

## A PROOFS FOR $l_2$ NORM

In this section, we present the full proofs for Theorem. 4.1. The main tool for our proofs is the Neyman-Pearson lemma for two variables, which we establish in Appendix A.1. Based on this lemma, we obtain the certified result in Appendix A.2. Finally, we provide the details of the linear transformation used for certification in Appendix A.3.

### A.1 Neyman-Pearson for Two Variables

LEMMA A.1 (NEYMAN-PEARSON FOR TWO VARIABLES). *Let  $X_1$  and  $Y_1$  be random variables in  $\mathbb{R}^d$  with densities  $\mu_{X_1}$  and  $\mu_{Y_1}$ . Then, let  $X_2$  and  $Y_2$  be random variables in  $\mathbb{R}^d$  with densities  $\mu_{X_2}$  and  $\mu_{Y_2}$ . Let  $h : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \{0, 1\}$  be any deterministic or random function with an input pair. Then:*

1. *If  $\mathcal{S}_1 \times \mathcal{S}_2 = \left\{z_1 \in \mathbb{R}^d, z_2 \in \mathbb{R}^d : \frac{\mu_{Y_1}(z_1)\mu_{Y_2}(z_2)}{\mu_{X_1}(z_1)\mu_{X_2}(z_2)} \leq t\right\}$  for some  $t > 0$  and  $P(h(X_1, X_2) = 1) \geq P((X_1, X_2) \in \mathcal{S}_1 \times \mathcal{S}_2)$ , then  $P(h(Y_1, Y_2) = 1) \geq P((Y_1, Y_2) \in \mathcal{S}_1 \times \mathcal{S}_2)$ .*
2. *If  $\mathcal{S}_1 \times \mathcal{S}_2 = \left\{z_1 \in \mathbb{R}^d, z_2 \in \mathbb{R}^d : \frac{\mu_{Y_1}(z_1)\mu_{Y_2}(z_2)}{\mu_{X_1}(z_1)\mu_{X_2}(z_2)} \geq t\right\}$  for some  $t > 0$  and  $P(h(X_1, X_2) = 1) \leq P((X_1, X_2) \in \mathcal{S}_1 \times \mathcal{S}_2)$ , then  $P(h(Y_1, Y_2) = 1) \leq P((Y_1, Y_2) \in \mathcal{S}_1 \times \mathcal{S}_2)$ .*

PROOF. We denote the complement of  $\mathcal{S}_1 \times \mathcal{S}_2$  as  $\mathcal{S}^c$ .

$$\begin{aligned}
 & P(h(Y_1, Y_2) = 1) - P((Y_1, Y_2) \in \mathcal{S}_1 \times \mathcal{S}_2) \\
 &= \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} h(1 \mid z_1, z_2) \mu_{Y_1}(z_1) \mu_{Y_2}(z_2) dz_1 dz_2 - \int \int_{\mathcal{S}_1 \times \mathcal{S}_2} \mu_{Y_1}(z_1) \mu_{Y_2}(z_2) dz_1 dz_2 \\
 &= \left[ \int \int_{\mathcal{S}_1 \times \mathcal{S}_2} h(1 \mid z_1, z_2) \mu_{Y_1}(z_1) \mu_{Y_2}(z_2) dz_1 dz_2 + \int \int_{\mathcal{S}^c} h(1 \mid z_1, z_2) \mu_{Y_1}(z_1) \mu_{Y_2}(z_2) dz_1 dz_2 \right] \\
 &\quad - \left[ \int \int_{\mathcal{S}_1 \times \mathcal{S}_2} h(1 \mid z_1, z_2) \mu_{Y_1}(z_1) \mu_{Y_2}(z_2) dz_1 dz_2 + \int \int_{\mathcal{S}_1 \times \mathcal{S}_2} h(0 \mid z_1, z_2) \mu_{Y_1}(z_1) \mu_{Y_2}(z_2) dz_1 dz_2 \right] \\
 &= \int \int_{\mathcal{S}^c} h(1 \mid z_1, z_2) \mu_{Y_1}(z_1) \mu_{Y_2}(z_2) dz_1 dz_2 - \int \int_{\mathcal{S}_1 \times \mathcal{S}_2} h(0 \mid z_1, z_2) \mu_{Y_1}(z_1) \mu_{Y_2}(z_2) dz_1 dz_2 \\
 &\geq t \left[ \int \int_{\mathcal{S}^c} h(1 \mid z_1, z_2) \mu_{X_1}(z_1) \mu_{X_2}(z_2) dz_1 dz_2 - \int \int_{\mathcal{S}_1 \times \mathcal{S}_2} h(0 \mid z_1, z_2) \mu_{X_1}(z_1) \mu_{X_2}(z_2) dz_1 dz_2 \right] \\
 &= t \left[ \int \int_{\mathcal{S}^c} h(1 \mid z_1, z_2) \mu_{X_1}(z_1) \mu_{X_2}(z_2) dz_1 dz_2 + \int \int_{\mathcal{S}_1 \times \mathcal{S}_2} h(1 \mid z_1, z_2) \mu_{X_1}(z_1) \mu_{X_2}(z_2) dz_1 dz_2 \right. \\
 &\quad \left. - \int \int_{\mathcal{S}_1 \times \mathcal{S}_2} h(1 \mid z_1, z_2) \mu_{X_1}(z_1) \mu_{X_2}(z_2) dz_1 dz_2 - \int \int_{\mathcal{S}_1 \times \mathcal{S}_2} h(0 \mid z_1, z_2) \mu_{X_1}(z_1) \mu_{X_2}(z_2) dz_1 dz_2 \right] \\
 &= t \left[ \int \int_{\mathbb{R}^d} h(1 \mid z_1, z_2) \mu_{X_1}(z_1) \mu_{X_2}(z_2) dz_1 dz_2 - \int \int_{\mathcal{S}_1 \times \mathcal{S}_2} \mu_{X_1}(z_1) \mu_{X_2}(z_2) dz_1 dz_2 \right] \\
 &= t [P(h(X_1, X_2) = 1) - P((X_1, X_2) \in \mathcal{S}_1 \times \mathcal{S}_2)] \\
 &\geq 0
 \end{aligned}$$

