

Problem Set 9, Nov 13, 2025

(Adversarial Robustness)

Security and robustness of machine learning models remain critical concerns in modern AI systems. Despite significant advances in model performance, neural networks are still vulnerable to adversarial attacks—carefully crafted perturbations that can cause models to misclassify inputs (e.g., images) while remaining imperceptible to humans. Understanding and defending against such attacks is essential both for deploying ML systems in real-world applications and for scientific understanding.

Goals. The goals of this exercise are to:

- Understand how to generate adversarial examples in practice.
- Use adversarial examples in adversarial training to obtain more robust models.
- Understand what adversarial examples correspond to in the simple case of linear models.

Problem 1 (Adversarial training for linear models):

It can be often very insightful to analyze what a method corresponds to in a simple setting of linear models.

Assume we have input points $\mathbf{x}_i \in \mathbb{R}^d$ and binary labels $y_i \in \{-1, 1\}$. Let ℓ be a monotonically decreasing margin-based loss function, for example the hinge loss $\ell(z) = \max\{0, 1 - z\}$ or logistic loss $\ell(z) = \log(1 + \exp(-z))$ that you have seen before.

Consider the adversarial training objective for a linear model $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$ with respect to ℓ_2 adversarial perturbations:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \max_{\|\hat{\mathbf{x}}_i - \mathbf{x}_i\|_2 \leq \varepsilon} \ell(y_i \cdot \mathbf{w}^\top \hat{\mathbf{x}}_i).$$

- Find a closed-form solution of the inner maximization problem $\max_{\|\hat{\mathbf{x}}_i - \mathbf{x}_i\|_2 \leq \varepsilon} \ell(y_i \cdot \mathbf{w}^\top \hat{\mathbf{x}}_i)$ and the minimizer $\hat{\mathbf{x}}_i^*$.

Solution:

$$\begin{aligned} \max_{\|\hat{\mathbf{x}}_i - \mathbf{x}_i\|_2 \leq \varepsilon} \ell(y_i \cdot \mathbf{w}^\top \hat{\mathbf{x}}_i) &= \max_{\|\delta\|_2 \leq \varepsilon} \ell(y_i \cdot \mathbf{w}^\top (\mathbf{x}_i + \delta)) = \ell\left(\min_{\|\delta\|_2 \leq \varepsilon} y_i \cdot \mathbf{w}^\top (\mathbf{x}_i + \delta)\right) \\ &= \ell\left(y_i \cdot \mathbf{w}^\top \mathbf{x}_i + \min_{\|\delta\|_2 \leq \varepsilon} y_i \mathbf{w}^\top \delta\right) = \ell\left(y_i \cdot \mathbf{w}^\top \mathbf{x}_i - \varepsilon \|\mathbf{w}\|_2\right) \end{aligned}$$

where we used in the second equality the fact that the loss is monotonically decreasing in its margin, and thus flips the max to the min. Moreover, in the last equality we used the Cauchy-Schwartz inequality to solve the inner minimization problem of a linear function over the ℓ_2 -ball:

$$y_i \mathbf{w}^\top \delta \geq -\|\mathbf{w}\|_2 \|\delta\|_2 = -\|\mathbf{w}\|_2 \|\delta\|_2 \geq -\|\mathbf{w}\|_2 \varepsilon.$$

And then the equalities are attained for $\delta^* = -y_i \varepsilon \frac{\mathbf{w}}{\|\mathbf{w}\|_2}$, i.e. δ^* is colinear with \mathbf{w} and scaled by an appropriate constant. Thus, δ^* is a minimizer of $\min_{\|\delta\|_2 \leq \varepsilon} y_i \mathbf{w}^\top \delta$. And equivalently we have that $\hat{\mathbf{x}}_i^* = \mathbf{x}_i - y_i \varepsilon \frac{\mathbf{w}}{\|\mathbf{w}\|_2}$.

This makes the overall adversarial training have the following form:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(y_i \cdot \mathbf{w}^\top \mathbf{x}_i - \varepsilon \|\mathbf{w}\|_2).$$

Note that adversarial training for linear models boils down to a convex optimization problem with respect to the weights \mathbf{w} , so it can be efficiently solved. This is unlike for neural networks where both inner maximization and outer minimization problems are computationally hard (but can be solved approximately, often with a lot of empirical success).

- In case of the hinge loss, $\ell(z) = \max\{0, 1 - z\}$, what is the connection between ℓ_2 adversarial training and the primal formulation of the soft-margin SVM?

Solution: For the hinge loss, we have the following adversarial training objective:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \max \{0, 1 - y_i \cdot \mathbf{w}^\top \mathbf{x} + \varepsilon \|\mathbf{w}\|_2\}.$$

And for the soft-margin SVM we have a quite similar objective:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \max \{0, 1 - y_i \cdot \mathbf{w}^\top \mathbf{x}\} + \varepsilon \|\mathbf{w}\|_2^2,$$

with the only difference that the ℓ_2 -regularization term is contained inside of the loss function and it is not squared (i.e. $\|\mathbf{w}\|_2$ instead of $\|\mathbf{w}\|_2^2$). Moreover, note that we have the following upper bound:

$$\frac{1}{n} \sum_{i=1}^n \max \{0, 1 - y_i \cdot \mathbf{w}^\top \mathbf{x} + \varepsilon \|\mathbf{w}\|_2\} \leq \frac{1}{n} \sum_{i=1}^n \max \{0, 1 - y_i \cdot \mathbf{w}^\top \mathbf{x}\} + \varepsilon \|\mathbf{w}\|_2$$

Thus, we can see that performing ℓ_2 -regularized standard training leads to a similar effect as the ℓ_2 adversarial training. In particular, a small norm of the weights \mathbf{w} is sufficient to have small adversarial loss.

