

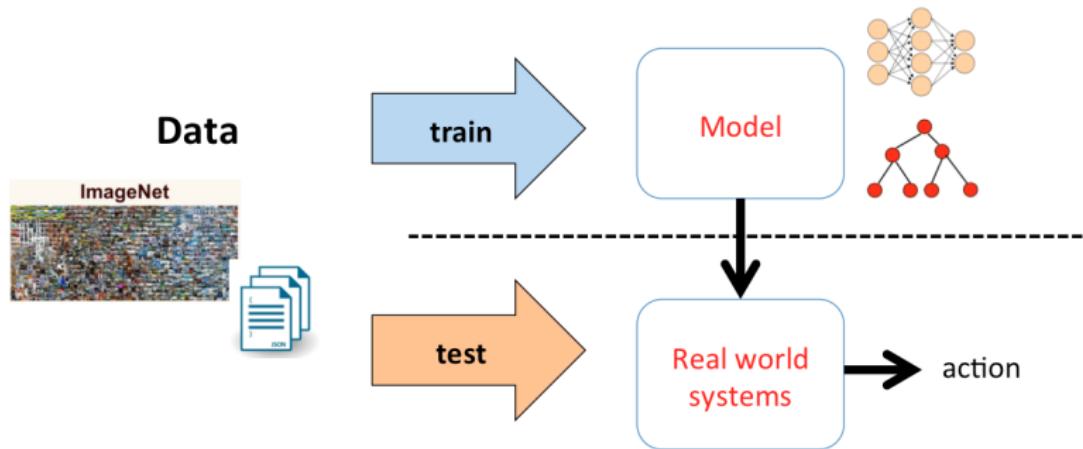
STOR566: Introduction to Deep Learning

Lecture 14: Adversarial Attack

Yao Li
UNC Chapel Hill

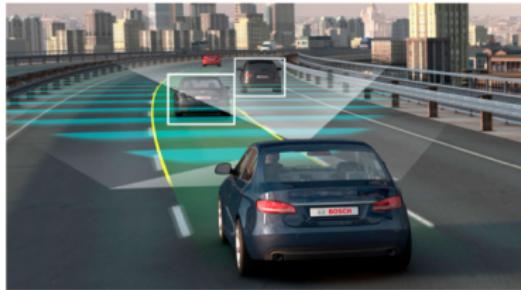
Oct 22, 2024

Machine Learning Systems



Testing: Robustness and Safety

ML systems need to interact with real world



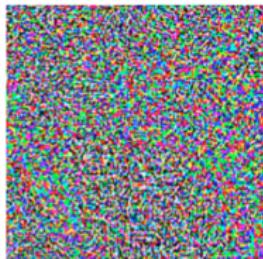
- Robustness and Safety

Adversarial Examples

Adversarial Examples



+ 0.001 ×



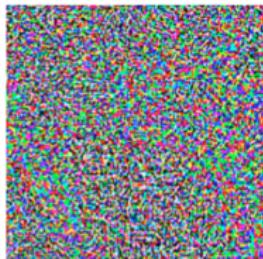
=



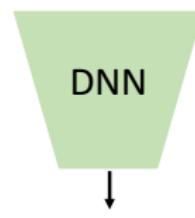
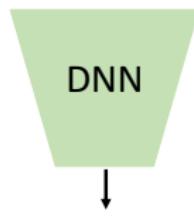
Adversarial Examples



