

Reliable and Interpretable Artificial Intelligence

Lecture 3: Adversarial Attacks II

Martin Vechev
ETH Zurich

Fall 2020

Recall: Our Optimization Problem

Two steps:

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

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

Step 2: Solve the following optimization problem:

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

Hard Box Constraint

Dealing with Constraints

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

Given \mathbf{x} is constant, this is the same as enforcing $\eta_i \in [-x_i, 1 - x_i]$ for every η_i . We can then use either of these two methods:

Projected gradient descent (PGD)

“Fit” all coordinates to be within the box

$$\text{project}((\eta_1, \dots, \eta_n)) = (\text{clip}_1(\eta_1), \dots, \text{clip}_n(\eta_n))$$

$$\text{clip}_i(\eta_i) = \begin{cases} -x_i & \text{if } \eta_i < -x_i \\ \eta_i & \text{if } \eta_i \in [-x_i, 1 - x_i] \\ 1 - x_i & \text{if } \eta_i > 1 - x_i \end{cases}$$

LBFGS-B optimizer:

Used by Carlini & Wagner

pass each $\eta_i \in [-x_i, 1 - x_i]$

separately to the optimizer.





































































































“-B” stands for box constraints

Note: if we also want $\|\boldsymbol{\eta}\|_\infty < e$ then we can also add the box constraints $\eta_i \in [-e, e]$

With this approach we get

Target label

Initial label

	0	1	2	3	4	5	6	7	8	9
0										
1										
2										
3										
4										
5										
6										
7										
8										
9										

What we see is that on the MNIST (digit recognition) data set **it is not difficult** to get a realistic looking image that fools the neural network classifier...

DeepSpeech Attack: more technically

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

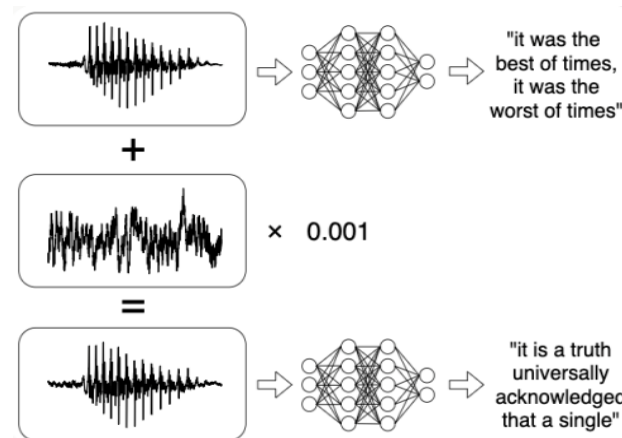
Instantiate
to audio

easier to optimize
(smooth function)

audio specific
loss

find η
 minimize $\|\eta\|_2^2 + c \cdot \text{loss}_t(x + \eta)$
 such that $\eta' \in [0, k]^n$
 where $\eta' = 20 * \log_{10}(\eta)$

change to decibels



find η
 minimize $\|\eta\|_2^2 + c \cdot \text{loss}_t(x + \eta)$
 such that $\eta \in [1, 10^{k/20}]^n$

A **re-write** so we work with η only.
 Now we can project on η as usual.

Another attack...often used during training

- So far, we looked at FGSM as well as an attack to minimize the distance to the original input (e.g., image, audio)
- Now, we illustrate another attack, a variant of FGSM applied iteratively **with projection**.
- The attack uses Projected Gradient Descent (PGD) and is referred to as a **PGD attack**.
- This is a commonly used attack for **adversarial training: training the network to be robust**.

