

---

# Pointwise Binary Classification with Pairwise Confidence Comparisons: Appendix

---

## A. Proof of Theorem 1

It is clear that each pair of examples  $(\mathbf{x}, \mathbf{x}')$  is independently drawn from the following data distribution:

$$\tilde{p}(\mathbf{x}, \mathbf{x}') = p((\mathbf{x}, \mathbf{x}') \mid (y, y') \in \tilde{\mathcal{Y}}) = \frac{p((\mathbf{x}, \mathbf{x}'), (y, y') \in \tilde{\mathcal{Y}})}{p((y, y') \in \tilde{\mathcal{Y}})},$$

where  $p((y, y') \in \tilde{\mathcal{Y}}) = \pi_+^2 + \pi_-^2 + \pi_+\pi_-$  and

$$\begin{aligned} p(\mathbf{x}, \mathbf{x}', (y, y') \in \tilde{\mathcal{Y}}) &= \sum_{(y, y') \in \tilde{\mathcal{Y}}} p(\mathbf{x}, \mathbf{x}' \mid (y, y')) \cdot p(y, y') \\ &= \pi_+^2 p_+(\mathbf{x}) p_+(\mathbf{x}') + \pi_-^2 p_-(\mathbf{x}) p_-(\mathbf{x}') + \pi_+\pi_- p_+(\mathbf{x}) p_-(\mathbf{x}'). \end{aligned}$$

Finally, let  $\tilde{p}(\mathbf{x}, \mathbf{x}') = p((\mathbf{x}, \mathbf{x}') \mid (y, y') \in \tilde{\mathcal{Y}})$ , the proof is completed.  $\square$

## B. Proof of Theorem 2

In order to decompose the pairwise comparison data distribution into pointwise distribution, we marginalize  $\tilde{p}(\mathbf{x}, \mathbf{x}')$  with respect to  $\mathbf{x}$  or  $\mathbf{x}'$ . Then we can obtain

$$\begin{aligned} \int \tilde{p}(\mathbf{x}, \mathbf{x}') d\mathbf{x}' &= \frac{1}{\tilde{\pi}} \left( \pi_+^2 p_+(\mathbf{x}) + \pi_-^2 p_-(\mathbf{x}) + \pi_+\pi_- p_+(\mathbf{x}) \right) \\ &= \frac{\pi_+}{\pi_-^2 + \pi_+} p_+(\mathbf{x}) + \frac{\pi_-^2}{\pi_-^2 + \pi_+} p_-(\mathbf{x}) \\ &= \tilde{p}_+(\mathbf{x}), \end{aligned}$$

and

$$\begin{aligned} \int \tilde{p}(\mathbf{x}, \mathbf{x}') d\mathbf{x} &= \frac{1}{\tilde{\pi}} \left( \pi_+^2 p_+(\mathbf{x}') + \pi_-^2 p_-(\mathbf{x}') + \pi_+\pi_- p_-(\mathbf{x}') \right) \\ &= \frac{\pi_+^2}{\pi_+^2 + \pi_-} p_+(\mathbf{x}') + \frac{\pi_-}{\pi_+^2 + \pi_-} p_-(\mathbf{x}') \\ &= \tilde{p}_-(\mathbf{x}'), \end{aligned}$$

which concludes the proof of Theorem 2.  $\square$

## C. Proof of Lemma 1

Based on Theorem 2, we can obtain the following linear equation:

$$\begin{bmatrix} \tilde{p}_+(\mathbf{x}) \\ \tilde{p}_-(\mathbf{x}) \end{bmatrix} = \frac{1}{\tilde{\pi}} \begin{bmatrix} \pi_+ & \pi_-^2 \\ \pi_+^2 & \pi_- \end{bmatrix} \begin{bmatrix} p_+(\mathbf{x}) \\ p_-(\mathbf{x}) \end{bmatrix}.$$