□

Next, we prove Lemma A.2, which is a special case of Lemma A.1 and states the Neyman-Pearson lemma for two joint Gaussian noise variables.

LEMMA A.2 (NEYMAN-PEARSON FOR TWO JOINT GAUSSIAN NOISE). *Let  $X_1 \sim \mathcal{N}(x_1, \Sigma_1)$ ,  $X_2 \sim \mathcal{N}(x_2, \Sigma_2)$  and  $Y_1 \sim \mathcal{N}(x_1 + \delta_1, \Sigma_1)$ ,  $Y_2 \sim \mathcal{N}(x_2 + \delta_2, \Sigma_2)$ . Let  $h : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \{0, 1\}$  be any deterministic or random function. Then:*

1. *If  $\mathcal{S}_1 \times \mathcal{S}_2 = \left\{z_1 \in \mathbb{R}^d, z_2 \in \mathbb{R}^d : \delta_1^\top \Sigma_1^{-1} z_1 + \delta_2^\top \Sigma_2^{-1} z_2 \leq \beta\right\}$  for some  $\beta$  and  $P(h(X_1, X_2) = 1) \geq P((X_1, X_2) \in \mathcal{S}_1 \times \mathcal{S}_2)$ , then  $P(h(Y_1, Y_2) = 1) \geq P((Y_1, Y_2) \in \mathcal{S}_1 \times \mathcal{S}_2)$ .*
2. *If  $\mathcal{S}_1 \times \mathcal{S}_2 = \left\{z_1 \in \mathbb{R}^d, z_2 \in \mathbb{R}^d : \delta_1^\top \Sigma_1^{-1} z_1 + \delta_2^\top \Sigma_2^{-1} z_2 \geq \beta\right\}$  for some  $\beta$  and  $P(h(X_1, X_2) = 1) \leq P((X_1, X_2) \in \mathcal{S}_1 \times \mathcal{S}_2)$ , then  $P(h(Y_1, Y_2) = 1) \leq P((Y_1, Y_2) \in \mathcal{S}_1 \times \mathcal{S}_2)$ .*

PROOF. This lemma is the special case of Neyman-Pearson for two variables when  $X_1, X_2, Y_1$ , and  $Y_2$  are joint Gaussian noises. It suffices to simply show that for any  $\beta$ , there is some  $t > 0$  for which:

$$\begin{aligned}
 \{z_1, z_2 : \delta_1^\top \Sigma_1^{-1} z_1 + \delta_2^\top \Sigma_2^{-1} z_2 \leq \beta\} &= \left\{z_1, z_2 : \frac{\mu_{Y_1}(z_1)\mu_{Y_2}(z_2)}{\mu_{X_1}(z_1)\mu_{X_2}(z_2)} \leq t\right\}, \\
 \{z_1, z_2 : \delta_1^\top \Sigma_1^{-1} z_1 + \delta_2^\top \Sigma_2^{-1} z_2 \geq \beta\} &= \left\{z_1, z_2 : \frac{\mu_{Y_1}(z_1)\mu_{Y_2}(z_2)}{\mu_{X_1}(z_1)\mu_{X_2}(z_2)} \geq t\right\}.
 \end{aligned} \tag{18}$$

