

## Appendix A Proof of Theorem 1

**Theorem 1.** Given a base classifier  $f$  and its smoothed version  $g: \mathbb{R}^d \rightarrow \mathcal{Y}$ , an example  $\mathbf{x}$ , a masked noise distribution  $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{M})$  with a binary mask  $\mathbf{M} \in \mathbb{R}^d$  where  $\mathbf{M}$  is all zero except that the selected part is one,  $c_A, c_B \in \mathcal{Y}$  and  $\underline{p}_A, \overline{p}_B \in [0, 1]$  that satisfy:

$$\mathbb{P}(f(\mathbf{x} + \epsilon) = c_A) \geq \underline{p}_A \geq \overline{p}_B \geq \max_{c \neq c_A} \mathbb{P}(f(\mathbf{x} + \epsilon) = c). \quad (\text{A1})$$

Then we have:

$$g(\mathbf{x} + \delta) = c_A, \quad \forall \|\delta\|_2 < R \quad (\text{A2})$$

where

$$R = \frac{\sigma}{2} (\Phi^{-1}(\underline{p}_A) - \Phi^{-1}(\overline{p}_B)). \quad (\text{A3})$$

*Proof.* Define two variables

$$\begin{aligned} X &:= \mathbf{x} + \epsilon = \mathcal{N}(\mathbf{x}, \sigma^2 \mathbf{M}) \\ Y &:= \mathbf{x} + \delta + \epsilon = \mathcal{N}(\mathbf{x} + \delta, \sigma^2 \mathbf{M}). \end{aligned}$$

We know from [1] [2] that if  $\exists \beta$  satisfying

$$\Omega = \left\{ \delta^T \mathbf{v} \leq \beta, \mathbf{v} \in \mathbb{R}^d \right\}$$

and

$$\mathbb{P}(f(X) = c_A) \geq \mathbb{P}(X \in \Omega),$$

then

$$\mathbb{P}(f(Y) = c_A) \geq \mathbb{P}(Y \in \Omega) \quad (\text{A4})$$

in which the inequality sign is reversed and  $c_A$  is replaced by  $c_B$ , and the lemma still holds. Define the half-spaces

$$\begin{aligned} \mathbb{A} &= \left\{ \mathbf{v} : \delta^T (\mathbf{v} - \mathbf{x}) \leq \sigma \|\delta\|_2 \Phi^{-1}(\underline{p}_A) \right\} \\ \mathbb{B} &= \left\{ \mathbf{v} : \delta^T (\mathbf{v} - \mathbf{x}) \geq \sigma \|\delta\|_2 \Phi^{-1}(1 - \overline{p}_B) \right\}, \end{aligned}$$

then we have

$$\begin{aligned} \mathbb{P}(X \in A) &= \mathbb{P}\left(\delta^T (X - \mathbf{x}) \leq \sigma \|\delta\|_2 \Phi^{-1}(\underline{p}_A)\right) \\ &= \mathbb{P}\left(\delta^T \mathcal{N}(0, \sigma^2 \mathbf{M}) \leq \sigma \|\delta\|_2 \Phi^{-1}(\underline{p}_A)\right) \\ &= \mathbb{P}\left(\sigma \|\mathbf{M} \odot \delta\|_2 \mathbb{Q} \leq \sigma \|\delta\|_2 \Phi^{-1}(\underline{p}_A)\right) \\ &= \Phi\left(\frac{\|\delta\|_2 \Phi^{-1}(\underline{p}_A)}{\|\mathbf{M} \odot \delta\|_2}\right), \mathbb{Q} \sim \mathcal{N}(0, 1) \end{aligned} \quad (\text{A5})$$

$$\begin{aligned} \mathbb{P}(X \in B) &= \mathbb{P}\left(\delta^T (X - \mathbf{x}) \geq \sigma \|\delta\|_2 \Phi^{-1}(1 - \overline{p}_B)\right) \\ &= \mathbb{P}\left(\delta^T \mathcal{N}(0, \sigma^2 \mathbf{M}) \geq \sigma \|\delta\|_2 \Phi^{-1}(1 - \overline{p}_B)\right) \\ &= \mathbb{P}\left(\sigma \|\mathbf{M} \odot \delta\|_2 \mathbb{Q} \geq \sigma \|\delta\|_2 \Phi^{-1}(1 - \overline{p}_B)\right) \\ &= \Phi\left(\frac{\|\delta\|_2 \Phi^{-1}(\overline{p}_B)}{\|\mathbf{M} \odot \delta\|_2}\right), \end{aligned} \quad (\text{A6})$$