Illustrating the PGD attack

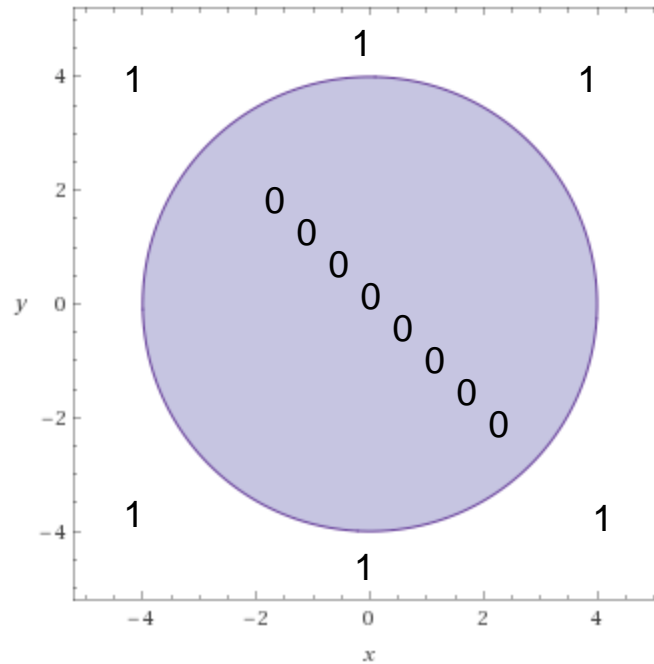
Given a **dataset** of points (x, y) where label is:

0 if $x^2+y^2 < 16$

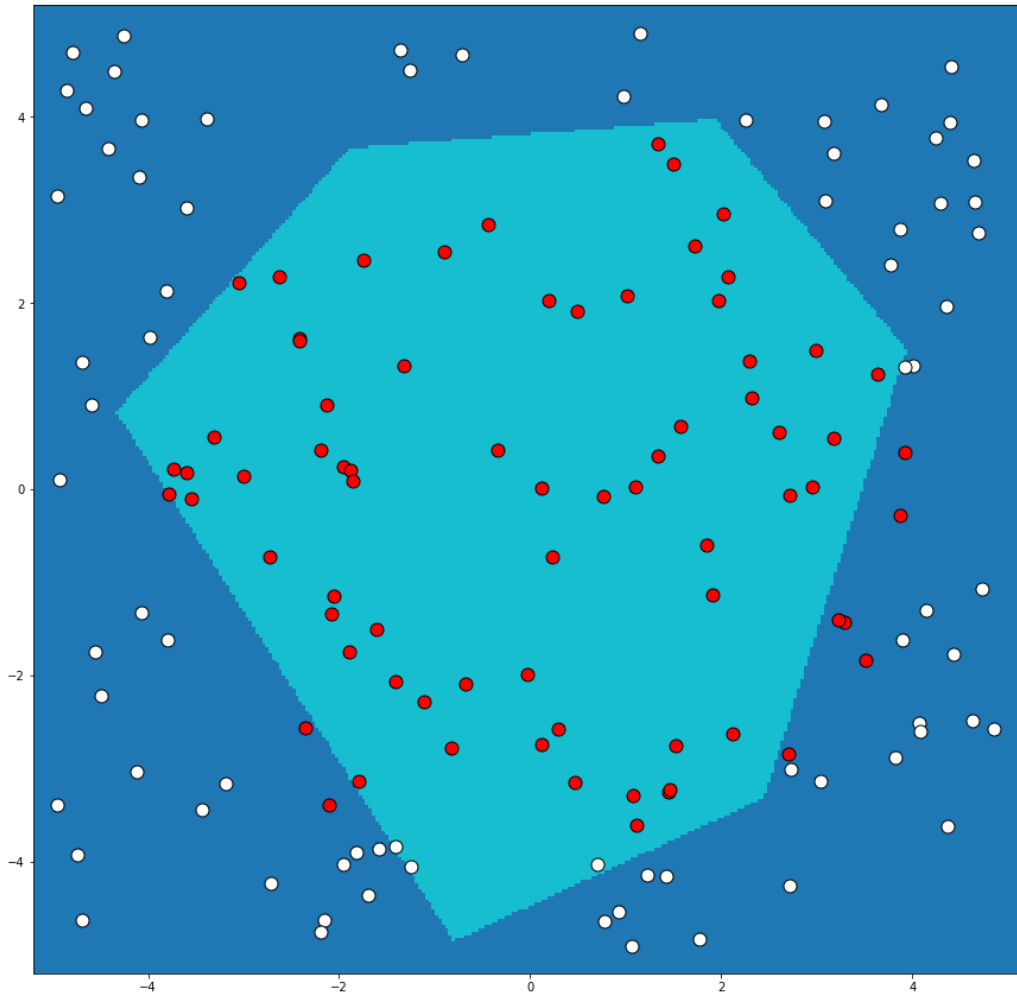
1 otherwise

train a neural network to classify the points correctly

Illustrating the PGD attack



After training we get the classifier:



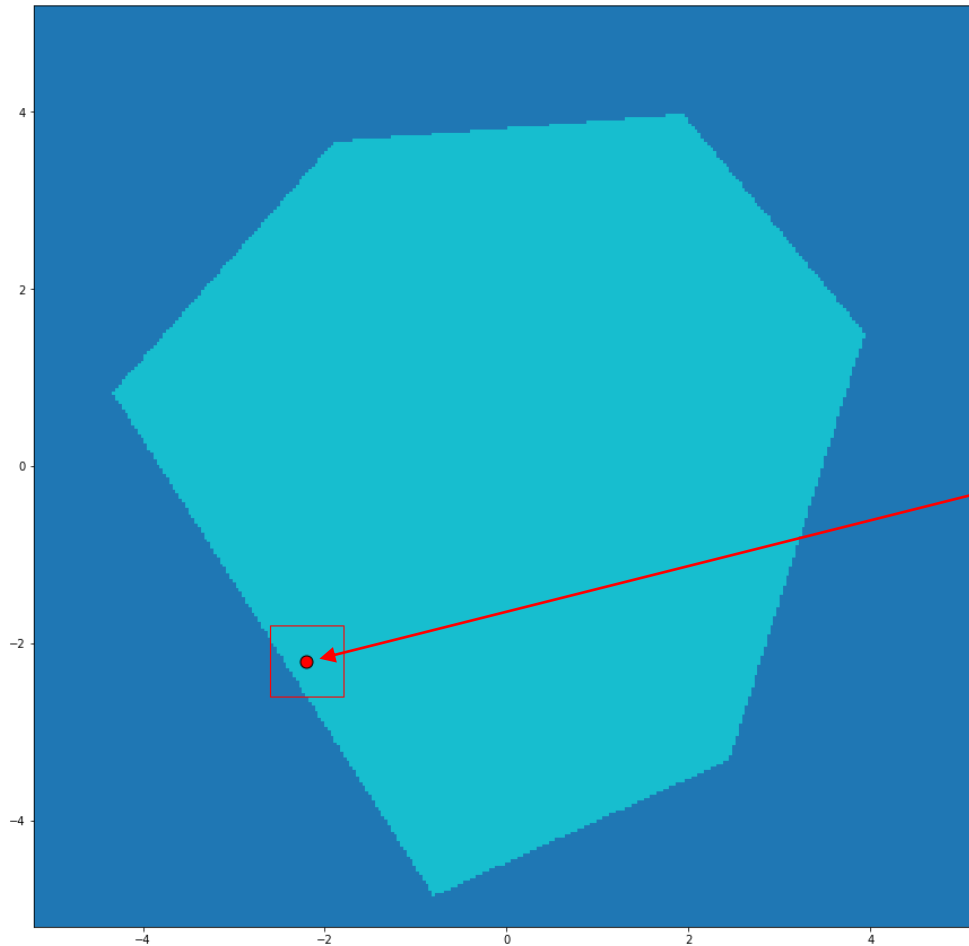
Dark blue – neural network predicts 1 (property does not hold)

Light blue – neural network predicts 0 (property holds)

Red dots – those where property actually holds

White dots – those where property actually does not hold

Lets pick a point...



Goal:

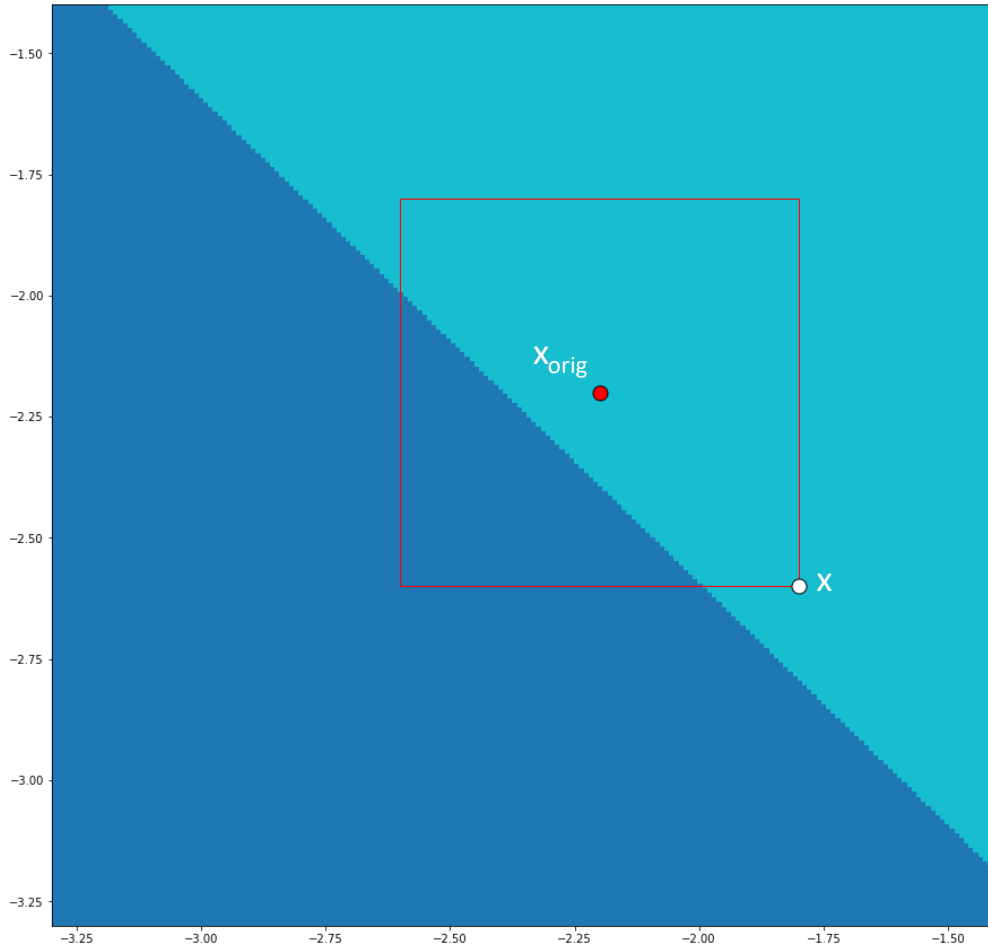
Find adversarial input in

L_{inf} ball around:

$x_{\text{orig}} = (-2.2, -2.2)$
(red point)

with $\epsilon=0.4$

Lets Zoom in a bit...

