

Exercise 12 - Solution

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 σ . 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.

Solution 1. 1. We start by applying the definitions and obtain

$$\mathbb{P}(\mathbf{X} \in A) = \mathbb{P}(\boldsymbol{\delta}^T (\mathbf{X} - \mathbf{x}) \leq \sigma \|\boldsymbol{\delta}\| \Phi^{-1}(\underline{p}_A)).$$

Next we replace $\mathbf{X} - \mathbf{x} \sim \mathcal{N}(0, \sigma^2 I)$ by $\mathbf{Z} \sim \mathcal{N}(0, \sigma^2 I)$ and get

$$\mathbb{P}(\mathbf{X} \in A) = \mathbb{P}(\boldsymbol{\delta}^T \mathbf{Z} \leq \sigma \|\boldsymbol{\delta}\| \Phi^{-1}(\underline{p}_A)).$$

Here, $\boldsymbol{\delta}^T \mathbf{Z} = \sum_i \boldsymbol{\delta}_i z_i$ where $z_i \sim \mathcal{N}(0, \sigma^2)$. Using the standard rule for multiplying the normally distributed gaussian variable z_i by a constant $\boldsymbol{\delta}_i$, we get $\boldsymbol{\delta}_i z_i \sim \mathcal{N}(0, \boldsymbol{\delta}_i^2 \sigma^2)$. By applying the sum rule to $\boldsymbol{\delta}_i z_i$ we get $\sum_i \boldsymbol{\delta}_i z_i \sim \mathcal{N}(0, \sigma^2 \sum_i \boldsymbol{\delta}_i^2)$. Further, with $\mathcal{N}(0, \sigma^2) = \sigma \mathcal{N}(0, 1)$ and $\sum_i \boldsymbol{\delta}_i^2 = \|\boldsymbol{\delta}\|^2$ we get

$$\begin{aligned} \mathbb{P}(\mathbf{X} \in A) &= \mathbb{P}(\sigma \|\boldsymbol{\delta}\| z \leq \sigma \|\boldsymbol{\delta}\| \Phi^{-1}(\underline{p}_A)) && (z \sim \mathcal{N}(0, 1)) \\ &= \mathbb{P}(z \leq \Phi^{-1}(\underline{p}_A)) \\ &= \Phi(\Phi^{-1}(\underline{p}_A)) \\ &= \underline{p}_A. \end{aligned}$$

Similarly for $\mathbb{P}(\mathbf{X} \in B) = \Phi\left(\Phi^{-1}(\overline{p}_B) + \frac{\|\boldsymbol{\delta}\|}{\sigma}\right)$.

2. Using the definitions for \mathbf{X} and \mathbf{Y} together with (Eq. (1)), we see that

$$\mathbb{P}(f(\mathbf{X}) = c_A) \geq \underline{p}_A \quad \text{and} \quad \mathbb{P}(f(\mathbf{X}) = c_B) \leq \overline{p}_B.$$

Applying Lemma 1 directly yields

$$\begin{aligned} \mathbb{P}(f(\mathbf{Y}) = c_A) &\geq \mathbb{P}(\mathbf{Y} \in A), \\ \mathbb{P}(f(\mathbf{Y}) = c_B) &\leq \mathbb{P}(\mathbf{Y} \in B). \end{aligned}$$

3. The calculation is similar to the one executed in sub-task 1:

$$\mathbb{P}(\mathbf{Y} \in A) = \mathbb{P}(\boldsymbol{\delta}^T(\mathbf{Y} - \mathbf{x}) \leq \sigma \|\boldsymbol{\delta}\| \Phi^{-1}(\underline{p}_A)).$$

Here we replace $\mathbf{Y} - \mathbf{x} \sim \mathcal{N}(\boldsymbol{\delta}, \sigma^2 I)$ by $\mathbf{Z} + \boldsymbol{\delta}$, where $\mathbf{Z} \sim \mathcal{N}(0, \sigma^2 I)$. Thus $\boldsymbol{\delta}^T(\mathbf{Y} - \mathbf{x}) = \boldsymbol{\delta}^T \mathbf{Z} + \|\boldsymbol{\delta}\|^2$:

$$\begin{aligned} \mathbb{P}(\mathbf{Y} \in A) &= \mathbb{P}(\boldsymbol{\delta}^T \mathbf{Z} + \|\boldsymbol{\delta}\|^2 \leq \sigma \|\boldsymbol{\delta}\| \Phi^{-1}(\underline{p}_A)) \\ &= \mathbb{P}(\boldsymbol{\delta}^T \mathbf{Z} \leq \sigma \|\boldsymbol{\delta}\| \Phi^{-1}(\underline{p}_A) - \|\boldsymbol{\delta}\|^2). \end{aligned}$$

