

ECE595 / STAT598: Machine Learning I

Lecture 34 Min-Distance Attacks

Spring 2020

Stanley Chan

School of Electrical and Computer Engineering
Purdue University



Today's Agenda

- Last lecture we have learned the basic terminologies of adversarial attack.
- In today's and the next lectures, we will go into the details of how to attack.
- We will discuss three forms of attacks
 - Min-distance attack
 - Max-loss attack
 - Regularized attack
- We will discuss everything for the **linear model**.
- And then we will talk about **deep models**.
- You are only required to know how to attack the linear model.
- For deep models, you probably need to have some prior experience with deep neural networks in order to understand what we are going to discuss.

Outline

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

Today's Lecture

- **Linear models**
 - **Definition**
 - **Geometry**
 - **Optimization solutions**
- Deep models
 - Deep fool
 - l_∞ case

Minimum Distance Attack

Definition (Minimum Distance Attack)

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

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \|\mathbf{x} - \mathbf{x}_0\| \\ & \text{subject to} && \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0, \end{aligned} \tag{1}$$

where $\|\cdot\|$ can be any norm specified by the user.

- I want to make you to class \mathcal{C}_t .
- So the constraint needs to be satisfied.
- But I also want to minimize the attack strength. This gives the objective.

Geometry: Attack as a Projection

What is the Geometry of the Attack?

- Claim: Attacking a data point = projecting it onto the decision boundary
- Let us look at ℓ_2 minimum distance attack

Theorem (Minimum-Distance Attack as a Projection)

The minimum-distance attack via ℓ_2 is equivalent to the projection

$$\begin{aligned} \mathbf{x}^* &= \underset{\mathbf{x} \in \Omega}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{x}_0\|^2, \quad \text{where } \Omega = \{\mathbf{x} \mid \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0\}, \\ &= \mathcal{P}_\Omega(\mathbf{x}_0), \end{aligned}$$

where $\mathcal{P}_\Omega(\cdot)$ denotes the projection onto the set Ω .

Geometry: Attack as a Projection

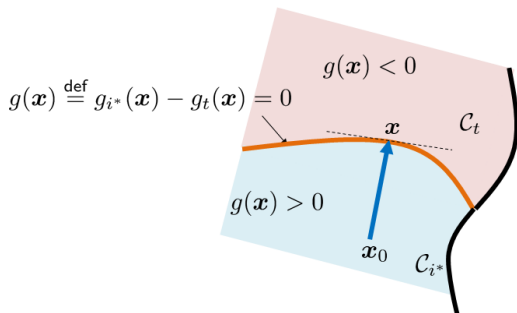


Figure: Geometry: Given an input data point \mathbf{x}_0 , our goal is to send \mathbf{x}_0 to a targeted class C_t by minimizing the distance between \mathbf{x} and \mathbf{x}_0 . The decision boundary is characterized by $g(\mathbf{x}) = g_{i^*}(\mathbf{x}) - g_t(\mathbf{x})$. The optimal solution is the projection of \mathbf{x}_0 onto the decision boundary.

Geometry: Overshoot

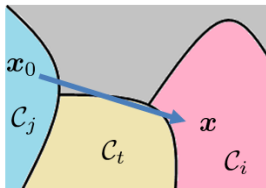
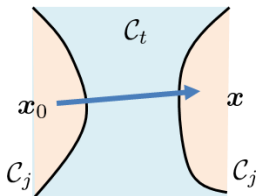
- What if you move along the attack direction but overshoot?

- Define

$$\mathbf{x} = \mathbf{x}_0 + \alpha(\mathcal{P}_\Omega(\mathbf{x}_0) - \mathbf{x}_0).$$

- Three cases:

- You overshoot but you still stay in the target class.
- You overshoot and you go back to the original class.
- You overshoot and you go to another class.



