

Trustworthy AI

Spring 2024

Yuan Tian

#2: Adversarial Machine Learning and Defenses

Pointing the need for broader view on AI

Human-Level Intelligence or Animal-Like Abilities?

Communications of the ACM, Oct 2018

Adnan Darwiche

<https://cacm.acm.org/magazines/2018/10/231373-human-level-intelligence-or-animal-like-abilities/fulltext>

“...We need a new generation of AI researchers who are well versed in and appreciate **classical AI, machine learning**, and computer science more broadly while also being informed about AI history...”

Today: Adversarial Examples

What are these? More examples in various domains

Why do they exist?

How to generate adversarial examples?

(techniques also used for adversarial training, logic, etc...discussed in later lectures)

“Adversarial examples are inputs to machine learning models that an attacker has intentionally designed to cause the model to make a mistake”

(Goodfellow et al 2017)

Adversarial Examples

Noisy attack: vision system thinks we now have a gibbon...

x
 “panda”
 57.7% confidence

$+ .007 \times$

$\text{sign}(\nabla_x J(\theta, x, y))$
 “nematode”
 8.2% confidence

$=$

$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$
 “gibbon”
 99.3 % confidence

Explaining and Harnessing Adversarial Examples, ICLR '15

Tape pieces make network predict a 45mph sign



Robust Physical-World Attacks on Deep Learning Visual Classification, CVPR'18

Self-driving car: in each picture one of the 3 networks makes a mistake...



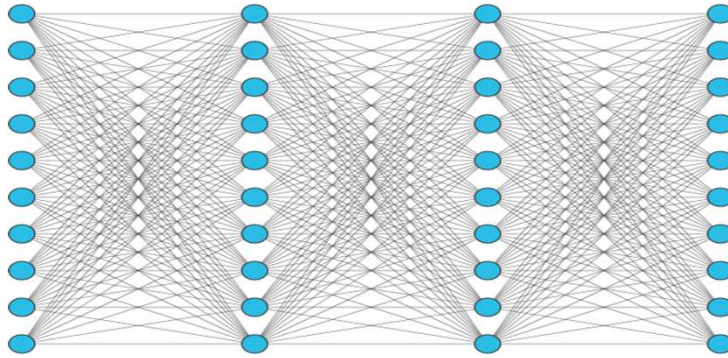
DRV_C1: right DRV_C2: right DRV_C3: right

DeepXplore: Automated Whitebox Testing of Deep Learning Systems, SOSP'17

Adversarial Geometric Perturbations



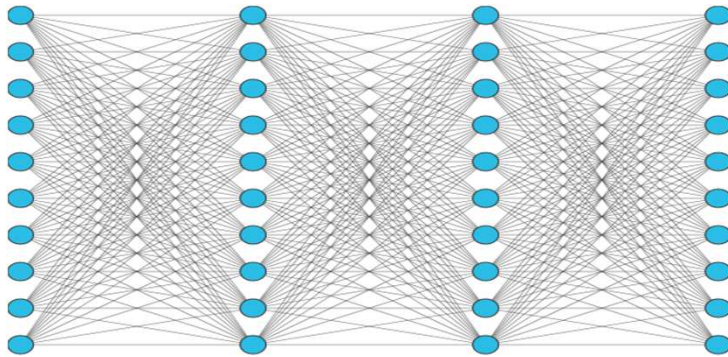
I_0



7



$I = rotate(I_0, -35)$



3

Adversarial Examples (more)