Following the same steps for the term $\boldsymbol{\delta}^T \mathbf{Z}$ as in sub-task 1, we get

$$\begin{aligned} \mathbb{P}(\mathbf{Y} \in A) &= \mathbb{P}(\sigma \|\boldsymbol{\delta}\| z \leq \sigma \|\boldsymbol{\delta}\| \Phi^{-1}(\underline{p}_A) - \|\boldsymbol{\delta}\|^2) && (z \sim \mathcal{N}(0, 1)) \\ &= \mathbb{P}\left(z \leq \Phi^{-1}(\underline{p}_A) - \frac{\|\boldsymbol{\delta}\|}{\sigma}\right) \\ &= \Phi\left(\Phi^{-1}(\underline{p}_A) - \frac{\|\boldsymbol{\delta}\|}{\sigma}\right). \end{aligned}$$

Similarly for $\mathbb{P}(\mathbf{Y} \in B) = \Phi\left(\Phi^{-1}(\overline{p}_B) + \frac{\|\boldsymbol{\delta}\|}{\sigma}\right)$.

4. Finally, algebra shows that $\mathbb{P}(\mathbf{Y} \in A) > \mathbb{P}(\mathbf{Y} \in B)$ if and only if:

$$\begin{aligned} &\mathbb{P}(\mathbf{Y} \in A) > \mathbb{P}(\mathbf{Y} \in B) \\ \iff &\Phi\left(\Phi^{-1}(\underline{p}_A) - \frac{\|\boldsymbol{\delta}\|}{\sigma}\right) > \Phi\left(\Phi^{-1}(\overline{p}_B) + \frac{\|\boldsymbol{\delta}\|}{\sigma}\right) \\ \iff &\Phi^{-1}(\underline{p}_A) - \frac{\|\boldsymbol{\delta}\|}{\sigma} > \Phi^{-1}(\overline{p}_B) + \frac{\|\boldsymbol{\delta}\|}{\sigma} \\ \iff &\Phi^{-1}(\underline{p}_A) > \Phi^{-1}(\overline{p}_B) + \frac{2\|\boldsymbol{\delta}\|}{\sigma} \\ \iff &\|\boldsymbol{\delta}\| < \frac{\sigma}{2}(\Phi^{-1}(\underline{p}_A) - \Phi^{-1}(\overline{p}_B)) \end{aligned}$$

which recovers the theorem statement.

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

$$\mathbb{R}^d = A \cup A^C. \quad (3)$$

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 $S := \{z \in \mathbb{R}^d : m(\mathbf{z}) \leq t'\}$ 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?

Solution 2.

1. Without loss of generality we can assume $\boldsymbol{\mu} = 0$, to simplify the notation.

$$\begin{aligned}
m(\mathbf{z}) &:= \frac{g_{\mathbf{Y}}(\mathbf{z})}{g_{\mathbf{X}}(\mathbf{z})} \\
&= \frac{\frac{1}{\sqrt{(2\pi\sigma^2)^d}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{z} - (\boldsymbol{\mu} + \boldsymbol{\delta}))^T(\mathbf{z} - (\boldsymbol{\mu} + \boldsymbol{\delta}))\right)}{\frac{1}{\sqrt{(2\pi\sigma^2)^d}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{z} - \boldsymbol{\mu})^T(\mathbf{z} - \boldsymbol{\mu})\right)} \\
&= \frac{\exp\left(-\frac{1}{2\sigma^2}(\mathbf{z} - \boldsymbol{\delta})^T(\mathbf{z} - \boldsymbol{\delta})\right)}{\exp\left(-\frac{1}{2\sigma^2}\mathbf{z}^T\mathbf{z}\right)} \\
&= \exp\left(-\frac{1}{2\sigma^2}\left((\mathbf{z} - \boldsymbol{\delta})^T(\mathbf{z} - \boldsymbol{\delta}) - \mathbf{z}^T\mathbf{z}\right)\right) \\
&= \exp\left(-\frac{1}{2\sigma^2}\left(\mathbf{z}^T\mathbf{z} - 2\mathbf{z}^T\boldsymbol{\delta} + \boldsymbol{\delta}^T\boldsymbol{\delta} - \mathbf{z}^T\mathbf{z}\right)\right) \\
&= \exp\left(-\frac{1}{2\sigma^2}\left(-2\mathbf{z}^T\boldsymbol{\delta} + \boldsymbol{\delta}^T\boldsymbol{\delta}\right)\right) \\
&= \exp\left(\frac{1}{\sigma^2}\mathbf{z}^T\boldsymbol{\delta} - \frac{\boldsymbol{\delta}^T\boldsymbol{\delta}}{2\sigma^2}\right) \tag{4}
\end{aligned}$$