$$\begin{aligned} \mathbb{P}(Y \in A) &= \mathbb{P}\left(\delta^T (Y - \mathbf{x}) \leq \sigma \|\delta\|_2 \Phi^{-1}(\underline{p}_A)\right) \\ &= \mathbb{P}\left(\delta^T \mathcal{N}(0, \sigma^2 \mathbf{M}) + \|\delta\|_2^2 \leq \sigma \|\delta\|_2 \Phi^{-1}(\underline{p}_A)\right) \\ &= \mathbb{P}\left(\sigma \|\mathbf{M} \odot \delta\|_2 \mathbb{Q} \leq \sigma \|\delta\|_2 \Phi^{-1}(\underline{p}_A) - \|\delta\|_2^2\right) \\ &= \mathbb{P}\left(\mathbb{Q} \leq \frac{\sigma \|\delta\|_2 \Phi^{-1}(\underline{p}_A) - \|\delta\|_2^2}{\sigma \|\mathbf{M} \odot \delta\|_2}\right) \\ &= \Phi\left(\frac{\sigma \|\delta\|_2 \Phi^{-1}(\underline{p}_A) - \|\delta\|_2^2}{\sigma \|\mathbf{M} \odot \delta\|_2}\right), \end{aligned} \quad (\text{A7})$$

$$\begin{aligned}
 \mathbb{P}(Y \in B) &= \mathbb{P}\left(\delta^T(Y - \mathbf{x}) \geq \sigma \|\delta\|_2 \Phi^{-1}(1 - \overline{p_B})\right) \\
 &= \mathbb{P}\left(\delta^T \mathcal{N}\left(0, \sigma^2 \mathbf{M}\right) + \|\delta\|_2^2 \geq \sigma \|\delta\|_2 \Phi^{-1}(1 - \overline{p_B})\right) \\
 &= \mathbb{P}\left(\sigma \|\mathbf{M} \odot \delta\|_2 \mathbb{Q} \geq \sigma \|\delta\|_2 \Phi^{-1}(1 - \overline{p_B}) - \|\delta\|_2^2\right) \\
 &= \mathbb{P}\left(\mathbb{Q} \geq \frac{\sigma \|\delta\|_2 \Phi^{-1}(1 - \overline{p_B}) - \|\delta\|_2^2}{\sigma \|\mathbf{M} \odot \delta\|_2}\right) \\
 &= \mathbb{P}\left(\mathbb{Q} \leq \frac{\sigma \|\delta\|_2 \Phi^{-1}(\overline{p_B}) + \|\delta\|_2^2}{\sigma \|\mathbf{M} \odot \delta\|_2}\right) \\
 &= \Phi\left(\frac{\sigma \|\delta\|_2 \Phi^{-1}(\overline{p_B}) + \|\delta\|_2^2}{\sigma \|\mathbf{M} \odot \delta\|_2}\right).
 \end{aligned} \tag{A8}$$

The mask  $\mathbf{M}$  covers the patch, so we have

$$\|\delta\|_2 = \|\mathbf{M} \odot \delta\|_2, \tag{A9}$$

thus we have

$$\mathbb{P}(X \in A) = \Phi\left(\Phi^{-1}(\underline{p_A})\right) = \underline{p_A} \tag{A10}$$

$$\mathbb{P}(X \in B) = \Phi\left(\Phi^{-1}(\overline{p_B})\right) = \overline{p_B} \tag{A11}$$

$$\mathbb{P}(f(X) = c_A) \geq \mathbb{P}(X \in A) = \underline{p_A} \tag{A12}$$

$$\mathbb{P}(f(X) = c_B) \leq \mathbb{P}(X \in B) = \overline{p_B}. \tag{A13}$$

From A4, we have

$$\mathbb{P}(f(Y) = c_A) \geq \mathbb{P}(Y \in A) \tag{A14}$$

$$\mathbb{P}(f(Y) = c_B) \leq \mathbb{P}(Y \in B). \tag{A15}$$

If we want to get  $\mathbb{P}(f(Y) = c_A) > \mathbb{P}(f(Y) = c_B)$ , just let  $\mathbb{P}(Y \in A) > \mathbb{P}(Y \in B)$ , we can get that  $\mathbb{P}(Y \in A) > \mathbb{P}(Y \in B)$  if and only if

$$R < \frac{\sigma}{2}(\Phi^{-1}(\underline{p_A}) - \Phi^{-1}(\overline{p_B})). \tag{A16}$$

## Appendix B Pseudocode for Practical Localization

---

### Algorithm B1 LOCALIZATION BASED ON JOINT VOTING

---

**Require:** Base classifier  $f$ , adversarial image  $\mathbf{x}$ , block set  $\mathbb{B}^{r \times r}$ , cross-shaped filter  $\mathbb{T}$ , threshold  $\tau$ .

