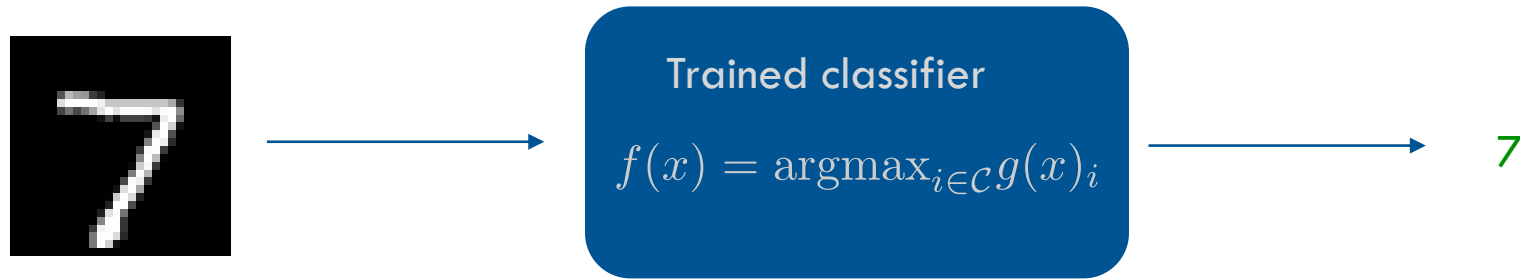# Generating Adversarial Examples using Gradient Descent

Arjun Bhagoji

University of Chicago

# Standard Classification



Trained classifier

$$f(x) = \mathrm{argmax}_{i \in \mathcal{C}} g(x)_i$$
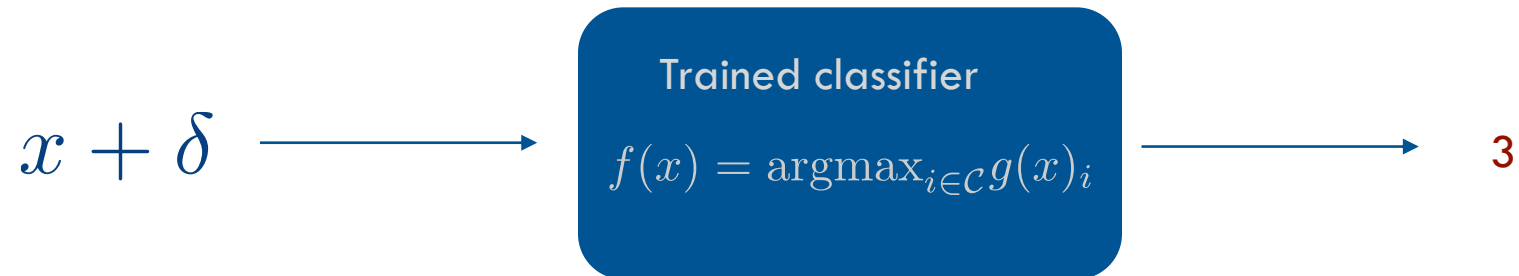
7

- Classifier chooses class with highest 'score'

- Examples of possible classifiers:
        - Logistic Regression
        - Support Vector Machines
        - Neural Networks

- **Q**: Can we modify this input so the trained classifier provides the wrong output?
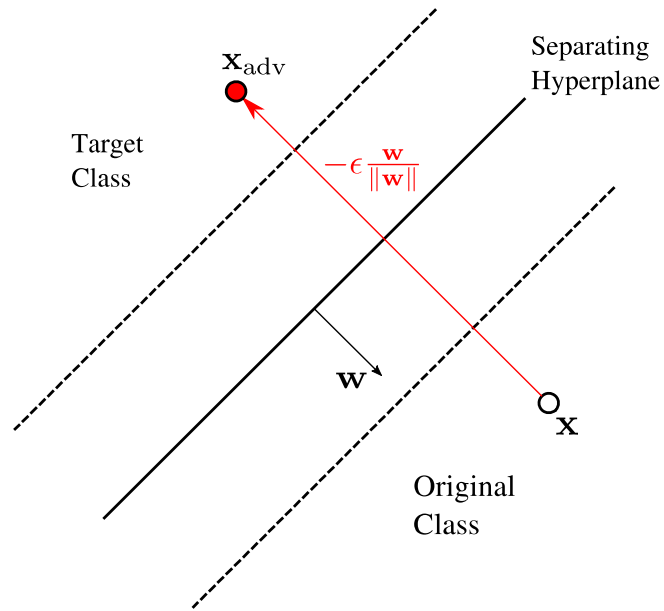
# Adversarial Examples

