

Chapter 3

Adversarial Attack

Consider a data point $\mathbf{x}_0 \in \mathbb{R}^d$ belonging to class \mathcal{C}_i . Adversarial attack is a malicious attempt which tries to perturb \mathbf{x}_0 to a new data point \mathbf{x} such that \mathbf{x} is misclassified by the classifier. Goodfellow et al. made this explicit by showing examples perceptually indistinguishable but sufficient to fool a classifier (See Figure 3.1). This raises many important questions such as vulnerability of machine learning algorithms and the associated security threats. While many people believe that adversarial attack is unique to deep neural networks as these classifiers tend to overfit, adversarial attack is in fact an intrinsic property of all classifier. In this chapter, we will study the fundamentals of adversarial attack through various examples of linear classifiers. Our goal is to understand the source of these attacks, the geometry, and possible defensive strategies.

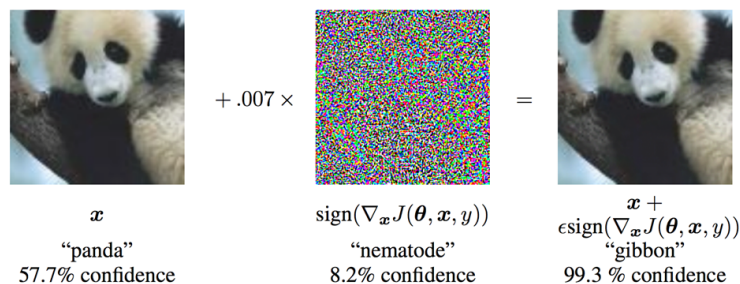


Figure 3.1: A classical example by Goodfellow et al 2014. By perturbing the input data \mathbf{x}_0 with a perceptually undistinguishable “noise”, the perturbed data is misclassified with high confidence. We will discuss how to generate these “noise” pattern, and their geometric meanings.

3.1 Problem Formulation

What is Adversarial Attack?

Adversarial attack is an malicious attempt to perturb a data point $\mathbf{x}_0 \in \mathbb{R}^d$ to another point $\mathbf{x} \in \mathbb{R}^d$ such that \mathbf{x} belongs to certain **target** adversarial class. For example if \mathbf{x}_0 is

a feature vector of a **cat** image, by adversarial attack we meant to create another feature vector \mathbf{x} which will be classifier as a **dog** (or another class specified by the attacker). In some scenarios, the goal may not be to push \mathbf{x}_0 to a specific target class \mathcal{C}_t , but just push it away from its original class \mathcal{C}_i . Attacks of this type is known as the **untargeted** attack. For the time being, let us focus on targeted attack first.

There are various approaches to move a data point \mathbf{x}_0 from \mathcal{C}_i to \mathcal{C}_t . The most general definition of such mapping is to consider an operator $\mathcal{A} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that $\mathbf{x} = \mathcal{A}(\mathbf{x}_0)$ is the perturbed data.

Definition 1 (Adversarial Attack). *Let $\mathbf{x}_0 \in \mathbb{R}^d$ be a data point belong to class \mathcal{C}_i . Define a target class \mathcal{C}_t . An **adversarial attack** is a mapping $\mathcal{A} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that the perturbed data*

$$\mathbf{x} = \mathcal{A}(\mathbf{x}_0)$$

is misclassified as \mathcal{C}_t .

Among many adversarial attack models, the most commonly used one is the additive model, where we define \mathcal{A} as a linear operator that adds perturbation to the input.

Definition 2 (Additive Adversarial Attack). *Let $\mathbf{x}_0 \in \mathbb{R}^d$ be a data point belong to class \mathcal{C}_i . Define a target class \mathcal{C}_t . An **additive** adversarial attack is an addition of a perturbation $\mathbf{r} \in \mathbb{R}^d$ such that the perturbed data*

$$\mathbf{x} = \mathbf{x}_0 + \mathbf{r}$$

is misclassified as \mathcal{C}_t .

Additive attacks are appealing for at least two reasons. First, an additive attack guarantees the input space remains unchanged. For example, if the input \mathbf{x}_0 is an image in \mathbb{R}^d , then the perturbed data \mathbf{x} remains in \mathbb{R}^d and still represents an image. A general mapping $\mathcal{A}(\cdot)$ may not preserve such relationship, e.g., \mathcal{A} is a sub-sampling operator or a dense matrix that transforms the input space to a feature space. Second, additive attack allows interpretable analysis with simple geometry. A general mapping $\mathcal{A}(\cdot)$ could be very complicated to analyze. Therefore, unless specified, in the rest of this chapter we will refer adversarial attack exclusively to the additive attacks.

Formulating Adversarial Attack

Since the primary goal of adversarial attack is to perturb a data point for misclassification, we need to first consider the input class \mathcal{C}_i and the target class \mathcal{C}_t . The input class \mathcal{C}_i is the class where \mathbf{x}_0 belongs to, and \mathcal{C}_t is the class which we wish the attack data \mathbf{x} to be.

Consider a multi-class scenario where we have classes $\mathcal{C}_1, \dots, \mathcal{C}_k$. The decision boundaries of these k classes are specified by the k discriminant functions $g_i(\cdot)$ for $i = 1, \dots, k$. If we want \mathbf{x} to be classified as class \mathcal{C}_t , then the discriminant functions should satisfy the condition that

$$g_t(\mathbf{x}) \geq g_j(\mathbf{x}), \quad \text{for all } j \neq t, \quad (3.1)$$

so that the target class \mathcal{C}_t will have a discriminant value $g_t(\mathbf{x})$ greater than all other classes in the classifier. Rewriting the $k - 1$ inequalities, we can equivalently express them as

$$g_t(\mathbf{x}) \geq \max_{j \neq t} \{g_j(\mathbf{x})\} \iff \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0. \quad (3.2)$$

Therefore, the goal of adversarial attack is to find \mathbf{x} such that the inequality in Equation (3.2) is satisfied.

The inequality in Equation (3.2) defines a set of \mathbf{x} that are allowed to be potential attacks. If this set is non-empty, then any element in the set could be used as the attack. However, some attacks are closer to the input \mathbf{x}_0 than others. In order to pick a unique attack, it is often necessary to post additional criteria in terms of the adversarial magnitude. That is, we want the perturbation as small as possible, or the perturbed data should stay as close to \mathbf{x}_0 as possible. Once such criteria are added to the problem, we will reach an optimization formulation. The following three definitions are the most commonly used definitions of attacks in the literature.

Definition 3 (Minimum Norm Attack). *The **minimum norm attack** finds a perturbed data \mathbf{x} by solving the optimization*

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \|\mathbf{x} - \mathbf{x}_0\| \\ & \text{subject to} && \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0, \end{aligned} \quad (3.3)$$

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

As shown in the definition, the goal of the minimum norm attack is to minimize the perturbation magnitude while ensuring the new data \mathbf{x} is classified as \mathcal{C}_t . If the constraint set is empty, then we say that it is impossible to attack the classifier for a particular input \mathbf{x}_0 .

An alternative to the minimum norm attack is the maximum allowable attack.

Definition 4 (Maximum Allowable Attack). *The **maximum allowable attack** finds a perturbed data \mathbf{x} by solving the optimization*

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \\ & \text{subject to} && \|\mathbf{x} - \mathbf{x}_0\| \leq \eta, \end{aligned} \quad (3.4)$$

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

In maximum allowable attack, the magnitude of the attack is upper bounded by η . That is, the attack cannot go beyond the range $\|\mathbf{x} - \mathbf{x}_0\| \leq \eta$. The objective function seeks \mathbf{x} to make $\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x})$ as negative as possible. If the best \mathbf{x} still cannot make the objective value negative, then we say that it is impossible to attack the classifier for a particular input \mathbf{x}_0 .

The third formulation is an unconstrained optimization based on regularization.

Definition 5 (Regularization-based Attack). *The **regularization-based attack** finds a perturbed data \mathbf{x} by solving the optimization*

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\| + \lambda (\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x})) \quad (3.5)$$

where $\|\cdot\|$ can be any norm specified by the user, and $\lambda > 0$ is a regularization parameter.

The parameter in the regularization-based attack is also known as the Lagrange multiplier. Intuitively, Equation (3.5) tries to simultaneously minimize two conflicting objectives $\|\mathbf{x} - \mathbf{x}_0\|$ and $\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x})$. Depending on the choice of λ , relative emphasis on the two terms can be controlled.

Which Optimization? While Equation (3.3), Equation (3.4) and Equation (3.5) offer three different ways of determining the adversarial attack (and hence three different solutions), it can be shown that for appropriately chosen η and λ , the three optimizations are in fact equivalent in the sense that the solutions are identical. In other words, if we can solve one problem, we will be able to solve the other two problems provided we have the right η and λ . Because of the equivalence between the optimization problems, we will spend most of the time on the **minimum norm attack**. There are two reasons.

- **No Parameter.** Unlike maximum allowable attack and regularization-based attack, **minimum norm attack does not require user defined parameters λ and η** . This allows us to determine unique solutions independent of the parameters.
- **Simple geometry.** It is easier to interpret $\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x})$ as a constraint set rather an objective function, for minimizing $\|\mathbf{x} - \mathbf{x}_0\|$ over a set can be treated as a **projection**.

Geometry of the Objective Function

The norm $\|\mathbf{x} - \mathbf{x}_0\|$ is measuring the distance between \mathbf{x} and \mathbf{x}_0 . There are various options, with geometric interpretation shown in Figure 3.2.

- **ℓ_0 -norm:** $\varphi(\mathbf{x}) = \|\mathbf{x} - \mathbf{x}_0\|_0$, which gives the most sparse solution. Useful when we want to limit the number of attack pixels.
- **ℓ_1 -norm:** $\varphi(\mathbf{x}) = \|\mathbf{x} - \mathbf{x}_0\|_1$, which is a convex surrogate of the ℓ_0 -norm.

- ℓ_∞ -norm: $\varphi(\mathbf{x}) = \|\mathbf{x} - \mathbf{x}_0\|_\infty$, which minimizes the maximum element of the perturbation.