Targeted VS Untargeted Attack

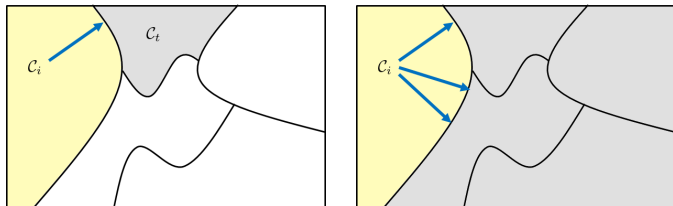


Figure: [Left] Targeted attack: The attack has to be specific from C_i to C_t .
[Right] Untargeted attack: The attack vector can point to anywhere outside C_i .

- Targeted attack: The constraint set Ω is

$$\Omega = \{\mathbf{x} \mid \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0\}$$

- Untargeted attack: The constraint set Ω is

$$\Omega = \{\mathbf{x} \mid g_i(\mathbf{x}) - \min_{j \neq i} \{g_j(\mathbf{x})\} \leq 0\}$$

White-box VS Black-box Attack

- **White-box:** You know everything about the classifier.
- So you know all g_i 's, completely.
- The constraint set is

$$\Omega = \{\mathbf{x} \mid \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0\}$$

- **Black-box:** You can only probe the classifier finite times.
- So you only know $\{g_i(\mathbf{x}^{(1)}), g_i(\mathbf{x}^{(2)}), \dots, g_i(\mathbf{x}^{(M)})\}$.
- The constraint set is

$$\Omega = \{\mathbf{x} \mid \max_{j \neq t} \{\hat{g}_j(\mathbf{x})\} - \hat{g}_t(\mathbf{x}) \leq 0\},$$

where \hat{g} is the best approximation you can get from the finite observations.

Launching the Attack: Basic Principles

- **Principle 1:** You need to solve the optimization

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

or its variations.

- **Principle 2:** You do not need to solve inequality. Equality is enough.
 - You just need to hit the decision boundary.
 - Then you add a small ϵ to your step.
- **Principle 3:** You do not need to be optimal.
 - Optimal = The nastiest attack.
 - You can still fool the classifier with a less nasty attack.
- **Our Plan:** Look at **linear** classifiers, and **binary** classifiers only.

So, if we restrict ourselves to binary linear classifiers ...

The min-distance attack (ℓ_2 -norm)

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \|\mathbf{x} - \mathbf{x}_0\|^2 \\ & \text{subject to} && \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0, \end{aligned}$$

will become ...

- **Linear** classifiers, we have

$$g_i(\mathbf{x}) - g_t(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0.$$

- **Two** class: the constraint is simplified to

$$g_i(\mathbf{x}) - g_t(\mathbf{x}) \leq 0$$

- And we just need to hit the boundary. So the attack becomes

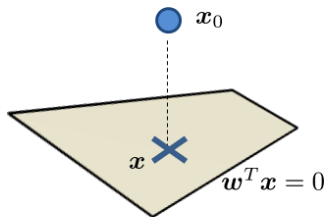
$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \|\mathbf{x} - \mathbf{x}_0\|^2 \\ & \text{subject to} && \mathbf{w}^T \mathbf{x} + w_0 = 0. \end{aligned}$$

Recall: Distance Between Point and Plane

What is the closest distance between a point and a plane?

- $\mathbf{w}^T \mathbf{x} = 0$ is a line.
- Find a point \mathbf{x} on the line that is closest to \mathbf{x}_0 .
- Solution is

$$\begin{aligned}\mathbf{x} &= \mathbf{x}_0 + \mathbf{w}(\mathbf{w}^T \mathbf{w})^{-1}(0 - \mathbf{w}^T \mathbf{x}_0) \\ &= \mathbf{x}_0 - \left(\frac{\mathbf{w}^T \mathbf{x}_0}{\|\mathbf{w}\|^2} \right)^T \mathbf{w}.\end{aligned}$$



Minimum-Distance Attack: Solving the Optimization

Theorem (Minimum ℓ_2 Norm Attack for Two-Class Linear Classifier)