For ease of representation, we use  $M_1 \in \mathbb{R}^{d \times d}$  (with element  $m_{1ij}$ ) instead of  $\Sigma_1^{-1}$  and  $M_2 \in \mathbb{R}^{d \times d}$  (with element  $m_{2ij}$ ) instead of  $\Sigma_2^{-1}$ . The likelihood ratio for this choice of  $X_1, X_2, Y_1$  and  $Y_2$  turns out to be:

$$\begin{aligned}
 & \frac{\mu_{Y_1}(z_1)\mu_{Y_2}(z_2)}{\mu_{X_1}(z_1)\mu_{X_2}(z_2)} \\
 &= \frac{\exp\left(-\frac{1}{2}(z_1 - (x_1 + \delta_1))^T \Sigma_1^{-1}(z_1 - (x_1 + \delta_1))\right)}{\exp\left(-\frac{1}{2}(z_1 - x_1)^T \Sigma_1^{-1}(z_1 - x_1)\right)} \times \frac{\exp\left(-\frac{1}{2}(z_2 - (x_2 + \delta_2))^T \Sigma_2^{-1}(z_2 - (x_2 + \delta_2))\right)}{\exp\left(-\frac{1}{2}(z_2 - x_2)^T \Sigma_2^{-1}(z_2 - x_2)\right)} \\
 &= \frac{\exp\left(-\frac{1}{2} \sum_i^d \sum_j^d (z_{1i} - (x_{1i} + \delta_{1i})) m_{1ij} (z_{1j} - (x_{1j} + \delta_{1j}))\right)}{\exp\left(-\frac{1}{2} \sum_i^d \sum_j^d (z_{1i} - x_{1i}) m_{1ij} (z_{1j} - x_{1j})\right)} \\
 &\quad \times \frac{\exp\left(-\frac{1}{2} \sum_i^d \sum_j^d (z_{2i} - (x_{2i} + \delta_{2i})) m_{2ij} (z_{2j} - (x_{2j} + \delta_{2j}))\right)}{\exp\left(-\frac{1}{2} \sum_i^d \sum_j^d (z_{2i} - x_{2i}) m_{2ij} (z_{2j} - x_{2j})\right)} \\
 &= \exp\left(\delta_1^T \Sigma_1^{-1} z_1 - \delta_1^T \Sigma_1^{-1} x_1 - \frac{1}{2} \delta_1^T \Sigma_1^{-1} \delta_1\right) \times \exp\left(\delta_2^T \Sigma_2^{-1} z_2 - \delta_2^T \Sigma_2^{-1} x_2 - \frac{1}{2} \delta_2^T \Sigma_2^{-1} \delta_2\right) \\
 &= \exp\left(\delta_1^T \Sigma_1^{-1} z_1 + \delta_2^T \Sigma_2^{-1} z_2 - \delta_1^T \Sigma_1^{-1} x_1 - \frac{1}{2} \delta_1^T \Sigma_1^{-1} \delta_1 - \delta_2^T \Sigma_2^{-1} x_2 - \frac{1}{2} \delta_2^T \Sigma_2^{-1} \delta_2\right) \\
 &= \exp\left(\delta_1^T \Sigma_1^{-1} z_1 + \delta_2^T \Sigma_2^{-1} z_2 + b\right) \leq t,
 \end{aligned}$$

where  $b$  is a constant, specifically  $b = -\delta_1^T \Sigma_1^{-1} x_1 - \frac{1}{2} \delta_1^T \Sigma_1^{-1} \delta_1 - \delta_2^T \Sigma_2^{-1} x_2 - \frac{1}{2} \delta_2^T \Sigma_2^{-1} \delta_2$ . Therefore given any  $\beta$ , we may take  $t = \exp(\beta + b)$  and obtain this correlation:

$$\begin{aligned}
 \delta_1^T \Sigma_1^{-1} z_1 + \delta_2^T \Sigma_2^{-1} z_2 \leq \beta &\iff \exp\left(\delta_1^T \Sigma_1^{-1} z_1 + \delta_2^T \Sigma_2^{-1} z_2 + b\right) \leq t, \\
 \delta_1^T \Sigma_1^{-1} z_1 + \delta_2^T \Sigma_2^{-1} z_2 \geq \beta &\iff \exp\left(\delta_1^T \Sigma_1^{-1} z_1 + \delta_2^T \Sigma_2^{-1} z_2 + b\right) \geq t.
 \end{aligned} \tag{19}$$

□

## A.2 Proof of the Certified Robustness

This subsection presents the logic for proving robustness guarantees and derives the certified spaces for these guarantees in Eq. 9.

To show that  $g_0(\mathbf{c}^1, \mathbf{c}^2, \mathbf{z}^1 + \delta_1, \mathbf{z}^2 + \delta_2) = 1$ , it follows from the definition of  $g_0$  that we need to show that:

$$P(f(\mathbf{c}^1, \mathbf{c}^2, \mathbf{z}^1 + \varepsilon_1 + \delta_1, \mathbf{z}^2 + \varepsilon_2 + \delta_2) \in \mathcal{X}') \geq P(f(\mathbf{c}^1, \mathbf{c}^2, \mathbf{z}^1 + \varepsilon_1 + \delta_1, \mathbf{z}^2 + \varepsilon_2 + \delta_2) \notin \mathcal{X}').$$

We define two random variables:

$$\begin{aligned}
 I &:= (\mathbf{c}^1, \mathbf{c}^2, \mathbf{z}^1 + \varepsilon_1, \mathbf{z}^2 + \varepsilon_2) = (\mathbf{c}^1, \mathbf{c}^2, \mathcal{N}(\mathbf{z}^1, \Sigma_1), \mathcal{N}(\mathbf{z}^2, \Sigma_2)) \\
 O &:= (\mathbf{c}^1, \mathbf{c}^2, \mathbf{z}^1 + \varepsilon_1 + \delta_1, \mathbf{z}^2 + \varepsilon_2 + \delta_2) = (\mathbf{c}^1, \mathbf{c}^2, \mathcal{N}(\mathbf{z}^1 + \delta_1, \Sigma_1), \mathcal{N}(\mathbf{z}^2 + \delta_2, \Sigma_2)).
 \end{aligned}$$