+ 0.001 ×



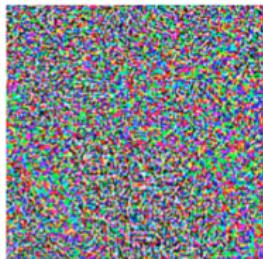
=



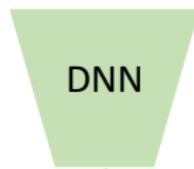
Adversarial Examples



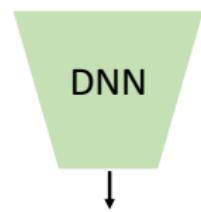
+ 0.001 ×



=



Bagel

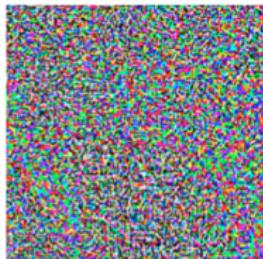


Bagel

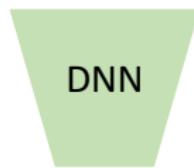
Adversarial Examples



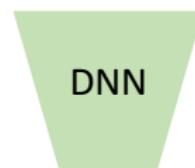
+ 0.001 ×



=



Bagel

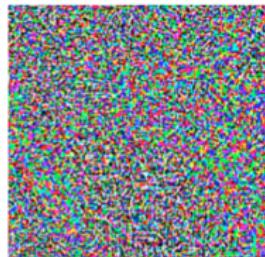


Piano

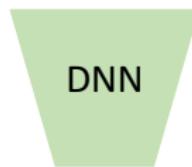
Adversarial Examples



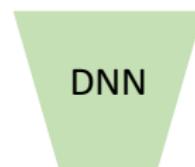
+ 0.001 ×



=



Bagel

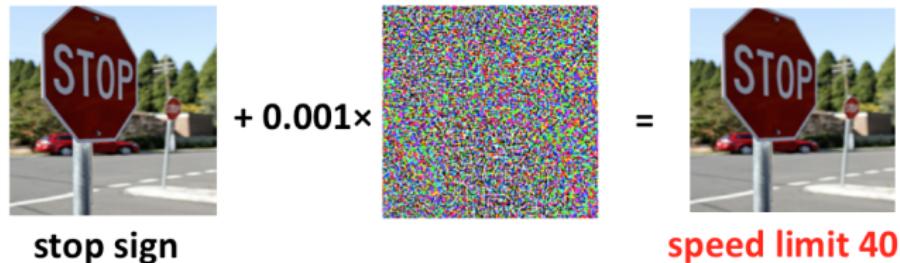


Piano

A carefully crafted **adversarial** example can easily fool a deep network

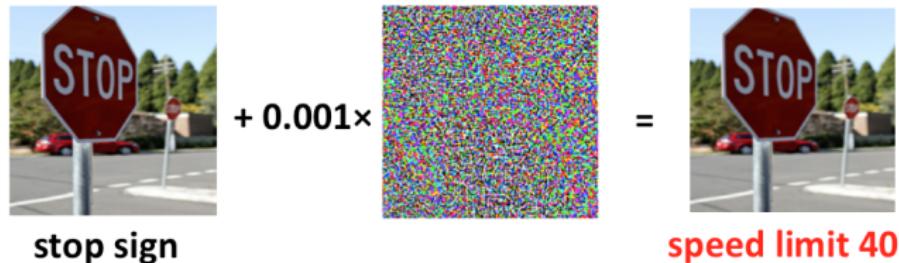
Adversarial Examples

- Robustness is critical in real systems



Adversarial Examples

- Robustness is critical in real systems



Not safe!

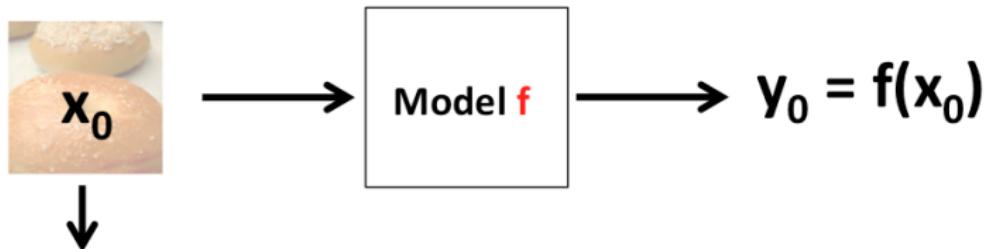
Adversarial Example in NLP

BAE attack on IMDB	
Groundtruth	Label changed : Positive → Negative
Original	This film offers many delights and surprises.
Attacked	This beautiful movie offers many pleasant de-lights and surprises.