2. We define $a := \frac{1}{\sigma^2}$ and $b := -\frac{(\boldsymbol{\delta}^T\boldsymbol{\delta})}{2\sigma^2}$ and thus can rewrite Eq. (4) as $\exp(a\boldsymbol{\delta}^T\mathbf{z} + b)$. Thus if we know $\boldsymbol{\delta}^T\mathbf{z} \leq t$ we can write $m(\mathbf{z}) = \exp(a\boldsymbol{\delta}^T\mathbf{z} + b) \leq \exp(at + b) =: t'$. Thus we obtain $t' = \exp(at + b)$.
3. By definition of $m(\mathbf{z})$ and S we know that $g_{\mathbf{Y}}(\mathbf{z}) \leq t'g_{\mathbf{X}}(\mathbf{z}) \quad \forall \mathbf{z} \in S$ and therefore $\int_S g_{\mathbf{Y}}(\mathbf{z})d\mathbf{z} \leq t' \int_S g_{\mathbf{X}}(\mathbf{z})d\mathbf{z}$. Similarly $g_{\mathbf{Y}}(\mathbf{z}) > t'g_{\mathbf{X}}(\mathbf{z}) \quad \forall \mathbf{z} \in S^C$ and thus $\int_{S^C} g_{\mathbf{Y}}(\mathbf{z})d\mathbf{z} > t' \int_{S^C} g_{\mathbf{X}}(\mathbf{z})d\mathbf{z}$. Combining both facts accounting for the sign the statement follows immediately.

By denoting $F := \{\mathbf{z} \in \mathbb{R}^d : [f(\mathbf{z}) = c]\}$ we can analogously conclude

$$\begin{aligned}
\int_S [f(\mathbf{z}) = c] \cdot g_{\mathbf{Y}}(\mathbf{z}) &= \int_{S \cap F} g_{\mathbf{Y}}(\mathbf{z})d\mathbf{z} \\
&\leq t' \int_{S \cap F} g_{\mathbf{X}}(\mathbf{z})d\mathbf{z} \\
&= \int_S [f(\mathbf{z}) = c] \cdot g_{\mathbf{X}}(\mathbf{z})d\mathbf{z}.
\end{aligned}$$

Likewise we can show $\int_S [f(\mathbf{z}) \neq c] \cdot g_{\mathbf{Y}}(\mathbf{z})d\mathbf{z} > t' \int_S [f(\mathbf{z}) \neq c] \cdot g_{\mathbf{X}}(\mathbf{z})d\mathbf{z}$ and finally

$$\begin{aligned} & \left(\int_{S^c} [f(\mathbf{z}) = c] \cdot g_{\mathbf{Y}}(\mathbf{z}) d\mathbf{z} - \int_S [f(\mathbf{z}) \neq c] \cdot g_{\mathbf{Y}}(\mathbf{z}) d\mathbf{z} \right) \\ & \geq t' \left(\int_{S^c} [f(\mathbf{z}) = c] \cdot g_{\mathbf{X}}(\mathbf{z}) d\mathbf{z} - \int_S [f(\mathbf{z}) \neq c] \cdot g_{\mathbf{X}}(\mathbf{z}) d\mathbf{z} \right). \end{aligned}$$

4.