By solving the above equation, we obtain

$$\begin{aligned} p_+(\mathbf{x}) &= \frac{1}{\pi_+ - \pi_- \pi_+^2} (\tilde{\pi} \cdot \tilde{p}_+(\mathbf{x}) - \pi_- \tilde{\pi} \cdot \tilde{p}_-(\mathbf{x})) = \frac{1}{\pi_+} (\tilde{p}_+(\mathbf{x}) - \pi_- \tilde{p}_-(\mathbf{x})), \\ p_-(\mathbf{x}) &= \frac{1}{\pi_- - \pi_+ \pi_-^2} (\tilde{\pi} \cdot \tilde{p}_-(\mathbf{x}) - \pi_+ \tilde{\pi} \cdot \tilde{p}_+(\mathbf{x})) = \frac{1}{\pi_-} (\tilde{p}_-(\mathbf{x}) - \pi_+ \tilde{p}_+(\mathbf{x})), \end{aligned}$$

which concludes the proof of Lemma 1.  $\square$

## D. Proof of Theorem 3

It is quite intuitive to derive

$$\begin{aligned}
 R(f) &= \mathbb{E}_{p(\mathbf{x}, y)} [\ell(f(\mathbf{x}), y)] \\
 &= \pi_+ \mathbb{E}_{p_+(\mathbf{x})} [\ell(f(\mathbf{x}), +1)] + \pi_- \mathbb{E}_{p_-(\mathbf{x})} [\ell(f(\mathbf{x}), -1)] \\
 &= \frac{\pi_+ \tilde{\pi}}{\pi_+ - \pi_- \pi_+^2} \mathbb{E}_{\tilde{p}_+(\mathbf{x})} [\ell(f(\mathbf{x}), +1)] - \frac{\pi_+ \pi_- \tilde{\pi}}{\pi_+ - \pi_- \pi_+^2} \mathbb{E}_{\tilde{p}_-(\mathbf{x}') } [\ell(f(\mathbf{x}), +1)] \quad (\text{Lemma 1}) \\
 &\quad + \frac{\pi_- \tilde{\pi}}{\pi_- - \pi_+ \pi_-^2} \mathbb{E}_{\tilde{p}_-(\mathbf{x}') } [\ell(f(\mathbf{x}), -1)] - \frac{\pi_+ \pi_- \tilde{\pi}}{\pi_- - \pi_+ \pi_-^2} \mathbb{E}_{\tilde{p}_+(\mathbf{x})} [\ell(f(\mathbf{x}), -1)] \\
 &= \mathbb{E}_{\tilde{p}_+(\mathbf{x})} [\ell(f(\mathbf{x}), +1) - \pi_+ \ell(f(\mathbf{x}), -1)] + \mathbb{E}_{\tilde{p}_-(\mathbf{x}') } [\ell(f(\mathbf{x}), -1) - \pi_- \ell(f(\mathbf{x}), +1)] \\
 &= R_{\text{PC}}(f),
 \end{aligned}$$

which concludes the proof of Theorem 3.  $\square$

## E. Proof of Theorem 4

First of all, we introduce the following notations:

$$\begin{aligned}
 R_{\text{PC}}^+(f) &= \mathbb{E}_{\tilde{p}_+(\mathbf{x})} [\ell(f(\mathbf{x}), +1) - \pi_+ \ell(f(\mathbf{x}), -1)], \\
 \widehat{R}_{\text{PC}}^+(f) &= \frac{1}{n} \sum_{i=1}^n \left( \ell(f(\mathbf{x}_i), +1) - \pi_+ \ell(f(\mathbf{x}_i), -1) \right), \\
 R_{\text{PC}}^-(f) &= \mathbb{E}_{\tilde{p}_-(\mathbf{x}') } [\ell(f(\mathbf{x}'), -1) - \pi_- \ell(f(\mathbf{x}'), +1)], \\
 \widehat{R}_{\text{PC}}^-(f) &= \frac{1}{n} \sum_{i=1}^n \left( \ell(f(\mathbf{x}'_i), -1) - \pi_- \ell(f(\mathbf{x}'_i), +1) \right).
 \end{aligned}$$

In this way, we could simply represent  $R_{\text{PC}}(f)$  and  $\widehat{R}_{\text{PC}}(f)$  as

$$R_{\text{PC}}(f) = R_{\text{PC}}^+(f) + R_{\text{PC}}^-(f), \quad \widehat{R}_{\text{PC}}(f) = \widehat{R}_{\text{PC}}^+(f) + \widehat{R}_{\text{PC}}^-(f).$$

Then we have the following lemma.