We know that:

$$P(f(I) \in \mathcal{X}') \geq \underline{p}. \tag{20}$$

Our goal is to show that

$$P(f(O) \in \mathcal{X}') > P(f(O) \notin \mathcal{X}'). \tag{21}$$

According to lemma A.2, we can define the half-spaces:

$$\begin{aligned}
 \mathcal{A} &= \left\{z_1, z_2 : \delta_1^T \Sigma_1^{-1}(z_1 - \mathbf{z}^1) + \delta_2^T \Sigma_2^{-1}(z_2 - \mathbf{z}^2) \leq \|\delta_1^T \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^T \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1}(\underline{p})\right\}, \\
 \mathcal{B} &= \left\{z_1, z_2 : \delta_1^T \Sigma_1^{-1}(z_1 - \mathbf{z}^1) + \delta_2^T \Sigma_2^{-1}(z_2 - \mathbf{z}^2) \geq \|\delta_1^T \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^T \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1}(\underline{p})\right\}.
 \end{aligned}$$

Claim 1 shows that  $P(I \in \mathcal{A}) = \underline{p}$ , therefore we can obtain  $P(f(I) \in \mathcal{X}') \geq P(I \in \mathcal{A})$ . Hence we may apply Lemma A.2 to conclude:

$$P(f(O) \in \mathcal{X}') \geq P(O \in \mathcal{A}). \tag{22}$$

Similarly, we obtain  $P(f(I) \notin \mathcal{X}') \leq P(I \in \mathcal{B})$ . Hence we may apply Lemma A.2 to conclude:

$$P(f(O) \notin \mathcal{X}') \leq P(O \in \mathcal{B}). \tag{23}$$

Combining Eq. 22 and 23, we can obtain the conditions of Eq. 21:

$$P(f(O) \in \mathcal{X}') \geq P(O \in \mathcal{A}) > P(O \in \mathcal{B}) \geq P(f(O) \notin \mathcal{X}'). \tag{24}$$

According to Claim 3 and Claim 4, we can obtain  $P(O \in \mathcal{A})$  and  $P(O \in \mathcal{B})$  as:

$$\begin{aligned} P(O \in \mathcal{A}) &= \Phi \left( \Phi^{-1} \left( \underline{p} \right) - \frac{\delta_1^\top \Sigma_1^{-1} \delta_1 + \delta_2^\top \Sigma_2^{-1} \delta_2}{\|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\|} \right), \\ P(O \in \mathcal{B}) &= \Phi \left( -\Phi^{-1} \left( \underline{p} \right) + \frac{\delta_1^\top \Sigma_1^{-1} \delta_1 + \delta_2^\top \Sigma_2^{-1} \delta_2}{\|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\|} \right). \end{aligned} \quad (25)$$

Finally, we obtain that  $P(O \in \mathcal{A}) > P(O \in \mathcal{B})$  if and only if:

$$\frac{\delta_1^\top \Sigma_1^{-1} \delta_1 + \delta_2^\top \Sigma_2^{-1} \delta_2}{\|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\|} < \Phi^{-1} \left( \underline{p} \right).$$

### A.3 Linear Transformation and Derivation

This subsection begins with Lemma A.3, which is the main tool for deriving all claims. Then, we present the proof process of claims, which is applied in Sec. A.2.

**LEMMA A.3 (JOINT GAUSSIAN DISTRIBUTION).** *If there is a random matrix  $X \sim \mathcal{N}(\mu, \Sigma)$ , where  $\mu \in \mathbb{R}^n$  is the mean matrix. A positive semi-definite real symmetric matrix  $\Sigma \in \mathbb{S}_{++}^{n \times n}$  is the covariance matrix of  $X$ . There is a full rank matrix  $\mathbf{B} \in \mathbb{R}^{n \times n}$ , which makes  $X = \mathbf{B}Z + \mu$ ,  $Z \sim \mathcal{N}(\mathbf{0}, I)$  and  $\mathbf{B}^\top \mathbf{B} = \Sigma$ .*

We obtain four claims based on linear transformation:

**Claim 1.**  $P(I \in \mathcal{A}) = \underline{p}$

**PROOF.** Recall that  $\mathcal{A} = \left\{ z_1, z_2 : \delta_1^\top \Sigma_1^{-1} (z_1 - z^1) + \delta_2^\top \Sigma_2^{-1} (z_2 - z^2) \leq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1} \left( \underline{p} \right) \right\}$ , according to lemma A.3, we can obtain:

$$\begin{aligned} P(I \in \mathcal{A}) &= P \left( \delta_1^\top \Sigma_1^{-1} (\mathcal{N}(z^1, \Sigma_1) - z^1) + \delta_2^\top \Sigma_2^{-1} (\mathcal{N}(z^2, \Sigma_2) - z^2) \leq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1} \left( \underline{p} \right) \right) \\ &= P \left( \delta_1^\top \Sigma_1^{-1} \mathcal{N}(0, \Sigma_1) + \delta_2^\top \Sigma_2^{-1} \mathcal{N}(0, \Sigma_2) \leq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1} \left( \underline{p} \right) \right) \\ &= P \left( \delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 \mathcal{N}(0, I) + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2 \mathcal{N}(0, I) \leq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1} \left( \underline{p} \right) \right) \\ &= P \left( \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \mathcal{N}(0, 1) \leq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1} \left( \underline{p} \right) \right) \\ &= \Phi \left( \Phi^{-1} \left( \underline{p} \right) \right) \\ &= \underline{p}. \end{aligned}$$

□

**Claim 2.**  $P(I \in \mathcal{B}) = 1 - \underline{p}$

**PROOF.** Recall that  $\mathcal{B} = \left\{ z_1, z_2 : \delta_1^\top \Sigma_1^{-1} (z_1 - z^1) + \delta_2^\top \Sigma_2^{-1} (z_2 - z^2) \geq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1} \left( \underline{p} \right) \right\}$ , according to lemma A.3, we can obtain:

$$\begin{aligned} P(I \in \mathcal{B}) &= P \left( \delta_1^\top \Sigma_1^{-1} (\mathcal{N}(z^1, \Sigma_1) - z^1) + \delta_2^\top \Sigma_2^{-1} (\mathcal{N}(z^2, \Sigma_2) - z^2) \geq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1} \left( \underline{p} \right) \right) \\ &= P \left( \delta_1^\top \Sigma_1^{-1} \mathcal{N}(0, \Sigma_1) + \delta_2^\top \Sigma_2^{-1} \mathcal{N}(0, \Sigma_2) \geq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1} \left( \underline{p} \right) \right) \\ &= P \left( \delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 \mathcal{N}(0, I) + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2 \mathcal{N}(0, I) \geq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1} \left( \underline{p} \right) \right) \\ &= P \left( \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \mathcal{N}(0, 1) \geq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1} \left( \underline{p} \right) \right) \\ &= 1 - \Phi \left( \Phi^{-1} \left( \underline{p} \right) \right) \\ &= 1 - \underline{p}. \end{aligned}$$

□

**Claim 3.**  $P(O \in \mathcal{A}) = \Phi \left( \Phi^{-1} \left( \underline{p} \right) - \frac{\delta_1^\top \Sigma_1^{-1} \delta_1 + \delta_2^\top \Sigma_2^{-1} \delta_2}{\|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\|} \right)$

PROOF. Recall that  $\mathcal{A} = \{z_1, z_2 : \delta_1^\top \Sigma_1^{-1}(z_1 - z^1) + \delta_2^\top \Sigma_2^{-1}(z_2 - z^2) \leq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1}(\underline{p})\}$  and  $O \sim (\mathbf{c}^1, \mathbf{c}^2, \mathcal{N}(z^1 + \delta_1, \Sigma_1), \mathcal{N}(z^2 + \delta_2, \Sigma_2))$ , according to lemma A.3, we can obtain:

$$\begin{aligned}
P(O \in \mathcal{A}) &= P\left(\delta_1^\top \Sigma_1^{-1}(\mathcal{N}(z^1 + \delta_1, \Sigma_1) - z^1) + \delta_2^\top \Sigma_2^{-1}(\mathcal{N}(z^2 + \delta_2, \Sigma_2) - z^2) \leq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1}(\underline{p})\right) \\
&= P\left(\delta_1^\top \Sigma_1^{-1} \mathcal{N}(\delta_1, \Sigma_1) + \delta_2^\top \Sigma_2^{-1} \mathcal{N}(\delta_2, \Sigma_2) \leq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1}(\underline{p})\right) \\
&= P\left(\delta_1^\top \Sigma_1^{-1}(\mathbf{B}_1 \mathcal{N}(0, I) + \delta_1) + \delta_2^\top \Sigma_2^{-1}(\mathbf{B}_2 \mathcal{N}(0, I) + \delta_2) \leq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1}(\underline{p})\right) \\
&= P\left(\|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \mathcal{N}(0, 1) + \delta_1^\top \Sigma_1^{-1} \delta_1 + \delta_2^\top \Sigma_2^{-1} \delta_2 \leq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1}(\underline{p})\right) \\
&= P\left(\mathcal{N}(0, 1) \leq \Phi^{-1}(\underline{p}) - \frac{\delta_1^\top \Sigma_1^{-1} \delta_1 + \delta_2^\top \Sigma_2^{-1} \delta_2}{\|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\|}\right) \\
&= \Phi\left(\Phi^{-1}(\underline{p}) - \frac{\delta_1^\top \Sigma_1^{-1} \delta_1 + \delta_2^\top \Sigma_2^{-1} \delta_2}{\|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\|}\right).
\end{aligned}$$

