

# Exercise 12

## Randomized Smoothing

Reliable and Interpretable Artificial Intelligence  
ETH Zurich

**Problem 1.** In this problem, we want to proof the following theorem:

**Theorem 1** (From [1]). *Let  $f : \mathbb{R}^d \rightarrow \mathcal{Y}$  be any deterministic or random function. Let  $\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$ . Let  $g(\mathbf{x}) = \operatorname{argmax}_c \mathbb{P}(f(\mathbf{x} + \varepsilon) = c)$ . Suppose that for a specific  $\mathbf{x} \in \mathbb{R}^d$ , there exist  $c_A \in \mathcal{Y}$  and  $\underline{p}_A, \overline{p}_B \in [0, 1]$  such that:*

$$\mathbb{P}(f(\mathbf{x} + \varepsilon) = c_A) \geq \underline{p}_A \geq \overline{p}_B \geq \max_{c \neq c_A} \mathbb{P}(f(\mathbf{x} + \varepsilon) = c) \quad (1)$$

Then  $g(\mathbf{x} + \boldsymbol{\delta}) = c_A$  for all  $\|\boldsymbol{\delta}\|_2 < R$ , where

$$R = \frac{\sigma}{2} (\Phi^{-1}(\underline{p}_A) - \Phi^{-1}(\overline{p}_B)). \quad (2)$$

The proof is broken down into the Steps (1) - (4). We (1) decompose the input space into two half spaces,  $A$  and  $B$ , such that the probabilities for the samples from the non-displaced gaussian  $\mathbf{X} \sim \mathcal{N}(\mathbf{x}, \sigma^2 I)$  to lie in  $A$  or  $B$  are  $\underline{p}_A$  and  $\overline{p}_B$  respectively. This induces a linear separation between the two gaussians (the worst case).

We (2) use Lemma 1 to relate the probabilities of the displaced gaussian  $\mathbf{Y} \sim \mathcal{N}(\mathbf{x} + \boldsymbol{\delta}, \sigma^2 I)$  to observe class  $c_A$  or  $c_B$  to the probabilities that  $Y \in A$  or  $Y \in B$  respectively.

Then in (3), we show how these probabilities relate to  $\boldsymbol{\delta}$  and  $\sigma$ . Finally, in (4) we obtain a condition on  $\|\boldsymbol{\delta}\|$  such that the classification is robust for all  $\boldsymbol{\delta}$  satisfying  $\|\boldsymbol{\delta}\| < R$ .

1. Show that  $\mathbb{P}(\mathbf{X} \in A) = \underline{p}_A$  and  $\mathbb{P}(\mathbf{X} \in B) = \overline{p}_B$ , where  $\mathbf{X} := \mathbf{x} + \varepsilon \sim \mathcal{N}(\mathbf{x}, \sigma^2 I)$  and

$$\begin{aligned} A &:= \{\mathbf{z} \in \mathbb{R}^d : \boldsymbol{\delta}^T (\mathbf{z} - \mathbf{x}) \leq \sigma \|\boldsymbol{\delta}\| \Phi^{-1}(\underline{p}_A)\} \\ B &:= \{\mathbf{z} \in \mathbb{R}^d : \boldsymbol{\delta}^T (\mathbf{z} - \mathbf{x}) \geq \sigma \|\boldsymbol{\delta}\| \Phi^{-1}(1 - \overline{p}_B)\}. \end{aligned}$$

*Hint:* Let  $x \sim \mathcal{N}(\mu_x, \sigma_x^2)$  and  $y \sim \mathcal{N}(\mu_y, \sigma_y^2)$ . Then  $x + y \sim \mathcal{N}(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$  and  $cx \sim \mathcal{N}(c\mu_x, c^2\sigma_x^2)$ .

2. Use Lemma 1 (see the bonus Problem 2), the results from sub-task 1 and the assumptions of the theorem to show

$$\mathbb{P}(f(\mathbf{Y}) = c_A) \geq \mathbb{P}(\mathbf{Y} \in A) \quad \text{and} \quad \mathbb{P}(f(\mathbf{Y}) = c_B) \leq \mathbb{P}(\mathbf{Y} \in B),$$

where  $\mathbf{Y} := (\mathbf{x} + \boldsymbol{\delta} + \boldsymbol{\varepsilon}) \sim \mathcal{N}(\mathbf{x} + \boldsymbol{\delta}, \sigma^2 I)$ .

