

## A Proofs

**Lemma 1.** *For any normalized probability distributions  $\mathbf{x}, \mathbf{x}' \in [0, 1]^{n \times m}$ , there exists at least one  $\delta$  such that  $\mathbf{x}' = \Delta(\mathbf{x}, \delta)$ . Furthermore:*

$$\min_{\delta: \mathbf{x}' = \Delta(\mathbf{x}, \delta)} \|\delta\|_1 = W_1(\mathbf{x}, \mathbf{x}') \quad (22)$$

Where  $W_1$  denotes the 1-Wasserstein metric, using  $L_1$  distance as the underlying distance metric.

*Proof.* We first show the equivalence of the above minimization problem with the linear program proposed by Ling and Okada (2007), restated here:

$$W_1(\mathbf{x}, \mathbf{x}') = \min_{\mathbf{g}} \sum_{(i,j)} \sum_{(i',j') \in \mathcal{N}((i,j))} \mathbf{g}_{(i,j),(i',j')} \quad (23)$$

where  $\mathbf{g} \geq 0$  and  $\forall (i, j)$ ,

$$\sum_{(i',j') \in \mathcal{N}((i,j))} \mathbf{g}_{(i,j),(i',j')} - \mathbf{g}_{(i',j'),(i,j)} = \mathbf{x}'_{i,j} - \mathbf{x}_{i,j}$$

It suffices to show that (1) there is a transformation from the variables  $\mathbf{g}$  in Equation 23 to the variables  $\delta$  in Equation 22, such that all points which are feasible in Equation 23 are feasible in 22 and the minimization objective in Equation 22 is less than or equal to the minimization objective in Equation 23, and (2) there is a transformation from the variables  $\delta$  in Equation 22 to the variables  $\mathbf{g}$  in Equation 23, such that all points which are feasible in Equation 22 are feasible in Equation 23 and the minimization objective in Equation 23 is less than or equal to the minimization objective in Equation 22.

We start with (1). We give the transformation as:

$$\begin{aligned} \delta_{i,j}^{\text{vert.}} &:= \mathbf{g}_{(i,j),(i+1,j)} - \mathbf{g}_{(i+1,j),(i,j)} \\ \delta_{i,j}^{\text{horiz.}} &:= \mathbf{g}_{(i,j),(i,j+1)} - \mathbf{g}_{(i,j+1),(i,j)} \end{aligned} \quad (24)$$

Where we let  $\mathbf{g}_{(n,j),(n+1,j)} = \mathbf{g}_{(n+1,j),(n,j)} = \mathbf{g}_{(i,m+1),(i,m)} = \mathbf{g}_{(i,m),(i,m+1)} = 0$ . To show feasibility, we write out fully the flow constraint of Equation 23:

$$\begin{aligned} &\mathbf{g}_{(i,j),(i+1,j)} - \mathbf{g}_{(i+1,j),(i,j)} + \\ &\mathbf{g}_{(i,j),(i-1,j)} - \mathbf{g}_{(i-1,j),(i,j)} + \\ &\mathbf{g}_{(i,j),(i,j+1)} - \mathbf{g}_{(i,j+1),(i,j)} + \\ &\mathbf{g}_{(i,j),(i,j-1)} - \mathbf{g}_{(i,j-1),(i,j)} = \mathbf{x}'_{i,j} - \mathbf{x}_{i,j} \end{aligned} \quad (25)$$

Substituting in Equation 24:

$$\delta_{i,j}^{\text{vert.}} - \delta_{i-1,j}^{\text{vert.}} + \delta_{i,j}^{\text{horiz.}} - \delta_{i,j-1}^{\text{horiz.}} = \mathbf{x}'_{i,j} - \mathbf{x}_{i,j} \quad (26)$$

But by Definition 3.1, this is exactly:

$$\Delta(\mathbf{x}, \delta)_{i,j} = \mathbf{x}'_{i,j} \quad (27)$$

Which is the sole constraint in Equation 22: then any solution which is feasible in Equation 23 is feasible in Equation 22. Also note that:

$$\begin{aligned} \|\delta\|_1 &= \sum_{i,j} |\delta_{i,j}^{\text{vert.}}| + |\delta_{i,j}^{\text{horiz.}}| \\ &\leq \sum_{i,j} |\mathbf{g}_{(i,j),(i+1,j)}| + |\mathbf{g}_{(i+1,j),(i,j)}| \\ &\quad + |\mathbf{g}_{(i,j),(i,j+1)}| + |\mathbf{g}_{(i,j+1),(i,j)}| \\ &= \sum_{i,j} \mathbf{g}_{(i,j),(i+1,j)} + \mathbf{g}_{(i+1,j),(i,j)} \\ &\quad + \mathbf{g}_{(i,j),(i,j+1)} + \mathbf{g}_{(i,j+1),(i,j)} \\ &= \sum_{i,j} \mathbf{g}_{(i,j),(i+1,j)} + \mathbf{g}_{(i,j),(i,j+1)} \\ &\quad + \sum_{i,j} \mathbf{g}_{(i+1,j),(i,j)} + \mathbf{g}_{(i,j+1),(i,j)} \\ &= \sum_{i,j} \mathbf{g}_{(i,j),(i+1,j)} + \mathbf{g}_{(i,j),(i,j+1)} \\ &\quad + \sum_{i,j} \mathbf{g}_{(i,j),(i-1,j)} + \mathbf{g}_{(i,j),(i,j-1)} \\ &= \sum_{(i,j)} \sum_{(i',j') \in \mathcal{N}((i,j))} \mathbf{g}_{(i,j),(i',j')} \end{aligned} \quad (28)$$

Where the inequality follows from triangle inequality applied to Equation 24, and in the second sum in the fourth line, we exploit the fact that  $\mathbf{g}_{(n,j),(n+1,j)} = \mathbf{g}_{(n+1,j),(n,j)} = \mathbf{g}_{(i,m+1),(i,m)} = \mathbf{g}_{(i,m),(i,m+1)} = 0$  to shift indices. This shows that the minimization objective in Equation 22 is less than or equal to the minimization objective in Equation 23.

Moving on to (2), we give the transformation as:

$$\begin{aligned} \mathbf{g}_{(i,j),(i+1,j)} &:= \max(\delta_{i,j}^{\text{vert.}}, 0) \\ \mathbf{g}_{(i,j),(i-1,j)} &:= \max(-\delta_{i-1,j}^{\text{vert.}}, 0) \\ \mathbf{g}_{(i,j),(i,j+1)} &:= \max(\delta_{i,j}^{\text{horiz.}}, 0) \\ \mathbf{g}_{(i,j),(i,j-1)} &:= \max(-\delta_{i,j-1}^{\text{horiz.}}, 0) \end{aligned} \quad (29)$$

Note that the non-negativity constraint of Equation 23 is automatically satisfied by the form of these definitions. Shifting indices, we also have:

$$\begin{aligned} \mathbf{g}_{(i-1,j),(i,j)} &= \max(\delta_{i-1,j}^{\text{vert.}}, 0) \\ \mathbf{g}_{(i+1,j),(i,j)} &= \max(-\delta_{i,j}^{\text{vert.}}, 0) \\ \mathbf{g}_{(i,j-1),(i,j)} &= \max(\delta_{i,j-1}^{\text{horiz.}}, 0) \\ \mathbf{g}_{(i,j+1),(i,j)} &= \max(-\delta_{i,j}^{\text{horiz.}}, 0) \end{aligned} \quad (30)$$

From the constraint on Equation 22, we have:

$$\begin{aligned}
 \mathbf{x}'_{i,j} - \mathbf{x}_{i,j} &= \delta_{i,j}^{\text{vert.}} + \\
 &\quad - \delta_{i-1,j}^{\text{vert.}} + \\
 &\quad \delta_{i,j}^{\text{horiz.}} + \\
 &\quad - \delta_{i,j-1}^{\text{horiz.}} \\
 &= \max(\delta_{i,j}^{\text{vert.}}, 0) - \max(-\delta_{i,j}^{\text{vert.}}, 0) + \\
 &\quad \max(-\delta_{i-1,j}^{\text{vert.}} - \max(\delta_{i-1,j}^{\text{vert.}}, 0), 0) + \\
 &\quad \max(\delta_{i,j}^{\text{horiz.}}, 0) - \max(-\delta_{i,j}^{\text{horiz.}}, 0) + \\
 &\quad \max(-\delta_{i,j-1}^{\text{horiz.}}, 0) - \max(\delta_{i,j-1}^{\text{horiz.}}, 0) \\
 &= \mathbf{g}(i,j),(i+1,j) - \mathbf{g}(i+1,j),(i,j) + \\
 &\quad \mathbf{g}(i,j),(i-1,j) - \mathbf{g}(i-1,j),(i,j) + \\
 &\quad \mathbf{g}(i,j),(i,j+1) - \mathbf{g}(i,j+1),(i,j) + \\
 &\quad \mathbf{g}(i,j),(i,j-1) - \mathbf{g}(i,j-1),(i,j)
 \end{aligned} \tag{31}$$