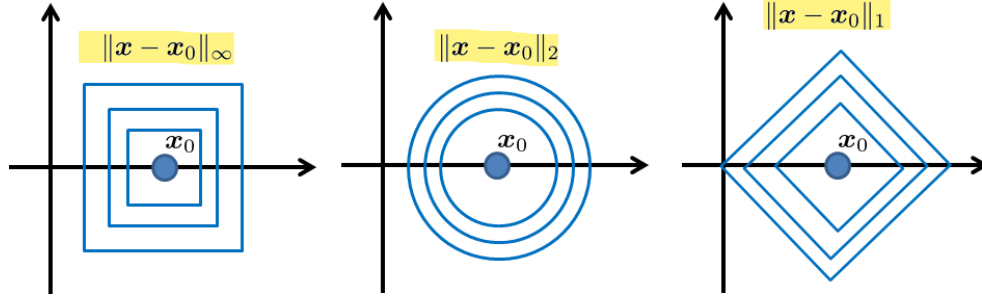


Figure 3.2: Geometry of different objective functions.

Besides changing the norm of the objective function, we can also post constraints to the objective. For example, in some applications we may want to limit the attack at certain spatial location. This can be achieved by adding a masking $\mathbf{M} \in \mathbb{R}^{m \times d}$ so that the distance becomes

$$\varphi(\mathbf{x}) = \|\mathbf{M}(\mathbf{x} - \mathbf{x}_0)\|^2. \quad \text{M can be 0 or 1} \quad (3.6)$$

In other examples, the objective function may contain a regularization term $h(\mathbf{x})$ so that φ becomes

$$\varphi(\mathbf{x}) = \|\mathbf{x} - \mathbf{x}_0\|^2 + \beta h(\mathbf{x}). \quad (3.7)$$

Geometry of the Constraints

The constraint in Equation (3.3) is a compact form of Equation (3.1). If we do not mind the tediousness of the constraints, we can show that $\Omega = \{\mathbf{x} \mid \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0\}$ is equivalent to

$$\Omega = \left\{ \mathbf{x} \mid \begin{array}{l} g_1(\mathbf{x}) - g_t(\mathbf{x}) \leq 0 \\ g_2(\mathbf{x}) - g_t(\mathbf{x}) \leq 0 \\ \vdots \\ g_k(\mathbf{x}) - g_t(\mathbf{x}) \leq 0 \end{array} \right\} \quad (3.8)$$

Depending on the nature of the discriminant function $g_i(\mathbf{x})$, the geometry of Ω could be polygon, convex, concave, or arbitrary.

To understand the behavior of the constraint, we show in Figure 3.3 a typical example of the constraint set. Given the discriminant function $g_i(\mathbf{x})$ for $i = 1, \dots, k$, the decision boundary between class \mathcal{C}_i and \mathcal{C}_t is $g_i(\mathbf{x}) - g_t(\mathbf{x}) = 0$. Depending on where \mathbf{x}_0 is located, the decision boundaries will change and hence the inequalities in Ω change with \mathbf{x}_0 . If a class \mathcal{C}_i does not have any direct contact with \mathcal{C}_t , a straight inequality will hold.

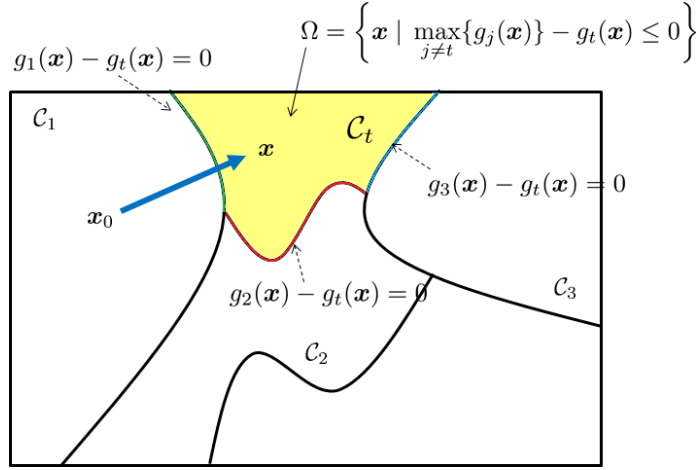


Figure 3.3: Geometry of the constraint set Ω . The decision boundaries are given by $g_i(\mathbf{x}) - g_t(\mathbf{x}) = 0$. The set Ω constitute all \mathbf{x} such that every $\mathbf{x} \in \Omega$ belongs to \mathcal{C}_t .

Caution: Note that we when say $\mathbf{x} \in \Omega$, we are referring to a point which simultaneously satisfies all the decision boundary constraints. It is therefore incorrect to say pick an index i^* which maximizes $g_i(\mathbf{x})$ for a particular \mathbf{x} , and simply study $g_{i^*}(\mathbf{x}) - g_t(\mathbf{x}) \leq 0$. In general, the index i^* must be a function of \mathbf{x} , i.e., $i^*(\mathbf{x})$. Therefore, the correct set is

$$\Omega = \{\mathbf{x} \mid g_{i^*(\mathbf{x})}(\mathbf{x}) - g_t(\mathbf{x}) \leq 0\}.$$

The reason for i^* to depend on \mathbf{x} is that the maximum is point-wise for every \mathbf{x} . As \mathbf{x} changes, the index corresponding to the maximizer also changes. The only exception is the linear classifier for two classes, because the decision boundary is a straight line. In this case, $\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0$ is simplified to $g_{i^*}(\mathbf{x}) - g_t(\mathbf{x}) \leq 0$ because for a two-class problem the index is either i^* or t . For non-linear classifiers, or linear classifiers with more than two classes, such simplification generally does not hold except for some special cases.

Geometry of the Constraints for Linear Classifiers

To gain insights about these constraints, we consider a k -class linear classifier. In this case, each discriminant function takes the form

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i,0}. \quad (3.9)$$

The decision boundary between the i -th class and the t -th class is therefore

$$g(\mathbf{x}) = (\mathbf{w}_i - \mathbf{w}_t)^T \mathbf{x} + w_{i,0} - w_{t,0} = 0. \quad (3.10)$$

Since we want the perturbed data \mathbf{x} to live in \mathcal{C}_t , the constraint set becomes

$$\begin{bmatrix} \mathbf{w}_1^T - \mathbf{w}_t^T \\ \vdots \\ \mathbf{w}_{t-1}^T - \mathbf{w}_t^T \\ \mathbf{w}_{t+1}^T - \mathbf{w}_t^T \\ \vdots \\ \mathbf{w}_k^T - \mathbf{w}_t^T \end{bmatrix} \mathbf{x} + \begin{bmatrix} w_{1,0} - w_{t,0} \\ \vdots \\ w_{t-1,0} - w_{t,0} \\ w_{t+1,0} - w_{t,0} \\ \vdots \\ w_{k,0} - w_{t,0} \end{bmatrix} \leq \mathbf{0} \Leftrightarrow \mathbf{A}^T \mathbf{x} \leq \mathbf{b} \quad (3.11)$$

where $\mathbf{A} = [\mathbf{w}_1 - \mathbf{w}_t, \dots, \mathbf{w}_k - \mathbf{w}_t] \in \mathbb{R}^{d \times (k-1)}$, and $\mathbf{b} = [w_{t,0} - w_{1,0}, \dots, w_{t,0} - w_{k,0}]^T \in \mathbb{R}^{k-1}$. This is summarized in the following Lemma.

Lemma 1 (Constraint Set of Linear Classifier). *Consider a k -class linear classifier with discriminant function $g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$ for $i = 1, \dots, k$. Define*

$$\mathbf{A} = [\mathbf{w}_1 - \mathbf{w}_t, \dots, \mathbf{w}_k - \mathbf{w}_t] \in \mathbb{R}^{d \times (k-1)}, \quad \text{and} \quad \mathbf{b} = [w_{t,0} - w_{1,0}, \dots, w_{t,0} - w_{k,0}]^T \in \mathbb{R}^{k-1},$$

Then, the constraint set is

$$\Omega = \{\mathbf{x} \mid \mathbf{A}^T \mathbf{x} \leq \mathbf{b}\}. \quad (3.12)$$

For linear classifiers, the constraint set Ω defines a d -dimensional polytope. Figure 3.4 shows one example. The normal vectors of each decision boundary is the column of the \mathbf{A} matrix, and the i -th decision boundary is

$$\mathbf{a}_i^T \mathbf{x} = b_i. \quad (3.13)$$

Since any polytope is convex, the constraint set Ω is convex. This can also be proved easily by letting $\mathbf{x}_1 \in \Omega$ and $\mathbf{x}_2 \in \Omega$. For any $0 \leq \lambda \leq 1$, we have

$$\mathbf{A}^T(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2) = \lambda \mathbf{A}^T \mathbf{x}_1 + (1 - \lambda) \mathbf{A}^T \mathbf{x}_2 \leq \lambda \mathbf{b} + (1 - \lambda) \mathbf{b} = \mathbf{b}.$$

Thus, $\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2 \in \Omega$ as well.

The set Ω is always non-empty if the classifier is a **linear machine**. A linear machine defines the decision based on Equation (3.1), i.e., decide class \mathcal{C}_t if $g_t(\mathbf{x}) \geq g_j(\mathbf{x})$ for every $j \neq t$. Therefore, the entire input space is **partitioned** into n decision boundaries, and each partition can contain up to $k - 1$ decision boundaries. Since every partition is not empty by definition, Ω must also be non-empty.

An immediate corollary of the lemma is that for linear classifiers, the adversarial attack problem is essentially a quadratic minimization

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 \quad \text{subject to} \quad \mathbf{A}^T \mathbf{x} \leq \mathbf{b}, \quad (3.14)$$

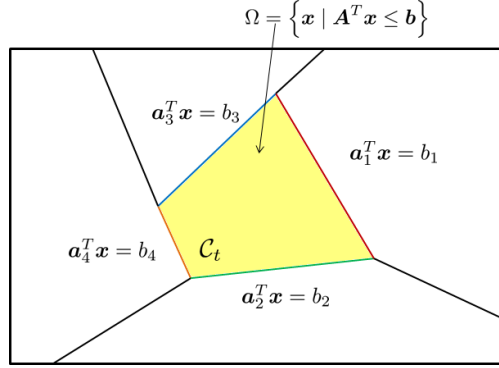


Figure 3.4: Geometry of the constraint set Ω for a linear classifier. The constraint set Ω is now a polygon with decision boundaries defined by $\mathbf{a}_i^T \mathbf{x} = b_i$, where $\mathbf{a}_i = \mathbf{w}_i - \mathbf{w}_t$ and $b_i = w_{i0} - w_{t0}$.

which can be solved easily using convex programming techniques. The following examples illustrates how to formulate an ℓ_1 -norm attack via linear programming.

