

ECE595 / STAT598: Machine Learning I

Lecture 35 Max-Loss Attacks and Regularized Attacks

Spring 2020

Stanley Chan

School of Electrical and Computer Engineering
Purdue University



Agenda

- Last lecture we have seen min-distance attack
- In linear case, there is a very simple geometry
- Today we are going to consider two of its variations
 - Max-loss attack
 - Regularized attack
- We will again talk about their geometry using **linear** models.
- And then we will link the results to deep models.
- You will see that some of the most popular deep attack models out there are based on one of the three formulations we discuss here

Outline

- Lecture 33 Overview
- Lecture 34 Min-distance attack
- **Lecture 35 Max-loss attack and regularized attack**

Today's Lecture

- **Max-loss attack**
 - **Linear models**
 - **Deep models: FGSM and PGD**
- Regularized attack
 - Linear models
 - CW attack

Maximum Loss Attack

Definition (Maximum Loss Attack)

The **maximum loss attack** finds a perturbed data \mathbf{x} by solving the optimization

$$\begin{aligned} & \underset{\mathbf{x}}{\text{maximize}} && g_t(\mathbf{x}) - \max_{j \neq t} \{g_j(\mathbf{x})\} \\ & \text{subject to} && \|\mathbf{x} - \mathbf{x}_0\| \leq \eta, \end{aligned} \tag{1}$$

where $\|\cdot\|$ can be any norm specified by the user, and $\eta > 0$ denotes the attack strength.

- I want to bound my attack $\|\mathbf{x} - \mathbf{x}_0\| \leq \eta$
- I want to make $g_t(\mathbf{x})$ as big as possible
- So I want to maximize $g_t(\mathbf{x}) - \max_{j \neq t} \{g_j(\mathbf{x})\}$
- This is equivalent to

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \\ & \text{subject to} && \|\mathbf{x} - \mathbf{x}_0\| \leq \eta, \end{aligned}$$

If you restrict yourself to two classes only ...

- The problem is

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \\ & \text{subject to} && \|\mathbf{x} - \mathbf{x}_0\| \leq \eta, \end{aligned}$$

- η is the maximum loss attack strength
- Want $g_t(\mathbf{x})$ to override $\max_{j \neq t} \{g_j(\mathbf{x})\}$
- So maximize $g_t(\mathbf{x})$
- If you restrict to linear, and only two classes, then

$$\underset{\mathbf{x}}{\text{minimize}} \quad \mathbf{w}^T \mathbf{x} + w_0 \quad \text{subject to} \quad \|\mathbf{x} - \mathbf{x}_0\| \leq \eta.$$

- Solvable in closed-form.

Max-Loss Attack using ℓ_2 -norm

- The problem is

$$\underset{\mathbf{r}}{\text{minimize}} \quad \mathbf{w}^T \mathbf{r} + b_0 \quad \text{subject to} \quad \|\mathbf{r}\|_2 \leq \eta.$$

- Cauchy inequality:

$$\mathbf{w}^T \mathbf{r} \geq -\|\mathbf{w}\|_2 \|\mathbf{r}\|_2 \geq -\eta \|\mathbf{w}\|_2.$$

- Claim: Lower bound of $\mathbf{w}^T \mathbf{r}$ is attained when $\mathbf{r} = -\eta \mathbf{w} / \|\mathbf{w}\|_2$:

$$\begin{aligned} \mathbf{w}^T \mathbf{r} &= \mathbf{w}^T \left(-\frac{\eta \mathbf{w}}{\|\mathbf{w}\|_2} \right) \\ &= -\eta \|\mathbf{w}\|_2. \end{aligned}$$

- So the solution is $\mathbf{r} = -\eta \mathbf{w} / \|\mathbf{w}\|_2$.