□

**Claim 4.**  $P(O \in \mathcal{B}) = \Phi\left(-\Phi^{-1}(\underline{p}) + \frac{\delta_1^\top \Sigma_1^{-1} \delta_1 + \delta_2^\top \Sigma_2^{-1} \delta_2}{\|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\|}\right)$

PROOF. Recall that  $\mathcal{B} = \{z_1, z_2 : \delta_1^\top \Sigma_1^{-1}(z_1 - z^1) + \delta_2^\top \Sigma_2^{-1}(z_2 - z^2) \geq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1}(\underline{p})\}$  and  $O \sim (\mathbf{c}^1, \mathbf{c}^2, \mathcal{N}(z^1 + \delta_1, \Sigma_1), \mathcal{N}(z^2 + \delta_2, \Sigma_2))$ , according to lemma A.3, we can obtain:

$$\begin{aligned}
P(O \in \mathcal{B}) &= P\left(\delta_1^\top \Sigma_1^{-1}(\mathcal{N}(z^1 + \delta_1, \Sigma_1) - z^1) + \delta_2^\top \Sigma_2^{-1}(\mathcal{N}(z^2 + \delta_2, \Sigma_2) - z^2) \geq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1}(\underline{p})\right) \\
&= P\left(\delta_1^\top \Sigma_1^{-1} \mathcal{N}(\delta_1, \Sigma_1) + \delta_2^\top \Sigma_2^{-1} \mathcal{N}(\delta_2, \Sigma_2) \geq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1}(\underline{p})\right) \\
&= P\left(\delta_1^\top \Sigma_1^{-1}(\mathbf{B}_1 \mathcal{N}(0, I) + \delta_1) + \delta_2^\top \Sigma_2^{-1}(\mathbf{B}_2 \mathcal{N}(0, I) + \delta_2) \geq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1}(\underline{p})\right) \\
&= P\left(\|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \mathcal{N}(0, 1) + \delta_1^\top \Sigma_1^{-1} \delta_1 + \delta_2^\top \Sigma_2^{-1} \delta_2 \geq \|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\| \Phi^{-1}(\underline{p})\right) \\
&= P\left(\mathcal{N}(0, 1) \geq \Phi^{-1}(\underline{p}) - \frac{\delta_1^\top \Sigma_1^{-1} \delta_1 + \delta_2^\top \Sigma_2^{-1} \delta_2}{\|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\|}\right) \\
&= \Phi\left(-\Phi^{-1}(\underline{p}) + \frac{\delta_1^\top \Sigma_1^{-1} \delta_1 + \delta_2^\top \Sigma_2^{-1} \delta_2}{\|\delta_1^\top \Sigma_1^{-1} \mathbf{B}_1 + \delta_2^\top \Sigma_2^{-1} \mathbf{B}_2\|}\right).
\end{aligned}$$

□

## B PROOF FOR $l_1$ NORM

In this section, we present the full proofs for the robustness guarantee for  $l_1$  norm. The main tool for our proofs is the Neyman-Pearson lemma for two variables, which we establish in Lemma A.1. Next, we prove Lemma B.1, which is a special case of Lemma A.1 and states the Neyman-Pearson lemma for two Laplace noise variables. Based on this lemma, we obtain the certified result in Appendix B.1.

**LEMMA B.1 (NEYMAN-PEARSON FOR TWO LAPLACE NOISE).** *Let  $X_1 \sim x_1 + \mathcal{L}(\lambda_1)$ ,  $X_2 \sim x_2 + \mathcal{L}(\lambda_2)$  and  $Y_1 \sim x_1 + \mathcal{L}(\lambda_1) + \delta_1$ ,  $Y_2 \sim x_2 + \mathcal{L}(\lambda_2) + \delta_2$ . Let  $h : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \{0, 1\}$  be any deterministic or random function. Then:*

1. *If  $\mathcal{S}_1 \times \mathcal{S}_2 = \left\{z_1 \in \mathbb{R}^d, z_2 \in \mathbb{R}^d : \frac{1}{\lambda_1}(\|z_1 - \delta_1\|_1 - \|z_1\|_1) + \frac{1}{\lambda_2}(\|z_2 - \delta_2\|_1 - \|z_2\|_1)\right\} \geq \beta$  for some  $\beta$  and  $P(h(X_1, X_2) = 1) \geq P((X_1, X_2) \in \mathcal{S}_1 \times \mathcal{S}_2)$ , then  $P(h(Y_1, Y_2) = 1) \geq P((Y_1, Y_2) \in \mathcal{S}_1 \times \mathcal{S}_2)$ .*
2. *If  $\mathcal{S}_1 \times \mathcal{S}_2 = \left\{z_1 \in \mathbb{R}^d, z_2 \in \mathbb{R}^d : \frac{1}{\lambda_1}(\|z_1 - \delta_1\|_1 - \|z_1\|_1) + \frac{1}{\lambda_2}(\|z_2 - \delta_2\|_1 - \|z_2\|_1)\right\} \leq \beta$  for some  $\beta$  and  $P(h(X_1, X_2) = 1) \leq P((X_1, X_2) \in \mathcal{S}_1 \times \mathcal{S}_2)$ , then  $P(h(Y_1, Y_2) = 1) \leq P((Y_1, Y_2) \in \mathcal{S}_1 \times \mathcal{S}_2)$ .*