Example. (Linear programming formulation for ℓ_1 -norm attack.) If we modify the ℓ_2 -norm objective to an ℓ_1 -norm, i.e.,

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|_1 \quad \text{subject to} \quad \mathbf{A}^T \mathbf{x} \leq \mathbf{b}, \quad (3.15)$$

then the problem can be formulated via a linear programming. To see this, we first rewrite Equation (3.15) by letting $\mathbf{r} = \mathbf{x} - \mathbf{x}_0$ so that we have

$$\underset{\mathbf{r}}{\text{minimize}} \quad \|\mathbf{r}\|_1 \quad \text{subject to} \quad \mathbf{A}^T \mathbf{r} \leq \tilde{\mathbf{b}},$$

where $\tilde{\mathbf{b}} = \mathbf{b} - \mathbf{A}^T \mathbf{x}_0$. Now, define \mathbf{r}_+ and \mathbf{r}_- be the positive and negative parts of \mathbf{r} such that $\mathbf{r} = \mathbf{r}_+ - \mathbf{r}_-$. Then, $\|\mathbf{r}\|_1 = \mathbf{r}_+ + \mathbf{r}_-$, and so the optimization becomes

$$\begin{aligned} &\underset{\mathbf{r}_+, \mathbf{r}_-}{\text{minimize}} && \mathbf{r}_+ + \mathbf{r}_- \\ &\text{subject to} && [\mathbf{A}^T \quad -\mathbf{A}^T] \begin{bmatrix} \mathbf{r}_+ \\ \mathbf{r}_- \end{bmatrix} \leq \tilde{\mathbf{b}}, \quad \mathbf{r}_+ \geq 0, \quad \text{and} \quad \mathbf{r}_- \geq 0. \end{aligned} \quad (3.16)$$

Note that this is a standard linear programming problem, which can be solved efficiently using the Simplex method.

Caution. The constraint set Ω could be empty if the classifier is not a linear machine. For example, if the decision rule is based on checking if $g_i(\mathbf{x}) > 0$, or if pairwise $g_i(\mathbf{x}) > g_j(\mathbf{x})$, then it is possible to have unclassified regions. Linear machines do not have the issue because it checks $g_i(\mathbf{x}) > g_j(\mathbf{x})$ for every $j \neq i$. See Duda-Hart-Stork Section 5.2.2 for discussion.

3.2 Geometry of the Attack

A Projection Perspective

The geometry of adversarial attack is specified by two factors. First is the distance metric, which defines whether we like to use a ball (ℓ_2 -norm), a diamond (ℓ_1 -norm), or square (ℓ_∞ -norm). Second is the feasible set Ω which specifies the decision boundary between the original class where \mathbf{x}_0 is positioned at, and the targeted class where \mathbf{x} should be sent to. If the objective is ℓ_2 , then we can re-formulate the minimum-norm attack problem as a projection.

Theorem 1 (Minimum-Norm Attack as a Projection). *The adversarial attack formulated by Equation (3.3) is equivalent to the projection*

$$\begin{aligned} \mathbf{x}^* &= \underset{\mathbf{x} \in \Omega}{\operatorname{argmin}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2, \quad \text{where } \Omega = \{\mathbf{x} \mid \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0\}, \\ &= \mathcal{P}_\Omega(\mathbf{x}_0), \end{aligned} \quad (3.17)$$

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

Equation (3.17) defines a projection of \mathbf{x}_0 onto the set Ω . Figure 3.5 illustrates the idea. Suppose we are originally at \mathbf{x}_0 and we like to move to class \mathcal{C}_t . In this particular illustration, the class being considered is \mathcal{C}_{i^*} . By solving the minimization problem, we effectively project \mathbf{x}_0 onto the decision boundary between \mathcal{C}_t and \mathcal{C}_{i^*} . The direction of the perturbation is the normal vector of the decision boundary at \mathbf{x} . The magnitude of the perturbation is the shortest distance between the input \mathbf{x}_0 and the target set \mathcal{C}_t . Note that because of the projection geometry, the optimal solution \mathbf{x}^* has to stay on the decision boundary and cannot be inside the interior.

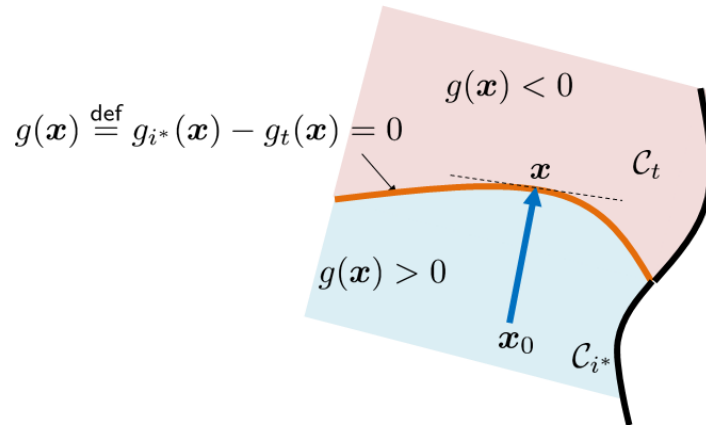


Figure 3.5: Geometry of Equation (3.3): Given an input data point \mathbf{x}_0 , our goal is to send \mathbf{x}_0 to a targeted class \mathcal{C}_t by minimizing the distance between \mathbf{x} and \mathbf{x}_0 . The decision boundary is characterized by $g(\mathbf{x}) = g_{i^*}(\mathbf{x}) - g_t(\mathbf{x})$. The optimal solution is the projection of \mathbf{x}_0 onto the decision boundary.

What if we insist of making a bigger step than needed so that \mathbf{x} will enter the interior of the target class? In theory this is possible if we make

$$\mathbf{x} = \mathbf{x}_0 + \alpha(\mathcal{P}_\Omega(\mathbf{x}_0) - \mathbf{x}_0), \quad (3.18)$$

for some step size $\alpha > 0$, where the residue vector $\mathbf{r} = \mathcal{P}_\Omega(\mathbf{x}_0) - \mathbf{x}_0$ defines the direction of perturbation. Geometrically, the new data point \mathbf{x} can be thought of moving the input \mathbf{x}_0 along the direction \mathbf{r} with certain step size α .

For complex geometry, attack based on Equation (3.18) could potentially cause problems if α is inappropriately chosen, e.g., landing on a class different from \mathcal{C}_t . Figure 3.6 illustrates two common examples. In the first example, the decision boundary is a hyperbola with decision regions \mathcal{C}_i on both sides of the space. For excessively large α , the perturbation will bypass the target class \mathcal{C}_t and land on the original class \mathcal{C}_i (though on the other side of the hyperbola). The second example shows a multi-class situation. Again, for excessively large α , we may shoot over the target class and land on a third party class.

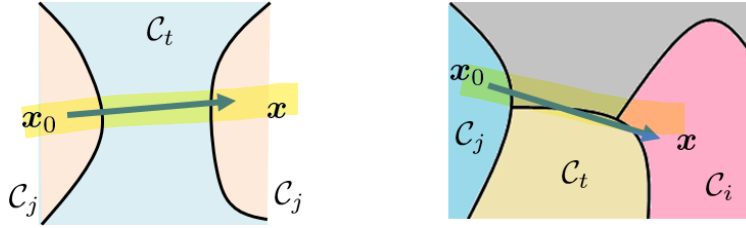


Figure 3.6: For inappropriately chosen step size α , the data point \mathbf{x}_0 can be sent to a wrong class. [Left] The decision boundary is such that \mathcal{C}_j occupies both sides of the space whereas \mathcal{C}_t is in the middle. \mathbf{x} is classified as the original class \mathcal{C}_j . [Right] In a multi-class decision, \mathbf{x}_0 is perturbed to \mathcal{C}_i . While this is not the targeted class \mathcal{C}_t , the new class \mathcal{C}_i is still a wrong class.

Targeted and Untargeted Attack

When considering adversarial attacks, there often two types of attack objectives: targeted or untargeted. By **targeted attack** we meant to move a data point \mathbf{x}_0 to the target class \mathcal{C}_t . This, however, does not require a specific class \mathcal{C}_i for \mathbf{x}_0 . Or in other words, as long as we can move \mathbf{x}_0 to \mathcal{C}_t , we have achieved the goal. We do not care about whether \mathbf{x}_0 is originally classified correctly, although in most cases we would wish \mathbf{x}_0 to be classified correctly.

Since targeted attack only concerns about the destination, the constraint is given by

$$\Omega = \{\mathbf{x} \mid \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0\}, \quad (3.19)$$

which is essentially the set of \mathbf{x} living in \mathcal{C}_t .

An **untargeted attack** is opposite to targeted attack in the sense that we like to move \mathbf{x}_0 away from its current class. That is, if $\mathbf{x}_0 \in \mathcal{C}_i$, then we want $\mathbf{x} \notin \mathcal{C}_i$. Same as targeted

attack, whether \mathbf{x}_0 is classified correctly is not a concern, although we would wish this to happen. The constraint set of an untargeted attack is given by

$$\Omega = \{\mathbf{x} \mid g_i(\mathbf{x}) - \min_{j \neq i} \{g_j(\mathbf{x})\} \leq 0\}, \quad (3.20)$$

which says that $g_i(\mathbf{x}) \leq g_j(\mathbf{x})$ for every $j \neq i$, i.e., \mathbf{x} cannot be classified as \mathcal{C}_i .

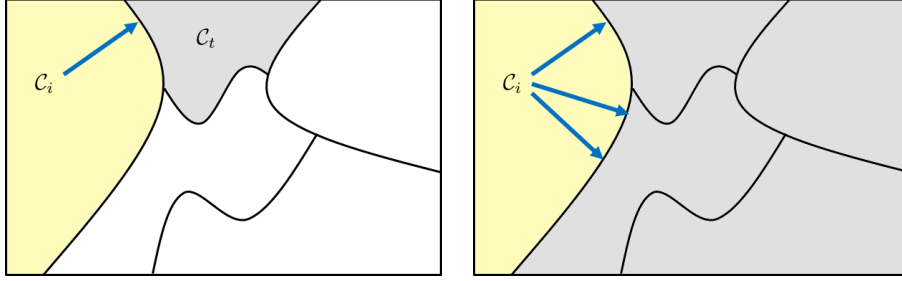


Figure 3.7: [Left] Targeted attack: The attack has to be specific from \mathcal{C}_i to \mathcal{C}_t . [Right] Untargeted attack: The attack vector can point to anywhere outside \mathcal{C}_i .