Max-Loss Attack using ℓ_∞ -norm

- Goal: Want to solve

$$\underset{\mathbf{x}}{\text{minimize}} \quad \mathbf{w}^T \mathbf{x} + w_0 \quad \text{subject to} \quad \|\mathbf{x} - \mathbf{x}_0\| \leq \eta.$$

- Define $\mathbf{x} = \mathbf{x}_0 + \mathbf{r}$. Then

$$\begin{aligned} \mathbf{w}^T \mathbf{x} + w_0 &= \mathbf{w}^T (\mathbf{x}_0 + \mathbf{r}) + w_0 \\ &= \mathbf{w}^T \mathbf{x}_0 + \mathbf{w}^T \mathbf{r} + w_0 \\ &= \mathbf{w}^T \mathbf{r} + \underbrace{\mathbf{w}^T \mathbf{x}_0 + w_0}_{=b_0} \end{aligned}$$

- Define $b_0 = (\mathbf{w}^T \mathbf{x}_0 + w_0)$. The optimization can be rewritten as

$$\underset{\mathbf{r}}{\text{minimize}} \quad \mathbf{w}^T \mathbf{r} + b_0 \quad \text{subject to} \quad \|\mathbf{r}\|_\infty \leq \eta.$$

Solution to Max-Loss Attack (ℓ_∞ -norm)

- Holder's inequality (the negative side):

$$\mathbf{w}^T \mathbf{r} \geq -\|\mathbf{r}\|_\infty \|\mathbf{w}\|_1 \geq -\eta \|\mathbf{w}\|_1.$$

- Claim: Lower bound of $\mathbf{w}^T \mathbf{r}$ is attained when $\mathbf{r} = -\eta \cdot \text{sign}(\mathbf{w})$

$$\begin{aligned}\mathbf{w}^T \mathbf{r} &= -\eta \mathbf{w}^T \text{sign}(\mathbf{w}) \\ &= -\eta \sum_{i=1}^d w_i \text{sign}(w_i) \\ &= -\eta \sum_{i=1}^d |w_i| \\ &= -\eta \|\mathbf{w}\|_1.\end{aligned}$$

- So the solution is $\mathbf{r} = -\eta \cdot \text{sign}(\mathbf{w})$.

To Summarize the Attack

Theorem (Maximum Loss ℓ_∞ Attack of Two-Class Linear Classifier)

The max-loss ℓ_∞ norm attack for a two-class linear classifier, i.e.,

$$\underset{\mathbf{x}}{\text{minimize}} \quad \mathbf{w}^T \mathbf{x} + w_0 \quad \text{subject to} \quad \|\mathbf{x} - \mathbf{x}_0\|_\infty \leq \eta.$$

is given by

$$\mathbf{x} = \mathbf{x}_0 - \eta \cdot \text{sign}(\mathbf{w}).$$

- Compare to minimum-distance attack:

$$\mathbf{x} = \mathbf{x}_0 - \left(\frac{\mathbf{w}^T \mathbf{x}_0 + w_0}{\|\mathbf{w}\|_1} \right) \cdot \text{sign}(\mathbf{w}).$$

- η is now a free variable. You need to pick.

FGSM (Goodfellow et al., NeurIPS 2014)

- Define training loss as

$$\begin{aligned} J(\mathbf{x}, \mathbf{w}) &= g_t(\mathbf{x}) - \max_{i \neq t} \{g_i(\mathbf{x})\} \\ &= - \left(\max_{i \neq t} \{g_i(\mathbf{x})\} - g_t(\mathbf{x}) \right). \end{aligned}$$

- Then max-loss attack is

$$\underset{\mathbf{x}}{\text{maximize}} \quad J(\mathbf{x}, \mathbf{w}) \quad \text{subject to} \quad \|\mathbf{x} - \mathbf{x}_0\|_\infty \leq \eta.$$

- Training: Minimize $J(\mathbf{x}, \mathbf{w})$ by finding a good \mathbf{w} .
- Attack: Maximize $J(\mathbf{x}, \mathbf{w})$ by finding a nasty \mathbf{x} .
- For neural networks, $J(\mathbf{x}, \mathbf{w})$ can be very general.