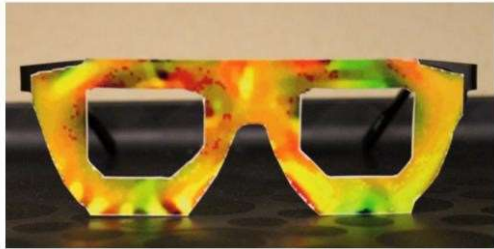


Reese Witherspoon

Russel Crowe

Real World Impersonation/Dodging Attacks

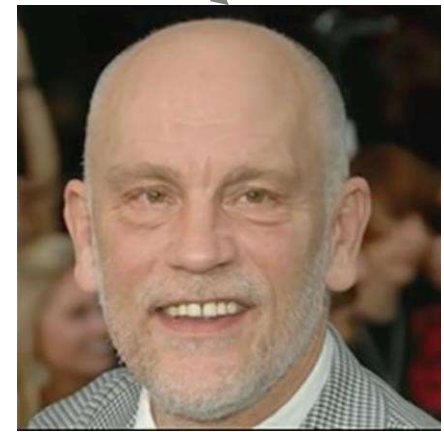
Real glasses



Lujo Bauer



=

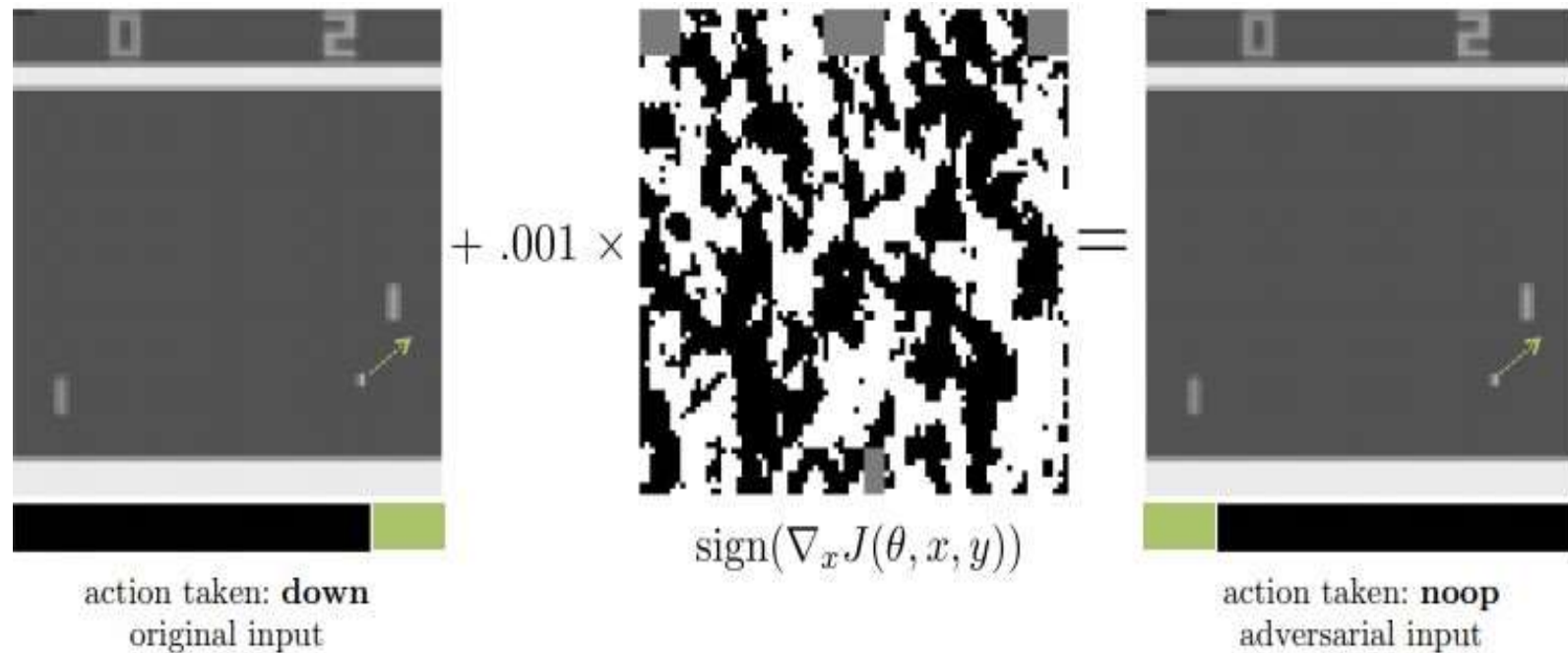


John Malkovich

100% success



Adversarial Examples in Reinforcement Learning



An agent (Deep Q Network) plays the game by selecting actions from a given state (image) that the game produces.

An attacker can perturb the image slightly so that the DQN agent chooses the wrong action: here, it wrongly picks noop (do nothing) in the right image, instead of moving the paddle down (left image).

Adversarial Examples in NLP

Article: Super Bowl 50

Paragraph: *“Peyton Manning became the first quarterback ever to lead two different teams to multiple Super Bowls. He is also the oldest quarterback ever to play in a Super Bowl at age 39. The past record was held by John Elway, who led the Broncos to victory in Super Bowl XXXIII at age 38 and is currently Denver’s Executive Vice President of Football Operations and General Manager. Quarterback Jeff Dean had jersey number 37 in Champ Bowl XXXIV.”*

Question: *“What is the name of the quarterback who was 38 in Super Bowl XXXIII?”*

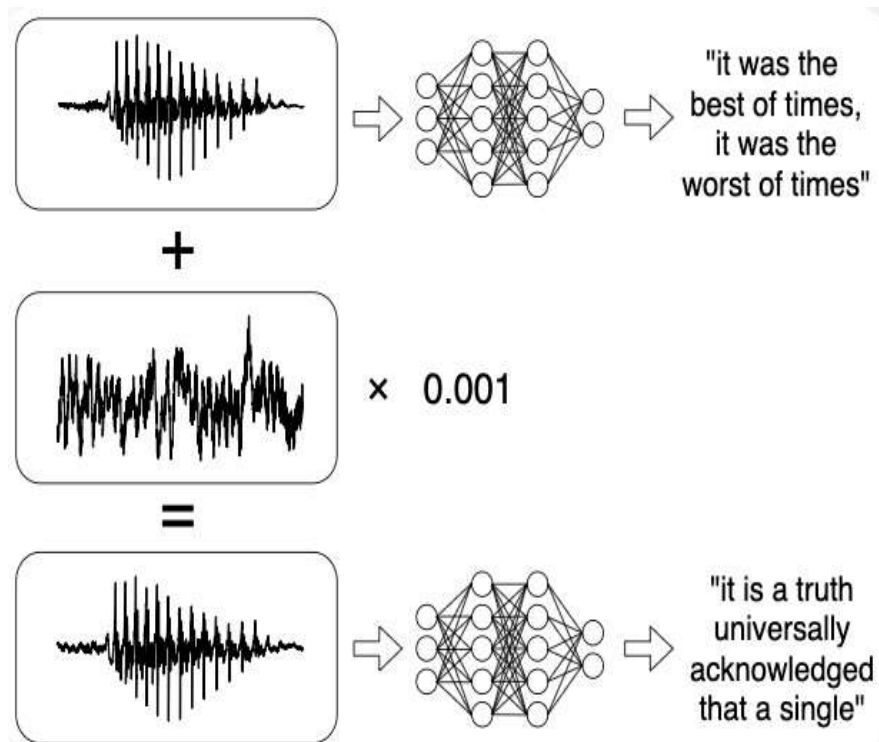
Original Prediction: John Elway

Prediction under adversary: Jeff Dean

The Ensemble model is fooled by the addition of an adversarial distracting sentence in blue.