White Box and Black Box Attack

The majority of the attacks we see are **white box** attacks. These attacks assume complete knowledge about the classifier, i.e., we know exactly the discriminant functions $g_i(\mathbf{x})$ for every i . For linear classifiers, this requires knowing all the decision boundaries \mathbf{w}_i and w_{i0} . For deep neural networks, this requires knowing all the weight values and network structure. Therefore, white box attack can be considered as the best possible attack we can ever design, or the worst scenario that a classifier should face.

The opposite to white box attacks is the **black box** attacks. However, if we know absolutely nothing about the classifier then it would be nonsense to ask for a meaningful attack except using random noise. Therefore, by black box we meant that we are able to probe the classifier for a fixed number of trials, say M times. In each of the m -th trial, we input a sample $\mathbf{x}^{(m)}$ and record its classification result. In this case, the constraint set becomes

$$\Omega = \{\mathbf{x} \mid \max_{j \neq t} \{\hat{g}_j(\mathbf{x})\} - \hat{g}_t(\mathbf{x}) \leq 0\}, \quad (3.21)$$

where \hat{g}_i has been evaluated at $\{\hat{g}_i(\mathbf{x}^{(1)}), \hat{g}_i(\mathbf{x}^{(2)}), \dots, \hat{g}_i(\mathbf{x}^{(M)})\}$.

3.3 Generating Attacks

A. Minimum-Norm Attack

In order to gain insight of the adversarial attack, it is useful to consider a simple linear classifier with only two classes. A linear classifier has a discriminant function

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}, \quad (3.22)$$

and therefore by defining $\mathbf{w} \stackrel{\text{def}}{=} \mathbf{w}_i - \mathbf{w}_t$ and $w_0 \stackrel{\text{def}}{=} w_{i0} - w_{t0}$, we can show that

$$g_i(\mathbf{x}) - g_t(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0. \quad (3.23)$$

Recalling Equation (3.3), the adversarial attack can be formulated as

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \|\mathbf{x} - \mathbf{x}_0\|^2 \\ & \text{subject to} && g_i(\mathbf{x}) - g_t(\mathbf{x}) = 0. \end{aligned} \quad (3.24)$$

In this optimization, we simplify $\max_{j \neq t} \{g_j(\mathbf{x})\}$ to $g_i(\mathbf{x})$ because there are only two classes in our problem. If the index is not t , then it has to be i . We also replaced the inequality $g_i(\mathbf{x}) - g_t(\mathbf{x}) \leq 0$ as an equality, because the projection theorem in Equation (3.17) suggests that the solution must be on the decision boundary.

Theorem 2 (Minimum ℓ_2 Norm Attack for Two-Class Linear Classifier). *The adversarial attack to a two-class linear classifier is the solution of*

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 \quad \text{subject to} \quad \mathbf{w}^T \mathbf{x} + w_0 = 0, \quad (3.25)$$

which is given by

$$\mathbf{x}^* = \mathbf{x}_0 - \left(\frac{\mathbf{w}^T \mathbf{x}_0 + w_0}{\|\mathbf{w}\|_2} \right) \frac{\mathbf{w}}{\|\mathbf{w}\|_2}. \quad (3.26)$$

Proof. The Lagrange multiplier of the constrained optimization is given by

$$\mathcal{L}(\mathbf{x}, \lambda) = \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda(\mathbf{w}^T \mathbf{x} + w_0).$$

The solution of the optimization is the saddle point $(\mathbf{x}^*, \lambda^*)$ such that $\nabla_{\mathbf{x}} \mathcal{L} = 0$ and $\nabla_{\lambda} \mathcal{L} = 0$.

Taking derivative with respect to \mathbf{x} and λ yields

$$\begin{aligned} \nabla_{\mathbf{x}} \mathcal{L} &= \mathbf{x} - \mathbf{x}_0 + \lambda \mathbf{w} = 0, \\ \nabla_{\lambda} \mathcal{L} &= \mathbf{w}^T \mathbf{x} + w_0 = 0. \end{aligned}$$

Multiplying the first equation by \mathbf{w}^T yields

$$0 = \mathbf{w}^T \mathbf{x} - \mathbf{w}^T \mathbf{x}_0 + \lambda \|\mathbf{w}\|^2 = -w_0 - \mathbf{w}^T \mathbf{x}_0 + \lambda \|\mathbf{w}\|^2. \quad (3.27)$$

Thus, the optimal λ is

$$\lambda^* = (\mathbf{w}^T \mathbf{x}_0 + w_0) / \|\mathbf{w}\|^2. \quad (3.28)$$

Correspondingly, the optimal \mathbf{x} is

$$\mathbf{x}^* = \mathbf{x}_0 - \lambda^* \mathbf{w} = \mathbf{x}_0 - \left(\frac{\mathbf{w}^T \mathbf{x}_0 + w_0}{\|\mathbf{w}\|^2} \right) \mathbf{w}.$$

□

The result of this theorem provides many useful interpretation. First, the search direction is $-\mathbf{w}$ (or $-\mathbf{w}/\|\mathbf{w}\|_2$ for unit norm.) The negative sign is there because $\mathbf{w} = \mathbf{w}_i - \mathbf{w}_t$ is pointing from class \mathcal{C}_t to \mathcal{C}_i , which is opposite to the desired search direction. The search step $\frac{\mathbf{w}^T \mathbf{x}_0 + w_0}{\|\mathbf{w}\|_2}$ measures the distance from the input \mathbf{x}_0 to the target class \mathcal{C}_t . This is coherent to the projection perspective we presented earlier.

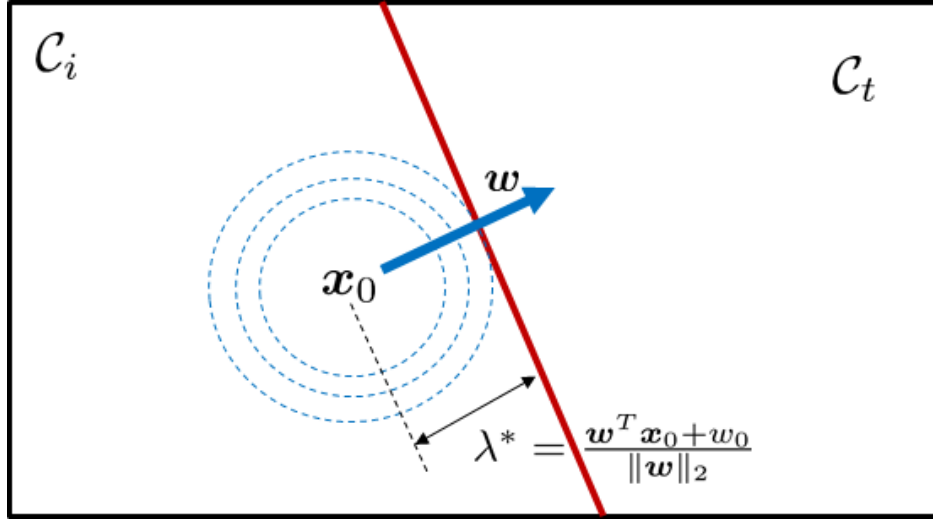


Figure 3.8: Geometry of minimum-norm attack for a two-class linear classifier with objective function $\|\mathbf{x} - \mathbf{x}_0\|^2$. The solution is a projection of the input \mathbf{x}_0 onto the separating hyperplane of the classifier.

DeepFool Attack by Moosavi-Dezfooli et al. 2016

The DeepFool attack by Moosavi-Dezfooli, Fawzi and Frossard (CVPR 2016) is a **generalization of the minimum ℓ_2 -norm attack** under general decision boundaries. Instead of having an affine constraint $\mathbf{w}^T \mathbf{x} + w_0 = 0$, the constraint is given by a nonlinear function $g(\mathbf{x}) = 0$. For a two-class problem, DeepFool formulates the optimization as follows.

Definition 6 (DeepFool Attack by Moosavi-Dezfooli et al. 2016). *The DeepFool attack for a two-class classification generates the attack by solving the optimization*

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 \quad \text{subject to} \quad g(\mathbf{x}) = 0, \quad (3.29)$$

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

Due to the nonlinearity of $g(\mathbf{x})$, it is generally very difficult to derive a closed-form solution. The numerical procedure to compute the solution can be derived by iteratively updating \mathbf{x} , where each iteration minimizes \mathbf{x} over the first order approximation of $g(\mathbf{x})$:

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

where $\mathbf{x}^{(k)}$ is the k -th iterate of the solution. In other words, we are defining a procedure

$$\mathbf{x}^{(k+1)} = \underset{\mathbf{x}}{\text{argmin}} \quad \|\mathbf{x} - \mathbf{x}^{(k)}\|^2 \quad \text{subject to} \quad g(\mathbf{x}^{(k)}) + \nabla_{\mathbf{x}} g(\mathbf{x}^{(k)})^T (\mathbf{x} - \mathbf{x}^{(k)}) = 0. \quad (3.30)$$

By identifying $\mathbf{w}^{(k)} = \nabla_{\mathbf{x}} g(\mathbf{x}^{(k)})$ and $w_0^{(k)} = g(\mathbf{x}^{(k)}) - \nabla_{\mathbf{x}} g(\mathbf{x}^{(k)})^T \mathbf{x}^{(k)}$, the linearized problem Equation (3.30) becomes

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

The solution can thus be found by using Equation (3.26), yielding

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

This result is summarized in the following corollary.