**Lemma 2.** *The following inequality holds:*

$$R(\widehat{f}_{\text{PC}}) - R(f^*) \leq 2 \sup_{f \in \mathcal{F}} \left| R_{\text{PC}}^+(f) - \widehat{R}_{\text{PC}}^+(f) \right| + 2 \sup_{f \in \mathcal{F}} \left| R_{\text{PC}}^-(f) - \widehat{R}_{\text{PC}}^-(f) \right|. \quad (9)$$

*Proof.* We could intuitively express  $R(\widehat{f}_{\text{PC}}) - R(f^*)$  as

$$\begin{aligned}
 R(\widehat{f}_{\text{PC}}) - R(f^*) &= R(\widehat{f}_{\text{PC}}) - \widehat{R}_{\text{PC}}(\widehat{f}_{\text{PC}}) + \widehat{R}_{\text{PC}}(\widehat{f}_{\text{PC}}) - \widehat{R}_{\text{PC}}(f^*) + \widehat{R}_{\text{PC}}(f^*) - R(f^*) \\
 &= R_{\text{PC}}(\widehat{f}_{\text{PC}}) - \widehat{R}_{\text{PC}}(\widehat{f}_{\text{PC}}) + \widehat{R}_{\text{PC}}(\widehat{f}_{\text{PC}}) - \widehat{R}_{\text{PC}}(f^*) + \widehat{R}_{\text{PC}}(f^*) - R_{\text{PC}}(f^*) \\
 &\leq \sup_{f \in \mathcal{F}} \left| R_{\text{PC}}(f) - \widehat{R}_{\text{PC}}(f) \right| + 0 + \sup_{f \in \mathcal{F}} \left| R_{\text{PC}}(f) - \widehat{R}_{\text{PC}}(f) \right| \\
 &= 2 \sup_{f \in \mathcal{F}} \left| R_{\text{PC}}(f) - \widehat{R}_{\text{PC}}(f) \right| \\
 &\leq 2 \sup_{f \in \mathcal{F}} \left| R_{\text{PC}}^+(f) - \widehat{R}_{\text{PC}}^+(f) \right| + 2 \sup_{f \in \mathcal{F}} \left| R_{\text{PC}}^-(f) - \widehat{R}_{\text{PC}}^-(f) \right|,
 \end{aligned}$$

where the second equality holds due to Theorem 3.  $\square$

As suggested by Lemma 2, we need to further upper bound the right hand size of Eq. (9). Before doing that, we introduce the *uniform deviation bound*, which is useful to derive estimation error bounds. The proof can be found in some textbooks such as (Mohri et al., 2012) (Theorem 3.1).

**Lemma 3.** Let  $Z$  be a random variable drawn from a probability distribution with density  $\mu$ ,  $\mathcal{H} = \{h : \mathcal{Z} \mapsto [0, M]\}$  ( $M > 0$ ) be a class of measurable functions,  $\{z_i\}_{i=1}^n$  be i.i.d. examples drawn from the distribution with density  $\mu$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ ,

$$\sup_{h \in \mathcal{H}} \left| \mathbb{E}_{Z \sim \mu} [h(Z)] - \frac{1}{n} \sum_{i=1}^n h(z_i) \right| \leq 2\mathfrak{R}_n(\mathcal{H}) + M \sqrt{\frac{\log \frac{2}{\delta}}{2n}},$$

where  $\mathfrak{R}_n(\mathcal{H})$  denotes the (expected) Rademacher complexity (Bartlett & Mendelson, 2002) of  $\mathcal{H}$  with sample size  $n$  over  $\mu$ .

**Lemma 4.** Suppose the loss function  $\ell$  is  $\rho$ -Lipschitz with respect to the first argument ( $0 < \rho < \infty$ ), and all the functions in the model class  $\mathcal{F}$  are bounded, i.e., there exists a constant  $C_b$  such that  $\|f\|_\infty \leq C_b$  for any  $f \in \mathcal{F}$ . Let  $C_\ell := \sup_{t=\pm 1} \ell(C_b, t)$ . For any  $\delta > 0$ , with probability  $1 - \delta$ ,

$$\sup_{f \in \mathcal{F}} \left| R_{\text{PC}}^+(f) - \widehat{R}_{\text{PC}}^+(f) \right| \leq (1 + \pi_+) 2\rho \widetilde{\mathfrak{R}}_n^+(\mathcal{F}) + (1 + \pi_+) C_\ell \sqrt{\frac{\log \frac{4}{\delta}}{2n}}.$$