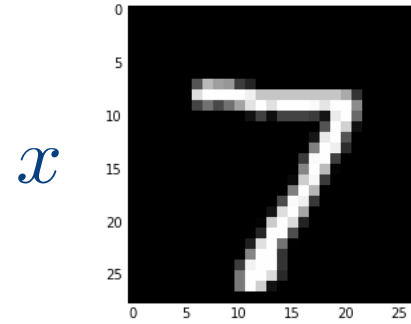Find the smallest perturbation $\delta$ such that:

$$f(x + \delta) \neq f(x)$$

$$x + \delta \longrightarrow$$

Trained classifier

$$f(x) = \text{argmax}_{i \in \mathcal{C}} g(x)_i$$

$$\longrightarrow 3$$

# Adv. Examples for Linear Classifiers



$$x_{\mathrm{adv}} = x - \epsilon \frac{\mathbf{w}}{\|\mathbf{w}\|_2}$$

$$\epsilon \in [0, \infty)$$

Attack on two-class Linear Classifier

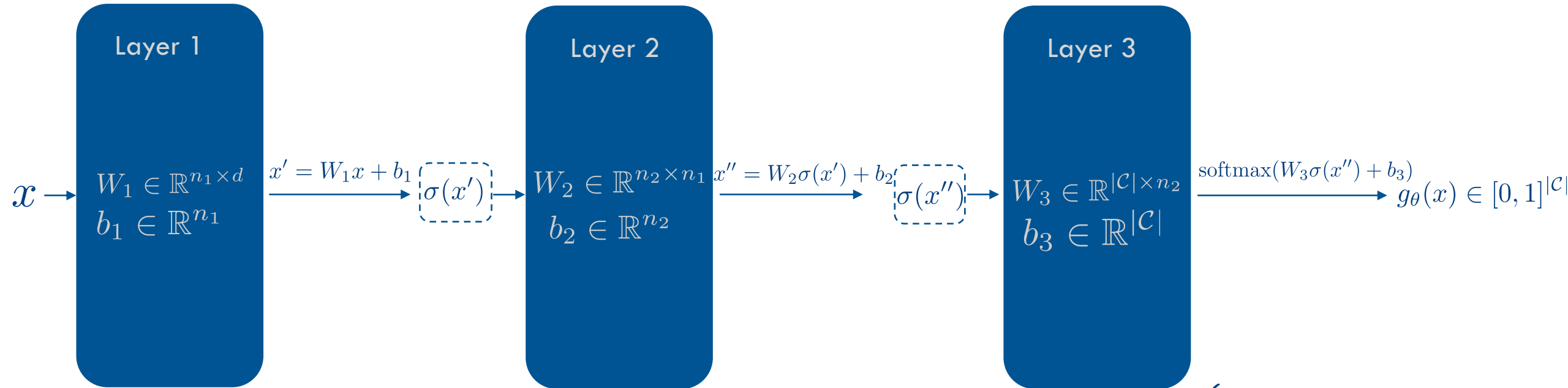$x$    Classified as 7

$x_{\mathrm{adv}}$    Classified as 3!

Adversarial image with $\epsilon$ =2.0.

Leads to 100% misclassification on test set.

$\epsilon$ is the *amount of perturbation* or *adversarial budget*

# Overview: Neural Networks



$x \longrightarrow$

**Layer 1**

$W_1 \in \mathbb{R}^{n_1 \times d}$
$b_1 \in \mathbb{R}^{n_1}$

$\xrightarrow{x' = W_1 x + b_1}$ $\sigma(x')$ $\rightarrow$

**Layer 2**

$W_2 \in \mathbb{R}^{n_2 \times n_1}$
$b_2 \in \mathbb{R}^{n_2}$

$\xrightarrow{x'' = W_2 \sigma(x') + b_2}$ $\sigma(x'')$ $\rightarrow$

**Layer 3**

$W_3 \in \mathbb{R}^{|\mathcal{C}| \times n_2}$
$b_3 \in \mathbb{R}^{|\mathcal{C}|}$

$\xrightarrow{\text{softmax}(W_3 \sigma(x'') + b_3)}$ $g_\theta(x) \in [0,1]^{|\mathcal{C}|}$

$\sigma(\cdot)$ is an element-wise non-linear function, e.g. Rectified Linear Unit

$$\sigma(a) = \begin{cases} 0 \; , a < 0, \\ a \; , a \geq 0 \end{cases}$$

$$\text{softmax}(z) = \frac{e^{z_i}}{\sum_i e^{z_i}}$$

$$J(g_{\{W_1, b_1, W_2, b_2, W_3, b_3\}}) = J(g_\theta) = \text{loss}(g_\theta(x), y)$$

# Remember: Stochastic Gradient Descent

Each update step to the **parameters** to minimize loss is:

$$\theta_t = \theta_{t-1} - \alpha \nabla_\theta J(g_\theta)$$

$$= \theta_{t-1} - \alpha \nabla_\theta \text{loss}(g_\theta(x), y))$$

Q: Can we use the same idea to maximize loss over an input?