3. Show that  $\mathbb{P}(\mathbf{Y} \in A) = \Phi\left(\Phi^{-1}(\underline{p}_A) - \frac{\|\boldsymbol{\delta}\|}{\sigma}\right)$  and  $\mathbb{P}(\mathbf{Y} \in B) = \Phi\left(\Phi^{-1}(\overline{p}_B) + \frac{\|\boldsymbol{\delta}\|}{\sigma}\right)$ .

*Hint:* Let  $z \sim \mathcal{N}(0, \sigma^2)$ . Then  $(z + \mu) \sim \mathcal{N}(\mu, \sigma^2)$ .

4. Find the condition for  $\boldsymbol{\delta}$  such that  $\mathbb{P}(\mathbf{Y} \in A) > \mathbb{P}(\mathbf{Y} \in B)$  holds.

**Problem (opt.) 2.** In this task we will prove the following Lemma:

**Lemma 1 (Special case of Neyman-Pearson).** *Let  $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 I)$ ,  $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu} + \boldsymbol{\delta}, \sigma^2 I)$ ,  $f: \mathbb{R}^d \rightarrow \mathcal{Y}$  a deterministic or random function and  $c \in \mathcal{Y}$ . Then:*

1. If  $S = \{\mathbf{z} \in \mathbb{R}^d : \boldsymbol{\delta}^T \mathbf{z} \leq t\}$  for some  $t$  and  $\mathbb{P}(f(\mathbf{X}) = c) \geq \mathbb{P}(\mathbf{X} \in S)$ , then  $\mathbb{P}(f(\mathbf{Y}) = c) \geq \mathbb{P}(\mathbf{Y} \in S)$
2. If  $S = \{\mathbf{z} \in \mathbb{R}^d : \boldsymbol{\delta}^T \mathbf{z} \geq t\}$  for some  $t$  and  $\mathbb{P}(f(\mathbf{X}) = c) \leq \mathbb{P}(\mathbf{X} \in S)$ , then  $\mathbb{P}(f(\mathbf{Y}) = c) \leq \mathbb{P}(\mathbf{Y} \in S)$ .

Let

$$g(\mathbf{x}; \boldsymbol{\mu}, \sigma^2 I) := \frac{1}{\sqrt{(2\pi\sigma^2)^d}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{x} - \boldsymbol{\mu})^T(\mathbf{x} - \boldsymbol{\mu})\right)$$

denote the Gaussian probability density function (for mean  $\boldsymbol{\mu}$  and co-variance matrix  $\sigma^2 I$ ) evaluate at  $\mathbf{x}$ . For convenience we write  $g_{\mathbf{X}}(\mathbf{x}) := g(\mathbf{x}; \boldsymbol{\mu}_{\mathbf{X}}, \sigma_{\mathbf{X}}^2 I)$  for Gaussian Random Variables  $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{X}}, \sigma_{\mathbf{X}}^2 I)$ . Further, we let  $A^C$  denote the complement of a set  $A$  over  $\mathbb{R}^d$ ,  $A^C := \mathbb{R}^d \setminus A$ . Then we can trivially decompose

$$\mathbb{R}^d = A \cup A^C. \tag{3}$$

1. Compute and simplify  $m(\mathbf{z}) := \frac{g_{\mathbf{Y}}(\mathbf{z})}{g_{\mathbf{X}}(\mathbf{z})}$ .
2. Show that for any  $t$  there exists a  $t' > 0$  such that  $\{z \in \mathbb{R}^d : \boldsymbol{\delta}^T \mathbf{z} \leq t\} = \{z \in \mathbb{R}^d : m(\mathbf{z}) \leq t'\}$ .
3. Given  $S := \{z \in \mathbb{R}^d : m(\mathbf{z}) \leq t'\}$  show that

$$\left( \int_{S^c} g_{\mathbf{Y}}(\mathbf{z}) d\mathbf{z} - \int_S g_{\mathbf{Y}}(\mathbf{z}) d\mathbf{z} \right) \geq t' \left( \int_{S^c} g_{\mathbf{X}}(\mathbf{z}) d\mathbf{z} - \int_S g_{\mathbf{X}}(\mathbf{z}) d\mathbf{z} \right).$$