Corollary 1 (DeepFool Algorithm for Two-Class Problem). *An iterative procedure to obtain the DeepFool attack solution is*

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \underset{\mathbf{x}}{\text{argmin}} \quad \|\mathbf{x} - \mathbf{x}^{(k)}\|^2 \quad \text{subject to} \quad g(\mathbf{x}^{(k)}) + \nabla_{\mathbf{x}} g(\mathbf{x}^{(k)})^T (\mathbf{x} - \mathbf{x}^{(k)}) = 0 \\ &= \mathbf{x}^{(k)} - \left(\frac{g(\mathbf{x}^{(k)})}{\|\nabla_{\mathbf{x}} g(\mathbf{x}^{(k)})\|^2} \right) \nabla_{\mathbf{x}} g(\mathbf{x}^{(k)}). \end{aligned} \quad (3.31)$$

For multi-class problems, DeepFool has a simple solution if the classifier is **one-vs-all**. Here, by one-vs-all we meant that the classifier decides the class by considering n two-class problems, where n is the number of classes which is also the number of discriminant functions. Interested readers can consult the paper for details. However, we should remark that the one-vs-all approach does not apply to a **linear machine**, because one-vs-all essentially is handling a sequence of separating hyperplanes, whereas a linear machine considers a polytope.

B. Maximum-Allowable Attack

The minimum-norm attack using $\|\mathbf{x} - \mathbf{x}_0\|^2$ is a projection. What about maximum allowable attack? Before we answer this question, let us consider one more example of the minimum-norm attack but this time using $\|\mathbf{x} - \mathbf{x}_0\|_\infty$. That is, consider the optimization problem

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

Let $\mathbf{r} = \mathbf{x} - \mathbf{x}_0$, we can rewrite the problem as

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

where $b_0 = -(\mathbf{w}^T \mathbf{x}_0 + w_0)$. What is the solution to this optimization? One very useful technique in handling minimax optimization is the Holder's inequality.

Theorem 3 (Holder's Inequality). Let $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{y} \in \mathbb{R}^d$. Then,

$$-\|\mathbf{x}\|_p \|\mathbf{y}\|_q \leq \mathbf{x}^T \mathbf{y} \leq \|\mathbf{x}\|_p \|\mathbf{y}\|_q \quad (3.32)$$

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

Holder's inequality can be considered as a generalization of the Cauchy-Schwarz inequality, which states that $|\mathbf{x}^T \mathbf{y}| \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2$, i.e., $p = q = 2$.

If we let $p = 1$ and $q = \infty$, then we can show that $|\mathbf{x}^T \mathbf{y}| \leq \|\mathbf{x}\|_1 \|\mathbf{y}\|_\infty$. This then implies that

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

In other words, the objective function $\|\mathbf{r}\|_\infty$ is lower bounded by $\frac{|b_0|}{\|\mathbf{w}\|_1}$, which is a constant independent of \mathbf{r} . If this lower bound can be attained, then $\|\mathbf{r}\|_\infty$ is minimized. So for what \mathbf{r} can the lower bound be attained? We claim that the vector is $\mathbf{r} = \eta \cdot \text{sign}(\mathbf{w})$ for some constant η to be determined. For this choice of \mathbf{r} , we can show that

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

Therefore, if we let $\eta = b_0/\|\mathbf{w}\|_1$ (which is a constant independent of \mathbf{r}), then we will have $\|\mathbf{r}\|_\infty = |\eta| = \frac{|b_0|}{\|\mathbf{w}\|_1}$, i.e., the lower bound is attained. To summarize, we have the following theorem.

Theorem 4 (Minimum ℓ_∞ Norm Attack for Two-Class Linear Classifier). *The minimum ℓ_∞ norm attack for a two-class linear classifier, i.e.,*

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|_\infty \quad \text{subject to} \quad \mathbf{w}^T \mathbf{x} + w_0 = 0 \quad (3.33)$$

is given by

$$\mathbf{x} = \mathbf{x}_0 - \left(\frac{\mathbf{w}^T \mathbf{x}_0 + w_0}{\|\mathbf{w}\|_1} \right) \cdot \text{sign}(\mathbf{w}), \quad (3.34)$$

this direction is not perpendicular to line. it is in direction of vertex of square contours
vector of 1 or -1

Comparing the ℓ_2 attack and the ℓ_∞ attack, ℓ_2 attack gives the minimum distance attack as it is the projection of \mathbf{x}_0 onto the separating hyperplane. In contrast, the ℓ_∞ attack does not achieve the minimum distance but the attack direction is basically a vector of +1 and -1 (due to the sign operation). In fact, as shown in Figure 3.9, the attacked data \mathbf{x} must be one of the vertices of the square norm ball defined by the ℓ_∞ norm, which by construction is vector containing only +1 or -1.

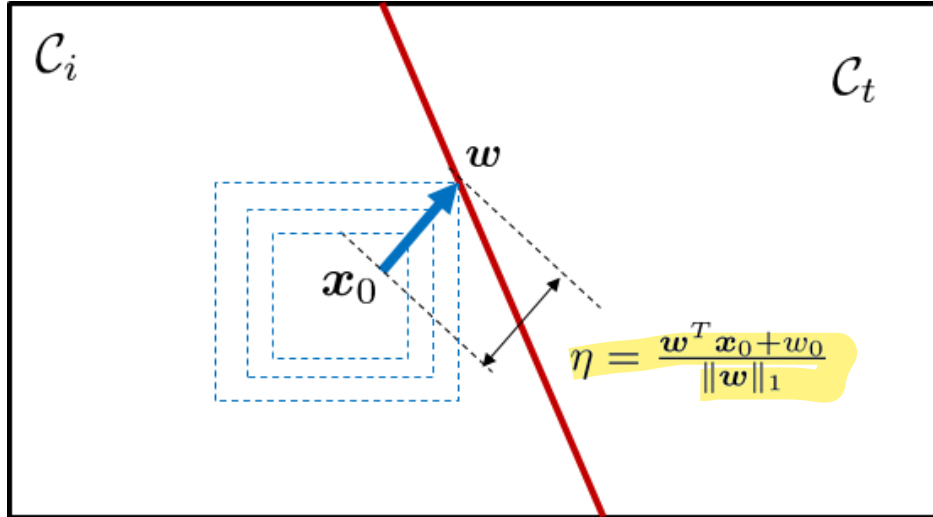


Figure 3.9: Geometry of minimum ℓ_∞ norm attack for a two-class linear classifier with objective function $\|\mathbf{x} - \mathbf{x}_0\|_\infty$. Because of the square geometry of the objective function, the perturbed data \mathbf{x} must be the vertex of the square. Correspondingly, the attack direction is $\text{sign}(\mathbf{w})$.

How is minimum ℓ_∞ norm attack related to the maximum allowable ℓ_∞ norm attack? Recall that the maximum allowable attack is given by

$$\underset{\mathbf{x}}{\text{minimize}} \quad \mathbf{w}^T \mathbf{x} + w_0 \quad \text{subject to} \quad \|\mathbf{x} - \mathbf{x}_0\|_\infty \leq \eta. \quad (3.35)$$

By defining $\mathbf{x} = \mathbf{x}_0 + \mathbf{r}$, we have $\mathbf{w}^T \mathbf{x} + w_0 = \mathbf{w}^T \mathbf{x}_0 + \mathbf{w}^T \mathbf{r} + w_0$. By defining $b_0 = -(\mathbf{w}^T \mathbf{x}_0 + w_0)$, the optimization can be rewritten as

$$\underset{\mathbf{r}}{\text{minimize}} \quad \mathbf{w}^T \mathbf{r} - b_0 \quad \text{subject to} \quad \|\mathbf{r}\|_\infty \leq \eta. \quad (3.36)$$

Using the Holder's inequality (the negative side), we can show that

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

The lower bound of $\mathbf{w}^T \mathbf{r}$ is attained when $\mathbf{r} = -\eta \cdot \text{sign}(\mathbf{w})$, because

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

Therefore, the attacked data point \mathbf{x} is $\mathbf{x} = \mathbf{x}_0 - \eta \cdot \text{sign}(\mathbf{w})$. This leads to the following theorem.

Theorem 5 (Maximum Allowable ℓ_∞ Norm Attack of Two-Class Linear Classifier). *The maximum allowable ℓ_∞ norm attack for a two-class linear classifier, i.e.,*

$$\underset{\mathbf{x}}{\text{minimize}} \quad \mathbf{w}^T \mathbf{x} + w_0 \quad \text{subject to} \quad \|\mathbf{x} - \mathbf{x}_0\|_\infty \leq \eta. \quad (3.37)$$

is given by

$$\mathbf{x} = \mathbf{x}_0 - \eta \cdot \text{sign}(\mathbf{w}). \quad (3.38)$$

The result of Equation (3.38) is exactly the same as Equation (3.34), except that the allowable attack magnitude η is unspecified. This should not be a surprise, as minimum norm attack does not involve internal parameters whereas maximum allowable attack requires η .

Fast Gradient Sign Method (FGSM)

The result of Equation (3.38) can also be considered as a special case of the **Fast Gradient Sign Method** (FGSM) by Goodfellow et al. (NIPS 2014). The FGSM was originally designed for attacking deep neural networks. The idea is to maximize certain loss function $J(\mathbf{x}; \mathbf{w})$, subject to an upper bound on the perturbation, e.g., $\|\mathbf{x} - \mathbf{x}_0\|_\infty \leq \eta$, in the original FGSM paper by Goodfellow et al. The loss function $J(\mathbf{x}; \mathbf{w})$ can be treated as the loss incurred by evaluating a data point \mathbf{x} to a classifier parameterized by \mathbf{w} . Therefore, by maximizing the loss we aim to make the classifier perform as worse as possible. The definition of the FGSM for perturbing a data point \mathbf{x}_0 with respect to a loss function $J(\mathbf{x}; \mathbf{w})$ is stated as follows.

Definition 7 (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*

$$\mathbf{x} = \mathbf{x}_0 + \eta \cdot \text{sign}(\nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})). \quad (3.39)$$

Note that in this definition, the gradient is taken with respect to the input \mathbf{x} , not the parameter \mathbf{w} . Therefore, $\nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})$ should be interpreted as the gradient of J with respect to \mathbf{x} evaluated at \mathbf{x}_0 . The definition of FGSM also explains its name: It is the gradient of the loss function, and because of the ℓ_∞ bound on the perturbation magnitude, the perturbation direction is the sign of the gradient. It is a fast method, because the gradient is typically not difficult to compute even for deep neural networks (by means of back propagation, which is available in most programming platforms.)

How do we specify the loss function? For most of the modern classifiers the loss function is also the training loss. But for simple toy problems, it is possible to explicitly derive the expression. Below is an example for the two-class linear classifier.