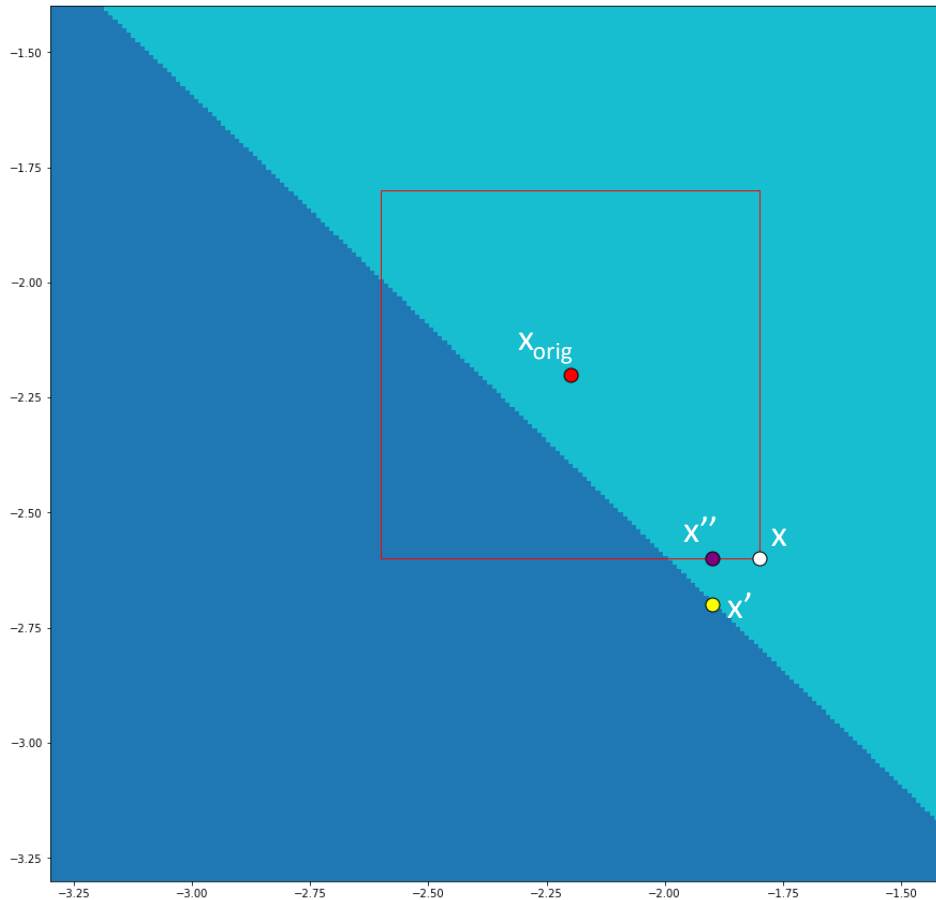


Initialize PGD with:

$$x = (-1.8, -2.6)$$

Note: this is just for the example to illustrate projection. In practice, one picks a point at random in the box

PGD Iteration 1



$$\text{NN}(x) = [0.5973, 0.4027]$$

$$\text{Loss}(x) = 0.5153$$

$$\nabla_x \text{Loss}(x) = [-0.852, -1.373]$$

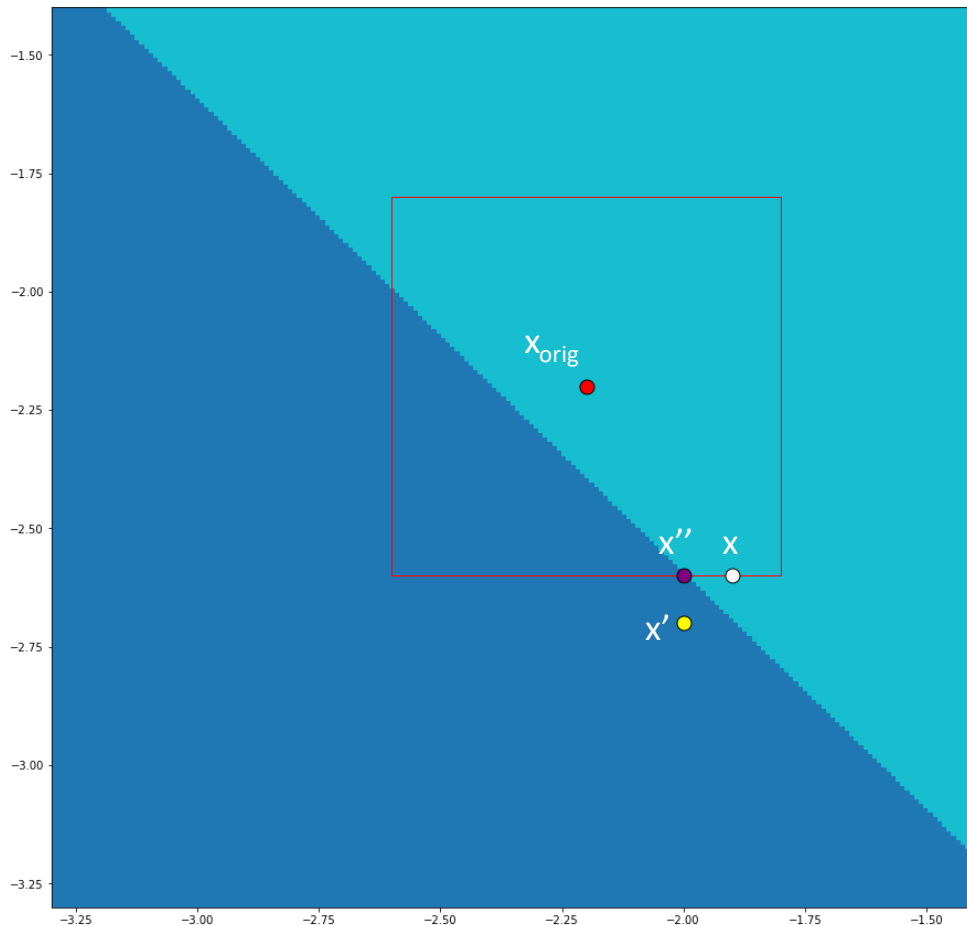
$$x' = x + \overbrace{0.1 * \text{sign}(\nabla_x \text{Loss}(x))}^{\text{Change } \Delta} \\ = [-1.9, -2.7] \text{ (yellow point)}$$