Which is exactly the second constraint of Equation 23: then any solution which is feasible in Equation 23 is feasible in Equation 22. Also note that:

$$\begin{aligned}
 &\sum_{(i,j)} \sum_{(i',j') \in \mathcal{N}((i,j))} \mathbf{g}(i,j),(i',j') \\
 &= \sum_{i,j} \mathbf{g}(i,j),(i+1,j) + \mathbf{g}(i,j),(i,j+1) \\
 &\quad + \sum_{i,j} \mathbf{g}(i,j),(i-1,j) + \mathbf{g}(i,j),(i,j-1) \\
 &= \sum_{i,j} \max(\delta_{i,j}^{\text{vert.}}, 0) + \max(\delta_{i,j}^{\text{horiz.}}, 0) \\
 &\quad + \sum_{i,j} \max(-\delta_{i-1,j}^{\text{vert.}}, 0) + \max(-\delta_{i,j-1}^{\text{horiz.}}, 0) \\
 &= \sum_{i,j} \max(\delta_{i,j}^{\text{vert.}}, 0) + \max(-\delta_{i,j}^{\text{vert.}}, 0) \\
 &\quad + \sum_{i,j} \max(\delta_{i,j}^{\text{horiz.}}, 0) + \max(-\delta_{i,j}^{\text{horiz.}}, 0) \\
 &= \sum_{i,j} |\delta_{i,j}^{\text{vert.}}| + |\delta_{i,j}^{\text{horiz.}}| \\
 &= \|\delta\|_1
 \end{aligned} \tag{32}$$

Where we again exploit the fact that  $\mathbf{g}(n,j),(n+1,j) = \mathbf{g}(n+1,j),(n,j) = \mathbf{g}(i,m+1),(i,m) = \mathbf{g}(i,m),(i,m+1) = 0$  to shift indices, in the fourth line. This shows that the minimization objective in Equation 23 is less than or equal to the minimization objective in Equation 22, completing (2).

Finally, now that we have shown that Equations 22 and 23 are in fact equivalent minimizations (i.e., we have proven Equation 22 correct), we would like to show that there is always a feasible solution to 22, as claimed. By the above transformations, it suffices to show that there is always a feasible solution to Equation 23. Ling and Okada (2007) show that any feasible solution the the general Wasserstein minimization LP (Definition 1) can be transformed into a solution to

Equation 23, so it suffices to show that the LP in Definition 1 always has a feasible solution. This is trivially satisfied by taking  $\Pi = \mathbf{x}(\mathbf{x}')^T$ , where we note that  $\mathbf{x}$ , a probability distribution, is non-negative.  $\square$

**Theorem 1.** Consider a normalized probability distribution  $\mathbf{x} \in [0, 1]^{n \times m}$ , and a classification score function  $\mathbf{f} : \mathbb{R}^{n \times m} \rightarrow [0, 1]^k$ . Let  $\bar{\mathbf{f}}$  refer to the Wasserstein-smoothed classification function:

$$\bar{\mathbf{f}}(\mathbf{x}) = \mathbb{E}_{\delta \sim \mathcal{L}(\sigma)} [\mathbf{f}(\Delta(\mathbf{x}, \delta))] \tag{33}$$

Let  $i$  be the class assignment of  $\mathbf{x}$  using the smoothed classifier  $\bar{\mathbf{f}}$  (i.e.  $i = \arg \max_{i'} \bar{\mathbf{f}}_{i'}(\mathbf{x})$ ). If

$$\bar{\mathbf{f}}_i(\mathbf{x}) \geq e^{2\sqrt{2}\rho/\sigma} \max_{i' \neq i} \bar{\mathbf{f}}_{i'}(\mathbf{x}) \tag{34}$$

Then for any perturbed probability distribution  $\tilde{\mathbf{x}}$  such that  $W_1(\mathbf{x}, \tilde{\mathbf{x}}) \leq \rho$ :

$$\bar{\mathbf{f}}_i(\tilde{\mathbf{x}}) \geq \max_{i' \neq i} \bar{\mathbf{f}}_{i'}(\tilde{\mathbf{x}}) \tag{35}$$

*Proof.* Let  $\mathbf{u}$  be the uniform probability vector. As a consequence of Lemma 1, for any distribution  $\mathbf{x}$ , there exists a nonempty set of local flow plans  $S_{\mathbf{x}}$ :

$$S_{\mathbf{x}} = \{\delta | \mathbf{x} = \Delta(\mathbf{u}, \delta)\} \tag{36}$$