Example. (FGSM Loss Function for a Two-Class Linear Classifier). Recall that the objective function in the maximum allowable attack is

$$\varphi(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = (\mathbf{w}_i^T \mathbf{x} + w_{i0}) - (\mathbf{w}_t^T \mathbf{x} + w_{t0}),$$

which is a measure of how far is \mathcal{C}_i from \mathcal{C}_t . If we can make $(\mathbf{w}_i^T \mathbf{x} + w_{i0}) - (\mathbf{w}_t^T \mathbf{x} + w_{t0}) \leq 0$, then we are guaranteed that \mathbf{x} is misclassified as \mathcal{C}_t . Putting this in another way, we can also define a loss function $J(\mathbf{x}) = (\mathbf{w}_t^T \mathbf{x} + w_{t0}) - (\mathbf{w}_i^T \mathbf{x} + w_{i0})$. If $J(\mathbf{x}) \geq 0$, then \mathbf{x} is misclassified and so we have a loss. By maximizing $J(\mathbf{x})$, we are guaranteed to make the classifier perform worse. Therefore, the loss function is

$$J(\mathbf{x}; \mathbf{w}) = -(\mathbf{w}^T \mathbf{x} + w_0). \quad (3.40)$$

Given a loss function $J(\mathbf{x}; \mathbf{w})$, is there an optimization that will result the FGSM attack? If we can answer this question, then perhaps we can generalize FGSM to a broader class of attacks. To this end, we notice that for general (possibly nonlinear) loss functions, FGSM can be interpreted by applying a first order approximation:

$$J(\mathbf{x}; \mathbf{w}) = J(\mathbf{x}_0 + \mathbf{r}; \mathbf{w}) \approx J(\mathbf{x}_0; \mathbf{w}) + \nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})^T \mathbf{r}. \quad (3.41)$$

Therefore, finding \mathbf{r} to maximize $J(\mathbf{x}_0 + \mathbf{r}; \mathbf{w})$ is approximately equivalent to finding \mathbf{r} which maximizes $J(\mathbf{x}_0; \mathbf{w}) + \nabla_{\mathbf{x}} J(\mathbf{x}_0; \mathbf{w})^T \mathbf{r}$. Hence, FGSM is the solution to

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

which can be simplified to

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

where we flipped the maximization to minimization by putting a negative sign to the gradient. Interestingly, this optimization is exactly the same as Equation (3.36) with \mathbf{w} and w_0 identified as above. Therefore, the solution is

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

and hence the perturbed data is $\mathbf{x} = \mathbf{x}_0 + \eta \cdot \text{sign}(\nabla_{\mathbf{x}}J(\mathbf{x}_0; \mathbf{w}))$. In other words, FGSM is the solution of the following maximum allowable attack problem.

Corollary 2 (FGSM as a Maximum Allowable 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

$$\mathbf{x} = \mathbf{x}_0 + \eta \cdot \text{sign}(\nabla_{\mathbf{x}}J(\mathbf{x}_0; \mathbf{w})). \quad (3.42)$$

Knowing that FGSM corresponds to a maximum allowable attack with ℓ_{∞} norm, we can easily generalize the attack to other ℓ_p norms.

Corollary 3 (FGSM with ℓ_2 Norm). *The FGSM attack with ℓ_2 -norm bound on the perturbation magnitude is*

$$\underset{\mathbf{r}}{\text{maximize}} \quad \nabla_{\mathbf{x}}J(\mathbf{x}_0; \mathbf{w})^T \mathbf{r} + J(\mathbf{x}_0; \mathbf{w}) \quad \text{subject to} \quad \|\mathbf{r}\|_2 \leq \eta,$$

of which the solution is given by

$$\mathbf{x} = \mathbf{x}_0 + \eta \cdot \frac{\nabla_{\mathbf{x}}J(\mathbf{x}_0; \mathbf{w})}{\|\nabla_{\mathbf{x}}J(\mathbf{x}_0; \mathbf{w})\|_2}. \quad (3.43)$$

Proof. The optimization problem can be rewritten as

$$\underset{\mathbf{r}}{\text{minimize}} \quad \mathbf{w}^T \mathbf{r} + w_0 \quad \text{subject to} \quad \|\mathbf{r}\|_2 \leq \eta,$$

where $\mathbf{w} = -\nabla_{\mathbf{x}}J(\mathbf{x}_0; \mathbf{w})$ and $w_0 = -J(\mathbf{x}_0; \mathbf{w})$. By Cauchy inequality, we have that $\mathbf{w}^T \mathbf{r} \geq -\|\mathbf{w}\|_2 \|\mathbf{r}\|_2 \geq -\eta \|\mathbf{w}\|_2$, where the second inequality holds because $\|\mathbf{r}\|_2 \leq \eta$. The lower bound of $\mathbf{w}^T \mathbf{r} \geq -\eta \|\mathbf{w}\|_2$ is attained when $\mathbf{r} = -\eta \mathbf{w} / \|\mathbf{w}\|_2$. Therefore, the perturbation is $\mathbf{x} = \mathbf{x}_0 + \eta \nabla_{\mathbf{x}}J(\mathbf{x}_0; \mathbf{w}) / \|\nabla_{\mathbf{x}}J(\mathbf{x}_0; \mathbf{w})\|_2$. \square

Iterative Gradient Sign Method by Kurakin et al. 2017

One obvious problem of the FGSM (in fact, all other attack strategies) is the potential unboundedness of the perturbed data \mathbf{x} . For many image classification problems, the unboundedness of \mathbf{x} is physically prohibited, e.g., a pixel can only take values between $[0, 255]$ or $[0, 1]$ for normalized intensity. The **Iterative Fast Gradient Sign Method** (I-FGSM) by Kurakin, Goodfellow and Bengio (ICLR 2017) addresses this problem by posing an bound constraint on \mathbf{x} . Taking FGSM as an example, this changes the optimization to

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

In addition to the bound constraint, **I-FGSM uses an iterative linearization** rather than the one-shot linearization in FGSM. Therefore, the algorithm of computing the attack becomes

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

By dropping the terms that are independent of \mathbf{x} , the optimization can be simplified to

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

where $\mathcal{P}_{[0,1]}(\mathbf{x})$ is a projection operator that elementwisely projects out of bound values to the bound $0 \leq \mathbf{x} \leq 1$.

Corollary 4 (I-FGSM Algorithm as Projected FGSM). *The Iterative FGSM by Kurakin et al. generates the attack by iteratively solving*

$$\mathbf{x}^{(k+1)} = \underset{0 \leq \mathbf{x} \leq 1}{\text{argmax}} \quad \nabla_{\mathbf{x}} J(\mathbf{x}^{(k)}; \mathbf{w})^T \mathbf{x} \quad \text{subject to} \quad \|\mathbf{x} - \mathbf{x}^{(k)}\|_{\infty} \leq \eta, \quad (3.44)$$

of which the per-iteration solution is given by

$$\mathbf{x}^{(k+1)} = \mathcal{P}_{[0,1]} \left\{ \mathbf{x}^{(k)} + \eta \cdot \text{sign}(\nabla_{\mathbf{x}} J(\mathbf{x}^{(k)}; \mathbf{w})) \right\}. \quad (3.45)$$

Remark. (Choice of η for I-FGSM.) Practically, when generating attack using I-FGSM, the magnitude upper bound η is usually reduced to a smaller constant so that the number of active attack pixels is limited.

C. Regularization-based Attack

In both minimum norm attack and maximum allowable attack, the optimizations are constrained. For certain simple cases, e.g., binary linear classifiers, closed-form solutions can be derived. However, for advanced classifiers such as deep neural networks, solving an optimization involving constraints are typically very difficult. The regularization-based attack alleviates the problem by considering

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

For two-class linear classifier, this is simplified to

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

Since the regularization-based attack is an unconstrained version of the minimum norm attack and the maximum allowable attack, one should expect to have a very similar solution.

Theorem 6 (Regularization-based Attack for Two-Class Linear Classifier). *The regularization-based attack for a two-class linear classifier generates the attack by solving*

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda(\mathbf{w}^T \mathbf{x} + w_0), \quad (3.46)$$

of which the solution is given by

$$\mathbf{x} = \mathbf{x}_0 - \lambda \mathbf{w}. \quad (3.47)$$

Proof. Let $\varphi(\mathbf{x}) = \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda(\mathbf{w}^T \mathbf{x} + w_0)$. Taking the first order derivative and sending to zero yields

$$0 = \nabla \varphi(\mathbf{x}) = (\mathbf{x} - \mathbf{x}_0) + \lambda \mathbf{w}. \quad (3.48)$$

Therefore, we have $\mathbf{x} = \mathbf{x}_0 - \lambda \mathbf{w}$. □

In a regularization based attack, the attack magnitude λ is specified by the user. It is not difficult to see that for linear classifiers, the optimal λ and η are both $\lambda = \eta = \frac{\mathbf{w}^T \mathbf{x}_0 + w_0}{\|\mathbf{w}\|^2}$. However, in minimum norm attack, this magnitude is automatically guaranteed by the algorithm. Maximum allowable attack is slightly more complicated, but it is still more interpretable than regularization based attack as defining η is equivalent to controlling the maximum perturbation.

Another drawback of the regularization based attack is that the objective value could potentially be unbounded. See the example below.

Example. (Unboundedness of Regularization-based Attack). Consider an ℓ_1 -norm 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 stated by the subgradient:

$$0 \in \partial \|\mathbf{r}\|_1 + \lambda \mathbf{w}, \quad \implies \quad 0 \in \text{sign}(\mathbf{r}) + \lambda \mathbf{w},$$

where $\text{sign}(r_i) = [-1, 1]$ if $r_i = 0$, $\text{sign}(r_i) = \{+1\}$ if $r_i > 0$, and $\text{sign}(r_i) = \{-1\}$ if $r_i < 0$. Therefore, the optimality condition holds only when $\lambda w_i \in [-1, 1]$, i.e., $|\lambda \mathbf{w}| \leq \mathbf{1}$. When this happens, we must have $\mathbf{r} = \mathbf{0}$ and so the optimal value is 0. However, if $|\lambda \mathbf{w}| > \mathbf{1}$, then it is impossible to have $\mathbf{0} \in \text{sign}(\mathbf{r}) + \lambda \mathbf{w}$ regardless of how we choose \mathbf{r} . This means that the objective function is unbounded below. In fact, one can show that the 1D function $f(x) = |x| + \alpha x \rightarrow -\infty$ as $x \rightarrow -\infty$ if $\alpha > 1$, or as $x \rightarrow \infty$ if $\alpha < -1$.