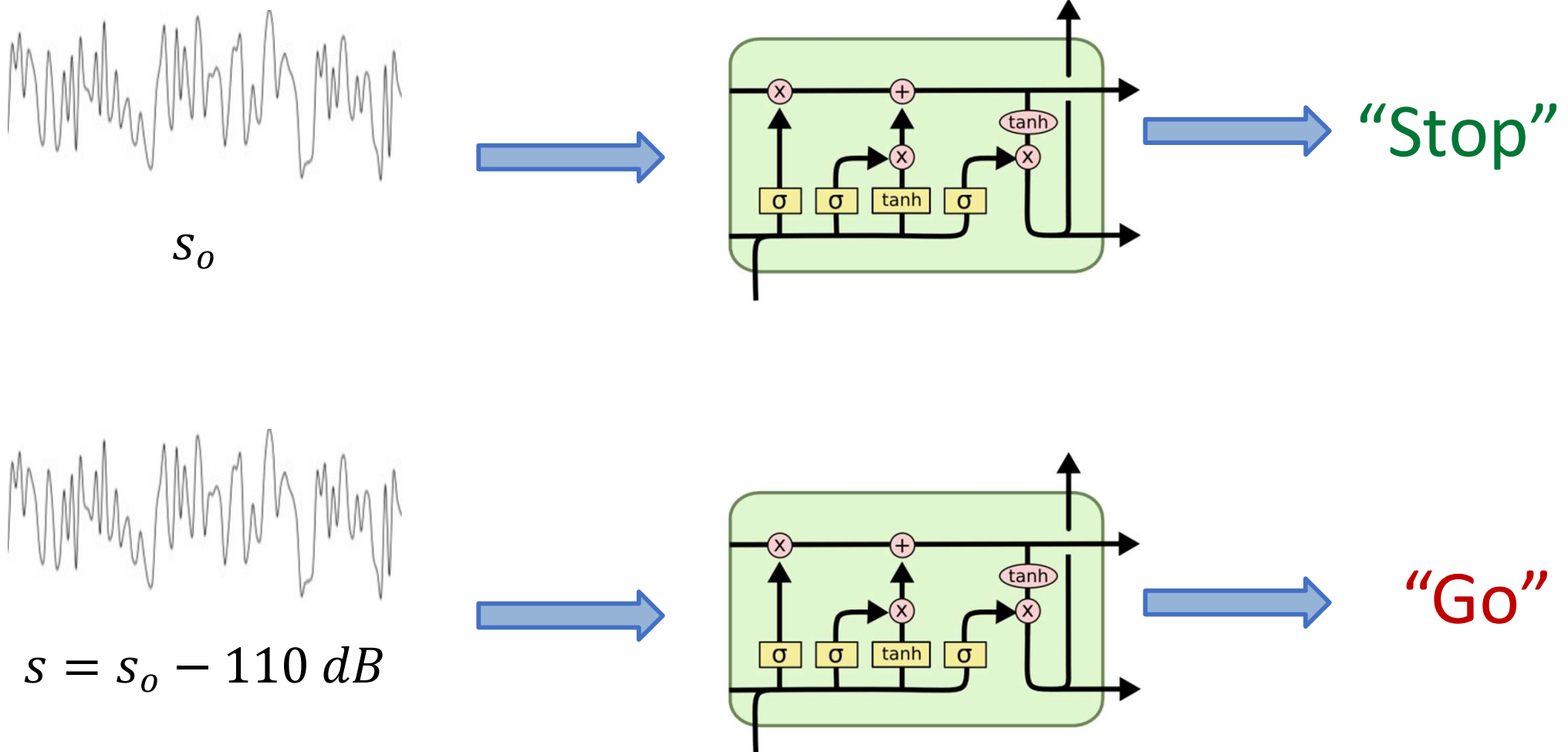
Adversarial Examples in Audio Processing: Speech to Text

An attack on DeepSpeech:



Adding small noise to the input audio makes the network transcribe any arbitrary phrase

Adversarial Examples in Audio Processing: Text classification



Adding small noise to the input audio changes the classification of the output

Adversarial Examples: Some History

2006: Deep learning models gain renewed interest

2012: Multiple works showed that deep networks can achieve near-human performance (sometimes even better)

2013: Research in understanding neural networks behavior becomes critical with society implications beyond computer science

2014: While trying to understand decision making in neural networks, Szegedy et al. discovered adversarial examples

2015-on: Finding adversarial examples and proving their absence becomes an active research area...

2023-on : Jail break LLMs, the alignment problem...

Robustness

Robustness: A network is **robust** if it returns correct output on all inputs

Impractical: the input space is too large to be covered

Local Robustness (informal): A learning model is **locally-robust** if it returns the correct output on inputs *similar* to inputs in the training set

This was believed to be evident by having high accuracy on the test set

Why is High Accuracy Not Enough?

Inputs in the training and test set are taken from a given distribution

Neural networks aim to achieve **high accuracy** on test sets drawn **from the given distribution**

There are still many similar inputs that are never tested (and have low-probability for the given distribution)

The Story of Clever Hans



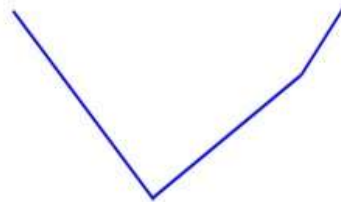
Why Do Adversarial Examples Exist?