Also, we may define a version of the classifier  $\mathbf{f}$  on the local flow plan domain:

$$\mathbf{f}^{\text{flow}}(\delta) = \mathbf{f}(\Delta(\mathbf{u}, \delta)) \tag{37}$$

Let  $\delta_{\mathbf{x}}$  be an arbitrary element in  $S_{\mathbf{x}}$ , and consider any perturbed  $\tilde{\mathbf{x}}$  such that  $W_1(\mathbf{x}, \tilde{\mathbf{x}}) \leq \rho$ . By Theorem 1:

$$\min_{\delta: \tilde{\mathbf{x}} = \Delta(\mathbf{x}, \delta)} \|\delta\|_1 = W_1(\mathbf{x}, \tilde{\mathbf{x}}) \tag{38}$$

Then, using Equation 6:

$$\min_{\delta: \tilde{\mathbf{x}} = \Delta(\mathbf{u}, \delta_{\mathbf{x}} + \delta)} \|\delta\|_1 = W_1(\mathbf{x}, \tilde{\mathbf{x}}) \tag{39}$$

Let the minimum be achieved at  $\delta^*$ . Making a change of variables ( $\delta_{\tilde{\mathbf{x}}} = \delta^* + \delta_{\mathbf{x}}$ ), we have:

$$\|\delta_{\tilde{\mathbf{x}}} - \delta_{\mathbf{x}}\|_1 = W_1(\mathbf{x}, \tilde{\mathbf{x}}) \quad \text{where } \tilde{\mathbf{x}} = \Delta(\mathbf{u}, \delta_{\tilde{\mathbf{x}}}) \tag{40}$$

Note that for any  $\mathbf{x}'$  (for  $\delta' \sim \mathcal{L}(\sigma)$ ):

$$\begin{aligned}
 \bar{\mathbf{f}}(\mathbf{x}') &= \mathbb{E} [\mathbf{f}(\Delta(\mathbf{x}', \delta'))] \\
 &= \mathbb{E} [\mathbf{f}(\Delta(\mathbf{u}, \delta_{\mathbf{x}'} + \delta'))] \\
 &= \mathbb{E} [\mathbf{f}^{\text{flow}}(\delta_{\mathbf{x}'} + \delta')]
 \end{aligned} \tag{41}$$

We can now apply Proposition 1 from Lecuyer et al. (2019), restated here:

**Proposition.** Consider a vector  $\mathbf{v} \in \mathbb{R}^d$ , and a classification score function  $\mathbf{h} : \mathbb{R}^d \rightarrow [0, 1]^k$ . Let  $\epsilon \sim \text{Laplace}(0, \sigma)^d$ , and let  $i$  be the class assignment of  $\mathbf{v}$  using a Laplace-smoothed version of the classifier  $\mathbf{h}$ :

$$i = \arg \max_{i'} \mathbb{E}_{\epsilon} [\mathbf{h}_{i'}(\mathbf{v} + \epsilon)] \quad (42)$$

If:

$$\mathbb{E}_{\epsilon} [\mathbf{h}_i(\mathbf{v} + \epsilon)] \geq e^{2\sqrt{2}\rho/\sigma} \max_{i' \neq i} \mathbb{E}_{\epsilon} [\mathbf{h}_{i'}(\mathbf{v} + \epsilon)] \quad (43)$$

Then for any perturbed probability distribution  $\tilde{\mathbf{v}}$  such that  $\|\mathbf{v} - \tilde{\mathbf{v}}\|_1 \leq \rho$ :

$$\mathbb{E}_{\epsilon} [\mathbf{h}_i(\tilde{\mathbf{v}} + \epsilon)] \geq \max_{i' \neq i} \mathbb{E}_{\epsilon} [\mathbf{h}_{i'}(\tilde{\mathbf{v}} + \epsilon)] \quad (44)$$

