

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 $1 - \sigma(z) = \sigma(-z)$.

$$\begin{aligned} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, y) &= \frac{-y}{\sigma(z)} \frac{\partial \sigma(z)}{\partial z} \nabla_{\mathbf{x}} z - \frac{1-y}{\sigma(-z)} \frac{\partial \sigma(-z)}{\partial z} \nabla_{\mathbf{x}} z \\ &= \frac{-y}{\sigma(z)} \sigma(z) \sigma(-z) \mathbf{w} + \frac{1-y}{\sigma(-z)} \sigma(-z) \sigma(z) \mathbf{w} \\ &= -y \sigma(-z) \mathbf{w} + (1-y) \sigma(z) \mathbf{w} \end{aligned}$$

The gradient is orthogonal to the decision boundary and points in the direction of the wrong class, depending on y .

- 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.

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

Since the loss is convex w.r.t. the data, taking a gradient step of magnitude ϵ towards the wrong class will result in the maximum increase in loss:

$$\begin{aligned} \tilde{\mathbf{x}}^* &= \mathbf{x} - \epsilon \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \text{ if } y = 1 \\ \tilde{\mathbf{x}}^* &= \mathbf{x} + \epsilon \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \text{ if } y = 0 \end{aligned}$$

- c) What is the smallest value of ϵ for which the sample \mathbf{x} is misclassified (assuming it was correctly classified before)?
-

For the sample to change classification we need to have $\sigma(z) = 0.5 \Leftrightarrow \mathbf{w}^T \tilde{\mathbf{x}} + b = 0$. Plugging in the perturbation we get for $y = 1$:

$$\begin{aligned}\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \epsilon \frac{\mathbf{w}}{\|\mathbf{w}\|_2} + b &= 0 \\ \mathbf{w}^T \mathbf{x} - \epsilon \|\mathbf{w}\|_2 + b &= 0 \\ \frac{1}{\|\mathbf{w}\|_2} (\mathbf{w}^T \mathbf{x} + b) &= \epsilon\end{aligned}$$

Thus, for a misclassification we need $\epsilon > \frac{1}{\|\mathbf{w}\|_2} (\mathbf{w}^T \mathbf{x} + b)$.

Analogously for $y = 0$ we obtain $\epsilon > \frac{1}{\|\mathbf{w}\|_2} (-\mathbf{w}^T \mathbf{x} - b)$

- 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?

$$\begin{aligned}\hat{\mathcal{L}}(\mathbf{x}, y) &= \max_{\|\tilde{\mathbf{x}} - \mathbf{x}\|_2 \leq \epsilon} \mathcal{L}(\tilde{\mathbf{x}}, y) \\ &= \mathcal{L}(\tilde{\mathbf{x}}^*, y) \\ &= -y \log(\sigma(\mathbf{w}^T \mathbf{x} - \epsilon \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|_2} + b)) - (1 - y) \log(\sigma(-\mathbf{w}^T \mathbf{x} - \epsilon \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|_2} - b)) \\ &= -y \log(\sigma(\mathbf{w}^T \mathbf{x} - \epsilon \|\mathbf{w}\|_2 + b)) - (1 - y) \log(\sigma(-\mathbf{w}^T \mathbf{x} - \epsilon \|\mathbf{w}\|_2 - b))\end{aligned}$$

Consider the case $y = 1$ ($y = 0$ follows symmetrically). The input to the sigmoid function is shifted to the left (i.e. negative direction) by $\epsilon \|\mathbf{w}\|_2$, reducing the predicted probability of the sample \mathbf{x} belonging to class 1. Thus, only if $\mathbf{w}^T \mathbf{x} + b \geq \epsilon \|\mathbf{w}\|_2$ the sample will be classified as belonging to class 1. We can interpret this as trying to enforce that each sample has at least a distance of $\epsilon \|\mathbf{w}\|_2$ to the decision boundary. Moreover, this margin is proportional to the norm of the weight vector, so simply increasing the norm of \mathbf{w} does not lead to the desired outcome, since we can move $\epsilon \|\mathbf{w}\|_2$ units towards the decision boundary for a unit norm change on the sample \mathbf{x} . Note that, in contrast to support vector machines (SVMs), even when the samples have a margin of at least $\epsilon \|\mathbf{w}\|_2$ to the decision boundary, we have non-zero loss and continue training.

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):

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

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:

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

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.

We first define $u_{\max} := \max_k \mathbf{u}_k$, i.e. the largest upper bound.

Now we introduce the following constraints:

$$y \leq \mathbf{x}_k + (1 - b_k)(u_{\max} - \mathbf{l}_k) \quad \forall 1 \leq k \leq K-1 \quad (1)$$

$$y \geq \mathbf{x}_i \quad \forall 1 \leq k \leq K-1 \quad (2)$$

$$\mathbf{b}_k \in \{0, 1\} \quad \forall 1 \leq k \leq K-1 \quad (3)$$

$$\sum_{k=1}^{K-1} \mathbf{b}_k = 1 \quad (4)$$

The last constraint (4) simply ensures that only one element in \mathbf{b} is 1 and all others are zero.

The only valid assignment of \mathbf{b} is to have $\mathbf{b}_k = 1$ for the (unique) maximum value $\mathbf{x}_k = \max(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{K-1})$. To see this, consider the case that $\mathbf{b}_k = 1$ but \mathbf{x}_k is not the maximum value. Then, (1) resolves to $y \leq \mathbf{x}_k$. However, for the maximum value $\mathbf{x}_{\max} > \mathbf{x}_k$ we have from (2) $y \geq \mathbf{x}_{\max}$, leads to a contradiction.

Consider the case $\mathbf{b}_k = 1$ and the corresponding value \mathbf{x}_k is indeed the (unique) maximum. (1) and (2) imply that $y = \mathbf{x}_k$. The remaining values \mathbf{b}_i are zero, and in this case we need to show that (1) and (2) are never binding, regardless of the values \mathbf{x}_i . (2) is not binding since \mathbf{x}_i is not the maximum value. (1) is not binding because we have that $\mathbf{x}_i + u_{\max} - \mathbf{l}_i \geq u_{\max} \geq y$.