Neural Networks are too linear

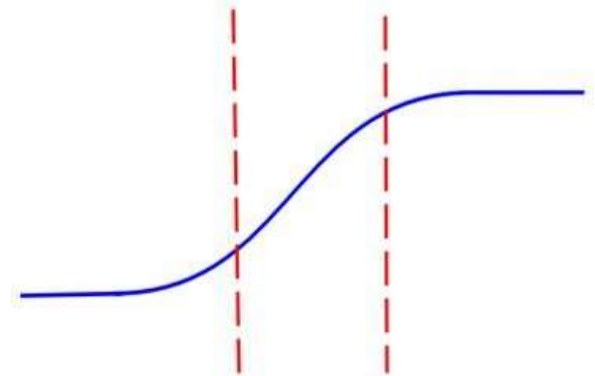
Rectified linear unit



Maxout



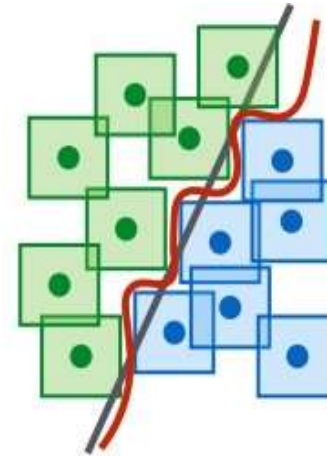
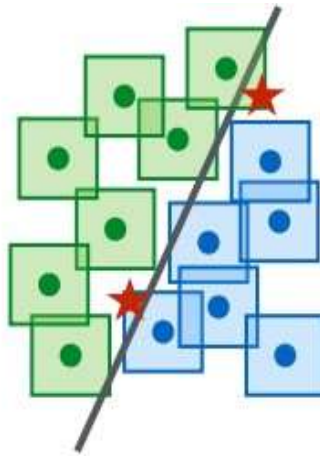
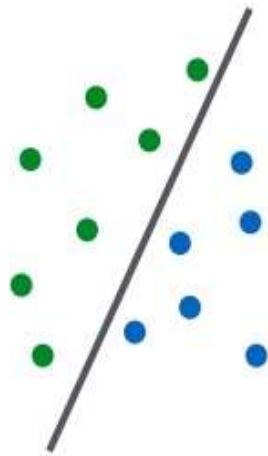
Carefully tuned sigmoid



Why are they designed to be linear?

Linear functions are easy to optimize!

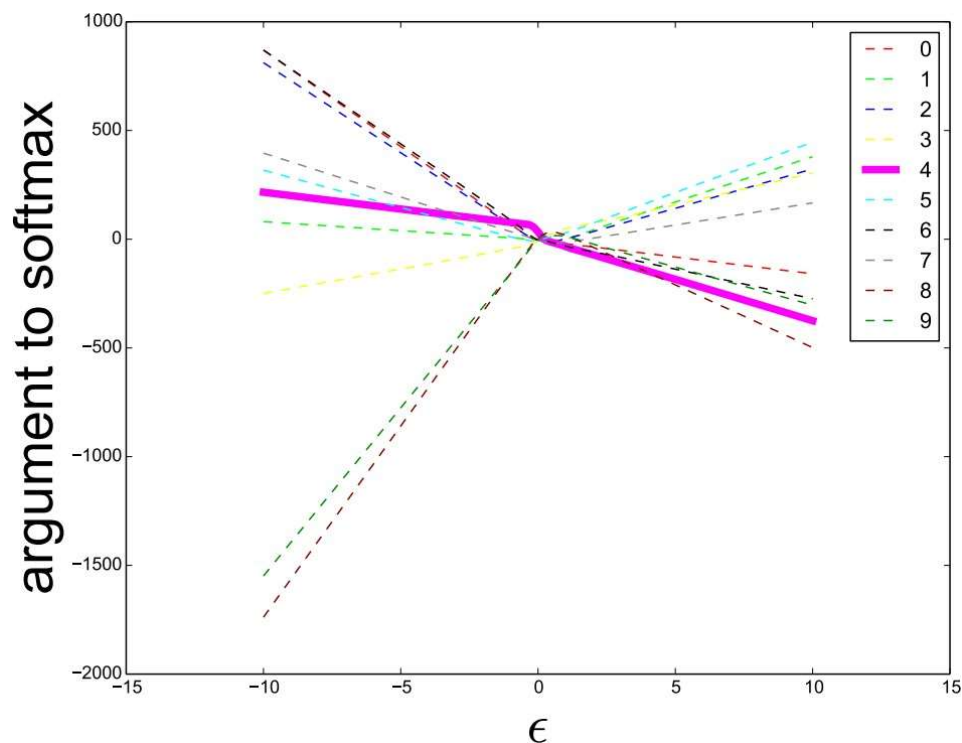
Neural Networks are too linear



Learned Model not powerful enough to fit data properly

Learned Model trained on adversarial examples now more powerful (less linear)

Experimental Linearity of Perturbations



- We have 1 image and 10 classes (0 to 9). The correct classification is 4.
- x-axis is perturbation (chose direction and perturb the image by that direction).
- y-axis are logits (values before calling softmax), unnormalized probabilities.

We see that the function which computes the particular logit, is basically almost (piece wise) linear in the perturbation.

Only around 0 (no perturbation) does the function behave in a non-linear manner and where classification is correct (i.e. 4).

Generating Adversarial Examples

(somewhat possible due to these mostly linear properties)

Targeted vs. Untargeted Attacks

Targeted Attack – aims to misclassify the input (e.g., image) to a **specific label** (e.g. panda to gibbon)

Untargeted Attack – aims to misclassify the input to **any wrong label** (e.g. panda to any other animal)

Formulated as a slightly different optimization problem

Targeted Attack: Problem Statement

Input:

- neural network $f: X \rightarrow \mathcal{C}$
- input $x \in X$
- target label $t \in \mathcal{C}$, such that $f(x) \neq t$

Output:

- A perturbation η such that $f(x + \eta) = t$

Adversarial example

$$x' = x + \eta$$

Untargeted Attack: Problem Statement

Input:

- neural network $f: X \rightarrow \mathcal{C}$
- input $x \in X$

Output:

- A perturbation η such that $f(x + \eta) \neq f(x)$

Adversarial example

$$x' = x + \eta$$

Types of Attacks

White box attacker: the attacker knows the model, the parameters, and the architecture

Black box attacker: the attacker knows the architecture (e.g., the layers) but not its parameters (e.g., weights)

Note: it was found adversarial examples are **transferrable**, hence given the same training data as the original network, an attacker can train their own **mirror network** of the black box original network and then attack the mirror network with white-box techniques. If attack on mirror network succeeds, it will likely succeed on the original.