We apply this proposition to  $\mathbf{f}^{\text{flow}}$ , noting that  $\|\delta_{\tilde{\mathbf{x}}}\|_1 = W_1(\mathbf{x}, \tilde{\mathbf{x}}) \leq \rho$ :

$$\begin{aligned} \mathbb{E}_{\delta'} [\mathbf{f}_i^{\text{flow}}(\delta_{\mathbf{x}} + \delta')] &\geq e^{2\sqrt{2}\rho/\sigma} \max_{i' \neq i} \mathbb{E}_{\delta'} [\mathbf{f}_{i'}^{\text{flow}}(\delta_{\mathbf{x}} + \delta')] \\ \implies \mathbb{E}_{\delta'} [\mathbf{f}_i^{\text{flow}}(\delta_{\tilde{\mathbf{x}}} + \delta')] &\geq \max_{i' \neq i} \mathbb{E}_{\delta'} [\mathbf{f}_{i'}^{\text{flow}}(\delta_{\tilde{\mathbf{x}}} + \delta')] \end{aligned} \quad (45)$$

Then, using Equation 41:

$$\begin{aligned} \bar{\mathbf{f}}_i(\mathbf{x}) &\geq e^{2\sqrt{2}\rho/\sigma} \max_{i' \neq i} \bar{\mathbf{f}}_{i'}(\mathbf{x}) \implies \\ \bar{\mathbf{f}}_i(\tilde{\mathbf{x}}) &\geq \max_{i' \neq i} \bar{\mathbf{f}}_{i'}(\tilde{\mathbf{x}}) \end{aligned} \quad (46)$$

Which was to be proven.  $\square$

**Corollary 1.** For any normalized probability distributions  $\mathbf{x}, \mathbf{x}' \in [0, 1]^{n \times m}$ , if  $W_1(\mathbf{x}, \mathbf{x}') \leq \rho/2$ , then  $\|\mathbf{x} - \mathbf{x}'\|_1 \leq \rho$ , where  $W_1$  is the 1-Wasserstein metric using any  $L_p$  norm as the underlying distance metric. Furthermore, there exist distributions where these inequalities are tight.

*Proof.* Let  $\Pi$  indicate the optimal transport plan between  $\mathbf{x}$  and  $\mathbf{x}'$ . From Definition 1, we have  $\Pi \mathbf{1} = \mathbf{x}$  and  $\Pi^T \mathbf{1} = \mathbf{x}'$ . Then:

$$(\Pi^T - \Pi) \mathbf{1} = \mathbf{x}' - \mathbf{x} \quad (47)$$

Let  $\Pi'$  represent a modified version of  $\Pi$ , with the diagonal elements set to zero. Note that  $\langle \Pi', C \rangle = \langle \Pi, C \rangle$  and  $\Pi^T - \Pi = (\Pi')^T - \Pi'$ . Then, using triangle inequality:

$$\begin{aligned} &\|(\Pi')^T \mathbf{1}\|_1 + \|(\Pi') \mathbf{1}\|_1 \\ &\geq \|((\Pi')^T - \Pi') \mathbf{1}\|_1 \\ &= \|\mathbf{x}' - \mathbf{x}\|_1 \end{aligned} \quad (48)$$

Because the elements of  $\Pi'$  are non-negative, this is simply:

$$2 \sum_{i,j} \Pi'_{i,j} \geq \|((\Pi')^T - \Pi') \mathbf{1}\|_1 = \|\mathbf{x}' - \mathbf{x}\|_1 \quad (49)$$

Then, because the (non-diagonal) elements of  $C$  are at least 1 for any  $L_p$  norm, we have,

$$2 \langle \Pi', C \rangle \geq 2 \sum_{i,j} \Pi'_{i,j} \geq \|\mathbf{x}' - \mathbf{x}\|_1 \quad (50)$$

Because  $\langle \Pi', C \rangle = \langle \Pi, C \rangle = W_1(\mathbf{x}, \mathbf{x}')$ , this means that  $\|\mathbf{x}' - \mathbf{x}\|_1 \leq 2W_1(\mathbf{x}, \mathbf{x}') \leq \rho$ , which was to be proven. Note that this inequality can be tight. For example, let  $\mathbf{x}$  be the distribution where the entire probability mass is at position  $(i, j)$ , and  $\mathbf{x}'$  be the distribution where the probability mass is equally split between at positions  $(i, j)$  and  $(i+1, j)$ . (In other words,  $\mathbf{x}_{(i,j)} = 1, \mathbf{x}'_{(i,j)} = .5, \mathbf{x}'_{(i+1,j)} = .5$ ). In this case,  $\|\mathbf{x}' - \mathbf{x}\|_1 = 1, W_1(\mathbf{x}, \mathbf{x}') = .5$ .  $\square$

