

Exercise 02

Adversarial Examples, Defenses and Box Verification

Reliable and Trustworthy Artificial Intelligence
ETH Zurich

Problem 1 (Projection onto ℓ_p -balls). The Euclidean projection $\mathbf{z} \in \mathbb{R}^n$ of a point $\mathbf{y} \in \mathbb{R}^n$ onto the ϵ -sized ℓ_p -ball around $\mathbf{x} \in \mathbb{R}^n$ is defined as (note the ℓ_2 and ℓ_p norms):

$$\mathbf{z} = \underset{\mathbf{x}' \text{ s.t. } \|\mathbf{x}' - \mathbf{x}\|_p \leq \epsilon}{\operatorname{arg\,min}} \|\mathbf{x}' - \mathbf{y}\|_2 \quad (1)$$

In general, this is a hard problem and exact closed-form solutions are only known for few p . In this task, we investigate this problem for different p .

a) Consider $p = \infty$. Derive a closed-form formula for projection onto the ϵ -sized ℓ_∞ -ball around $\mathbf{x} \in \mathbb{R}^n$.

b) Consider $p = 2$.

(i) Derive a closed-form formula for the projection $\mathbf{z} \in \mathbb{R}^n$ of a point $\mathbf{y} \in \mathbb{R}^n$ onto the ϵ -sized ℓ_2 -ball around $\mathbf{x} \in \mathbb{R}^n$.

(ii) Prove that for $n = 2$, your closed-form solution is correct, i.e., show that there exists no $\mathbf{q} \neq \mathbf{z}$ in the ϵ -sized ℓ_2 -ball around \mathbf{x} that is closer to \mathbf{y} than \mathbf{z} .

Hint: Assume for the sake of contradiction that there exists such a point \mathbf{q} . Use the triangle inequality.

c) Consider a general $p \geq 1$. Instead of finding an exact solution, we are only looking for an *approximate* projection.

(i) Assume $\|\mathbf{x} - \mathbf{y}\|_p > \epsilon$. Like for $p = 2$, we can try to move \mathbf{y} closer to \mathbf{x} along a direct line by shifting the former along the vector $\mathbf{y} - \mathbf{x}$ rescaled according to p . Use this idea to derive a closed-form formula for approximate projection, where the result \mathbf{z} is guaranteed to lie in the ϵ -sized ℓ_p -ball around \mathbf{x} (but is not necessarily the closest point).

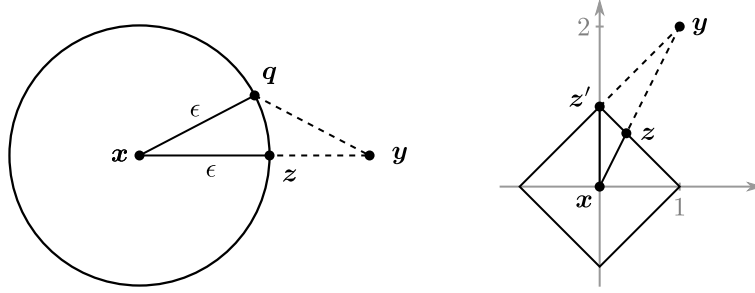


Figure 1: Geometric interpretation for solution 1.

- (ii) Construct a concrete counter-example for $n = 2$ showing that for $p = 1$, your formula is only approximate and does not solve (1).

Note: For general p , we can alternatively try to solve (1) using expensive iterative gradient-based optimization algorithms.

Solution 1.

- a) To project onto the ϵ -sized ℓ_∞ -ball around \mathbf{x} , we clamp \mathbf{y} coordinate-wise:

$$\forall i. \quad z_i = \max(\min(y_i, x_i + \epsilon), x_i - \epsilon).$$

- b) (i) The closed-form solution is

$$\mathbf{z} = \mathbf{x} + \frac{\mathbf{y} - \mathbf{x}}{\max(1, \|\mathbf{y} - \mathbf{x}\|_2 / \epsilon)}. \quad (2)$$

- (ii) *Proof:* We distinguish two cases.

