

Machine Learning for Graphs and Sequential Data Exercise Sheet 04

Robustness of Machine Learning Models I

Problem 1: Suppose we have a trained binary logistic regression classifier with weight vector $\mathbf{w} \in \mathbb{R}^d$ and bias $b \in \mathbb{R}$. Given a sample $\mathbf{x} \in \mathbb{R}^d$ we want to construct an adversarial example via gradient descent on the binary cross entropy loss:

$$\mathcal{L}(\mathbf{x}, y) = -y \log(\sigma(z)) - (1 - y) \log(1 - \sigma(z)),$$

where $\sigma(z) = \frac{1}{1+e^{-z}}$ is the logistic sigmoid function, $z = \mathbf{w}^T \mathbf{x} + b$, and $y \in \{0, 1\}$ is the class label of the sample at hand.

- a) Derive the gradient $\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, y)$. How do you interpret the result?

Hint: You may use the relation $1 - \sigma(z) = \sigma(-z)$.

- b) Provide a closed-form expression for the worst-case perturbed instance $\tilde{\mathbf{x}}^*$ (measured by the loss \mathcal{L}) for the perturbation set $\mathcal{P}(\mathbf{x}) = \{\tilde{\mathbf{x}} : \|\tilde{\mathbf{x}} - \mathbf{x}\|_2 \leq \epsilon\}$, i.e.

$$\tilde{\mathbf{x}}^* = \arg \max_{\|\tilde{\mathbf{x}} - \mathbf{x}\|_2 \leq \epsilon} \mathcal{L}(\tilde{\mathbf{x}}, y)$$

- c) What is the smallest value of ϵ for which the sample \mathbf{x} is misclassified (assuming it was correctly classified before)?
- d) We would now like to perform adversarial training. Provide a closed-form expression of the worst-case loss

$$\hat{\mathcal{L}}(\mathbf{x}, y) = \max_{\|\tilde{\mathbf{x}} - \mathbf{x}\|_2 \leq \epsilon} \mathcal{L}(\tilde{\mathbf{x}}, y)$$

as a function of \mathbf{x} and \mathbf{w} . How do you interpret the results?

Problem 2: In the lecture on exact certification of neural network robustness we have considered $K - 1$ optimization problems (one for each incorrect class) of the form (c.f. slide 42):

$$m_t^* = \min_{\tilde{\mathbf{x}}, \mathbf{y}^{(t)}, \hat{\mathbf{x}}^{(t)}, \mathbf{a}^{(t)}} [\hat{\mathbf{x}}^{(L)}]_{c^*} - [\hat{\mathbf{x}}^{(L)}]_t \quad \text{subject to MILP constraints.}$$

That is, for each class $t \neq c^*$, we optimize for the **worst-case margin** m_t^* , and conclude that the classifier is robust if and only if

$$\min_{t \neq c^*} m_t^* \geq 0.$$

However, we can equivalently solve the following single optimization problem:

$$m^* = \min_{\tilde{\mathbf{x}}, \mathbf{y}^{(t)}, \hat{\mathbf{x}}^{(t)}, \mathbf{a}^{(t)}} \left([\hat{\mathbf{x}}^{(L)}]_{c^*} - y \right) \quad \text{subject to } y = \max_{t \neq c^*} [\hat{\mathbf{x}}^{(L)}]_t \wedge \text{MILP constraints,}$$

where we have introduced a new variable y into the objective function.

Express the equality constraint

$$y = \max(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{K-1})$$

using only linear and integer constraints. To simplify notation, here $\mathbf{x}_k \in \mathbb{R}$ denotes the logit corresponding to the k -th incorrect class, and \mathbf{l}_k and \mathbf{u}_k its corresponding lower and upper bound.

Hint: You might want to introduce binary variables to indicate which logit is the maximum.
