

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 $[x]_+$ is shorthand for $\max(x, 0)$ 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 .

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 $\mathbf{a}_i \in [0, 1]$.