We will look at **white box attacks first**

Targeted Fast Gradient Sign Method

1. Compute perturbation:

$\eta = \epsilon \cdot \text{sign}(\nabla_x \text{loss}_t(x))$, where

$$\nabla_x \text{loss}_t = \left(\frac{\partial \text{loss}_t}{\partial x_1}, \dots, \frac{\partial \text{loss}_t}{\partial x_n} \right) \quad \text{sign}(g) = \begin{cases} -1, & \text{if } g < 0 \\ 0, & \text{if } g = 0 \\ 1, & \text{if } g > 0 \end{cases}$$

2. Perturb the input:

$$x' = x - \eta$$

3. Check if:

$$f(x') = t$$

- Here, each x_i is a pixel
- ϵ is a very small constant (e.g., 0.007)
- As FGSM is 1-step, x' is **guaranteed** to stay inside the box $[x - \epsilon, x + \epsilon]$, so no need to project.
- t is the target, **bad label**
- loss_t is the loss w.r.t target label
- FGSM was designed to be fast, not optimal (may not compute minimal perturbation)

Untargeted version of FGSM

1. Compute perturbation:

$$\eta = \epsilon \cdot \text{sign}(\nabla_x \text{loss}_s(x))$$

2. Perturb the input:

$$x' = x + \eta$$

3. Check if:

$$f(x') \neq s$$

- With untargeted FGSM, we do not know what the target (bad) label is that we want.
- We just want **some label** different than the correct label s .
- So we try to “get away” from the correct label by maximizing the value of the loss

FGSM



$$\epsilon$$
$$+ .007 \times$$



=



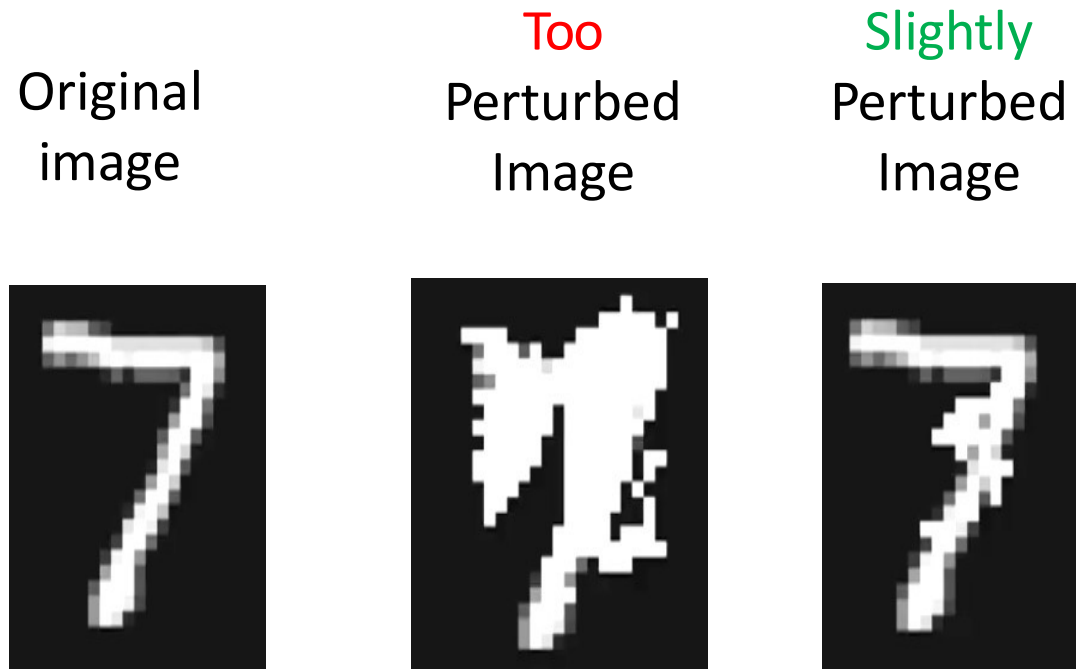
57.7% panda

“noise”

93.3% gibbon

$$\text{sign}(\nabla_x \text{loss}_t(x))$$

Importance of Small Perturbations



We need some notion of distance....

Norm: Notion of Distance

Similarity of $x \sim x'$ is usually captured by an l_p norm:

$$x \sim x' \text{ iff } \|x - x'\|_p < \epsilon,$$

$$\text{where } \|x - x'\|_p = \left((|x_1 - x'_1|)^p + \dots + (|x_n - x'_n|)^p \right)^{\frac{1}{p}}$$

l_0 (when $0^0 = 0$ and we get rid of $1/p$ root) captures the number of changed pixels.

l_2 captures the **Euclidian distance** between x and x' . It can remain small if there are many small changes to many pixels.

l_∞ captures **maximum noise (change)** added to any coordinate. It is the maximum of the absolute values of the entries:

$$\|x - x'\|_\infty = \max(|x_1 - x'_1|, \dots, |x_n - x'_n|)$$

This is the most common norm used for adversarial example generation and it is argued that it most naturally captures human vision.

To derive the max equation, see:

<https://math.stackexchange.com/questions/3099179/proving-the-infinity-norm-is-equal-to-the-maximum-value-of-the-vector&sa=D&source=hangouts&ust=1600861806197000&usg=AFQjCNFLvo4hpAfiNBK06EwAZVFkuRiDNw>

Targeted Attack with Small Changes

Input:

- neural network $f: X \rightarrow \mathcal{C}$
- input $x \in X$
- target label $t \in \mathcal{C}$, such that $f(x) \neq t$

Output:

- A perturbation η such that $f(x + \eta) = t$
- $\|\eta\|_p$ is minimized

Optimization Problem

The problem of generating small perturbations can be phrased as an **optimization problem**:

$$p \in \{0, 2, \infty\}$$

find η
minimize $\|\eta\|_p$
such that $f(x + \eta) = t$
 $x + \eta \in [0, 1]^n$

This is a **hard discrete constraint** which is difficult to optimize for with gradient methods.