*Proof.* By the definition of  $R_{\text{PC}}^+(f)$  and  $\widehat{R}_{\text{PC}}^+(f)$ , we can obtain

$$\begin{aligned} \sup_{f \in \mathcal{F}} \left| R_{\text{PC}}^+(f) - \widehat{R}_{\text{PC}}^+(f) \right| &\leq \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{\tilde{p}_+(\mathbf{x})} [\ell(f(\mathbf{x}), +1)] - \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}), +1) \right| \\ &\quad + \pi_+ \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{\tilde{p}_+(\mathbf{x})} [\ell(f(\mathbf{x}), -1)] - \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}), -1) \right|. \end{aligned} \quad (10)$$

By applying Lemma 3, we have for any  $\delta > 0$ , with probability  $1 - \delta$ ,

$$\sup_{f \in \mathcal{F}} \left| \mathbb{E}_{\tilde{p}_+(\mathbf{x})} [\ell(f(\mathbf{x}), +1)] - \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}), +1) \right| \leq 2\widetilde{\mathfrak{R}}_n^+(\ell \circ \mathcal{F}) + C_\ell \sqrt{\frac{\log \frac{2}{\delta}}{2n}}, \quad (11)$$

and for any  $\delta > 0$ , with probability  $1 - \delta$ ,

$$\sup_{f \in \mathcal{F}} \left| \mathbb{E}_{\tilde{p}_+(\mathbf{x})} [\ell(f(\mathbf{x}), -1)] - \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}), -1) \right| \leq 2\widetilde{\mathfrak{R}}_n^+(\ell \circ \mathcal{F}) + C_\ell \sqrt{\frac{\log \frac{2}{\delta}}{2n}}, \quad (12)$$

where  $\ell \circ \mathcal{F}$  means  $\{\ell \circ f \mid f \in \mathcal{F}\}$ . By Talagrand's lemma (Lemma 4.2 in (Mohri et al., 2012)),

$$\widetilde{\mathfrak{R}}_n^+(\ell \circ \mathcal{F}) \leq \rho \widetilde{\mathfrak{R}}_n^+(\mathcal{F}). \quad (13)$$

Finally, by combing Eqs. (10), (11), (12), and (13), we have for any  $\delta > 0$ , with probability at least  $1 - \delta$ ,

$$\sup_{f \in \mathcal{F}} \left| R_{\text{PC}}^+(f) - \widehat{R}_{\text{PC}}^+(f) \right| \leq (1 + \pi_+) 2\rho \widetilde{\mathfrak{R}}_n^+(\mathcal{F}) + (1 + \pi_+) C_\ell \sqrt{\frac{\log \frac{4}{\delta}}{2n}}, \quad (14)$$

which concludes the proof of Lemma 4.  $\square$

**Lemma 5.** Suppose the loss function  $\ell$  is  $\rho$ -Lipschitz with respect to the first argument ( $0 < \rho < \infty$ ), and all the functions in the model class  $\mathcal{F}$  are bounded, i.e., there exists a constant  $C_b$  such that  $\|f\|_\infty \leq C_b$  for any  $f \in \mathcal{F}$ . Let  $C_\ell := \sup_{t=\pm 1} \ell(C_b, t)$ . For any  $\delta > 0$ , with probability  $1 - \delta$ ,

$$\sup_{f \in \mathcal{F}} \left| R_{\text{PC}}^-(f) - \widehat{R}_{\text{PC}}^-(f) \right| \leq (1 + \pi_-) 2\rho \widetilde{\mathfrak{R}}_n^-(\mathcal{F}) + (1 + \pi_-) C_\ell \sqrt{\frac{\log \frac{4}{\delta}}{2n}}.$$