Extend this to  $(\int_{S^c} [f(\mathbf{z}) = c] \cdot g_{\mathbf{Y}}(\mathbf{z}) d\mathbf{z} - \int_S [f(\mathbf{z}) \neq c] \cdot g_{\mathbf{Y}}(\mathbf{z}) d\mathbf{z})$   
 $\geq t' (\int_{S^c} [f(\mathbf{z}) = c] \cdot g_{\mathbf{X}}(\mathbf{z}) d\mathbf{z} - \int_S [f(\mathbf{z}) \neq c] \cdot g_{\mathbf{X}}(\mathbf{z}) d\mathbf{z}).$

4. Let  $f : \mathbb{R}^d \rightarrow Y$  be a classifier (function) and  $c \in Y$  a class. Show that for  $S := \{\mathbf{z} \in \mathbb{R}^d : m(\mathbf{z}) \leq t'\}$  for a given  $t' > 0$  and  $\mathbb{P}(f(\mathbf{X}) = c) \geq \mathbb{P}(\mathbf{X} \in S)$ , then  $\mathbb{P}(f(\mathbf{Y}) = c) \geq \mathbb{P}(\mathbf{Y} \in S)$ . Hint: Show that  $\mathbb{P}(f(\mathbf{Y}) = c) - \mathbb{P}(\mathbf{Y} \in S) \leq 0$  and use the results from the previous tasks.
5. Putting the previous tasks together gives you the proof for the first part of Lemma 1. What changes are required for the second part?

**Problem 3.** Randomized smoothing currently is usually formulated for the  $\ell^1$  or  $\ell^2$ -norm. However, well-known equalities can be used to bound norms other than the one guaranteed by the method used; for an example see [2] which uses the  $\ell_2$ -norm to obtain  $\ell^\infty$ -bounds. In the following we will show different useful inequalities. Show the inequality and provide the tightest value of  $c$  you can find. (You don't need to prove the tightness, although you easily can through an example). Hint: To obtain the tightest bounds you might need to use an additional theorem such as the subadditivity of the square root function or the Cauchy-Schwartz inequality. As a reminder:

$$\|\mathbf{x}\|_1 = \sum_{i=1}^d |x_i| \quad \|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^d |x_i|^2} \quad \|\mathbf{x}\|_\infty = \max_{i \in \{1, \dots, d\}} |x_i| \quad \text{for } \mathbf{x} \in \mathbb{R}^d$$

1. Show  $\|\mathbf{x}\|_\infty \leq c_1 \|\mathbf{x}\|_1$ .
2. Show  $\|\mathbf{x}\|_1 \leq c_2 \|\mathbf{x}\|_\infty$ .
3. Show  $\|\mathbf{x}\|_\infty \leq c_3 \|\mathbf{x}\|_2$ .
4. Show  $\|\mathbf{x}\|_2 \leq c_4 \|\mathbf{x}\|_\infty$ .
5. Show  $\|\mathbf{x}\|_2 \leq c_5 \|\mathbf{x}\|_1$ .
6. Show  $\|\mathbf{x}\|_1 \leq c_6 \|\mathbf{x}\|_2$ .
7. Let  $\mathbb{B}_\epsilon^p := \{x \in \mathbb{R}^d \mid \|x\|_p \leq \epsilon\}$  denote the  $l_p$ -norm ball of size  $\epsilon$ . Order  $\mathbb{B}_1^1, \mathbb{B}_1^2, \mathbb{B}_1^\infty, \mathbb{B}_d^1, \mathbb{B}_{\sqrt{d}}^2$  with respect to the inclusion relation  $\subseteq$ .

## References

- [1] Jeremy M. Cohen, Elan Rosenfeld, and J. Zico Kolter. “Certified Adversarial Robustness via Randomized Smoothing”. In: *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*. Ed. by Kamalika Chaudhuri and Ruslan Salakhutdinov. Vol. 97. Proceedings of Machine Learning Research. PMLR, 2019, pp. 1310–1320. URL: <http://proceedings.mlr.press/v97/cohen19c.html>.
- [2] Hadi Salman et al. “Provably Robust Deep Learning via Adversarially Trained Smoothed Classifiers”. In: *Advances in Neural Information Processing Systems 32*. Ed. by H. Wallach et al. Curran Associates, Inc., 2019, pp. 11289–11300. URL: <http://papers.nips.cc/paper/9307-provably-robust-deep-learning-via-adversarially-trained-smoothed-classifiers.pdf>.