Case 1. $\|\mathbf{y} - \mathbf{x}\|_2 \leq \epsilon$. The solution $\mathbf{z} = \mathbf{y}$ is correct since the constraint $\|\mathbf{z} - \mathbf{x}\|_2 \leq \epsilon$ holds and $\|\mathbf{x}' - \mathbf{y}\|_2$ is minimized for $\mathbf{x}' = \mathbf{y}$.

Case 2. $\|\mathbf{y} - \mathbf{x}\|_2 > \epsilon$. In this case, (2) induces that:

$$\|\mathbf{x} - \mathbf{z}\|_2 = \epsilon \quad (3)$$

Assume there exists a point $\mathbf{q} \neq \mathbf{z}$ that is closer to \mathbf{y} than \mathbf{z} and lies in the ϵ -sized ℓ_2 -ball around \mathbf{x} , as visualized in Fig. 1, left. Formally, for $d(\mathbf{a}, \mathbf{b}) =$

$\|\mathbf{a} - \mathbf{b}\|_2$ being the ℓ_2 -distance between \mathbf{a} and \mathbf{b} :

$$d(\mathbf{q}, \mathbf{y}) < d(\mathbf{z}, \mathbf{y}) \quad \wedge \quad d(\mathbf{x}, \mathbf{q}) \leq \epsilon \quad (4)$$

The triangle inequality in 2 dimensions states

$$d(\mathbf{x}, \mathbf{q}) + d(\mathbf{q}, \mathbf{y}) \geq d(\mathbf{x}, \mathbf{y}) \stackrel{(*)}{=} d(\mathbf{x}, \mathbf{z}) + d(\mathbf{z}, \mathbf{y})$$

where in (*) we used the fact that \mathbf{x} , \mathbf{z} , \mathbf{y} lie along a line in two dimensions according to (2) (see the triangle \mathbf{xqy} in Fig. 1, left). We can instantiate the inequalities from (4) and obtain

$$\begin{aligned} \epsilon + d(\mathbf{z}, \mathbf{y}) &> d(\mathbf{x}, \mathbf{z}) + d(\mathbf{z}, \mathbf{y}) \\ \epsilon &> d(\mathbf{x}, \mathbf{z}), \end{aligned}$$

which is a contradiction to (3).

c) (i) Analogously to $p = 2$, we can compute

$$\mathbf{z} = \mathbf{x} + \frac{\mathbf{y} - \mathbf{x}}{\max(1, \|\mathbf{y} - \mathbf{x}\|_p / \epsilon)}. \quad (5)$$

Note that it is $\|\mathbf{z} - \mathbf{x}\|_p \leq \epsilon$.

(ii) For example, consider $\epsilon = 1$, $\mathbf{x} = (0, 0)^\top$ and $\mathbf{y} = (1, 2)^\top$ as visualized in Fig. 1, right. The approximate projection according to (5) is $\mathbf{z} = (1/3, 2/3)^\top$. However, this is not the closest point in the ℓ_1 -ball around \mathbf{x} according to the ℓ_2 -norm: In particular, the point $\mathbf{z}' = (0, 1)^\top$ (see Fig. 1, right) has distance $d(\mathbf{y}, \mathbf{z}') = \sqrt{2} \approx 1.414$, which is less than $d(\mathbf{y}, \mathbf{z}) = \sqrt{4/9 + 16/9} \approx 1.491$.

