

Machine Learning for Graphs and Sequential Data Exercise Sheet 05

Robustness of Machine Learning Models II

Problem 1: On slide 15 of the robustness chapter, we have defined an optimization problem for untargeted attacks, i.e. we aim to have the sample $\hat{\mathbf{x}}$ classified as **any** class other than the correct one:

$$\min_{\hat{\mathbf{x}}} \mathcal{D}(\mathbf{x}, \hat{\mathbf{x}}) + \lambda \cdot L(\hat{\mathbf{x}}, y)$$

The loss function is defined as:

$$L(\hat{\mathbf{x}}, y) = \left[Z(\hat{\mathbf{x}})_y - \max_{i \neq y} Z(\hat{\mathbf{x}})_i \right]_+,$$

where $[\mathbf{x}]_+$ is shorthand for $\max(\mathbf{x}, 0)$ and $Z(\mathbf{x})_i = \log f(\mathbf{x})_i$ (i.e. log probability of class i). Here, $L(\hat{\mathbf{x}}, y)$ is positive if $\hat{\mathbf{x}}$ is classified correctly and 0 otherwise.

Provide an alternative loss function to turn this attack into a targeted attack, i.e. we aim to have the sample \mathbf{x} classified as a *specific* target class t .

$$L(\hat{\mathbf{x}}, t) = \left[\max_{i \neq t} Z(\hat{\mathbf{x}})_i - Z(\hat{\mathbf{x}})_t \right]_+$$

This loss is positive if $\hat{\mathbf{x}}$ is classified as a class that is **not** t , and is zero otherwise.

Problem 2: Recall from slide 41 the MILP constraints expressing the ReLU activation function. Show that a continuous relaxation on \mathbf{a} leads to the convex relaxation constraints on slide 54. That is, we replace the constraint $\mathbf{a}_i \in \{0, 1\}$ with $\mathbf{a}_i \in [0, 1]$.

We can combine the first two constraints on slide 41:

$$\begin{aligned} \mathbf{y}_i &\leq \mathbf{x}_i - \mathbf{l}_i(1 - \mathbf{a}_i) \\ \mathbf{y}_i &\leq \mathbf{u}_i \cdot \mathbf{a}_i \end{aligned}$$

by expressing it as

$$\mathbf{y}_i \leq \min(\mathbf{x}_i - \mathbf{l}_i(1 - \mathbf{a}_i), \mathbf{u}_i \cdot \mathbf{a}_i)$$

Note that we are free to choose any value for \mathbf{a}_i between 0 and 1. We want to choose \mathbf{a}_i so that it leads to the loosest-possible constraint on \mathbf{y}_i , since this leads to the maximum ‘leeway’ to optimize the objective function. More formally,

$$\mathbf{y}_i \leq \max_{\mathbf{a}_i} \min(\mathbf{x}_i - \mathbf{l}_i(1 - \mathbf{a}_i), \mathbf{u}_i \cdot \mathbf{a}_i)$$

Further note that the two terms in the $\min(\cdot, \cdot)$ are two linear functions in \mathbf{a}_i . For $\mathbf{l}_i < 0$, the former is a function with negative slope in \mathbf{a}_i . We only need to consider the case $\mathbf{l}_i < 0$, since if $\mathbf{l}_i \geq 0$, we know that the unit is *stably active* and therefore linear.

The second term in the $\min(\cdot, \cdot)$ is a function with positive slope in \mathbf{a}_i if $\mathbf{u}_i > 0$. Again, we only need to consider this case since $\mathbf{u}_i \leq 0$ implies that the unit is *stably inactive* and therefore $\mathbf{y}_i = 0$.

Consequently, the function $\min(\mathbf{x}_i - \mathbf{l}_i(1 - \mathbf{a}_i), \mathbf{u}_i \cdot \mathbf{a}_i)$ is maximal at the intersection of the two linear functions. Solving for \mathbf{a}_i we get:

$$\begin{aligned}\mathbf{x}_i - (1 - \mathbf{a}_i)\mathbf{l}_i &= \mathbf{a}_i\mathbf{u}_i \\ \Leftrightarrow \mathbf{a}_i &= \frac{\mathbf{x}_i - \mathbf{l}_i}{\mathbf{u}_i - \mathbf{l}_i}\end{aligned}$$

Plugging the expression of \mathbf{a}_i into one of the original constraints, e.g. $\mathbf{y}_i \leq \mathbf{a}_i \cdot \mathbf{u}_i$ we get:

$$\begin{aligned}\mathbf{y}_i &\leq \frac{\mathbf{x}_i - \mathbf{l}_i}{\mathbf{u}_i - \mathbf{l}_i} \mathbf{u}_i \\ \Leftrightarrow \mathbf{y}_i(\mathbf{u}_i - \mathbf{l}_i) - \mathbf{u}_i\mathbf{x}_i &\leq \mathbf{u}_i\mathbf{l}_i,\end{aligned}$$

and therefore we have recovered the constraint of the convex relaxation.