**Corollary 2.** Consider a color image with three channels, denoted  $\mathbf{x} = [\mathbf{x}^R, \mathbf{x}^G, \mathbf{x}^B]$ , normalized such that  $\sum_{(i,j)} \mathbf{x}_{(i,j)}^R + \mathbf{x}_{(i,j)}^G + \mathbf{x}_{(i,j)}^B = 1$ . Consider a perturbed image  $\tilde{\mathbf{x}}$  such that  $\forall K \in \{R, G, B\}, \sum_{(i,j)} \mathbf{x}_{(i,j)}^K = \sum_{(i,j)} \tilde{\mathbf{x}}_{(i,j)}^K$ . Let  $W_1(\mathbf{x}, \tilde{\mathbf{x}})$  denote the 1-Wasserstein distance (with  $L_1$  distance metric) between  $\mathbf{x}$  and  $\tilde{\mathbf{x}}$ , where, when determining the minimum transport plan, transport between channels is not permitted. Using this definition, let  $W_1(\mathbf{x}, \tilde{\mathbf{x}}) \leq \rho$ . Define:

$$\begin{aligned} \delta &= \{\delta^R, \delta^G, \delta^B\} \\ \Delta(\mathbf{x}, \delta) &= \{\Delta(\mathbf{x}^R, \delta^R), \Delta(\mathbf{x}^G, \delta^G), \Delta(\mathbf{x}^B, \delta^B)\} \end{aligned} \quad (51)$$

and let  $\mathcal{L}^{\text{color}}(\sigma)$  represent independent draws of Laplace noise each with standard deviation  $\sigma$  in the shape of  $\delta$ . Then if

$$\bar{\mathbf{f}}_i(\mathbf{x}) \geq e^{2\sqrt{2}\rho/\sigma} \max_{i' \neq i} \bar{\mathbf{f}}_{i'}(\mathbf{x}) \quad (52)$$

then

$$\bar{\mathbf{f}}_i(\tilde{\mathbf{x}}) \geq \max_{i' \neq i} \bar{\mathbf{f}}_{i'}(\tilde{\mathbf{x}}). \quad (53)$$

*Proof.* Let the mass in each channel be denoted  $s_K$ :

$$s_K := \sum_{(i,j)} \mathbf{x}_{(i,j)}^K = \sum_{(i,j)} \tilde{\mathbf{x}}_{(i,j)}^K \quad (54)$$

Consider the formulation of Wasserstein distance given in Definition 1. If we represent the elements of  $\mathbf{x}$  as a vector by concatenating the elements of  $\delta^R, \delta^G$ , and  $\delta^B$ , then the restriction that there is no flow between channels amounts to the requirement that  $\Pi$  is block-diagonal:

$$\Pi = \begin{bmatrix} \Pi^R & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \Pi^G & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \Pi^B \end{bmatrix} \quad (55)$$

Let  $C^{1,1}$  represent the standard cost matrix for 1-Wasserstein transport (with  $L_1$  distance metric). Because the cost of transport within each channel is the

same for standard 1-Wasserstein transport (with  $L_1$  distance metric), we have:

$$C = \begin{bmatrix} C^{1,1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & C^{1,1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & C^{1,1} \end{bmatrix} \quad (56)$$

Then we have:

$$\langle \Pi, C \rangle = \langle \Pi^R, C^{1,1} \rangle + \langle \Pi^G, C^{1,1} \rangle + \langle \Pi^B, C^{1,1} \rangle \quad (57)$$

And by Equation 55, the constraints also factorize out:

$$\begin{aligned} \Pi^R \mathbf{1} &= \mathbf{x}^R, & (\Pi^R)^T \mathbf{1} &= \tilde{\mathbf{x}}^R, \\ \Pi^B \mathbf{1} &= \mathbf{x}^B, & (\Pi^B)^T \mathbf{1} &= \tilde{\mathbf{x}}^B, \\ \Pi^G \mathbf{1} &= \mathbf{x}^G, & (\Pi^G)^T \mathbf{1} &= \tilde{\mathbf{x}}^G \end{aligned} \quad (58)$$

Then the variables of each  $\Pi^K$  are separable (in that they appear together in the objective only in the sum and share no constraints). We can then factorize the minimization:

$$W_1(\mathbf{x}, \tilde{\mathbf{x}}) = \sum_K \min_{\Pi^K \in \mathbb{R}_+^{(n-m) \times (n-m)}} \langle \Pi^K, C^{(1,1)} \rangle, \quad (59)$$

$$\forall K, \Pi^K \mathbf{1} = \mathbf{x}^K, (\Pi^K)^T \mathbf{1} = \tilde{\mathbf{x}}^K \quad (60)$$

$$(61)$$

We can transform each  $x^K$  into a normalized probability distribution by scaling it by a factor of  $1/s_K$ . We similarly scale each  $\Pi^K$ :

$$\mathbf{x}_{sc.}^K := \frac{\mathbf{x}^K}{s_K}, \quad \Pi_{sc.}^K := \frac{\Pi^K}{s_K} \quad (62)$$