Problem 2. In the lecture, we discussed the Carlini-Wagner optimization problem [1]:

$$\begin{aligned} &\text{find } \boldsymbol{\eta} \\ &\text{minimize } \|\boldsymbol{\eta}\|_p + c \cdot \text{obj}(\mathbf{x} + \boldsymbol{\eta}) \\ &\text{such that } \mathbf{x} + \boldsymbol{\eta} \in [0, 1]^n \end{aligned} \quad (6)$$

Directly optimizing this objective with gradient-based methods can be problematic in the case of $p = \infty$. In this task, we will investigate this and discuss a surrogate term for $\|\cdot\|_\infty$. In the following, assume that $\mathbf{x}, \boldsymbol{\eta} \in \mathbb{R}^n$. We define $h(\boldsymbol{\eta}) = \|\boldsymbol{\eta}\|_\infty$ and $g_\tau(\boldsymbol{\eta}) := \sum_{i=0}^{n-1} \max(\boldsymbol{\eta}_i - \tau, 0)$ for $\tau \in \mathbb{R}$.

a) Calculate the partial derivatives $\frac{\partial}{\partial \boldsymbol{\eta}_i} h(\boldsymbol{\eta})$ and $\frac{\partial}{\partial \boldsymbol{\eta}_i} g_\tau(\boldsymbol{\eta})$.

b) Instantiate the above derivatives with $\tau = 0.9$ and $\tau = 2.0$ for

$$\boldsymbol{\eta} = (1.00005, 1.00004, 1.00003, 1.00002, 0.001, 0.001)^\top.$$

- c) Assume we are minimizing $h(\boldsymbol{\eta})$ using gradient descent with step size $\gamma = 0.01$. How many iterations are required until $h(\boldsymbol{\eta}) < 0.01$? What is the problem?
- d) When combining $h(\boldsymbol{\eta})$ with the objective function to obtain the full optimization problem in (6), optimizing the latter using gradient descent may actually result in an infinite number of iterations. Explain why this can happen.
- e) To circumvent the above problems, we can use g_τ as a surrogate for h . How does g_τ improve the situation?

Note: For the reasons discussed above, the authors of [1] suggest to use g_τ instead of h , where τ is lowered repeatedly during optimization whenever $\max_i(\boldsymbol{\eta}_i) \leq \tau$ in order to produce a solution $\boldsymbol{\eta}$ with small ℓ_∞ -norm.

Solution 2.

a)

$$\frac{\partial}{\partial \boldsymbol{\eta}_i} h(\boldsymbol{\eta}) = \begin{cases} \text{sign}(\boldsymbol{\eta}_i) & \text{if } i \in \text{argmax}_k(|\boldsymbol{\eta}_k|) \\ 0 & \text{else} \end{cases}$$

$$\frac{\partial}{\partial \boldsymbol{\eta}_i} g(\boldsymbol{\eta}) = \begin{cases} 1 & \text{if } \boldsymbol{\eta}_i > \tau \\ 0 & \text{else} \end{cases}$$

Note that for $\frac{\partial}{\partial \boldsymbol{\eta}_i} h(\boldsymbol{\eta})$, there could be multiple maxima $\boldsymbol{\eta}_k$. Mathematically, the derivative is not defined in this case. However, in automatic differentiation engines such as PyTorch, the derivative is 1 for all maxima. For this exercise, either definition is fine.

b) Instantiating the above derivatives for the given $\boldsymbol{\eta}$, we obtain:

$$\begin{aligned} \nabla_{\boldsymbol{\eta}} h(\boldsymbol{\eta}) &= (1, 0, 0, 0, 0, 0)^T \\ \nabla_{\boldsymbol{\eta}} g_{0.9}(\boldsymbol{\eta}) &= (1, 1, 1, 1, 0, 0)^T \\ \nabla_{\boldsymbol{\eta}} g_{2.0}(\boldsymbol{\eta}) &= (0, 0, 0, 0, 0, 0)^T. \end{aligned}$$

c) Recall that gradient descent with step size γ repeatedly updates

$$\boldsymbol{\eta}^{t+1} \leftarrow \boldsymbol{\eta}^t - \gamma \cdot \nabla_{\boldsymbol{\eta}} h(\boldsymbol{\eta}^t).$$

The first four iterations are as follows.