$$\begin{aligned} & \mathbb{P}(f(\mathbf{Y}) = c) - \mathbb{P}(\mathbf{Y} \in S) \\ & \stackrel{\text{definition}}{=} \int_{\mathbb{R}^d} [f(\mathbf{z}) = c] \cdot g_{\mathbf{Y}}(\mathbf{z}) d\mathbf{z} - \int_S g_{\mathbf{Y}}(\mathbf{z}) d\mathbf{z} \\ & \stackrel{\text{by Eq. (3)}}{=} \left[\int_{S^c} [f(\mathbf{z}) = c] \cdot g_{\mathbf{Y}}(\mathbf{z}) d\mathbf{z} + \int_S [f(\mathbf{z}) = c] \cdot g_{\mathbf{Y}}(\mathbf{z}) d\mathbf{z} \right] - \int_S g_{\mathbf{Y}}(\mathbf{z}) d\mathbf{z} \\ & = \left[\int_{S^c} [f(\mathbf{z}) = c] \cdot g_{\mathbf{Y}}(\mathbf{z}) d\mathbf{z} + \int_S [f(\mathbf{z}) = c] \cdot g_{\mathbf{Y}}(\mathbf{z}) d\mathbf{z} \right] \\ & \quad - \left[\int_S [f(\mathbf{z}) = c] \cdot g_{\mathbf{Y}}(\mathbf{z}) d\mathbf{z} + \int_S [f(\mathbf{z}) \neq c] \cdot g_{\mathbf{Y}}(\mathbf{z}) d\mathbf{z} \right] \\ & = \int_{S^c} [f(\mathbf{z}) = c] \cdot g_{\mathbf{Y}}(\mathbf{z}) d\mathbf{z} - \int_S [f(\mathbf{z}) \neq c] \cdot g_{\mathbf{Y}}(\mathbf{z}) d\mathbf{z} \\ & \stackrel{\text{sub-task 3}}{\geq} t' \left[\int_{S^c} [f(\mathbf{z}) = c] \cdot g_{\mathbf{X}}(\mathbf{z}) d\mathbf{z} - \int_S [f(\mathbf{z}) \neq c] \cdot g_{\mathbf{X}}(\mathbf{z}) d\mathbf{z} \right] \\ & \stackrel{\text{add 0}}{=} t' \left[\int_{S^c} [f(\mathbf{z}) = c] \cdot g_{\mathbf{X}}(\mathbf{z}) d\mathbf{z} - \int_S [f(\mathbf{z}) \neq c] \cdot g_{\mathbf{X}}(\mathbf{z}) d\mathbf{z} \right. \\ & \quad \left. + \int_S [f(\mathbf{z}) = c] \cdot g_{\mathbf{X}}(\mathbf{z}) d\mathbf{z} - \int_S [f(\mathbf{z}) = c] \cdot g_{\mathbf{X}}(\mathbf{z}) d\mathbf{z} \right] \\ & = t' \left[\int_{\mathbb{R}^d} [f(\mathbf{z}) = c] \cdot g_{\mathbf{X}}(\mathbf{z}) d\mathbf{z} - \int_S g_{\mathbf{X}}(\mathbf{z}) d\mathbf{z} \right] \\ & \stackrel{\text{definition}}{=} t' [\mathbb{P}(f(\mathbf{X}) = c) - \mathbb{P}(\mathbf{X} \in S)] \\ & \stackrel{\text{assumption}}{\geq} 0 \end{aligned}$$

5. The generalization for the second part is straight-forward and only requires to change the direction of some inequalities.

Problem 3. Randomized smoothing currently is usually formulated for the ℓ^1 or ℓ^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 l_2 -norm to obtain ℓ^∞ -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 $\mathbb{B}_e^p := \{\mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x}\|_p \leq e\}$ denote the l_p -norm ball of size e . 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 .

Solution 3.

1. $\|\mathbf{x}\|_\infty = \max_{i \in \{1, \dots, d\}} |x_i| \leq \sum_{i=1}^d |x_i| = \|\mathbf{x}\|_1$.
Thus $c_1 = 1$.
Example for tightness: $\|(\frac{1}{0})\|_\infty = 1 = \|(\frac{1}{0})\|_1$.
2. $\|\mathbf{x}\|_1 = \sum_{i=1}^d |x_i| \leq \sum_{i=1}^d \max_{i \in \{1, \dots, d\}} |x_i| = d \cdot \max_{i \in \{1, \dots, d\}} |x_i| = d \|\mathbf{x}\|_\infty$.
Thus $c_2 = d$.
Example for tightness: $\|(\frac{1}{1})\|_1 = 2 = 2 \cdot 1 = 2 \cdot \|(\frac{1}{1})\|_\infty$.
3. Let $i \in \operatorname{argmax}_{i \in \{1, \dots, d\}} |x_i|^2$. Then

$$\|\mathbf{x}\|_\infty^2 = |x_i|^2 \leq |x_i|^2 + \sum_{\substack{j=1 \\ j \neq i}}^d |x_j|^2 = \sum_{j=1}^d |x_j|^2 = \|\mathbf{x}\|_2^2.$$

The taking the square root yields $\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_2$.