- What if instead of ℓ_2 adversarial training, we performed ℓ_∞ adversarial training, how would the solution of the inner maximization problem change? Does the maximizer for ℓ_∞ -perturbations resemble the Fast Gradient Sign Method (FGSM)?

Solution: The only difference would be only in the step where we used the Cauchy-Schwartz inequality. Instead, we can use its generalization known as the *Hölder's inequality* which states that for ℓ_p -norms with $p \geq 1$ such that $\frac{1}{p} + \frac{1}{q} = 1$ we have $|\mathbf{w}^\top \delta| \leq \|\mathbf{w}\|_p \|\delta\|_q$ and the equality can be always attained for some choice of δ . Using this fact for $p = 1$ and $q = \infty$:

$$y_i \mathbf{w}^\top \delta \geq -\|\mathbf{w}\|_1 \|\delta\|_\infty = -\|\mathbf{w}\|_1 \|\delta\|_\infty \geq -\|\mathbf{w}\|_1 \varepsilon \implies \min_{\|\delta\|_\infty \leq \varepsilon} y_i \mathbf{w}^\top \delta = -\varepsilon \|\mathbf{w}\|_1.$$

Moreover, we can recover the minimizer δ^* by solving $y_i \mathbf{w}^\top \delta = -\varepsilon \|\mathbf{w}\|_1$. We get that $\sum_{j=1}^d y_i w_j \delta_j = \sum_{j=1}^d -\varepsilon |w_j|$ for which we can find a solution by equating each term in the sum. This gives us $y_i w_j \delta_j = -\varepsilon |w_j|$ which implies that $\delta_j = -y_i \varepsilon \text{ sign}(w_j)$ or in a vector form $\delta^* = -y_i \varepsilon \text{ sign}(\mathbf{w})$. This solution is optimal since it satisfies $y_i \mathbf{w}^\top \delta = -\varepsilon \|\mathbf{w}\|_1$ and is feasible $\|\delta\|_\infty$.

We also note that we can alternatively solve the optimization problem $\min_{\|\delta\|_\infty \leq \varepsilon} y_i \mathbf{w}^\top \delta$ directly, without relying on Hölder's inequality:

$$\begin{aligned} \min_{\|\delta\|_\infty \leq \varepsilon} y_i \mathbf{w}^\top \delta &= \min_{\forall j \ |\delta_j| \leq \varepsilon} \sum_{j=1}^d y_i w_j \delta_j \stackrel{(a)}{=} \sum_{j=1}^d \min_{|\delta_j| \leq \varepsilon} y_i w_j \delta_j \stackrel{(b)}{=} \sum_{j=1}^d \min_{|\delta_j| = \varepsilon} y_i w_j \delta_j \\ &= \sum_{j=1}^d \min\{-w_j \varepsilon, w_j \varepsilon\} = \sum_{j=1}^d -\varepsilon |w_j| = -\varepsilon \|\mathbf{w}\|_1, \end{aligned}$$

where in (a) we used the separability of the objective and constraints over dimensions and in (b) the fact that a minimum of a linear function is attained at a boundary of the feasible set.

Additionally, we note that for an arbitrary constrained convex optimization problem (that satisfies some mild technical conditions), we could solve the Karush-Kuhn-Tucker conditions to find the solution. In the case of $\min_{\|\delta\|_\infty \leq \varepsilon} y_i \mathbf{w}^\top \delta$, it would lead to the same closed-form solution as shown above.

Regarding the interpretation of δ^* , the choice $\delta^* = -y_i \varepsilon \text{ sign}(\mathbf{w})$ is very similar to the Fast Gradient Sign Method (FGSM), where instead of \mathbf{w} we have the negative gradient of the loss function $-\nabla_x \ell(\mathbf{x}, y)$. This

is motivated by the fact that we want to maximize the loss (i.e. to make the prediction on \mathbf{x} as dissimilar to the true label y as possible), and for this we try to maximize the *linear approximation* of the loss:

$$\max_{\|\boldsymbol{\delta}\|_\infty \leq \varepsilon} \ell(\mathbf{x} + \boldsymbol{\delta}, y) \approx \max_{\|\boldsymbol{\delta}\|_\infty \leq \varepsilon} [\ell(\mathbf{x}, y) + \nabla_x \ell(\mathbf{x}, y)^\top \boldsymbol{\delta}] = \max_{\|\boldsymbol{\delta}\|_\infty \leq \varepsilon} \nabla_x \ell(\mathbf{x}, y)^\top \boldsymbol{\delta},$$

where the maximizer is similarly given as $\boldsymbol{\delta}^* = \varepsilon \cdot \text{sign}(\nabla_x \ell(\mathbf{x}, y))$. In case of binary labels, monotonically decreasing ℓ , and a linear classifier, we get that $\boldsymbol{\delta}^* = \varepsilon \text{ sign}(y \ell'(\mathbf{y} \mathbf{w}^\top \mathbf{x}) \mathbf{w}) = \varepsilon \text{ sign}(-y \mathbf{w}) = -y \varepsilon \text{ sign}(\mathbf{w})$, i.e. we recover the expression above.