The adversarial attack to a two-class linear classifier is the solution of

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 \quad \text{subject to} \quad \mathbf{w}^T \mathbf{x} + w_0 = 0,$$

which is given by

$$\mathbf{x}^* = \mathbf{x}_0 - \left(\frac{\mathbf{w}^T \mathbf{x}_0 + w_0}{\|\mathbf{w}\|_2} \right) \frac{\mathbf{w}}{\|\mathbf{w}\|_2}.$$

- This is just finding the closest point to a hyperplane!
- $\mathbf{w}/\|\mathbf{w}\|_2$ is the normal direction = best attack angle.
- $\frac{\mathbf{w}^T \mathbf{x}_0 + w_0}{\|\mathbf{w}\|_2}$ is the step size.

Minimum-Distance Attack: Two-Class Linear Classifier

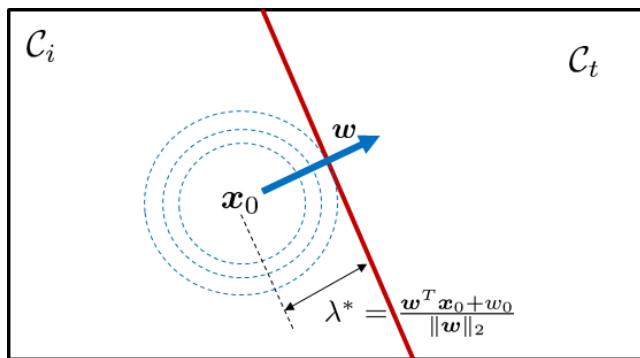


Figure: Geometry of minimum-distance attack for a two-class linear classifier with objective function $\|x - x_0\|^2$. The solution is a projection of the input x_0 onto the separating hyperplane of the classifier.

Outline

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

Today's Lecture

- Linear models
 - Definition
 - Geometry
 - Optimization solutions
- **Deep models**
 - **Deep fool**
 - l_∞ case

Deep-Fool (CVPR, 2016)

Let's Connect to the Real Problem.

- Proposed by Moosavi-Dezfooli, Fawzi and Frossard
- Generalize linear classifier to neural network

Definition (DeepFool Attack by Moosavi-Dezfooli et al. 2016)

The DeepFool attack for a two-class classification generates the attack by solving the optimization

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 \quad \text{subject to} \quad g(\mathbf{x}) = 0,$$

where $g(\mathbf{x}) = 0$ is the nonlinear decision boundary separating the two classes.

How to deal with non-linearity?

- First order approximation

$$g(\mathbf{x}) \approx g(\mathbf{x}^{(k)}) + \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})^T(\mathbf{x} - \mathbf{x}^{(k)}),$$

- Modify the problem (assume $\mathbf{x}^{(0)} = \mathbf{x}_0$)

$$\mathbf{x}^{(k+1)} = \underset{\mathbf{x}}{\operatorname{argmin}} \quad \|\mathbf{x} - \mathbf{x}^{(k)}\|^2 \quad \text{subject to } g(\mathbf{x}) = 0.$$

⋮

$$\mathbf{x}^{(k+1)} = \underset{\mathbf{x}}{\operatorname{argmin}} \quad \|\mathbf{x} - \mathbf{x}^{(k)}\|^2$$

subject to $g(\mathbf{x}^{(k)}) + \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})^T(\mathbf{x} - \mathbf{x}^{(k)}) = 0.$

- Now, rewrite

$$\begin{aligned} & g(\mathbf{x}^{(k)}) + \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})^T(\mathbf{x} - \mathbf{x}^{(k)}) \\ &= \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})^T\mathbf{x} + g(\mathbf{x}^{(k)}) - \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})^T\mathbf{x}^{(k)}. \end{aligned}$$

How to deal with non-linearity?