Thus $c_3 = 1$.

Example for tightness: $\|(\frac{1}{0})\|_\infty = 1 = \|(\frac{1}{0})\|_2$.

4. $\|\mathbf{x}\|_2^2 = \sum_{i=1}^d |x_i|^2 \leq d \cdot \max_{i \in \{1, \dots, d\}} |x_i|^2 = d \|\mathbf{x}\|_\infty^2$. Taking the square root of both sides we obtain $\|\mathbf{x}\|_2 \leq \sqrt{d} \|\mathbf{x}\|_\infty$.

Thus $c_4 = \sqrt{d}$.

Example for tightness: $\|(\frac{1}{1})\|_2 = \sqrt{2} = \sqrt{2} \cdot 1 = \sqrt{2} \cdot \|(\frac{1}{1})\|_\infty$.

5. $\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^d |x_i|^2} \leq \sum_{i=1}^d \sqrt{|x_i|^2} = \sum_{i=1}^d |x_i| = \|\mathbf{x}\|_1$ where the inequality holds due to the subadditivity of $\sqrt{\cdot}$ as

$$x_1 + x_2 \leq x_1 + 2\sqrt{x_1 x_2} + x_2 = (\sqrt{x_1} + \sqrt{x_2})^2 \implies \sqrt{x_1 + x_2} \leq \sqrt{x_1} + \sqrt{x_2}$$

which by induction holds for arbitrarily many $x_i \geq 0$.

Thus $c_5 = 1$.

Example for tightness: $\|(\frac{1}{0})\|_2 = 1 = \|(\frac{1}{0})\|_1$.

6. Let $\mathbf{z} \in \mathbb{R}^d; z_i = \text{sgn}(x_i)$. $\|\mathbf{x}\|_1 = \sum_{i=1}^d \mathbf{x}_i \mathbf{z}_i = \langle \mathbf{x}, \mathbf{z} \rangle \leq \|\mathbf{x}\|_2 \|\mathbf{z}\|_2 = \sqrt{d} \|\mathbf{x}\|_2$. Where the inequality holds due to the Cauchy-Schwarz inequality:

$$|\langle \mathbf{a}, \mathbf{b} \rangle|^2 \leq \langle \mathbf{a}, \mathbf{a} \rangle \langle \mathbf{b}, \mathbf{b} \rangle, \text{ which implies } |\langle \mathbf{a}, \mathbf{b} \rangle| \leq \|\mathbf{a}\|_2 \|\mathbf{b}\|_2$$

Thus $c_6 = \sqrt{d}$.

Example for tightness: $\|(\frac{1}{1})\|_1 = 2 = \sqrt{2} \cdot \sqrt{2} = \sqrt{2} \cdot \|(\frac{1}{1})\|_2$.

7. **Claim** $\mathbb{B}_\epsilon^\infty \subseteq \mathbb{B}_{\epsilon \cdot \sqrt{d}}^2$.

We can show the claim by taking the result from sub-task 4 and rearrange it to show $\frac{1}{\sqrt{d}} \|\mathbf{x}\|_2 \leq \|\mathbf{x}\|_\infty$. and we know that $\forall \mathbf{x} \in \mathbb{B}_\epsilon^\infty. \|\mathbf{x}\| \leq \epsilon$. Combining these two insights we obtain:

$$\begin{aligned} \forall \mathbf{x} \in \mathbb{B}_\epsilon^\infty. \quad \|\mathbf{x}\|_\infty &\leq \epsilon \\ \implies \forall \mathbf{x} \in \mathbb{B}_\epsilon^\infty. \quad \frac{1}{\sqrt{d}} \|\mathbf{x}\|_2 &\leq \|\mathbf{x}\|_\infty \leq \epsilon \\ \implies \forall \mathbf{x} \in \mathbb{B}_\epsilon^\infty. \quad \|\mathbf{x}\|_2 &\leq \epsilon \cdot \sqrt{d} \end{aligned}$$

Note: It is crucial here to use the result of sub-task 4 ($\|\mathbf{x}\|_2 \leq \sqrt{d} \|\mathbf{x}\|_\infty$) rather than the more intuitive seeming result from sub-task 3 ($\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_2$) as these are statements about individual vectors. So for example $\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_2$ **does not** imply $\mathbb{B}_1^\infty \subseteq \mathbb{B}_1^2$ as $\|(\frac{1}{1})\|_\infty = 1 \leq \sqrt{2} = \|(\frac{1}{1})\|_2$.