# Single Gradient Ascent Step

**Untargeted Adversarial Example:**

$$x_{\mathrm{adv}} = x + \epsilon \frac{\nabla_x \mathrm{loss}(g_\theta(x), y)}{\|\nabla_x \mathrm{loss}(g_\theta(x), y)\|}$$

**Intuition**: Move the point in the direction of the gradient, this locally increases the loss the most

**Q**: How do we move the example such that it is classified in some specific target class $T$ ?

**Targeted Adversarial Example:**

$$x_{\mathrm{adv}} = x - \epsilon \frac{\nabla_x \mathrm{loss}(g_\theta(x), T)}{\|\nabla_x \mathrm{loss}(g_\theta(x), T)\|}$$

# Multi-step Projected Gradient Descent

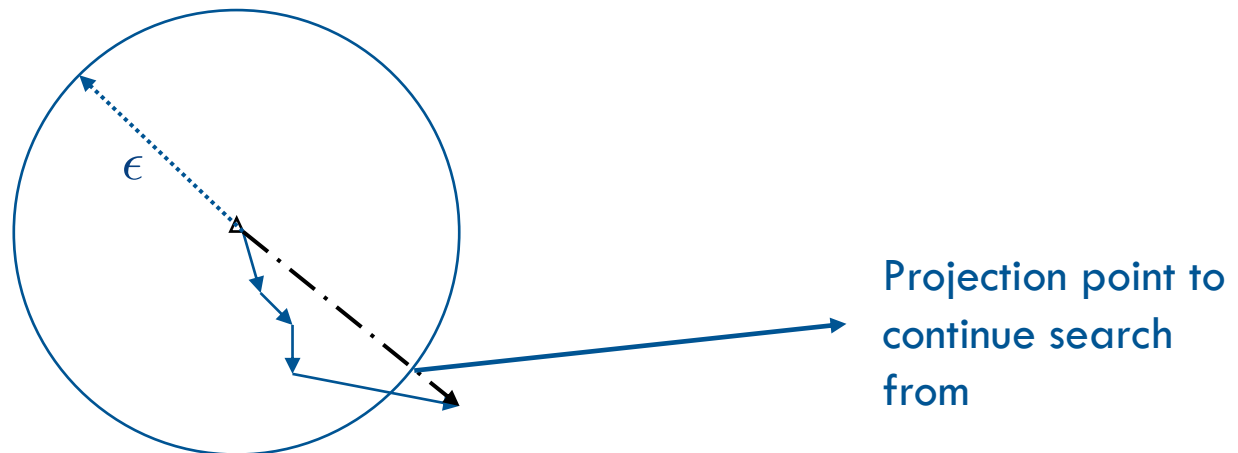**Multi-step, untargeted adversarial example:**

$$x_{\mathrm{adv}}^t = P_{x,\epsilon}\left(x_{\mathrm{adv}}^{t-1} + \alpha\frac{\nabla_{x_{\mathrm{adv}}^{t-1}}\mathrm{loss}(g_\theta(x_{\mathrm{adv}}^{t-1}), y)}{\|\nabla_{x_{\mathrm{adv}}^{t-1}}\mathrm{loss}(g_\theta(x_{\mathrm{adv}}^{t-1}), y)\|}\right)$$

**Multi-step, targeted adversarial example:**

$$x_{\mathrm{adv}}^t = P_{x,\epsilon}\left(x_{\mathrm{adv}}^{t-1} - \alpha\frac{\nabla_{x_{\mathrm{adv}}^{t-1}}\mathrm{loss}(g_\theta(x_{\mathrm{adv}}^{t-1}), T)}{\|\nabla_{x_{\mathrm{adv}}^{t-1}}\mathrm{loss}(g_\theta(x_{\mathrm{adv}}^{t-1}), T)\|}\right)$$

$P_{x,\epsilon}$ projects the adversarial example back onto an $\epsilon$- ball around the original input

**Intuition**: Search for empirically best perturbation using small steps



$\epsilon$

Projection point to continue search from

# Advanced: Robust Training

For e in [1,E]

    for samples i in [1,N]

        1. Compute adversarial example $x_{i,\mathrm{adv}}^T$ for current state $\theta_{t-1}$

        2. Compute adversarial loss $\mathrm{loss}(g_{\theta_{t-1}}(x_{i,\mathrm{adv}}^T), y)$

        3. Compute gradient $\nabla_\theta \mathrm{loss}(g_{\theta_{t-1}}(x_{i,\mathrm{adv}}^T), y)$

        4. Update gradient $\theta_t = \theta_{t-1} - \alpha \nabla_\theta \mathrm{loss}(g_{\theta_{t-1}}(x_{i,\mathrm{adv}}^T), y)$