- So here is our problem

$$\mathbf{x}^{(k+1)} = \underset{\mathbf{x}}{\operatorname{argmin}} \quad \|\mathbf{x} - \mathbf{x}^{(k)}\|^2$$

subject to $\mathbf{g}(\mathbf{x}^{(k)}) + \nabla_{\mathbf{x}}\mathbf{g}(\mathbf{x}^{(k)})^T(\mathbf{x} - \mathbf{x}^{(k)}) = 0.$

- Let $\mathbf{w}^{(k)} = \nabla_{\mathbf{x}}\mathbf{g}(\mathbf{x}^{(k)})$ and $w_0^{(k)} = \mathbf{g}(\mathbf{x}^{(k)}) - \nabla_{\mathbf{x}}\mathbf{g}(\mathbf{x}^{(k)})^T\mathbf{x}^{(k)}$
- Then equivalent to

$$\mathbf{x}^{(k+1)} = \underset{\mathbf{x}}{\operatorname{argmin}} \quad \|\mathbf{x} - \mathbf{x}^{(k)}\|^2 \quad \text{subject to} \quad (\mathbf{w}^{(k)})^T\mathbf{x} + w_0^{(k)} = 0$$

- This is just a linear problem!

How to deal with non-linearity?

- Here is the optimization

$$\mathbf{x}^{(k+1)} = \underset{\mathbf{x}}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{x}^{(k)}\|^2 \quad \text{subject to } (\mathbf{w}^{(k)})^T \mathbf{x} + w_0^{(k)} = 0$$

- So the solution is

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \left(\frac{(\mathbf{w}^{(k)})^T \mathbf{x}^{(k)} + w_0^{(k)}}{\|\mathbf{w}^{(k)}\|^2} \right) \mathbf{w}^{(k)} \\ &= \mathbf{x}^{(k)} - \left(\frac{g(\mathbf{x}^{(k)})}{\|\nabla_{\mathbf{x}} g(\mathbf{x}^{(k)})\|^2} \right) \nabla_{\mathbf{x}} g(\mathbf{x}^{(k)}). \end{aligned}$$

- How to evaluate the gradient?
- $\nabla_{\mathbf{x}} g(\mathbf{x}^{(k)})$ can be computed via back propagation.

How to deal with non-linearity?

- Now, for this attack

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \left(\frac{g(\mathbf{x}^{(k)})}{\|\nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})\|^2} \right) \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)}).$$

- You can control the perturbation magnitude:

$$\mathbf{x}^{(k+1)} = \mathcal{P}_{[0,1]} \left\{ \mathbf{x}^{(k)} - \left(\frac{g(\mathbf{x}^{(k)})}{\|\nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})\|^2} \right) \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)}) \right\}.$$

- $\mathcal{P}_{[0,1]}$: Projection onto a ball, e.g., $\mathcal{P}_{[0,1]}(\mathbf{x})$ clips \mathbf{x} to $[0, 1]$.

Deep-Fool (CVPR, 2016)

Corollary (DeepFool Algorithm for Two-Class Problem)

An iterative procedure to obtain the DeepFool attack solution is

$$\begin{aligned}\mathbf{x}^{(k+1)} &= \underset{\mathbf{x}}{\operatorname{argmin}} \quad \|\mathbf{x} - \mathbf{x}^{(k)}\|^2 \\ &\text{subject to } g(\mathbf{x}^{(k)}) + \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})^T(\mathbf{x} - \mathbf{x}^{(k)}) = 0 \\ &= \mathbf{x}^{(k)} - \left(\frac{g(\mathbf{x}^{(k)})}{\|\nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})\|^2} \right) \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)}),\end{aligned}$$

with $\mathbf{x}^{(0)} = \mathbf{x}_0$.

- This is not the complete Deep-fool.
- We assume two classes only.
- If you have multiple classes, you need to take care of “ $\max_{j \neq t} g_j(\mathbf{x})$ ”

The l_∞ Case

- How about we try to solve this?