To summarize, while regularization based attack leads to an unconstrained optimization which is attractive to many off-the-shelf optimization packages, the parameter tuning of λ and the potential unboundedness of the problem are its major drawbacks. In addition, for problems involving nonlinear regularization terms, the optimization may be difficult as well.

Carlini-Wagner Attack

The Carlini-Wagner attack (2016) is a regularization-based attack with some critical modifications which can resolve the unboundedness issue. The main idea is to realize that for general classifiers, the optimization can be written as

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 + \iota_\Omega(\mathbf{x}), \quad (3.49)$$

where

$$\iota_\Omega(\mathbf{x}) = \begin{cases} 0, & \text{if } \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0, \\ +\infty, & \text{otherwise.} \end{cases} \quad (3.50)$$

That is, when the constraint $\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0$ is satisfied, the optimization is simplified to $\min \|\mathbf{x} - \mathbf{x}_0\|^2$ which has a trivial solution $\mathbf{x} = \mathbf{x}_0$. In other words, when \mathbf{x}_0 is already in the targeted class \mathcal{C}_t (so that $\max_{j \neq t} \{g_j(\mathbf{x}_0)\} \leq g_t(\mathbf{x}_0)$), then there is no need to perturb the data point. Of course, in most of the cases $\mathbf{x}_0 \notin \mathcal{C}_t$, and so we need to iteratively check whether the current iterate satisfies the constraint. Yet, the introduction

of $\iota_\Omega(\mathbf{x})$ allows us to rewrite the constrained problem (i.e., the minimum-norm attack) into an unconstrained problem.

Unfortunately, the unconstrained problem in Equation (3.49) remains uneasy to solve, for ι_Ω is either 0 or $+\infty$. Carlini-Wagner attack relaxes ι_Ω by considering a rectifier function

$$\zeta(x) = \max(x, 0). \quad (3.51)$$

Putting into the optimization, we have

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda \cdot \zeta(\mathbf{x}), \quad (3.52)$$

for some constant $\lambda > 0$. Intuitively, $\zeta(\cdot)$ maps any positive x to itself, and any negative x to 0. It is a useful relaxation because when the constraint is satisfied, $\zeta(x)$ will be 0. When the constraint is not satisfied, $\eta(x)$ penalizes the objective linearly.

Using the rectifier function, the attack can now be written as

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda \cdot \max \left\{ \left(\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \right), 0 \right\}, \quad (3.53)$$

which is known as the Carlini-Wagner attack. In the original paper by Carlini and Wagner, it was shown that besides the rectifier function there are many other choices one can consider. Readers interested in the details can consult their paper.

The rectifier function is particularly useful for two reasons. The first reason is the computation, because the derivative of a rectifier function is a unit step function which is computationally very easy to implement. The second reason is its effect in handling the non-boundedness issue. Below is an example.

Example. (Carlini-Wagner Attack for ℓ_1 Problem.) Let us revisit the ℓ_1 minimization but using the rectifier function this time. Thus, the optimization becomes

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

The objective function is always non-negative, as $\|\mathbf{x} - \mathbf{x}_0\|_1 \geq 0$ and $\max(\mathbf{w}^T \mathbf{x} + w_0, 0) \geq 0$. Therefore, the optimization is always bounded below by 0.

Note that the lower bound is achieved when $\mathbf{x} = \mathbf{x}_0$ and $\mathbf{w}^T \mathbf{x}_0 + w_0 = 0$. This happens only when the attack solution is $\mathbf{x} = \mathbf{x}_0$ and \mathbf{x}_0 is right on the decision boundary. Clearly, this is unlikely to happen, and if this happens the optimal objective value is 0. Therefore, besides this trivial solution, the optimization is guarantee to attain a finite minimum objective.

Remark. (Convexity of CW attack.) Let $\varphi(\mathbf{x}) = \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x})$. It is not difficult to show that the function $h(\mathbf{x}) = \max(\varphi(\mathbf{x}), 0)$ is convex if φ is convex:

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

In particular, if g_i is linear, then φ is linear and so h is convex.

CW attack is appeal because it allows us to use simple algorithms to compute the attack. To do so, we first note that 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: $i^*(\mathbf{x}) = \underset{j \neq t}{\operatorname{argmax}} \{g_j(\mathbf{x})\}$. If, for the time being, we assume that the index i^* is independent of \mathbf{x} , then the gradient for the rectifier function applied to the constraint is

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

Letting $\varphi(\mathbf{x})$ be the overall objective function, i.e.,

$$\varphi(\mathbf{x}) = \|\mathbf{x} - \mathbf{x}_0\|^2 + \lambda \cdot \max \left\{ \left(\max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \right), 0 \right\}, \quad (3.54)$$

it follows that the gradient is

$$\nabla \varphi(\mathbf{x}; i^*) = 2(\mathbf{x} - \mathbf{x}_0) + \lambda \cdot \mathbb{I}\{g_{i^*}(\mathbf{x}) - g_t(\mathbf{x}) > 0\} \cdot (\nabla g_{i^*}(\mathbf{x}) - \nabla g_t(\mathbf{x})). \quad (3.55)$$

Here, we denote the gradient as $\nabla \varphi(\mathbf{x}; i^*)$ by parametrizing using i^* , as the gradient depends on how i^* is defined.

Applying gradient descent, the solution can be found by iteratively updating the estimate.

Algorithm. (CW Attack Gradient Descent.) The Gradient Descent algorithm for generating CW attack is given by the following iteration for $k = 1, 2, \dots$

$$\begin{aligned} i^* &= \underset{j \neq t}{\operatorname{argmax}} \{g_j(\mathbf{x}^k)\} \\ \mathbf{x}^{k+1} &= \mathbf{x}^k - \alpha \nabla \varphi(\mathbf{x}^k; i^*). \end{aligned} \quad (3.56)$$

As we can see in the above algorithm, the two parameters we need to tune are:

- α : Gradient descent step size, which controls the rate of convergence. For non-convex problems, it is advised to use smaller α so that a local minimum is not missed.
- λ : Regularization parameter, which controls the relative strength between the distance term and the constraint term. The general rule of thumb is either by trial-and-error or cross-validation if training set is available.

3.4 Can Random Noise Attack?

Adversarial attack works because the perturbation is carefully designed. What if we perturb the data by pure i.i.d. Gaussian noise?

Recall that when launching attacks for a linear classifier, the perturbation always takes the form of $\mathbf{x} = \mathbf{x}_0 + \lambda \mathbf{w}$. That is, we want to move \mathbf{x}_0 along \mathbf{w} by an amount λ so that \mathbf{x} is misclassified. Now, if we do not move along \mathbf{w} but along a random vector \mathbf{r} such that

$$\mathbf{x} = \mathbf{x}_0 + \sigma_r \mathbf{r}, \quad (3.57)$$

where $\mathbf{r} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, then we can ask if we can still misclassify the data point \mathbf{x} . Clearly, this requires us to check whether $\mathbf{w}^T \mathbf{r} > 0$. If $\mathbf{w}^T \mathbf{r} > 0$, then \mathbf{r} and \mathbf{w} will form an acute angle and so for sufficient step size we will be able to move \mathbf{x}_0 to another class. If $\mathbf{w}^T \mathbf{r} < 0$, then \mathbf{w} and \mathbf{r} are form an obtuse angle and so \mathbf{r} will move \mathbf{x}_0 to an opposite direction.

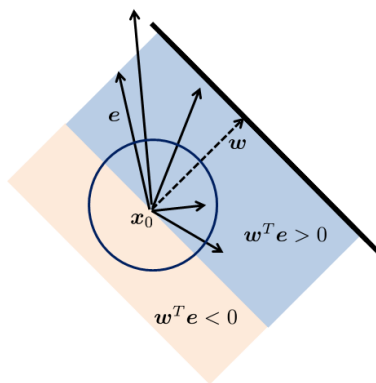


Figure 3.10: Attacking the linear classifier with i.i.d. noise is equivalent to putting an uncertainty circle around \mathbf{x}_0 with radius σ_r . The possible attack directions are those that form acute angle with the normal vector \mathbf{w} . Therefore, among all the possible \mathbf{r} 's, we are only interested in those such that $\mathbf{w}^T \mathbf{r} > 0$, which occupy half of the space in 2D.

Figure 3.10 provides a pictorial illustration of the noise attack. From this figure, we are tempted to think that i.i.d. noise can achieve 50% attack rate because $\mathbf{w}^T \mathbf{r} > 0$ occupies half of the space. The problem of such argument is that in high dimension, the probability of $\mathbf{w}^T \mathbf{r} > 0$ is diminishing very quickly as the dimensionality of \mathbf{r} grows. This is a well-known

phenomenon called **curse of dimensionality**. To illustrate the idea, let us evaluate the probability of $\mathbf{w}^T \mathbf{r} \geq \epsilon$ for some $\epsilon > 0$. To this end, let us consider

$$\mathbb{P} \left[\frac{1}{d} \mathbf{w}^T \mathbf{r} \geq \epsilon \right] = \mathbb{P} \left[\frac{1}{d} \sum_{j=1}^d w_j r_j \geq \epsilon \right],$$

where d is the dimensionality of \mathbf{w} , i.e., $\mathbf{w} \in \mathbb{R}^d$. The tolerance level ϵ is a small positive constant that stays away from 0.

Example. Before we discuss the theoretical results, it will be useful to do a quick numerical experiment. Consider a special case where $\mathbf{w} = \mathbf{1}_{d \times 1}$, i.e., a d -dimensional all-one vector, and $\mathbf{r} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. In this case, we define the average as

$$Y \stackrel{\text{def}}{=} \frac{1}{d} \sum_{j=1}^d r_j. \quad (3.58)$$

It follows that Y is a Gaussian random variable because linear combination of Gaussian remains a Gaussian. The mean and variance are

$$\mathbb{E}[Y] = 0, \quad \text{and} \quad \text{Var}[Y] = \frac{1}{d}. \quad (3.59)$$

Therefore, the probability of the event $\{Y > \epsilon\}$ is

$$\begin{aligned} \mathbb{P}[Y > \epsilon] &= \int_{\epsilon}^{\infty} \frac{1}{\sqrt{2\pi/d}} \exp \left\{ -\frac{t^2}{2/d} \right\} dt \\ &= \int_{\epsilon\sqrt{d/2}}^{\infty} \frac{1}{\sqrt{\pi}} \exp \{ -t^2 \} dt \\ &= \frac{1}{2} \text{erfc} \left(\epsilon\sqrt{d/2} \right), \end{aligned} \quad (3.60)$$

where erfc is the complementary error function defined as $\text{erfc}(z) \stackrel{\text{def}}{=} \frac{2}{\sqrt{\pi}} \int_z^{\infty} e^{-t^2} dt$. As we can see in Figure 3.11, the probability drops rapidly as d increases.

