{7}(x') = \operatorname{softplus}(\max_{i \neq t} (Z(x')_{i}) - Z(x')_{t}) - \log(2)$$

Box constraints

- Projected gradient descent
- Clipped gradient descent
 - Incorporate clipping into objective function

$$f(\min(\max(x+\delta,0),1))$$

• Change of variables

$$\delta_i = \frac{1}{2}(\tanh(w_i) + 1) - x_i$$

Examples



Evaluation metrics – what is a successful adversarial example

- Misclassification
 - *Targeted*: *C*(*x*')=*t*, an attacker-chosen target label
 - Untargeted: $C(x') \neq C(x)$
- Human perceives x' and x as the same
 - Hard to implement involves user studies
 - Approximate using L_p norm of noise

Other methods

- Beyond L_p norm
- Physically realizable adversarial examples





Beyond images

- Text
- Audio
- Video
- Software

Preserving semantics

C misclassifies x'

Formulation as optimization problem