$$\begin{aligned}\boldsymbol{\eta}^0 &= (1.00005, 1.00004, 1.00003, 1.00002, 0.001, 0.001)^\top \\ \boldsymbol{\eta}^1 &= (0.99005, 1.00004, 1.00003, 1.00002, 0.001, 0.001)^\top \\ \boldsymbol{\eta}^2 &= (0.99005, 0.99004, 1.00003, 1.00002, 0.001, 0.001)^\top \\ \boldsymbol{\eta}^3 &= (0.99005, 0.99004, 0.99003, 1.00002, 0.001, 0.001)^\top \\ \boldsymbol{\eta}^4 &= (0.99005, 0.99004, 0.99003, 0.99002, 0.001, 0.001)^\top\end{aligned}$$

In each iteration, only the coordinate with the maximum value is updated. Hence, we require 400 iterations until all coordinates are less than 0.01. As demonstrated in this example, one problem with optimizing the ℓ_∞ norm using gradient descent is slow convergence.

- d) As the gradient $\nabla_{\boldsymbol{\eta}}h(\boldsymbol{\eta})$ is zero everywhere except at the coordinate i with the maximum value, the complete objective function is not penalized for increasing any non-maximum coordinate $j \neq i$.

For example, consider $\mathbf{x} = (0, 0)^\top$ and $\boldsymbol{\eta} = (1, 0.9)^\top$. Assume the objective function has derivative $\nabla_{\boldsymbol{\eta}}\text{obj}(\mathbf{x} + \boldsymbol{\eta}) = (0, -1)^\top$, meaning that its output can be lowered by increasing the last coordinate. Hence, for $c = 1$, the gradient of the full objective becomes

$$\nabla_{\boldsymbol{\eta}}(h(\boldsymbol{\eta}) + \text{obj}(\mathbf{x} + \boldsymbol{\eta})) = (1, 0)^\top + (0, -1)^\top = (1, -1)^\top.$$

Note how the gradient of $h(\boldsymbol{\eta})$ does not penalize for increasing the last coordinate. With a step size of $\gamma = 0.1$, the gradient step updates $\boldsymbol{\eta}^0 = \boldsymbol{\eta}$ to $\boldsymbol{\eta}^1 = (0.9, 1)^\top$. In the next step, by symmetric reasoning, the optimizer may oscillate back to $\boldsymbol{\eta}^2 = (1, 0.9)^\top = \boldsymbol{\eta}^0$ and never terminate.

- e) Using $g_\tau(\boldsymbol{\eta})$, we can mitigate the problems above for a well-chosen τ . In the example of subtask (c), using $g_{0.9}$ will update multiple coordinates at the same time and hence lead to much faster convergence. In the example of subtask (d), picking a sufficiently large τ ensures that the overall objective is penalized for increasing the non-maximum coordinate, thereby breaking the cycle.

Note 1: Our surrogate function $g_\tau(\boldsymbol{\eta})$ only works in the case where we know the entries of $\boldsymbol{\eta}$ to be positive. More generally, we would instead use the function $g_\tau^*(\boldsymbol{\eta}) := \sum_{i=0}^{n-1} \max(|\boldsymbol{\eta}_i| - \tau, 0)$. However, for optimizing $\boldsymbol{\eta} \in [0, 1]^n$, this does not matter.

Note 2: From a theoretical point of view, this approach provides no guarantees. The authors of [1] observed that using g_τ improves convergence in practice. for the reasons intuitively presented here.

Problem 3 (Coding). In the ZIP file provided on the course webpage, you can find a python skeleton `task3.ipynb` along with a pre-trained MNIST classifier `model`.

Note: The skeleton is based on the PyTorch¹ framework. We strongly recommend that you familiarize yourself with PyTorch now, because the course project will rely heavily on PyTorch. This exercise allows you to gain some initial experience with PyTorch.

Re-using your code from last week, implement the (untargeted) PGD attack ([2]) discussed in the lecture (`pgd`). Here, `k` is the number of FGSM iterations with `eps_step` step size. Each FGSM iteration is projected back to the `eps`-sized ℓ_∞ -ball around `x`. Your implementation should clamp the adversarial example back to the image domain.