FGSM (Goodfellow et al., NeurIPS 2014)

- How to attack $J(\mathbf{x}, \mathbf{w})$?
- Linearize:

$$J(\mathbf{x}; \mathbf{w}) = J(\mathbf{x}_0 + \mathbf{r}; \mathbf{w}) \approx J(\mathbf{x}_0; \mathbf{w}) + \nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})^T \mathbf{r}.$$

- Then solve

$$\underset{\mathbf{r}}{\text{maximize}} \quad J(\mathbf{x}_0; \mathbf{w}) + \nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})^T \mathbf{r} \quad \text{subject to} \quad \|\mathbf{r}\|_{\infty} \leq \eta$$

- Equivalent to

$$\underset{\mathbf{r}}{\text{minimize}} \quad \underbrace{-\nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})^T \mathbf{r}}_{\mathbf{w}^T \mathbf{r}} - \underbrace{J(\mathbf{x}_0; \mathbf{w})}_{w_0} \quad \text{subject to} \quad \|\mathbf{r}\|_{\infty} \leq \eta$$

- Solution is

$$\mathbf{r} = \eta \cdot \text{sign}(-\nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w}))$$

FGSM (Goodfellow et al., NeurIPS 2014)

Definition (Fast Gradient Sign Method (FGSM) by Goodfellow et al 2014)

Given a loss function $J(\mathbf{x}; \mathbf{w})$, the FGSM creates an attack \mathbf{x} by

$$\mathbf{x} = \mathbf{x}_0 + \eta \cdot \text{sign}(\nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})). \quad (2)$$

Corollary (FGSM as a Max-Loss Attack Problem)

The FGSM attack can be formulated as the optimization with $J(\mathbf{x}; \mathbf{w})$ being the loss function:

$$\underset{\mathbf{r}}{\text{maximize}} \quad \nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})^T \mathbf{r} + J(\mathbf{x}_0; \mathbf{w}) \quad \text{subject to} \quad \|\mathbf{r}\|_{\infty} \leq \eta,$$

of which the solution is given by

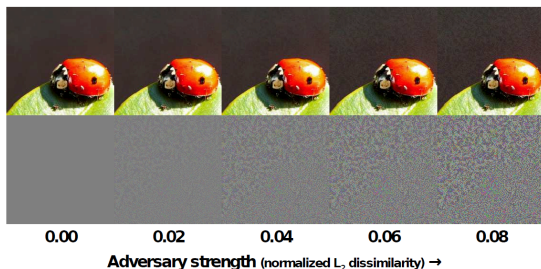
$$\mathbf{x} = \mathbf{x}_0 + \eta \cdot \text{sign}(\nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})). \quad (3)$$

FGSM (Goodfellow et al., NeurIPS 2014)

Definition (Fast Gradient Sign Method (FGSM) by Goodfellow et al 2014)

Given a loss function $J(\mathbf{x}; \mathbf{w})$, the FGSM creates an attack \mathbf{x} by

$$\mathbf{x} = \mathbf{x}_0 + \eta \cdot \text{sign}(\nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})). \quad (4)$$



<https://arxiv.org/pdf/1711.00117.pdf>

l_∞ and l_2 FGSM

Corollary (FGSM as a Max-Loss Attack)

The FGSM attack can be formulated as the optimization with $J(\mathbf{x}; \mathbf{w})$ being the loss function:

$$\underset{\mathbf{r}}{\text{maximize}} \quad \nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})^T \mathbf{r} + J(\mathbf{x}_0; \mathbf{w}) \quad \text{subject to} \quad \|\mathbf{r}\| \leq \eta,$$

of which the solution is given by

$$\mathbf{x} = \mathbf{x}_0 + \eta \cdot \text{sign}(\nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})) \quad (l_\infty\text{-norm})$$

and

$$\mathbf{x} = \mathbf{x}_0 + \eta \cdot \frac{\nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})}{\|\nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})\|_2} \quad (l_2\text{-norm})$$

Iterative Fast Gradient Sign Method

- By Kurakin, Goodfellow and Bengio (ICLR 2017)
- Recall this equation

$$\begin{aligned} J(\mathbf{x}; \mathbf{w}) &= J(\mathbf{x}_0 + \mathbf{r}; \mathbf{w}) \\ &\approx J(\mathbf{x}_0; \mathbf{w}) + \nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})^T \mathbf{r} \\ &= J(\mathbf{x}_0; \mathbf{w}) + \nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})^T (\mathbf{x} - \mathbf{x}_0) \\ &= J(\mathbf{x}_0; \mathbf{w}) + \nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})^T \mathbf{x} - \nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})^T \mathbf{x}_0. \end{aligned}$$