Let us discuss the general case. One classical result to derive the probability is the tail bound on Gaussian by Feller 1968.¹ If Y is a Gaussian random variable with $Y \sim \mathcal{N}(\mu, \sigma)$, then

$$\mathbb{P}[Y \geq \mu + \sigma\epsilon] \leq \frac{1}{\epsilon} \frac{e^{-\epsilon^2/2}}{\sqrt{2\pi}}. \quad (3.61)$$

¹<http://www.stat.yale.edu/~pollard/Books/Mini/Basic.pdf>

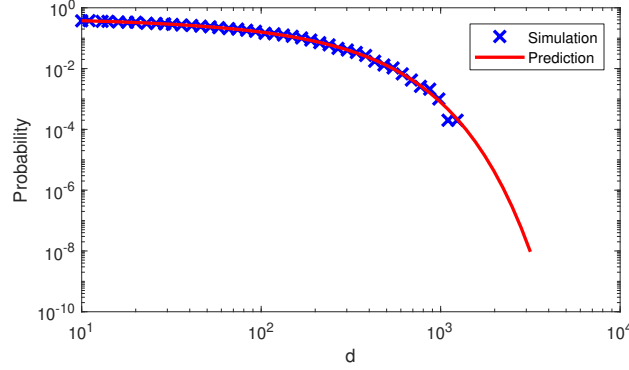


Figure 3.11: Empirical probability and the theoretically predicted value. In this experiment, we let $Y = \frac{1}{d} \sum_{j=1}^d r_j$ where $r_j \sim \mathcal{N}(0, 1)$ is an iid Gaussian random variable. We plot the probability $\mathbb{P}[Y > \epsilon]$ as a function of d . Note the tightness of the prediction, and the rapid decaying behavior of the probability. The result shows that when d grows, the value Y is concentrated at 0.

Now, let $Y = \frac{1}{d} \sum_{j=1}^d w_j r_j$. Since a linear combination of Gaussian remains a Gaussian, it holds that Y is Gaussian and

$$\mu = \mathbb{E}[Y] = 0, \quad \text{and} \quad \sigma^2 = \text{Var}[Y] = \frac{1}{d^2} \sum_{j=1}^d w_j^2 = \frac{\|\mathbf{w}\|^2}{d^2}. \quad (3.62)$$

Therefore, by substituting $\varepsilon = \sigma\epsilon$ we can show that

$$\mathbb{P} \left[\frac{1}{d} \sum_{j=1}^d w_j r_j \geq \varepsilon \right] = \mathbb{P}[Y \geq \varepsilon] \leq \frac{\sigma}{\varepsilon} \frac{e^{-\frac{\varepsilon^2}{2\sigma^2}}}{\sqrt{2\pi}} = \frac{\|\mathbf{w}\|}{\varepsilon d \sqrt{2\pi}} \exp \left\{ -d^2 \frac{\varepsilon^2}{2\|\mathbf{w}\|^2} \right\}. \quad (3.63)$$

As $d \rightarrow \infty$, it holds that $\mathbb{P} \left[\frac{1}{d} \sum_{j=1}^d w_j r_j \geq \varepsilon \right] \rightarrow 0$. That means, the probability of getting a “good attack direction” is diminishing to zero exponentially. Putting everything together we have the following theorem.

Theorem 7. Let $\mathbf{w}^T \mathbf{x} + w_0 = 0$ be the decision boundary of a linear classifier, and let $\mathbf{x}_0 \in \mathbb{R}^d$ be an input data point. Suppose we attack the classifier by adding i.i.d. Gaussian noise $\mathbf{r} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ to \mathbf{x}_0 . The probability of a successful attack at a tolerance level ε is $\mathbb{P} \left[\frac{1}{d} \sum_{j=1}^d w_j r_j \geq \varepsilon \right]$, and such probability is upper bounded by

$$\mathbb{P} \left[\frac{1}{d} \sum_{j=1}^d w_j r_j \geq \varepsilon \right] \leq \frac{\|\mathbf{w}\|}{\varepsilon d \sqrt{2\pi}} \exp \left\{ -d^2 \frac{\varepsilon^2}{2\|\mathbf{w}\|^2} \right\}. \quad (3.64)$$

Therefore, as $d \rightarrow \infty$ it becomes increasingly more difficult for i.i.d. Gaussian noise to succeed in attacking.

Why does the probability goes to zero so quickly? If we look at the inner product, we realize that

$$\mathbf{w}^T \mathbf{r} = \sum_{j=1}^d w_j r_j \quad (3.65)$$

follows Central Limit Theorem to a Gaussian. As d increases, the random estimate $\sum_{j=1}^d w_j r_j$ approaches its expectation (which is zero in this case because $\mathbb{E}[r_i] = 0$ for all i) exponentially according to the Gaussian tail probability. This can be seen from the fact that there are positive and negative values of r_i 's in the sum. The more r_i 's we have, the more likely that the terms will cancel out each other.

The concentration inequality also suggests a surprising picture in high-dimensional space. Instead of having a half space, the event $\left| \frac{1}{d} \sum_{j=1}^d w_j r_j \right| < \epsilon$ is actually concentrated at the equator of a high-dimensional sphere. See Figure 3.12 for an illustration. Therefore, as d increases, there is less and less probability that we can ever find a point NOT on the equator.

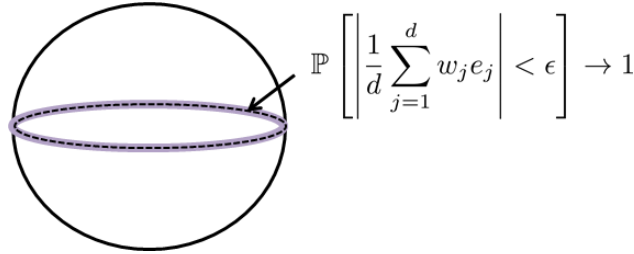


Figure 3.12: The magnitude of i.i.d. noise attack can be determined by the cosine angle, $\cos \theta = \mathbf{w}^T \mathbf{r} / \|\mathbf{w}\|_2 \|\mathbf{r}\|_2$, which is equivalent to λ^* / σ_e .

What about the attack magnitude σ_e ? Let $\cos \theta$ be the angle between \mathbf{w} and \mathbf{r} . That is,

$$\frac{\mathbf{w}^T \mathbf{r}}{\|\mathbf{w}\|_2 \|\mathbf{r}\|_2} = \cos \theta. \quad (3.66)$$

Since the shortest distance between \mathbf{x}_0 and the decision boundary $\mathbf{w}^T \mathbf{x} + w_0 = 0$ is λ^* , it follows that the ratio between λ^* and σ_e , the distance from \mathbf{x}_0 to the decision boundary along the direction \mathbf{r} , is given by

$$\cos \theta = \frac{\lambda^*}{\sigma_e}. \quad (3.67)$$

Equating the results yields the following corollary.

Corollary 5. Let $\mathbf{r} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ be a random attack direction, and \mathbf{w} be the linear discriminant function. The search step σ_r to move an input \mathbf{x}_0 to the decision boundary along the direction \mathbf{r} is

$$\sigma_r = \frac{\lambda^* \|\mathbf{w}\|_2 \|\mathbf{r}\|_2}{\mathbf{w}^T \mathbf{r}}, \quad (3.68)$$

where λ^* defined in Equation (3.28) is the minimum distance from \mathbf{x}_0 to the decision boundary.

Intuitively, Equation (3.68) provides a scaling to the shortest distance λ^* . The scaling constant $\frac{\|\mathbf{w}\|_2 \|\mathbf{r}\|_2}{\mathbf{w}^T \mathbf{r}}$ is the reciprocal of the cosine angle. By Cauchy inequality we can show a lower bound on the ratio: $\frac{\|\mathbf{w}\|_2 \|\mathbf{r}\|_2}{\mathbf{w}^T \mathbf{r}} \geq 1$. Therefore, if \mathbf{r} aligns well with \mathbf{w} then $\frac{\|\mathbf{w}\|_2 \|\mathbf{r}\|_2}{\mathbf{w}^T \mathbf{r}}$ will be close to 1. In contrast, if \mathbf{r} is almost perpendicular to \mathbf{w} , then $\frac{\|\mathbf{w}\|_2 \|\mathbf{r}\|_2}{\mathbf{w}^T \mathbf{r}}$ can become big.

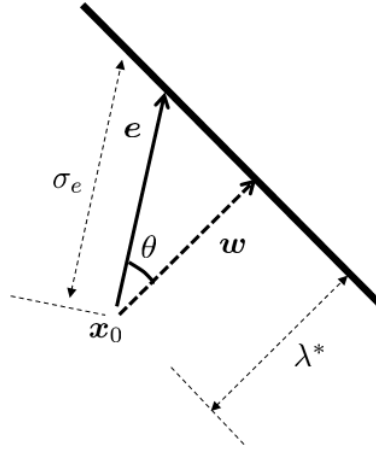


Figure 3.13: The magnitude of i.i.d. noise attack can be determined by the cosine angle, $\cos \theta = \mathbf{w}^T \mathbf{r} / \|\mathbf{w}\|_2 \|\mathbf{r}\|_2$, which is equivalent to λ^* / σ_r .

The concentration inequality predicts that for large d , the probability for $\mathbf{w}^T \mathbf{r}$ to stay away from 0 is exponentially vanishing. This suggests that as $d \rightarrow \infty$, $\cos \theta \propto \mathbf{w}^T \mathbf{r} \rightarrow 0$ exponentially. Therefore, in order to use i.i.d. noise as an attack, the magnitude of the attack σ_r has to grow exponentially with d .