

Problem Set 10, Nov 19, 2020 (Adversarial Robustness)

Problem 2 (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:

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

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:

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

And then the equalities are attained for $\delta^* = -y_i \varepsilon \frac{\mathbf{w}}{\|\mathbf{w}\|_2^2}$, 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 $\hat{\mathbf{x}}_i^* = \mathbf{x}_i - y_i \varepsilon \frac{\mathbf{w}}{\|\mathbf{w}\|_2^2}$.

This makes the overall adversarial training have the following form:

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

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) from the previous exercise?

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 $|\mathbf{w}^\top \boldsymbol{\delta}| \leq \|\mathbf{w}\|_p \|\boldsymbol{\delta}\|_q$ and the equality can be always attained for some $\mathbf{w}, \boldsymbol{\delta}$. Using this fact for $p = 1$ and $q = \infty$:

$$\min_{\|\boldsymbol{\delta}\|_\infty \leq \varepsilon} y_i \mathbf{w}^\top \boldsymbol{\delta} = -\varepsilon \|\mathbf{w}\|_1,$$

and moreover the minimizer is $\boldsymbol{\delta}^* = -\varepsilon \cdot \text{sign}(\mathbf{w})$, where the sign is taken elementwise.

This choice of $\boldsymbol{\delta}^*$ is very similar to that the Fast Gradient Sign Method, 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 minimizer is similarly given as $\boldsymbol{\delta}^* = \varepsilon \cdot \text{sign}(\nabla_x L(\mathbf{x}, y))$.