- Let us consider the problem

$$\begin{aligned} &\underset{\mathbf{x}}{\text{maximize}} \quad \cancel{J(\mathbf{x}_0; \mathbf{w})} + \nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})^T \mathbf{x} - \cancel{\nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})^T \mathbf{x}_0} \\ &\text{subject to} \quad \|\mathbf{x} - \mathbf{x}_0\| \leq \eta, \quad 0 \leq \mathbf{x} \leq 1. \end{aligned}$$

Iterative Gradient Sign Method

- Introduce iterative linearization

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \underset{\mathbf{x}}{\operatorname{argmax}} \quad \nabla_{\mathbf{x}} J(\mathbf{x}^{(k)}; \mathbf{w})^T \mathbf{x} \\ &\text{subject to } \|\mathbf{x} - \mathbf{x}^{(k)}\|_{\infty} \leq \eta, \quad 0 \leq \mathbf{x} \leq 1 \end{aligned}$$

- The optimization becomes

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \underset{\mathbf{x}}{\operatorname{argmax}} \quad \nabla_{\mathbf{x}} J(\mathbf{x}^{(k)}; \mathbf{w})^T \mathbf{x} \\ &\text{subject to } \|\mathbf{x} - \mathbf{x}^{(k)}\|_{\infty} \leq \eta, \quad 0 \leq \mathbf{x} \leq 1 \\ &= \mathcal{P}_{[0,1]} \left\{ \mathbf{x}^{(k)} + \eta \cdot \operatorname{sign}(\nabla_{\mathbf{x}} J(\mathbf{x}^{(k)}; \mathbf{w})) \right\}, \end{aligned}$$

- This is known as the projected gradient descent (PGD).
- Strongest first order attack, so far.
- You can add random noise to $\mathbf{x}^{(k)}$ to make it less predictable.

Outline

- Lecture 33 Overview
- Lecture 34 Min-distance attack
- **Lecture 35 Max-loss attack and regularized attack**

Today's Lecture

- Max-loss attack
 - Linear models
 - Deep models: FGSM and PGD
- **Regularized attack**
 - **Linear models**
 - **CW attack**

Two-Class Linear Classifier

- We want to study

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda \left(\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \right).$$

- If we restrict to two-class, linear classifier, then simplified to

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda \left((\mathbf{w}_j^T \mathbf{x} + w_{j,0}) - (\mathbf{w}_t^T \mathbf{x} + w_{t,0}) \right),$$

which is

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda(\mathbf{w}^T \mathbf{x} + w_0).$$

- Unconstrained minimization.
- Let $\varphi(\mathbf{x}) = \frac{1}{2}\|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda(\mathbf{w}^T \mathbf{x} + w_0)$. Then

$$0 = \nabla \varphi(\mathbf{x}) = (\mathbf{x} - \mathbf{x}_0) + \lambda \mathbf{w}.$$

- Solution is $\mathbf{x} = \mathbf{x}_0 - \lambda \mathbf{w}$.

Two-Class Linear Classifier

Theorem (Regularization-based Attack for Two-Class Linear Classifier)

The regularization-based attack for a two-class linear classifier generates the attack by solving

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda(\mathbf{w}^T \mathbf{x} + w_0),$$

of which the solution is given by

$$\mathbf{x} = \mathbf{x}_0 - \lambda \mathbf{w}.$$

- \mathbf{w} is search direction
- λ is step size
- You need to choose λ .

Unboundedness of ℓ_1 Attack

- Can we do ℓ_1 attack?

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|_1 + \lambda(\mathbf{w}^T \mathbf{x} + w_0),$$

which is equivalent to

$$\underset{\mathbf{r}}{\text{minimize}} \quad \|\mathbf{r}\|_1 + \lambda \mathbf{w}^T \mathbf{r}.$$

- The optimality condition is (sort of):

$$\text{sign}(r_i) + \lambda w_i = 0.$$

- This requires that

$$\lambda w_i = \begin{cases} \pm 1, & |r_i| > 0, \\ \in (-1, 1) & r_i = 0. \end{cases}$$

- So $|\lambda w_i|$ will never exceed 1.

Unboundedness of ℓ_1 Attack



$$\lambda w_i = \begin{cases} \pm 1, & |r_i| > 0, \\ \in (-1, 1) & r_i = 0. \end{cases}$$

- Therefore, if $|\lambda \mathbf{w}| > \mathbf{1}$, then the above equation is impossible to hold regardless of how we choose \mathbf{r} .
- This means that the optimization does not have a solution.
- You can show that the function

$$f(x) = |x| + \alpha x$$

goes to $-\infty$ as $x \rightarrow -\infty$ if $\alpha > 1$.