Solution 3. See `task3_sol.ipynb` in the provided ZIP file.

Problem 4 (Coding). In this task, you are going to implement adversarial training with PGD (originally introduced in [2]) and an alternative defense.

- a) Complete the provided code skeleton in `task4.py` to train the network `model` with PGD defense. That is, for data distribution D (the MNIST dataset in our case) and network parameters θ , optimize the following objective:

$$\min_{\theta} \mathbb{E}_{(x,y) \sim D} \left[\max_{x' \in \mathbb{B}_\epsilon(x)} L(\theta, x', y) \right]. \quad (7)$$

Here, $\mathbb{B}_\epsilon(x) := \{x' \mid \|x - x'\|_\infty \leq \epsilon\}$ denotes the ϵ -sized ℓ_∞ -ball around x . L is the usual classification loss $L(\theta, x', y) := \mathcal{H}(y, f_\theta(x'))$, where $f_\theta = (\text{softmax} \circ \text{model})$ denotes the output distribution of the neural network, and \mathcal{H} the cross entropy² between distributions (being a discrete value, we treat y as a one-hot distribution). In PyTorch, you can use `nn.CrossEntropyLoss`³ to implement L .

Use PGD to solve the inner optimization problem with $\epsilon = 0.1$, $k = 7$ steps, and $\epsilon_{\text{step}} = 2.5 \frac{\epsilon}{k}$. You can reuse your implementation of untargeted PGD from the previous exercise, or create a more efficient (batched) version better suited for training.

Compare the accuracy results with and without PGD training.

- b) The TRADES [3] algorithm minimizes the following objective (see [3] for details):

¹<https://pytorch.org>

²https://en.wikipedia.org/wiki/Cross_entropy

³<https://pytorch.org/docs/stable/generated/torch.nn.CrossEntropyLoss.html>

$$\min_{\theta} \mathbb{E}_{(x,y) \sim D} \left[\underbrace{L(\theta, x, y)}_{\text{for accuracy}} + \lambda \underbrace{\max_{x' \in \mathbb{B}_{\epsilon}(x)} L(\theta, x', f_{\theta}(x))}_{\text{regularization for robustness}} \right]. \quad (8)$$

Extend your implementation in `train.py` to the TRADES defense. Again, the inner optimization problem is solved using PGD. Use $\lambda = 1.0$ and the same parameters as in the previous task. Compare your results with the previous task.

Note: Here, $L(\theta, x', f_{\theta}(x))$ still denotes the cross entropy loss.

Solution 4. See `task4_sol.py`.

Defense	Clean Accuracy	Adversarial Accuracy
None	0.98	0.29
PGD	0.98	0.90
TRADES	0.98	0.89

Table 1: Test set accuracy.

The accuracy obtained with the default parameters and 50 (instead of 10) epochs is shown in table 1. Without any defense, the adversarial accuracy drops significantly, while it remains high with a defense. Note that the PGD defense outperformed TRADES here (higher adversarial accuracy at the same clean accuracy)—MNIST is a very “easy” task where adversarial training does not impact the clean accuracy too much. Note that the model and hyper-parameters used here are far from optimal.

References

- [1] Nicholas Carlini and David A. Wagner. “Towards Evaluating the Robustness of Neural Networks”. In: *2017 IEEE Symposium on Security and Privacy, SP 2017, San Jose, CA, USA, May 22-26, 2017*. IEEE Computer Society, 2017, pp. 39–57. DOI: 10.1109/SP.2017.49. URL: <https://arxiv.org/abs/1608.04644>.
- [2] Aleksander Madry et al. “Towards deep learning models resistant to adversarial attacks”. In: *arXiv preprint arXiv:1706.06083* (2017).
- [3] Hongyang Zhang et al. “Theoretically Principled Trade-off between Robustness and Accuracy”. In: *ICML*. Vol. 97. Proceedings of Machine Learning Research. PMLR, 2019, pp. 7472–7482.