*Proof.* Lemma 5 can be proved similarly to Lemma 4.  $\square$

By combining Lemma 2, Lemma 4, and Lemma 5, Theorem 4 is proved.  $\square$

## F. Proof of Theorem 5

Suppose there are  $n$  pairs of paired data points, which means there are in total  $2n$  data points. For our Pcomp classification problem, we could simply regard  $\mathbf{x}$  sampled from  $\tilde{p}_+(\mathbf{x})$  as (noisy) positive data and  $\mathbf{x}'$  sampled from  $\tilde{p}_-(\mathbf{x}')$  as (noisy) negative data. Given  $n$  pairs of examples  $\{(\mathbf{x}_i, \mathbf{x}'_i)\}_{i=1}^n$ , for the  $n$  observed positive examples, there are actually  $n \cdot p(y = +1 | \tilde{y} = +1)$  true positive examples; for the  $n$  observed negative examples, there are actually  $n \cdot p(y = -1 | \tilde{y} = -1)$  true negative examples. From our defined data generation process in Theorem 1, it is intuitive to obtain

$$p(y = +1 | \tilde{y} = +1) = \frac{\pi_+^2 + \pi_+ \pi_-}{\pi_+^2 + \pi_-^2 + \pi_+ \pi_-} = \frac{\pi_+}{\pi_+^2 + \pi_-^2 + \pi_+ \pi_-},$$

$$p(y = -1 | \tilde{y} = -1) = \frac{\pi_-^2 + \pi_+ \pi_-}{\pi_+^2 + \pi_-^2 + \pi_+ \pi_-} = \frac{\pi_-}{\pi_+^2 + \pi_-^2 + \pi_+ \pi_-}.$$

Since  $\phi_+ = p(y = -1 | \tilde{y} = +1) = 1 - p(y = +1 | \tilde{y} = +1)$  and  $\phi_- = p(y = +1 | \tilde{y} = -1) = 1 - p(y = -1 | \tilde{y} = -1)$ , we can obtain

$$\phi_+ = p(y = -1 | \tilde{y} = +1) = 1 - \frac{\pi_+}{\pi_+^2 + \pi_-^2 + \pi_+ \pi_-} = \frac{\pi_-^2}{\pi_+^2 + \pi_-^2 + \pi_+ \pi_-},$$

$$\phi_- = p(y = +1 | \tilde{y} = -1) = 1 - \frac{\pi_-}{\pi_+^2 + \pi_-^2 + \pi_+ \pi_-} = \frac{\pi_+^2}{\pi_+^2 + \pi_-^2 + \pi_+ \pi_-}.$$

In this way, we can further obtain the following noise transition ratios:

$$\rho_+ = p(\tilde{y} = -1 | y = +1) = \frac{p(y = +1 | \tilde{y} = -1)p(\tilde{y} = -1)}{p(y = +1 | \tilde{y} = -1)p(\tilde{y} = -1) + p(y = +1 | \tilde{y} = +1)p(\tilde{y} = +1)} = \frac{\pi_+}{1 + \pi_+},$$

$$\rho_- = p(\tilde{y} = +1 | y = -1) = \frac{p(y = -1 | \tilde{y} = +1)p(\tilde{y} = +1)}{p(y = -1 | \tilde{y} = +1)p(\tilde{y} = +1) + p(y = -1 | \tilde{y} = -1)p(\tilde{y} = -1)} = \frac{\pi_-}{1 + \pi_-},$$

where  $p(\tilde{y} = 1) = p(\tilde{y} = -1) = \frac{1}{2}$ , because we have the same number of observed positive examples and negative examples.

## G. Proof of Theorem 7

First of all, we introduce the following notations:

$$R_{\text{pPC}}^+(f) = \mathbb{E}_{\tilde{p}_+(\mathbf{x})}[\ell(f(\mathbf{x}), +1)\mathbb{I}[\mathbf{x} \in \text{PP}\tilde{]}],$$

$$\hat{R}_{\text{pPC}}^+(f) = \frac{1}{n} \sum_{i=1}^n (\ell(f(\mathbf{x}_i), +1)\mathbb{I}[\mathbf{x}_i \in \text{PP}\tilde{]}),$$

$$R_{\text{pPC}}^-(f) = \mathbb{E}_{\tilde{p}_-(\mathbf{x}') }[\ell(f(\mathbf{x}'), -1)\mathbb{I}[\mathbf{x}' \in \text{NN}\tilde{]}],$$

$$\hat{R}_{\text{pPC}}^-(f) = \frac{1}{n} \sum_{i=1}^n (\ell(f(\mathbf{x}'_i), -1)\mathbb{I}[\mathbf{x}'_i \in \text{NN}\tilde{]}).$$

In this way, we could simply represent  $R_{\text{pPC}}(f)$  and  $\hat{R}_{\text{pPC}}(f)$  as

$$R_{\text{pPC}}(f) = \frac{1}{1 - \rho_+} R_{\text{pPC}}^+(f) + \frac{1}{1 - \rho_-} R_{\text{pPC}}^-(f), \quad \hat{R}_{\text{pPC}}(f) = \frac{1}{1 - \rho_+} \hat{R}_{\text{pPC}}^+(f) + \frac{1}{1 - \rho_-} \hat{R}_{\text{pPC}}^-(f).$$

Then we have the following lemma.