Up-to-here, its just standard untargeted FGSM attack but with **smaller step-size** of 0.1 than ϵ which is 0.4.

But now we also **project**:

$$x'' = \text{project}(x', x_{\text{orig}}, \epsilon) \\ = [-1.9, -2.6] \text{ (purple point)}$$

PGD Iteration 2



x'' from before now named x :

$\text{NN}(x) = [0.5455, 0.4545]$
(so point $x = (-1.9, -2.6)$ is
not yet a counter example)

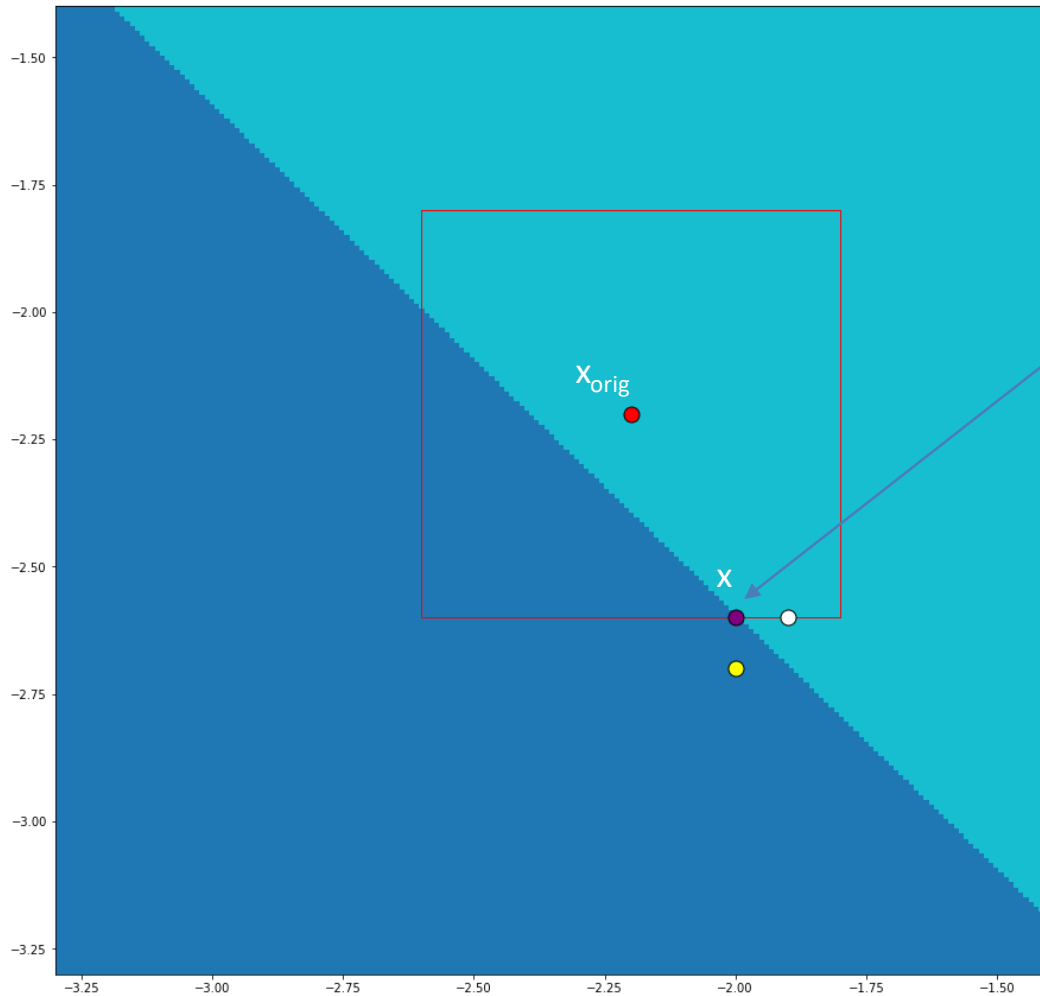
$\text{Loss}(x) = 0.6060$

$\nabla_x \text{Loss}(x) = [-0.9621, -1.5493]$

$$x' = x + \overbrace{0.1 * \text{sign}(\nabla_x \text{Loss}(x))}^{\text{Change } \Delta}$$
$$= [-2, -2.7]$$

$$x'' = \text{project}(x', x_{\text{orig}}, \epsilon)$$
$$= [-2, -2.6]$$

PGD Iteration 3

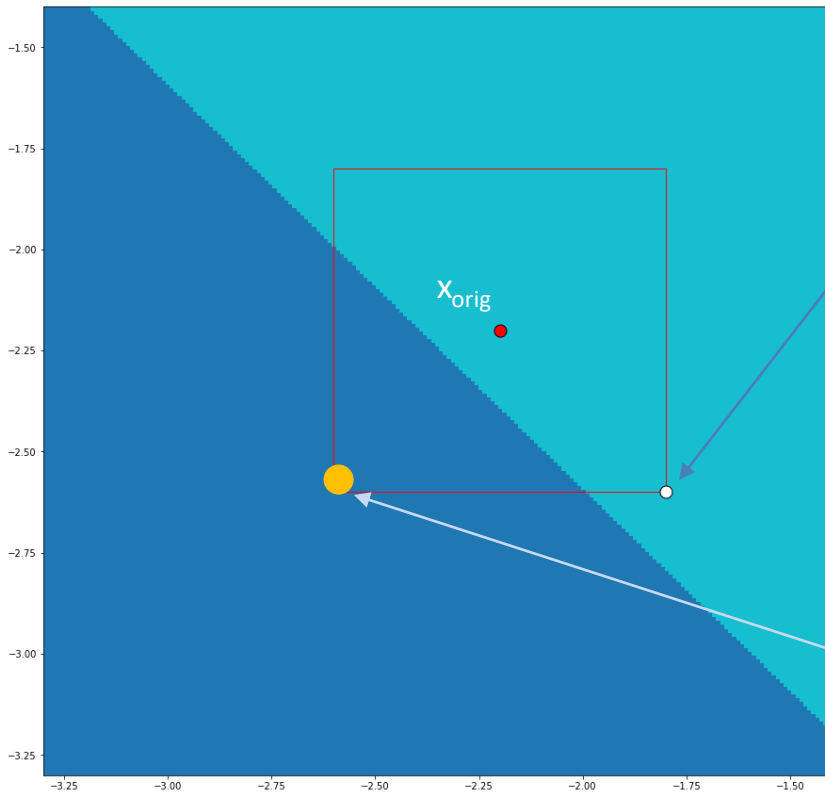


$\text{NN}(x) = [0.4927, 0.5073]$

found adversarial example
 $x = [-2, -2.6]$

Neural network predicts 1,
although $(-2)^2 + (-2.6)^2 < 16$
so it should have been
classified as 0

Some notes on PGD



- The goal of the PGD attack is to **find a point in the region which maximizes the loss** (it may still classify to the same label as x_{orig})
- For our example, we started at the corner. Typically one starts the search with a **random point inside the box**.
- One stops PGD after a **pre-defined number of iterations** (e.g., 10).
- In our example, we always stepped outside the box to illustrate projection, and then projected to the box. It is possible to never step outside the box and thus **projection will have no effect**.
- It is possible the final produced example **is inside the box**, and not on the boundary. However, when we project, if outside the box, we will end up on the boundary.
- In this example, loss is **likely to be highest** somewhere around the big orange point (typically far from the decision boundary). Of course, when we are searching, we **don't know the actual decision boundary**.
- One can implement PGD in **two ways**:
 - **a)** by projecting current **point** x' to the ϵ -box around x_{orig} as well as $[0,1]$ for each dimension, or
 - **b)** by projecting the **change** Δ to $[-\epsilon, +\epsilon]$ as well as to the constraints needed so each element in the resulting point is between $[0,1]$ (see slide 3 in this lecture)
- Step size (in our example 0.1) is **typically smaller than ϵ** (in FGSM it is ϵ).