**Ensure:** Patch location set  $\mathbb{L}$  or ABSTAIN.

```

1: initial  $\mathbb{L} \leftarrow \emptyset$ ,  $\mathbb{C}^{r \times r} \leftarrow \mathbf{0}$ , score set  $\mathbb{S} \leftarrow \emptyset$  and  $\text{pred} \leftarrow \mathbf{0}$ 
2: for  $b^{(i,j)}$  in  $\mathbb{B}$  do
3:    $\mathbf{x}^{(i,j)} \leftarrow \text{INPAINT}(\mathbf{x}, b^{(i,j)})$ 
4:    $\text{pred}(i, j) \leftarrow f(\mathbf{x}^{(i,j)})$ 
5: end for
6:  $\hat{l}_A \leftarrow$  top indices in  $\text{pred}$ 
7: for  $l^{(i,j)}$  in  $\text{pred}$  do
8:   if  $l^{(i,j)} \neq \hat{l}_A$  then
9:      $\mathbb{C}[i, j] \leftarrow 1$ 
10:   end if
11: end for
12: for  $\mathbb{C}[p, q] \in \mathbb{C}$  do
13:    $\mathbb{S}[p, q] \leftarrow \text{SUM}(\mathbb{T}_{\mathbb{C}[p, q]} \odot \mathbb{C})$ 
14: end for
15:  $s_A \leftarrow \max(\mathbb{S})$ 
16: for  $\mathbb{S}[m, n] \in \mathbb{S}$  do
17:   if  $\mathbb{S}[m, n] = s_A$  then
18:      $\mathbb{L} \leftarrow \mathbb{L} \cup \mathbb{B}[m, n]$ 
19:   end if
20: end for
21: for  $\mathbb{L}[m, n] \in \mathbb{L}$  do
22:   for  $\mathbb{C}[p, q] \in \mathbb{C}$  do
23:     if  $|p - m| < \tau \wedge |q - n| < \tau$  then
24:        $\mathbb{L} \leftarrow \mathbb{L} \cup \mathbb{B}[p, q]$ 
25:     end if
26:   end for
27: end for
28: return  $\mathbb{L}$ 

```

---

The pseudocode of localization is described in algorithm B1. Function  $\text{INPAINT}(\mathbf{x}, b_i)$  reconstructs block  $b_i$  using its neighbor pixels. The filter is centered on the currently calculated block and takes out the values in the filter. Different from the certification, the process of the prediction will perform a binomial test on the number of  $c_A$  and  $c_B$  to ensure that the majority voted category is correct with high probability.

## Appendix C Other Results

### Appendix C.1 More comparisons with DRS

For DRS, the certified condition is  $\text{Counts}[c] > \text{Counts}[c_{next}] + 2 \times \text{AffectedBands}$ , which implies that the larger the patch is, the harder the condition is reached. So we compare the performance to Column smoothing of DRS when assuming the patch size is  $10 \times 10$  and  $15 \times 15$  on CIFAR-10 datasets. Table C1 shows that our method is much better than DRS on large patches.

**Table C1** Certified accuracy of our method and DRS (Column smoothing) on CIFAR-10 when the patch size is  $10 \times 10$  and  $15 \times 15$ . We use 100-accuracy for the  $10 \times 10$  patch and 225-accuracy for the  $15 \times 15$  patch.

	( $10 \times 10$ )	( $15 \times 15$ )
DRS	29.0	1.9
Ours	60.0	32.7

### Appendix C.2 Accuracy of Localization

From Appendix A, we only need to ensure that the localization can completely cover the patch to calculate a credible certificate, so the localization accuracy (LA) is defined as the ratio of images where the patch is completely covered to the total test images. The coverage area is defined as the ratio of the mask covering the patch to the image, and we count the average coverage area (ACA) of the successfully localized images to help measure the effectiveness of the localization algorithm. The larger the image resolution, the more accurate the localization, in agreement with the experimental part. And as long as the localization algorithm provides the mask, the certificate provided by MRS can be obtained.

**Table C2** The localization accuracy (LA) and average coverage area (ACA) on CIFAR-10 and ImageNet. The larger the LA and the smaller the ACA the better the localization algorithm performs.

	LA	ACA
CIFAR-10	96.2	12.2
ImageNet	99.8	11.7

## References

- 1 Neyman J, Pearson E S. On the problem of the most efficient tests of statistical hypotheses. *Philosophical Transactions of the Royal Society A*, 1933, 231: 289-337
- 2 Cohen J, Rosenfeld E, Kolter Z. Certified adversarial robustness via randomized smoothing. In: *Proceedings of 36th International Conference on Machine Learning*, 2019. 1310-1320