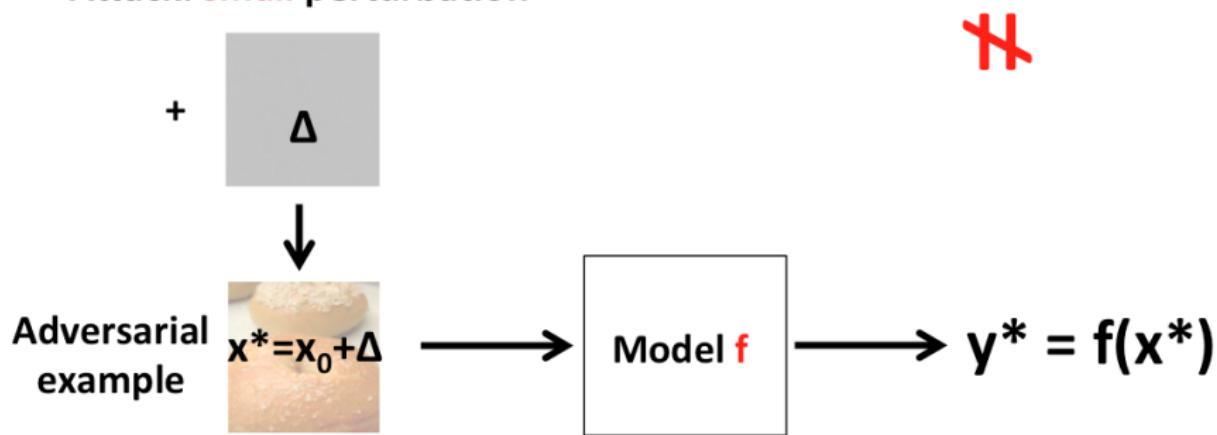
Table: BERT-based Adversarial Examples for Text Classification on IMDB dataset.
(Inserts: Red, Replacements: Blue)

Adversarial Attack

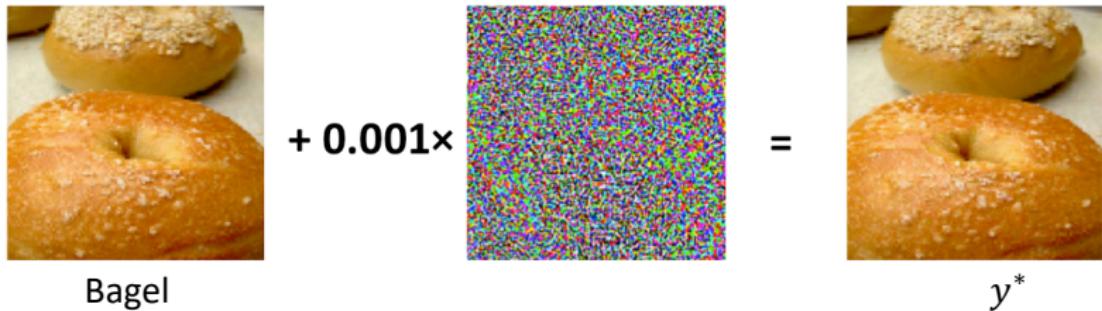
Notations and Attack Procedure



Attack: **small** perturbation



Type of Attacks



- ① Untargeted attack: $y^* \neq \text{Bagel}$
- ② Targeted attack: For target class $t = \text{Piano}$, the attacker wants $y^* = \text{Piano}$

Gradient-based Attack

Attack as Optimization

Attack as an **optimization** problem:

- Given prediction model with fixed parameter θ
- (x_0, y_0) : input image and label

$$\delta = \arg \max_{\delta \in \mathcal{S}} L(\theta, x_0 + \delta, y_0)$$

$$x^* = x_0 + \delta$$

L : loss function training the classifier

δ : adversarial perturbation

$\mathcal{S} \in \mathbb{R}^d$: allowed perturbation set, usually chosen to be $\{\delta | \|\delta\|_\infty \leq \epsilon\}$

FGSM and PGD

Fast gradient sign method (FGSM):

- One-step gradient ascent:

$$\mathbf{x}^* = \mathbf{x}_0 + \epsilon \cdot \text{sign}\left(\nabla_{\mathbf{x}} L(\theta, \mathbf{x}_0, y)\right) \quad (1)$$

Projected Gradient Descent Attack (PGD):

- Multiple-steps gradient ascent:

$$\mathbf{x}^{t+1} = \Pi_{\epsilon} \left\{ \mathbf{x}^t + \alpha \cdot \text{sign}\left(\nabla_{\mathbf{x}} L(\theta, \mathbf{x}^t, y)\right), \mathbf{x}_0 \right\} \quad (2)$$

$\Pi_{\epsilon}(\cdot, \mathbf{x}_0)$: projection to the set $\{\mathbf{x} | \|\mathbf{x} - \mathbf{x}_0\|_{\infty} \leq \epsilon\}$

α : step size

Carlini and Wagner Attack (C&W Attack)

Given model $f(\cdot)$ and input image and label (\mathbf{x}_0, y_0)

Craft adversarial example by solving

$$\arg \min_{\mathbf{x}} \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda \cdot g(\mathbf{x})$$

(Carlini et al., 2017)

C&W Attack

Given model $f(\cdot)$ and input image and label (\mathbf{x}_0, y_0)
Craft adversarial example by solving

$$\arg \min_{\mathbf{x}} \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda \cdot g(\mathbf{x})$$