Then we have:

$$W_1(\mathbf{x}, \tilde{\mathbf{x}}) = \sum_K s_K \cdot \min_{\Pi_{sc.}^K \in \mathbb{R}_+^{(n-m) \times (n-m)}} \langle \Pi_{sc.}^K, C^{(1,1)} \rangle, \quad (63)$$

$$\forall K, \Pi_{sc.}^K \mathbf{1} = \mathbf{x}_{sc.}^K, (\Pi_{sc.}^K)^T \mathbf{1} = \tilde{\mathbf{x}}_{sc.}^K \quad (64)$$

$$(65)$$

But note that this is simply:

$$W_1(\mathbf{x}, \tilde{\mathbf{x}}) = \sum_K s_K \cdot W_1(\mathbf{x}_{sc.}^K, \tilde{\mathbf{x}}_{sc.}^K) \quad (66)$$

By Lemma 1, this is:

$$W_1(\mathbf{x}, \tilde{\mathbf{x}}) = \sum_K s_K \cdot \min_{\delta_{sc.}^K: \tilde{\mathbf{x}}_{sc.}^K = \Delta(\mathbf{x}_{sc.}^K, \delta_{sc.}^K)} \|\delta_{sc.}^K\|_1 \quad (67)$$

By the linearity to scaling of  $\Delta$  and the  $L_1$  norm, this is simply:

$$W_1(\mathbf{x}, \tilde{\mathbf{x}}) = \sum_K \min_{\delta^K: \tilde{\mathbf{x}}^K = \Delta(\mathbf{x}^K, \delta^K)} \|\delta^K\|_1 \quad (68)$$

Which, by Equation 51, is simply,

$$W_1(\mathbf{x}, \tilde{\mathbf{x}}) = \min_{\delta: \tilde{\mathbf{x}} = \Delta(\mathbf{x}, \delta)} \|\delta\|_1 \quad (69)$$

Then all of the mechanics of the proof of Theorem 1 apply, and (avoiding unnecessary repetition), we conclude the result.  $\square$

**Corollary 3.** *Let  $W^1$  denote the  $L_1$  1-Wasserstein distance, and  $W^2$  denote the  $L_2$  1-Wasserstein distance. For a radius  $\rho_2$ , define  $\rho_1 := \sqrt{2}\rho_2$ . Then, for any classifier  $f$  and input  $\mathbf{x}$ , if there does not exist any adversarial example  $\tilde{\mathbf{x}}$  with  $W_1(\mathbf{x}, \tilde{\mathbf{x}}) \leq \rho_1$ , then there are also no adversarial examples  $\tilde{\mathbf{x}}'$  with  $W_2(\mathbf{x}, \tilde{\mathbf{x}}') \leq \rho_2$ .*

*Proof.* We show the contrapositive: If there is an adversarial example  $\tilde{\mathbf{x}}'$  with  $W^2(\mathbf{x}, \tilde{\mathbf{x}}') \leq \rho_2$ , then there is an adversarial example  $\tilde{\mathbf{x}}$  with  $W^1(\mathbf{x}, \tilde{\mathbf{x}}) \leq \rho_1$ . It is sufficient to show that for any arbitrary  $\mathbf{x}'$ , if  $W^2(\mathbf{x}, \mathbf{x}') \leq \rho_2$ , then  $W^1(\mathbf{x}, \mathbf{x}') \leq \rho_1$ . (The predicate is then satisfied with  $\tilde{\mathbf{x}}' = \tilde{\mathbf{x}} = \mathbf{x}'$ ). In other words, we need to show that

$$\sqrt{2}W^2(\mathbf{x}, \mathbf{x}') \geq W^1(\mathbf{x}, \mathbf{x}'), \quad \forall \mathbf{x}, \mathbf{x}'. \quad (70)$$