PROOF. This lemma is the special case of Neyman-Pearson for two variables when  $X_1, X_2, Y_1$ , and  $Y_2$  are Laplace noises. It suffices to simply show that for any  $\beta$ , there is some  $t > 0$  for which:

$$\begin{aligned} \left\{ z_1, z_2 : \frac{1}{\lambda_1} (\|z_1 - \delta_1\|_1 - \|z_1\|_1) + \frac{1}{\lambda_2} (\|z_2 - \delta_2\|_1 - \|z_2\|_1) \geq \beta \right\} &= \left\{ z_1, z_2 : \frac{\mu_{Y_1}(z_1)\mu_{Y_2}(z_2)}{\mu_{X_1}(z_1)\mu_{X_2}(z_2)} \leq t \right\}, \\ \left\{ z_1, z_2 : \frac{1}{\lambda_1} (\|z_1 - \delta_1\|_1 - \|z_1\|_1) + \frac{1}{\lambda_2} (\|z_2 - \delta_2\|_1 - \|z_2\|_1) \leq \beta \right\} &= \left\{ z_1, z_2 : \frac{\mu_{Y_1}(z_1)\mu_{Y_2}(z_2)}{\mu_{X_1}(z_1)\mu_{X_2}(z_2)} \geq t \right\}. \end{aligned} \quad (26)$$

$$\begin{aligned} &\frac{\mu_{Y_1}(z_1)\mu_{Y_2}(z_2)}{\mu_{X_1}(z_1)\mu_{X_2}(z_2)} \\ &= \frac{\exp\left(-\frac{1}{\lambda_1}\|z_1 - \delta_1\|_1\right) \exp\left(-\frac{1}{\lambda_2}\|z_2 - \delta_2\|_1\right)}{\exp\left(-\frac{1}{\lambda_1}\|z_1\|_1\right) \exp\left(-\frac{1}{\lambda_2}\|z_2\|_1\right)} \\ &= \exp\left(-\frac{1}{\lambda_1}(\|z_1 - \delta_1\|_1 - \|z_1\|_1) - \frac{1}{\lambda_2}(\|z_2 - \delta_2\|_1 - \|z_2\|_1)\right) \end{aligned}$$

By choosing  $\beta = -\log(t)$ , we can derive that

$$\begin{aligned} \frac{1}{\lambda_1} (\|z_1 - \delta_1\|_1 - \|z_1\|_1) + \frac{1}{\lambda_2} (\|z_2 - \delta_2\|_1 - \|z_2\|_1) \geq \beta &\iff \frac{\mu_{Y_1}(z_1)\mu_{Y_2}(z_2)}{\mu_{X_1}(z_1)\mu_{X_2}(z_2)} \leq t, \\ \frac{1}{\lambda_1} (\|z_1 - \delta_1\|_1 - \|z_1\|_1) + \frac{1}{\lambda_2} (\|z_2 - \delta_2\|_1 - \|z_2\|_1) \leq \beta &\iff \frac{\mu_{Y_1}(z_1)\mu_{Y_2}(z_2)}{\mu_{X_1}(z_1)\mu_{X_2}(z_2)} \geq t. \end{aligned}$$

□

## B.1 Proof of the Certified Robustness for $l_1$ norm

**THEOREM B.2 ( $\ell_1$  NORM CERTIFIED SPACE FOR VISUAL GM).** Let  $f$  be a matching function,  $f_0$  and  $g_0$  be defined as in Eq. 6 and Eq. 7,  $\varepsilon_1 \sim \mathcal{L}(\lambda_1)$ ,  $\varepsilon_2 \sim \mathcal{L}(\lambda_2)$ . Suppose  $\underline{p} \in (\frac{1}{2}, 1]$  satisfy:

$$\begin{aligned} P(f_0(\mathbf{c}^1, \mathbf{c}^2, \mathbf{z}^1 + \varepsilon_1, \mathbf{z}^2 + \varepsilon_2) = 1) &= \\ P(f(\mathbf{c}^1, \mathbf{c}^2, \mathbf{z}^1 + \varepsilon_1, \mathbf{z}^2 + \varepsilon_2) \in \mathcal{X}') &= p \geq \underline{p}. \end{aligned} \quad (27)$$

Then we obtain the  $\ell_1$  norm certified space for the perturbation pair  $(\delta_1, \delta_2)$ :

$$\frac{\|\delta_1\|_1}{\lambda_1} + \frac{\|\delta_2\|_1}{\lambda_2} \leq -\log \left[ 2 \left( 1 - \underline{p} \right) \right], \quad (28)$$

which guarantees  $g_0(\mathbf{c}^1, \mathbf{c}^2, \mathbf{z}^1 + \delta_1, \mathbf{z}^2 + \delta_2) = 1$ .