Similarly we can show $\mathbb{B}_\epsilon^1 \subseteq \mathbb{B}_\epsilon^2$ by sub-task 5, $\mathbb{B}_\epsilon^2 \subseteq \mathbb{B}_\epsilon^\infty$ by sub-task 3, $\mathbb{B}_\epsilon^\infty \subseteq \mathbb{B}_{\epsilon \cdot d}^1$ by sub-task 2 and $\mathbb{B}_\epsilon^2 \subseteq \mathbb{B}_{\epsilon \cdot \sqrt{d}}^1$ by sub-task 6.

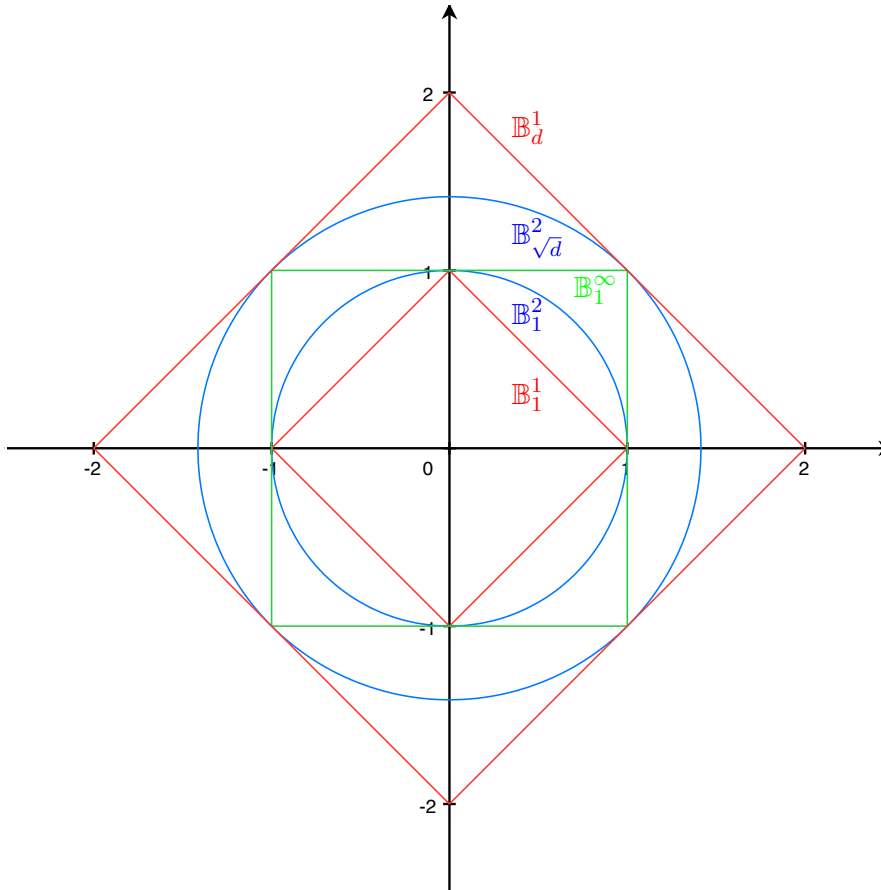


Figure 1: Boundaries of different normballs.

We thus obtain $\mathbb{B}_1^1 \subseteq \mathbb{B}_1^2 \subseteq \mathbb{B}_1^\infty \subseteq \mathbb{B}_{\sqrt{d}}^2 \subseteq \mathbb{B}_d^1$. Figure 1 shows this for $d = 2$.

Take-away message: If we prove a classifier to be safe for a l_2 -radius of ϵ that implies it is safe for a l_∞ radius of $\frac{\epsilon}{\sqrt{d}}$.

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>.