- and goes to $-\infty$ as $x \rightarrow +\infty$ if $\alpha > -1$.
- So unbounded below.

Carlini-Wagner Attack (2016)

- The idea is to solve

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\| + \lambda \cdot \max \left\{ \left(\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \right), 0 \right\},$$

- If $(\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x})) < 0$: Already misclassified. No action needed.
- If $(\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x})) > 0$: Not yet misclassified. Need action.
- Here we used the rectifier function

$$\zeta(x) = \max(x, 0).$$

- So the problem can be written as

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\| + \lambda \cdot \zeta \left(\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \right).$$

- The norm here can be ℓ_1 or ℓ_2 , or any other norm.

Comparing Regularized and Min-Norm

- Regularized attack is

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\| + \lambda \cdot \zeta \left(\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \right).$$

- Min-distance attack is

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\| + \iota_{\Omega}(\mathbf{x}),$$

where

$$\iota_{\Omega}(\mathbf{x}) = \begin{cases} 0, & \text{if } \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0, \\ +\infty, & \text{otherwise.} \end{cases}$$

- So the regularized attack (CW attack) is a soft-version of the min-distance attack.

CW Attack for ℓ_1 -norm

- We showed that this problem is unbounded below.

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|_1 + \lambda(\mathbf{w}^T \mathbf{x} + w_0),$$

- Now consider the CW attack:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|_1 + \lambda \max(\mathbf{w}^T \mathbf{x} + w_0, 0).$$

- The objective function is always non-negative: $\|\mathbf{x} - \mathbf{x}_0\|_1 \geq 0$ and $\max(\mathbf{w}^T \mathbf{x} + w_0, 0) \geq 0$.
- We are guaranteed to have a solution.
- Here is a trivial solution.
- Lower bound is achieved when $\mathbf{x} = \mathbf{x}_0$ and $\mathbf{w}^T \mathbf{x}_0 + w_0 = 0$.
- This happens when the attack solution is $\mathbf{x} = \mathbf{x}_0$ and \mathbf{x}_0 is on the decision boundary.
- Of course, the chance for this to happen is unlikely. So we can safely ignore this trivial case.

Convexity for Linear Classifier

- The function $h(\mathbf{x}) = \max(\varphi(\mathbf{x}), 0)$ is convex in \mathbf{x} if $\varphi(\mathbf{x})$ is convex.

-

$$\begin{aligned}h(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) &= \max(\varphi(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}), 0) \\ &\leq \max(\alpha \varphi(\mathbf{x}) + (1 - \alpha) \varphi(\mathbf{y}), 0) \\ &\leq \alpha \max(\varphi(\mathbf{x}), 0) + (1 - \alpha) \max(\varphi(\mathbf{y}), 0) \\ &= \alpha h(\mathbf{x}) + (1 - \alpha) h(\mathbf{y}).\end{aligned}$$

- Our $\varphi(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$. So φ is convex.
- So the overall optimization is convex

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\| + \lambda \max(\mathbf{w}^T \mathbf{x} + w_0, 0).$$

- That means you can solve using CVX.

General g

- In general, CW attack solves

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda \cdot \zeta \left(\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \right).$$

- We can use gradient algorithms.
- The gradient of $\zeta(\cdot)$ is

$$\frac{d}{ds} \zeta(s) = \mathbb{I}\{s > 0\} \stackrel{\text{def}}{=} \begin{cases} 1, & \text{if } s > 0, \\ 0, & \text{otherwise.} \end{cases}$$

- Let $i^*(\mathbf{x})$ be the index of the maximum response
- $i^*(\mathbf{x}) = \underset{j \neq t}{\operatorname{argmax}} \{g_j(\mathbf{x})\}$
- For the time being, let us assume that the index i^* is independent of \mathbf{x}
- Then, the gradient is