To show that  $g_0(\mathbf{c}^1, \mathbf{c}^2, \mathbf{z}^1 + \delta_1, \mathbf{z}^2 + \delta_2) = 1$ , it follows from the definition of  $g_0$  that we need to show that:

$$P(f(\mathbf{c}^1, \mathbf{c}^2, \mathbf{z}^1 + \varepsilon_1 + \delta_1, \mathbf{z}^2 + \varepsilon_2 + \delta_2) \in \mathcal{X}') \geq P(f(\mathbf{c}^1, \mathbf{c}^2, \mathbf{z}^1 + \varepsilon_1 + \delta_1, \mathbf{z}^2 + \varepsilon_2 + \delta_2) \notin \mathcal{X}').$$

We define two random variables:

$$\begin{aligned} I &:= (\mathbf{c}^1, \mathbf{c}^2, \mathbf{z}^1 + \varepsilon_1, \mathbf{z}^2 + \varepsilon_2) \\ O &:= (\mathbf{c}^1, \mathbf{c}^2, \mathbf{z}^1 + \varepsilon_1 + \delta_1, \mathbf{z}^2 + \varepsilon_2 + \delta_2). \end{aligned}$$

We know that:

$$P(f(I) \in \mathcal{X}') \geq \underline{p}. \quad (29)$$

Our goal is to show that

$$P(f(O) \in \mathcal{X}') > P(f(O) \notin \mathcal{X}'). \quad (30)$$

Denote  $T(\mathbf{z}^1, \mathbf{z}^2) = \frac{1}{\lambda_1} (\|\mathbf{z}^1 - \delta_1\|_1 - \|\mathbf{z}^1\|_1) + \frac{1}{\lambda_2} (\|\mathbf{z}^2 - \delta_2\|_1 - \|\mathbf{z}^2\|_1)$ . Use Triangle Inequality we can derive a bound for  $T(\mathbf{z}^1, \mathbf{z}^2)$  :

$$-\frac{\|\delta_1\|_1}{\lambda_1} - \frac{\|\delta_2\|_1}{\lambda_2} \leq T(\mathbf{z}^1, \mathbf{z}^2) \leq \frac{\|\delta_1\|_1}{\lambda_1} + \frac{\|\delta_2\|_1}{\lambda_2}. \quad (31)$$

Pick  $\beta'$  such that there exists  $B' \subseteq \{z_1, z_2 : T(z_1, z_2) = \beta'\}$ , and

$$P(I \in \{z_1, z_2 : T(z_1, z_2) < \beta'\} \cup B') = 1 - \underline{p} = P(f(I) \notin \mathcal{X}'). \quad (32)$$

Define

$$S := \{z_1, z_2 : T(z_1, z_2) < \beta'\} \cup B', \quad (33)$$

so we also have  $P(X \notin S) = p = P(f(I) \notin \mathcal{X}')$ . Plug into Lemma B.1, we can get

$$\begin{aligned} P(Y \notin S) &\leq P(f(O) \in \mathcal{X}'), \\ P(Y \in S) &\geq P(f(O) \notin \mathcal{X}'). \end{aligned} \tag{34}$$

Then we can obtain

$$\begin{aligned} \mathbb{P}(Y \in S) &= \int \int_S [2\lambda_1]^{-d} [2\lambda_2]^{-d} \exp\left(-\frac{\|z^1 - \delta_1\|_1}{\lambda_1}\right) \exp\left(-\frac{\|z^2 - \delta_2\|_1}{\lambda_2}\right) dz^1 dz^2 \\ &= \int \int_S [2\lambda_1]^{-d} [2\lambda_2]^{-d} \exp\left(-\frac{\|z^1\|_1}{\lambda_1}\right) \exp\left(-\frac{\|z^2\|_1}{\lambda_2}\right) \exp\left(-T(z^1, z^2)\right) dz^1 dz^2 \\ &\leq \exp\left(\frac{\|\delta_1\|_1}{\lambda_1} + \frac{\|\delta_2\|_1}{\lambda_2}\right) \int \int_S [2\lambda_1]^{-d} [2\lambda_2]^{-d} \exp\left(-\frac{\|z^1\|_1}{\lambda_1}\right) \exp\left(-\frac{\|z^2\|_1}{\lambda_2}\right) dz^1 dz^2 \\ &= \exp\left(\frac{\|\delta_1\|_1}{\lambda_1} + \frac{\|\delta_2\|_1}{\lambda_2}\right) (1 - \underline{p}). \end{aligned} \tag{35}$$

Thus, if  $\frac{\|\delta_1\|_1}{\lambda_1} + \frac{\|\delta_2\|_1}{\lambda_2} \leq -\log\left[2\left(1 - \underline{p}\right)\right]$ , it holds that

$$\begin{aligned} P(Y \in S) &\leq \exp\left(\frac{\|\delta_1\|_1}{\lambda_1} + \frac{\|\delta_2\|_1}{\lambda_2}\right) (1 - \underline{p}) \\ &\leq \exp\left(-\log\left[2\left(1 - \underline{p}\right)\right]\right) (1 - \underline{p}) \\ &= \frac{1}{2}. \end{aligned} \tag{36}$$