- $\|\mathbf{x} - \mathbf{x}_0\|^2$: the distortion

(Carlini et al., 2017)

C&W Attack

Given model $f(\cdot)$ and input image and label (\mathbf{x}_0, y_0)

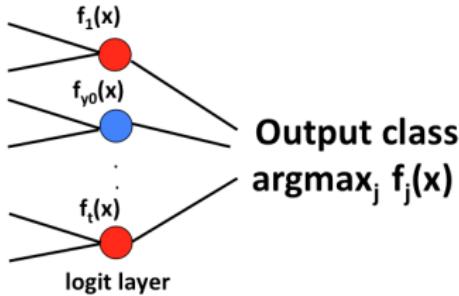
Craft adversarial example by solving

$$\arg \min_{\mathbf{x}} \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda \cdot g(\mathbf{x})$$

- $\|\mathbf{x} - \mathbf{x}_0\|_2^2$: the distortion
- $g(\mathbf{x})$: loss to measure the **successfulness** of attack
- $\lambda \geq 0$: controls the trade-off

(Carlini et al., 2017)

Untargeted attack

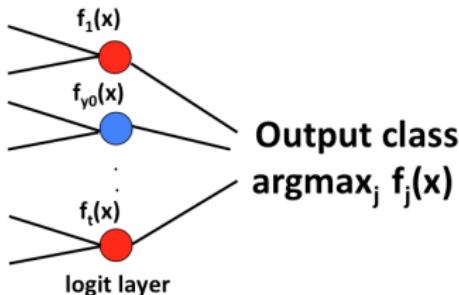


- $g(\mathbf{x})$: loss to measure the **successfulness** of attack

Untargeted attack: success if $\arg \max_j f_j(\mathbf{x}) \neq y_0$

$$g(\mathbf{x}) = \max\{f_{y_0}(\mathbf{x}) - \max_{j \neq y_0} f_j(\mathbf{x}), 0\}$$

Targeted attack



- $g(\mathbf{x})$: loss to measure the **successfulness** of attack

Targeted attack: success if $\arg \max_j f_j(\mathbf{x}) = t$

$$g(\mathbf{x}) = \max\{\max_{j \neq t} f_j(\mathbf{x}) - f_t(\mathbf{x}), 0\}$$

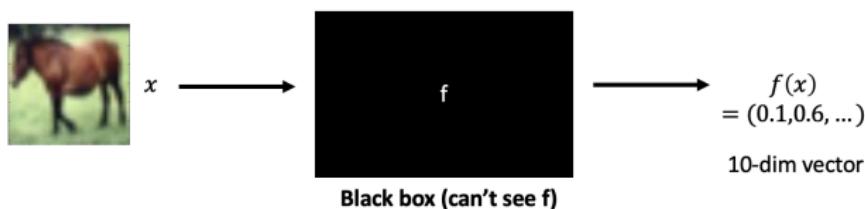
Gradient-Based (White-box) Setting

- Model (network structure and weights) is revealed to attacker
 - ⇒ gradient of $g(\mathbf{x})$ (or $\nabla_{\mathbf{x}} L(\theta, \mathbf{x}^t, y)$) can be computed by back-propagation
 - ⇒ attacker searches for \mathbf{x}^* by gradient descent
- Distortion can be measured by other norms:
(e.g., ℓ_∞ , Elastic net, ...)
- Black-box setting: Only part of the information is available

Score-based Attack

Score-based Setting

- In practice, the deep network parameters are not revealed to attackers
cannot compute gradient
- Score-based setting: Attacker can **query** the model and get the **score vector**



- Example: CIFAR10 image classification
 $f(\mathbf{x}) \in \mathbb{R}^{10}$, $f(\mathbf{x})_i$: predicted score of class i

Zeroth Order Optimization based Attack (ZOO)

- We can solve the following problem without access to parameters:

$$\arg \min_{\mathbf{x}} \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda \cdot g(\mathbf{x})$$

- Estimate gradient:

$$\hat{g}_i \approx \frac{\partial f(\mathbf{x})}{\partial x_i} \approx \frac{f(\mathbf{x} + \delta \mathbf{e}_i) - f(\mathbf{x} - \delta \mathbf{e}_i)}{2\delta}$$

Zeroth Order Optimization based Attack (ZOO)

- We can solve the following problem without access to parameters:

$$\arg \min_{\mathbf{x}} \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda \cdot g(\mathbf{x})$$