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|_\infty \quad \text{subject to} \quad \mathbf{w}^T \mathbf{x} + w_0 = 0.$$

- Not the l_2 -norm, but the l_∞ -norm.
- Let $\mathbf{r} = \mathbf{x} - \mathbf{x}_0$, $b_0 = -(\mathbf{w}^T \mathbf{x}_0 + w_0)$.
- Rewrite the problem as

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

- Setup Lagrangian function and take derivative?

$$\mathcal{L}(\mathbf{r}, \lambda) = \|\mathbf{r}\|_\infty + \lambda(b_0 - \mathbf{w}^T \mathbf{r}).$$

- Doesn't work because l_∞ is not differentiable.

Solving the ℓ_∞ -norm Problem

Theorem (Holder's Inequality)

Let $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{y} \in \mathbb{R}^d$. Then,

$$-\|\mathbf{x}\|_p \|\mathbf{y}\|_q \leq \mathbf{x}^T \mathbf{y} \leq \|\mathbf{x}\|_p \|\mathbf{y}\|_q$$

for any p and q such that $\frac{1}{p} + \frac{1}{q} = 1$, where $p \in [1, \infty]$.

- Let $p = 1$ and $q = \infty$
- Can show that $|\mathbf{x}^T \mathbf{y}| \leq \|\mathbf{x}\|_1 \|\mathbf{y}\|_\infty$
- Then

$$|b_0| = |\mathbf{w}^T \mathbf{r}| \leq \|\mathbf{w}\|_1 \|\mathbf{r}\|_\infty, \quad \implies \quad \|\mathbf{r}\|_\infty \geq \frac{|b_0|}{\|\mathbf{w}\|_1}.$$

- So $\|\mathbf{r}\|_\infty$ is lower bounded by a constant.
- If \mathbf{r}^* can reach this lower bound, then \mathbf{r}^* is the minimizer.

Solving the ℓ_∞ -norm Problem

- How about this candidate?

$$\mathbf{r} = \eta \cdot \text{sign}(\mathbf{w})$$

for some constant η to be determined.

- We can show that

$$\|\mathbf{r}\|_\infty = \max_i |\eta \cdot \text{sign}(w_i)| = |\eta|.$$

- So if we let $\eta = b_0 / \|\mathbf{w}\|_1$, then we will have

$$\|\mathbf{r}\|_\infty = |\eta| = \frac{|b_0|}{\|\mathbf{w}\|_1}.$$

- Lower bound achieved! So the solution is

$$\mathbf{r} = \frac{|b_0|}{\|\mathbf{w}\|_1} \cdot \text{sign}(\mathbf{w})$$

The ℓ_∞ Solution

Theorem (Minimum Distance ℓ_∞ Norm Attack for Two-Class Linear Classifier)

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

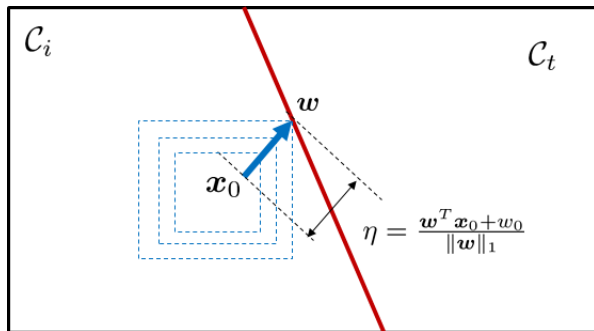
$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|_\infty \quad \text{subject to} \quad \mathbf{w}^T \mathbf{x} + w_0 = 0$$

is given by

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

- Search direction is $\text{sign}(\mathbf{w})$.
- This means ± 1 for every entry.
- In 2D, the search direction is $\pm 45^\circ$ or $\pm 135^\circ$.

The l_∞ Solution



- Is it the "optimal" direction? No.
- The fastest search direction is l_2 .
- Can it move x_0 to another class? Yes, if η is large enough.

Summary

Min-Distance Attack

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

- We have talked about the geometry.
- You can see that the geometry applies beyond linear models.
- For linear models, we can derive closed-form solutions.
- Deep models apply successive approximations.

Next Lecture

- **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**

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