CW Attack Algorithm

- The gradient is

$$\begin{aligned}\nabla_{\mathbf{x}}\zeta\left(\max_{j\neq t}\{g_j(\mathbf{x})\}-g_t(\mathbf{x})\right) &= \nabla_{\mathbf{x}}\zeta(\{g_{i^*}(\mathbf{x})\}-g_t(\mathbf{x})) \\ &= \begin{cases} \nabla_{\mathbf{x}}g_{i^*}(\mathbf{x})-\nabla_{\mathbf{x}}g_t(\mathbf{x}), & \text{if } g_{i^*}(\mathbf{x})-g_t(\mathbf{x}) > 0, \\ 0, & \text{otherwise.} \end{cases} \\ &= \mathbb{I}\{g_{i^*}(\mathbf{x})-g_t(\mathbf{x}) > 0\}\cdot(\nabla_{\mathbf{x}}g_{i^*}(\mathbf{x})-\nabla_{\mathbf{x}}g_j(\mathbf{x}))\end{aligned}$$

- Letting $\varphi(\mathbf{x})$ be the overall objective function

$$\varphi(\mathbf{x}) = \|\mathbf{x}-\mathbf{x}_0\|^2 + \lambda\cdot\max\left\{\left(\max_{j\neq t}\{g_j(\mathbf{x})\}-g_t(\mathbf{x})\right), 0\right\},$$

- The gradient is

$$\nabla\varphi(\mathbf{x}; i^*) = 2(\mathbf{x}-\mathbf{x}_0) + \lambda\cdot\mathbb{I}\{g_{i^*}(\mathbf{x})-g_t(\mathbf{x}) > 0\}\cdot(\nabla g_{i^*}(\mathbf{x})-\nabla g_j(\mathbf{x})).$$

CW Attack Algorithm

- Gradient is

$$\nabla\varphi(\mathbf{x}; i^*) = 2(\mathbf{x} - \mathbf{x}_0) + \lambda \cdot \mathbb{I} \{g_{i^*}(\mathbf{x}) - g_t(\mathbf{x}) > 0\} \cdot (\nabla g_{i^*}(\mathbf{x}) - \nabla g_j(\mathbf{x})).$$

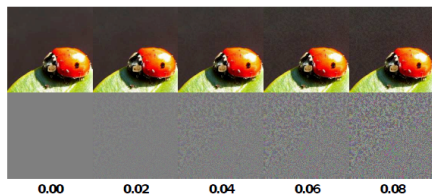
- The algorithm is
- For iteration $k = 1, 2, \dots$

$$i^* = \operatorname{argmax}_{j \neq t} \{g_j(\mathbf{x}^k)\}$$
$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha \nabla\varphi(\mathbf{x}^k; i^*).$$

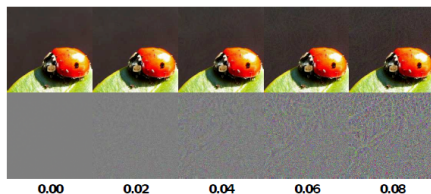
- α is gradient descent step size. You need to tune it.
- λ is regularization parameter. You need to tune it.

Comparison

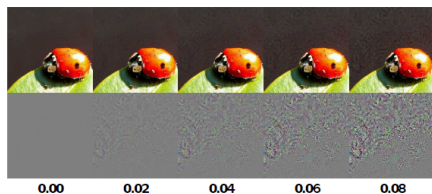
FGSM



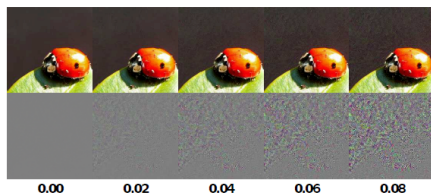
I-FGSM



DeepFool



Carlini-Wagner



<https://arxiv.org/pdf/1711.00117.pdf>

Summary

So we have discussed three forms of adversarial attacks.

Min-Distance Attack

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \|\mathbf{x} - \mathbf{x}_0\| \\ & \text{subject to} && \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0, \end{aligned}$$

Max-Loss Attack

$$\begin{aligned} & \underset{\mathbf{x}}{\text{maximize}} && g_t(\mathbf{x}) - \max_{j \neq t} \{g_j(\mathbf{x})\} \\ & \text{subject to} && \|\mathbf{x} - \mathbf{x}_0\| \leq \eta, \end{aligned}$$

Regularized Attack

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\| + \lambda (\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}))$$

- Next time we will talk about defense
- And then we will talk about fundamental trade off between robustness and accuracy