- Projection is linear-time in the dimension for L_∞ and L_2 norms.
- An **open problem**: finding efficient projections for various convex regions that are **more expressive than boxes** (e.g., convex polyhedral restrictions).

Another Attack Example: Diffing Networks

Finding a **differencing input**:

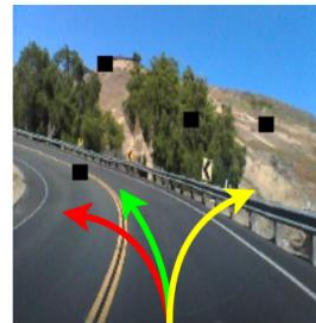
Given two neural networks f_1 and f_2 trained to learn the same function $f^*: X \rightarrow C$, find an input $x \in X$ such that $f_1(x) \neq f_2(x)$



DRV_C1: right



DRV_C2: right



DRV_C3: right

Another Attack Example: Diffing Networks

Define the following objective (good if we want $f_1(x)$ to classify x to t):

$$obj_t(x) = f_1(x)_t - f_2(x)_t$$

$f_i(x)_t$ returns the probability that f_i predicts x to be t


We can use absolute value loss if we just want to get a different classification by both (need not be t).

Select input $x \in X$ which classifies as t with both networks

while $class(f_1(x)) = class(f_2(x))$:

$$x = x + \epsilon \cdot \frac{\partial obj_t(x)}{\partial x}$$

return x



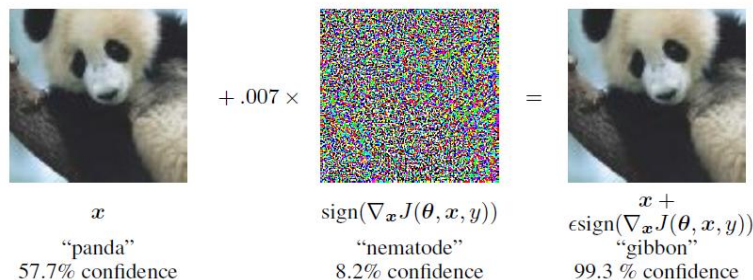
Maximize loss: make f_1
more confident about
 t while making f_2 less
confident about t

Summary of adversarial attacks

Attack Type	Region	Optimization	Outcome
FGSM (targeted, untargeted)	Change η fixed to $[-\epsilon, +\epsilon]$.	Take exactly one ϵ -sized step	Produced example will be on boundary of region.
PGD (typically untargeted, but can be targeted)	Can be instantiated with any region one can project to.	Take many steps. Uses projection to stay inside region. For special case of l_∞ , step size smaller than ϵ .	Result will be inside region. Tries to maximize loss.
C&W [Images] (presented as targeted)	No real restriction, except image has to be in $[0,1]$ (like all other methods). This restricts the region for the change η : η has to be bounded s.t. original image + η stays in $[0,1]$.	Aims to produce a change η with small l_∞ . Takes many steps, using LBFGS-B to ensure η stays in bounds.	Result will be inside $[0,1]$, with a hopefully small l_∞ distance from original image.
C&W [Audio] (presented as targeted)	A fixed region for η .	Aims to produce a change η with small squared l_2 . Takes many steps, using LBFGS-B to ensure η stays in bounds.	Result will be inside the fixed η region, with a hopefully small l_2 distance from original sound wave.
Diffing Networks (targeted)	Can work with a fixed region around change η .	Aims to produce a change η where one neural networks classifies image as t while the other as not t .	Ideally, a result where two networks disagree on their classification.

Lecture Summary

Deep Learning is susceptible to adversarial examples



Generating Adversarial examples (an optimization problem)

- FGSM
- C&W (minimize perturbation)
- PGD
- Diffing

An example of the PGD attack