**Lemma 6.** *The following inequality holds:*

$$R(\hat{f}_{\text{pPC}}) - R(f^*) \leq \frac{2}{1 - \rho_+} \sup_{f \in \mathcal{F}} |R_{\text{pPC}}^+(f) - \hat{R}_{\text{pPC}}^+(f)| + \frac{2}{1 - \rho_-} \sup_{f \in \mathcal{F}} |R_{\text{pPC}}^-(f) - \hat{R}_{\text{pPC}}^-(f)|. \quad (15)$$

*Proof.* We omit the proof of Lemma 6 since it is quite similar to that of Lemma 2.  $\square$

As suggested by Lemma 6, we need to further upper bound the right hand size of Eq. (15). According to Lemma 3, we have the following two lemmas.

**Lemma 7.** *Suppose the loss function  $\ell$  is  $\rho$ -Lipschitz with respect to the first argument ( $0 < \rho < \infty$ ), and all the functions in the model class  $\mathcal{F}$  are bounded, i.e., there exists a constant  $C_b$  such that  $\|f\|_\infty \leq C_b$  for any  $f \in \mathcal{F}$ . Let  $C_\ell := \sup_{z \leq C_b, t = \pm 1} \ell(z, t)$ . For any  $\delta > 0$ , with probability  $1 - \delta$ ,*

$$\sup_{f \in \mathcal{F}} \left| R_{\text{pPC}}^+(f) - \widehat{R}_{\text{pPC}}^+(f) \right| \leq 2\rho \widetilde{\mathfrak{R}}_n^+(\mathcal{F}) + C_\ell \sqrt{\frac{\log \frac{2}{\delta}}{2n}}.$$

**Lemma 8.** *Suppose the loss function  $\ell$  is  $\rho$ -Lipschitz with respect to the first argument ( $0 < \rho < \infty$ ), and all the functions in the model class  $\mathcal{F}$  are bounded, i.e., there exists a constant  $C_b$  such that  $\|f\|_\infty \leq C_b$  for any  $f \in \mathcal{F}$ . Let  $C_\ell := \sup_{z \leq C_b, t = \pm 1} \ell(z, t)$ . For any  $\delta > 0$ , with probability  $1 - \delta$ ,*

$$\sup_{f \in \mathcal{F}} \left| R_{\text{pPC}}^-(f) - \widehat{R}_{\text{pPC}}^-(f) \right| \leq 2\rho \widetilde{\mathfrak{R}}_n^-(\mathcal{F}) + C_\ell \sqrt{\frac{\log \frac{2}{\delta}}{2n}}.$$

We omit the proofs of Lemma 7 and Lemma 8 since they are similar to that of Lemma 4.

By combing Lemma 6, Lemma 7, and Lemma 8, Theorem 7 is proved.

## H. Supplementary Information of Experiments

We report the detailed information of the used datasets as follows.

**MNIST<sup>1</sup>** (LeCun et al., 1998). This is a grayscale image dataset composed of handwritten digits from 0 to 9 where the size of the each image is  $28 \times 28$ . It contains 60,000 training images and 10,000 test images. Because the original dataset has 10 classes, we regard the even digits as the positive class and the odd digits as the negative class. We generate 30,000 pointwise corrupted examples from MNIST for model training.

**Fashion-MNIST<sup>2</sup>** (Xiao et al., 2017). Similarly to MNIST, this is also a grayscale image dataset composed of fashion items (‘T-shirt’, ‘trouser’, ‘pullover’, ‘dress’, ‘sandal’, ‘coat’, ‘shirt’, ‘sneaker’, ‘bag’, and ‘ankle boot’). It contains 60,000 training examples and 10,000 test examples. It is converted into a binary classification dataset as follows:

- The positive class is formed by ‘T-shirt’, ‘pullover’, ‘coat’, ‘shirt’, and ‘bag’.
- The negative class is formed by ‘trouser’, ‘dress’, ‘sandal’, ‘sneaker’, and ‘ankle boot’.