Note: η can have negative components.

Key insight: Relaxation of the hard constraint

Optimization Problem

Two steps:

Step 1: Define an objective function obj_t such that:

if $\mathit{obj}_t(\mathbf{x} + \boldsymbol{\eta}) \leq \mathbf{0}$ then $f(\mathbf{x} + \boldsymbol{\eta}) = t$

Step 2: Solve the following optimization problem:

find	$\boldsymbol{\eta}$
minimize	$\ \boldsymbol{\eta}\ _p + c \cdot \mathit{obj}_t(\mathbf{x} + \boldsymbol{\eta})$
such that	$\mathbf{x} + \boldsymbol{\eta} \in [\mathbf{0}, \mathbf{1}]^n$

Optimization Problem

Two steps:

Step 1: Define an objective function \mathbf{obj}_t such that:

$$\text{if } \mathbf{obj}_t(\mathbf{x} + \boldsymbol{\eta}) \leq \mathbf{0} \text{ then } \mathbf{f}(\mathbf{x} + \boldsymbol{\eta}) = \mathbf{t}$$

What are examples of functions for \mathbf{obj} with the property of Step 1?

Choice I:

$$\mathbf{obj}_t(\mathbf{x}') = \text{loss}_t(\mathbf{x}') - 1$$

Lets take cross
entropy loss for loss_t

Choice II:

$$\mathbf{obj}_t(\mathbf{x}') = \max(0, 0.5 - \mathbf{p}_f(\mathbf{x}')_t)$$

$\mathbf{p}_f(\mathbf{x}')_t$ returns the probability of
class \mathbf{t} for input \mathbf{x}' on network \mathbf{f}

Choice I: $\mathit{obj}(\mathbf{x}) = \mathit{loss}_t(\mathbf{x}) - 1$

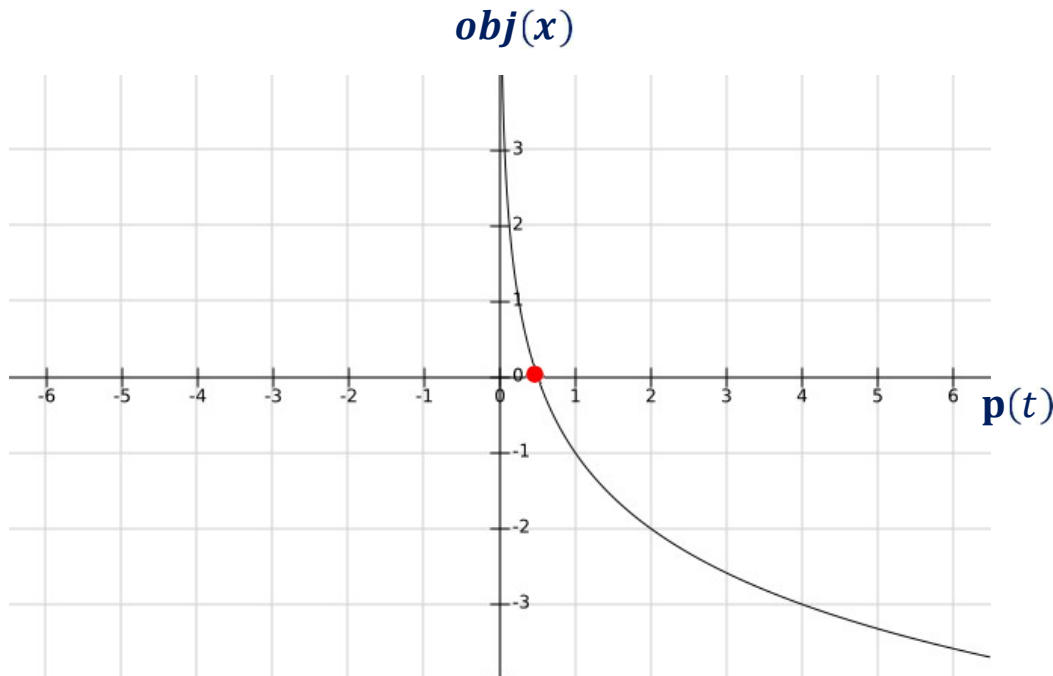
Choice I:

$$\mathit{obj}_t(\mathbf{x}) = \mathit{loss}_t(\mathbf{x}) - 1$$

$$= -\log_2(\mathbf{p}(t)) - 1$$

Plug in cross entropy loss for loss_t with logarithm base 2

Here, we use $\mathbf{p}(t)$ as a shortcut for $\mathbf{p}_f(\mathbf{x})_t$ so to avoid clutter



What we see here is that if the obj_t function is 0 or negative, then the probability $\mathbf{p}(t)$ is ≥ 0.5 (50%).