- Estimate gradient:

$$\hat{g}_i \approx \frac{\partial f(\mathbf{x})}{\partial x_i} \approx \frac{f(\mathbf{x} + \delta \mathbf{e}_i) - f(\mathbf{x} - \delta \mathbf{e}_i)}{2\delta}$$

- zero-order gradient descent

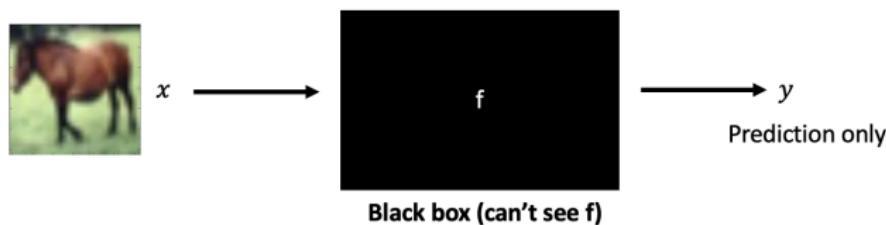
$$\mathbf{x} \leftarrow \mathbf{x} - \eta \begin{pmatrix} \hat{g}_1 \\ \vdots \\ \hat{g}_d \end{pmatrix}$$

(Chen et al., 2017)

Decision-based Attack

Decision-based Setting

- In practice, the attacker may only have access to the predicted label
- Decision-based setting: Attacker can **query** the model and get the **predicted label**



- Example: CIFAR10 image classification
 $y \in \{1, \dots, K\}$, K : total number of classes

Substitute Attack

- Intuition: adversarial examples transfer across models

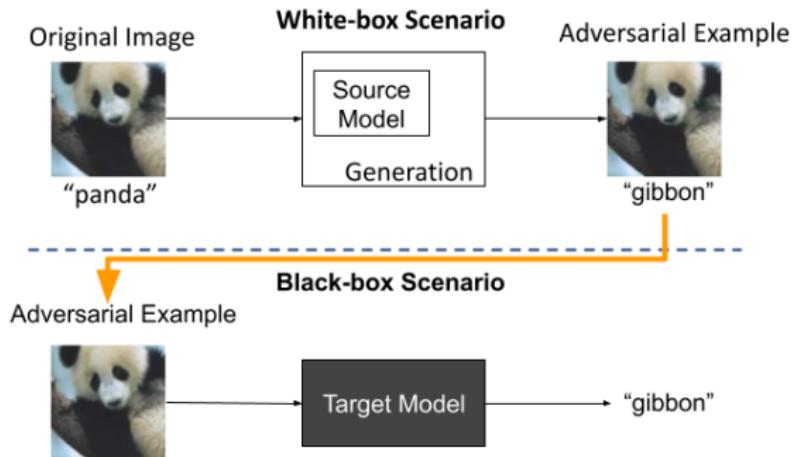
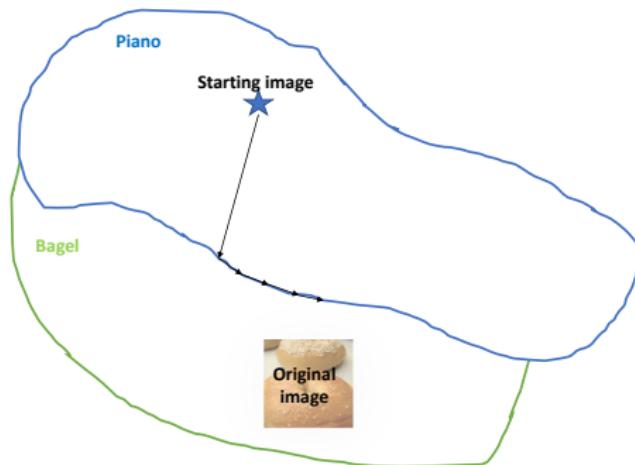


Figure from [link](#)

- Train a **substitute** model using a small amount of training data.
- Generate adversarial examples based on the substitute model.
(Papernot et al., 2017)

Boundary Attack

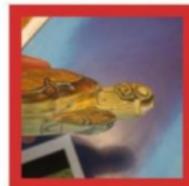


- Start as a random example from a target class.
- Perform rejection sampling along the boundary.
(Brendel et al., 2018)

Attack in Physical World



Attack in Physical World



■ classified as turtle

■ classified as rifle

■ classified as other

Conclusions

- Adversarial examples
- Different types of adversarial attacks

Questions?