We generate 30,000 pointwise corrupted examples from Fashion-MNIST for model training.

**Kuzushiji-MNIST<sup>3</sup>** (Netzer et al., 2011). This is another grayscale image dataset that is similar to MNIST. It is a 10-class dataset of cursive Japanese (‘‘Kuzushiji’’) characters. It consists of 60,000 training images and 10,000 test images. It is converted into a binary classification dataset as follows:

- The positive class is formed by ‘o’, ‘su’, ‘na’, ‘ma’, ‘re’.
- The negative class is formed by ‘ki’, ‘tsu’, ‘ha’, ‘ya’, ‘wo’.

We generate 30,000 pointwise corrupted examples from Kuzushiji-MNIST for model training.

**CIFAR-10<sup>4</sup>** (Krizhevsky et al., 2009). This is also a color image dataset of 10 different objects (‘airplane’, ‘bird’, ‘automobile’, ‘cat’, ‘deer’, ‘dog’, ‘frog’, ‘horse’, ‘ship’, and ‘truck’), where the size of each image is  $32 \times 32 \times 3$ . There are 5,000 training images and 1,000 test images per class. This dataset is converted into a binary classification dataset as follows:

<sup>1</sup><http://yann.lecun.com/exdb/mnist/>

<sup>2</sup><https://github.com/zalandoresearch/fashion-mnist>

<sup>3</sup><https://github.com/rois-codh/kmnist>

<sup>4</sup><https://www.cs.toronto.edu/~kriz/cifar.html>

- The positive class is formed by ‘bird’, ‘deer’, ‘dog’, ‘frog’, ‘cat’, and ‘horse’.
- The negative class is formed by ‘airplane’, ‘automobile’, ‘ship’, and ‘truck’.

We generate 30,000 pointwise corrupted examples from CIFAR-10 for model training.

**USPS, Pendigits, Otdigits.** These datasets are composed of handwritten digits from 0 to 9. Because each of the original datasets has 10 classes, we regard the even digits as the positive class and the odd digits as the negative class. We generate 4,000 pointwise corrupted examples from USPS (5,000 from Pendigits and 2,000 from Otdigits) for model training.

**CNAE-9.** This dataset contains 1,080 documents of free text business descriptions of Brazilian companies categorized into a subset of 9 categories cataloged in a table called National Classification of Economic Activities.

- The positive class is formed by ‘2’, ‘4’, ‘6’ and ‘8’.
- The negative class is formed by ‘1’, ‘3’, ‘5’, ‘7’ and ‘9’.

We generate 400 pointwise corrupted examples from CNAE-9 for model training.

For the four datasets, USPS can be downloaded from the website of the late Sam Roweis<sup>5</sup>, and the other three datasets can be downloaded from the UCI machine learning repository<sup>6</sup>.

For MNIST, Kuzushiji-MNIST, and Fashion-MNIST, we set learning rate to  $1e - 3$  and weight decay to  $1e - 5$ . For CIFAR-10, we set learning rate to  $1e - 3$  and weight decay to  $1e - 3$ . For the four datasets including USPS, Pendigits, Otdigits, and CNAE-9, we search learning rate and weight decay from  $\{1e - 5, 1e - 4, \dots, 1e - 1\}$  for all learning methods. For Pcomp-teacher, the regularization parameter is searched from  $\{1e - 3, 1e - 2, \dots, 1e + 3\}$  and the exponential moving average decay is fixed at 0.97. Hyper-parameters for all learning methods are selected so as to maximize the accuracy of five-fold cross validation on the training set.

## References

- Bartlett, P. L. and Mendelson, S. Rademacher and gaussian complexities: Risk bounds and structural results. *JMLR*, 3(11): 463–482, 2002.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- LeCun, Y., Bottou, L., Bengio, Y., Haffner, P., et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Mohri, M., Rostamizadeh, A., and Talwalkar, A. *Foundations of Machine Learning*. MIT Press, 2012.
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. Reading digits in natural images with unsupervised feature learning. In *NeurIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
- Xiao, H., Rasul, K., and Vollgraf, R. Fashion-mnist: A novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.

<sup>5</sup><http://cs.nyu.edu/~roweis/data.html>

<sup>6</sup><http://archive.ics.uci.edu/ml/>