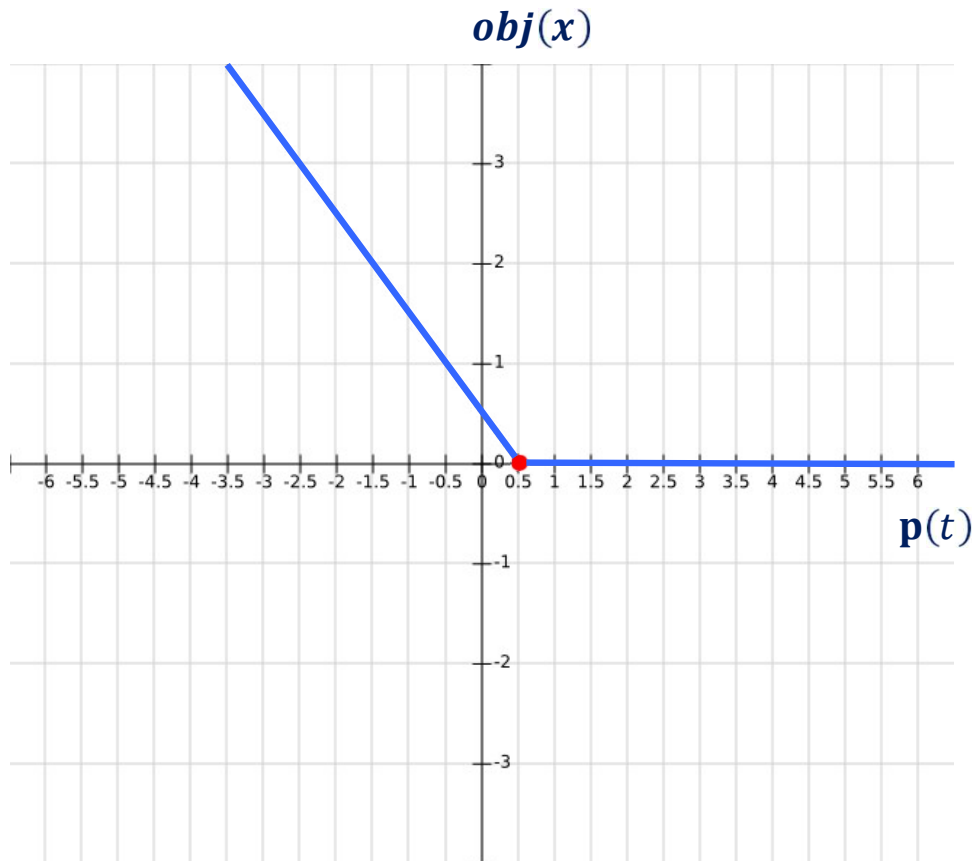
But if $\mathbf{p}(t)$ is ≥ 0.5 for the input \mathbf{x} , then f will return t as a classification for \mathbf{x} because this is the highest probability class. Hence, the desired property of Step 1 holds.

Choice II: $\max(0, 0.5 - \mathbf{p}_f(\mathbf{x})_t)$

Choice II:

$$\mathit{obj}_t(\mathbf{x}) = \max(0, 0.5 - \mathbf{p}(t))$$

← Again we use $\mathbf{p}(t)$ as a shortcut for $\mathbf{p}_f(\mathbf{x})_t$ so to avoid clutter



What we see here is that the obj_t function is always 0 or greater.

It is only 0 when $\mathbf{p}(t)$ is ≥ 0.5 for the input \mathbf{x} .

Again, then f will return t as a classification for \mathbf{x} because this is the highest probability class.

Hence, the desired property holds for Step 1.

Optimization Problem

Two steps:

Step 1: Define an objective function obj_t such that:

if $\mathit{obj}_t(\mathbf{x} + \boldsymbol{\eta}) \leq \mathbf{0}$ then $f(\mathbf{x} + \boldsymbol{\eta}) = t$

Step 2: Solve the following optimization problem:

find $\boldsymbol{\eta}$
minimize $\|\boldsymbol{\eta}\|_p + c \cdot \mathit{obj}_t(\mathbf{x} + \boldsymbol{\eta})$
such that $\mathbf{x} + \boldsymbol{\eta} \in [\mathbf{0}, \mathbf{1}]^n$

Optimization Problem

Two steps:

Step 1: Define an objective function obj_t such that:

if $obj_t(x + \eta) \leq 0$ then $f(x + \eta) = t$

Step 2: Solve the following optimization problem:

find η
minimize $\|\eta\|_\infty + c \cdot obj_t(x + \eta)$
such that $x + \eta \in [0, 1]^n$

This is a problem for
optimization

Lets take a closer look at $\|\boldsymbol{\eta}\|_{\infty}$

$\|\boldsymbol{\eta}\|_{\infty}$ computes the maximum change: it takes the absolute value of every coordinate in $\boldsymbol{\eta}$ and returns the maximum value.

$$\boldsymbol{\eta} = (0.5, 0.49, 0.48) \quad \frac{\partial \|\boldsymbol{\eta}\|_{\infty}}{\partial \eta_1} = 1 \quad \frac{\partial \|\boldsymbol{\eta}\|_{\infty}}{\partial \eta_2} = 0 \quad \frac{\partial \|\boldsymbol{\eta}\|_{\infty}}{\partial \eta_3} = 0$$

After one step we get: $\boldsymbol{\eta} = \boldsymbol{\eta} - \gamma \cdot (1,0,0) = \boldsymbol{\eta} - 0.03 \cdot (1,0,0) = (0.47, 0.49, 0.48)$

After two steps we get: $\boldsymbol{\eta} = \boldsymbol{\eta} - \gamma \cdot (0,1,0) = \boldsymbol{\eta} - 0.03 \cdot (0,1,0) = (0.47, 0.46, 0.48)$

After three steps we get: $\boldsymbol{\eta} = \boldsymbol{\eta} - \gamma \cdot (0,0,1) = \boldsymbol{\eta} - 0.03 \cdot (0,0,1) = (0.47, 0.46, 0.45)$

What we see is that because the gradient is **0 at all non-max locations**, the gradient does not impose a penalty on the optimizer increasing a little bit those locations (due to the $obj_t(\boldsymbol{x} + \boldsymbol{\eta})$ term in the optimization). Also, **only one entry is changed at a time**.