By the definition of 1-Wasserstein distance, can rewrite this goal as

$$\sqrt{2} \min_{\Pi} \langle \Pi, C^2 \rangle \geq \min_{\Pi} \langle \Pi, C^1 \rangle \quad (71)$$

where in both minimizations  $\Pi$  is non-negative and subject to  $\Pi \mathbf{1} = \mathbf{x}$ ,  $\Pi^T \mathbf{1} = \mathbf{x}'$ . Here,  $C^2$  and  $C^1$  are the weight matrices for  $L_2$  and  $L_1$  Wasserstein distances. Note that  $\Pi$  is subject to the same constraints in both minimizations: therefore any  $\Pi$  that is feasible in one is feasible in the other. Let  $\Pi_2^*$  be the minimum of the first ( $L_2$ ) minimization. Recall that

$$C_{(i,j),(i',j')}^2 = \sqrt{(i-i')^2 + (j-j')^2} \quad (72)$$

while

$$C_{(i,j),(i',j')}^1 = |i-i'| + |j-j'|. \quad (73)$$

By equivalence of norms, we have:

$$\sqrt{2}C_{(i,j),(i',j')}^2 \geq C_{(i,j),(i',j')}^1. \quad (74)$$

Then by linearity (and using  $\Pi$  non-negative),

$$\sqrt{2} \langle \Pi, C^2 \rangle \geq \langle \Pi, C^1 \rangle \quad \forall \Pi \geq 0 \quad (75)$$

So

$$\begin{aligned} \sqrt{2} \min_{\Pi} \langle \Pi, C^2 \rangle &= \sqrt{2} \langle \Pi_2^*, C^2 \rangle \\ &\geq \langle \Pi_2^*, C^1 \rangle \\ &\geq \min_{\Pi} \langle \Pi, C^1 \rangle, \end{aligned} \quad (76)$$

as desired.  $\square$

## B Training Parameters

In this paper, network architectures models used were identical to those used in Wong et al. (2019). Unless stated otherwise, all parameters of attacks are the same as used in that paper for each data set. For training smoothed models, we train the base classifier using standard cross-entropy loss on individual noised sample images, using the same noise distribution as used when performing smoothed classification. However, during training, rather than using the same image repeatedly while adding different noise (as at test time), we instead train with each image only once per epoch, with one noise draw. In fact, for computational efficiency and as suggested by Lecuyer et al. (2019), we re-use the same noise for each image in a batch. Training parameters are as follows (Tables 3, 4):

**Table 3** Training Parameters for MNIST Experiments

Training Epochs	200
Batch Size	128
Optimizer	Stochastic Gradient Descent with Momentum
Learning Rate	.001
Momentum	0.9
$L_2$ Weight Penalty	0.0005

**Table 4** Training Parameters for CIFAR-10 Experiments

Training Epochs	200
Batch Size	128
Training Set Preprocessing	Normalization, Random Cropping (Padding:4) and Random Horizontal Flip
Optimizer	Stochastic Gradient Descent with Momentum
Learning Rate	.01 (Epochs 1-200) .001 (Epochs 201-400)
Momentum	0.9
$L_2$ Weight Penalty	0.0005

## C Comparison to other Defenses in Wong et al. (2019)

In addition to proposing adversarial training as a defense against Wasserstein Adversarial attacks, Wong et al. (2019) also tests other defenses. On MNIST, binarization of the input and using a provably  $L_\infty$ -robust classifier were also tested as defenses: our randomized smoothing method is more effective than these methods at all attack magnitudes (see Figure 7). On CIFAR-10, Wong et al. (2019) only tested a provably

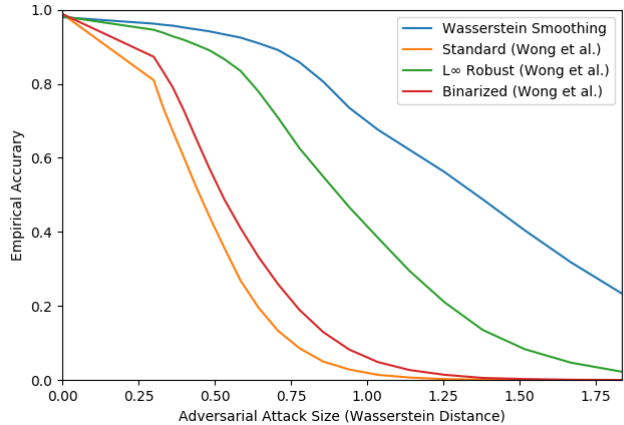


Figure 7: Comparison of empirical robustness on MNIST to additional defenses from (Wong et al., 2019), other than adversarial training. Randomized Smoothing shown here is Wasserstein smoothing with  $\sigma = 0.01$ . (This is the amount of noise which maximizes certified robustness, as seen in Table 1.)

$L_\infty$ -robust classifier as an additional defense: unfortunately, code was not provided for this model, so we did not attempt to replicate the results.