Lets take a closer look at $\|\eta\|_{\infty}$

Going back to the full optimization problem which also includes $obj(x + \eta)$

After one step we get: $\eta = \eta - \gamma \cdot (1,0,0) = \eta - 0.03 \cdot (1,0,0) = (0.47, 0.49, 0.48)$

Now the optimizer can slightly bump up the second location:

After one full step of optimizer, it may also bump up the 2nd location: $(0.47, 0.5, 0.48)$

After second full step of optimizer, we may get: $(0.5, 0.47, 0.48)$

Well, we are just **oscillating** now and bouncing around...turns out SGD may not be a good way to optimize $\|\eta\|_{\infty}$ especially with other terms

One approach to solving the issue

Replace $\|\boldsymbol{\eta}\|_{\infty}$ with other proxy functions that reflect the distance

One idea is to penalize large values in $\boldsymbol{\eta}$ via a term ϵ :

$$\text{Replace } \|\boldsymbol{\eta}\|_{\infty} \text{ with } \sum_i \max(0, (|\boldsymbol{\eta}_i| - \epsilon))$$

- is basically intended to capture the $\|\boldsymbol{\eta}\|_{\infty}$ bound when optimization finishes.
- will be continuously minimized. Initially, ϵ starts at 1
- is decreased with some factor (say 0.9) at every iteration if all $|\boldsymbol{\eta}_i|$ are less than ϵ (then, entire expression will be 0).

Note: an iteration consist of K small steps.

Note: when ϵ is large, gradient of $\sum_i \max(0, (|\boldsymbol{\eta}_i| - \epsilon))$ is similar to gradient of $\|\boldsymbol{\eta}\|_{\infty}$.

Example & Notes on Optimization

Let $L(\boldsymbol{\eta}) = \sigma_i \max(0, (|\boldsymbol{\eta}_i| - \cdot))$ and $\boldsymbol{\eta} = (0.47, 0.49, 0.48)$

Start with $\cdot = 1$, then $L(\boldsymbol{\eta}) = 0$

Next iteration: $\cdot = 0.9$, then $L(\boldsymbol{\eta}) = 0$

Next iteration: $\cdot = 0.81$, then $L(\boldsymbol{\eta}) = 0$

...

At some iteration: $\cdot = 0.478$, now via one step, we get:

$$L(\boldsymbol{\eta}) = \sigma_i \max(0, (|\boldsymbol{\eta}_i| - 0.478)) = \sigma_i (0, 0.012, 0.002) = \mathbf{0.014}$$

$$\square_{\boldsymbol{\eta}} L(\boldsymbol{\eta}) = (\mathbf{0}, \mathbf{1}, \mathbf{1})$$

We can then update $\boldsymbol{\eta}$ as usual, complete this step, and continue with the next step

Notes on optimization:

- There are K steps within an iteration, each updating $\boldsymbol{\eta}$.
- Entire optimization stops if after K steps $L(\boldsymbol{\eta}) \neq 0$. Otherwise, if $L(\boldsymbol{\eta}) = 0$, optimization continues with a new $\cdot = 0.9 * \text{previous } \cdot$.
- Entire optimization stops if it also reaches some pre-defined value of \cdot ($1/256$ for Carlini & Wagner).
- When optimization stops, we return $\boldsymbol{\eta}$ at the iteration before the last one. This means $\|\boldsymbol{\eta}\|_{\infty} \leq \cdot$ where \cdot is the one used at iteration before last.
- If $\boldsymbol{\eta}_{top2} < r < \boldsymbol{\eta}_{top1}$ where $\boldsymbol{\eta}_{top1}$ is the largest element and $\boldsymbol{\eta}_{top2}$ is second largest element in $\boldsymbol{\eta}$, then the gradient $\square_{\boldsymbol{\eta}} L(\boldsymbol{\eta})$ will be the same as $\square_{\boldsymbol{\eta}} \|\boldsymbol{\eta}\|_{\infty}$.

Here we only show the optimization of one term, namely $L(\boldsymbol{\eta})$, to illustrate the relationship with optimizing $\|\boldsymbol{\eta}\|_{\infty}$.

Summary: Optimization Problem

Two steps:

Step 1: Define an objective function obj_t such that:

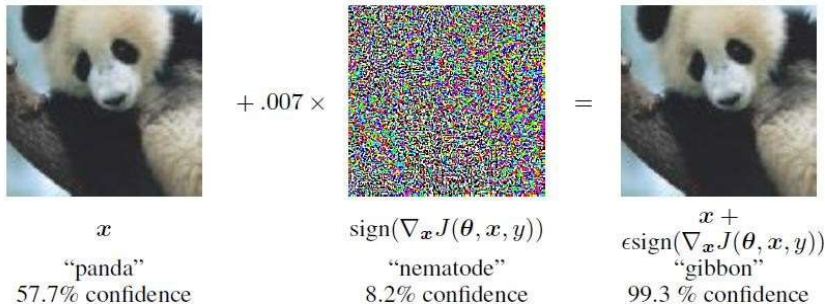
$$\text{if } \mathit{obj}_t(\mathbf{x} + \boldsymbol{\eta}) \leq \mathbf{0} \text{ then } \mathbf{f}(\mathbf{x} + \boldsymbol{\eta}) = t$$

Step 2: Solve the following optimization problem:

$$\begin{array}{ll} \text{find} & \boldsymbol{\eta} \\ \text{minimize} & \|\boldsymbol{\eta}\|_{\infty} + c \cdot \mathit{obj}_t(\mathbf{x} + \boldsymbol{\eta}) \\ \text{such that} & \mathbf{x} + \boldsymbol{\eta} \in [\mathbf{0}, \mathbf{1}]^n \end{array}$$

Lecture Summary

Deep Learning is susceptible to adversarial examples in various domains



Generating Adversarial examples
(basically, an optimization problem)

- FGSM: targeted and untargeted
- Small perturbation attacks
- Need suitable optimization problem

Next lecture: Dealing with Constraints and Adversarial defenses