

# Appendix

## 1 PROOFS

In this section, we provide proofs for all the lemmas and theorems in the main paper. We always assume that a class-conditional extension of the Classification Noise Process (CNP) (Angluin & Laird, 1988) maps true labels  $y$  to observed labels  $s$  such that each label in  $P$  is flipped independently with probability  $\rho_1$  and each label in  $N$  is flipped independently with probability  $\rho_0$  ( $s \leftarrow \text{CNP}(y, \rho_1, \rho_0)$ ), so that  $P(s = s|y = y, x) = P(s = s|y = y)$ . Remember that  $\rho_1 + \rho_0 < 1$  is a necessary condition of minimal information, otherwise we may learn opposite labels.

In Lemma 1, Theorem 2, Lemma 3 and Theorem 4, we assume that  $P$  and  $N$  have infinite number of examples so that they are the true, hidden distributions.

A fundamental equation we use in the proofs is the following lemma:

**Lemma S1** *When  $g$  is ideal, i.e.  $g(x) = g^*(x)$  and  $P$  and  $N$  have non-overlapping support, we have*

$$g(x) = (1 - \rho_1) \cdot \mathbb{1}[[y = 1]] + \rho_0 \cdot \mathbb{1}[[y = 0]] \quad (1)$$

**Proof:** Since  $g(x) = g^*(x)$  and  $P$  and  $N$  have non-overlapping support, we have

$$\begin{aligned} g(x) &= g^*(x) = P(s = 1|x) \\ &= P(s = 1|y = 1, x) \cdot P(y = 1|x) + P(s = 1|y = 0, x) \cdot P(y = 0|x) \\ &= P(s = 1|y = 1) \cdot P(y = 1|x) + P(s = 1|y = 0) \cdot P(y = 0|x) \\ &= (1 - \rho_1) \cdot \mathbb{1}[[y = 1]] + \rho_0 \cdot \mathbb{1}[[y = 0]] \end{aligned}$$

### 1.1 PROOF OF LEMMA 1

**Lemma 1** *When  $g$  is ideal, i.e.  $g(x) = g^*(x)$  and  $P$  and  $N$  have non-overlapping support, we have*

$$\begin{cases} \tilde{P}_{y=1} = \{x \in P | s = 1\}, \tilde{N}_{y=1} = \{x \in P | s = 0\} \\ \tilde{P}_{y=0} = \{x \in N | s = 1\}, \tilde{N}_{y=0} = \{x \in N | s = 0\} \end{cases} \quad (2)$$

**Proof:** Firstly, we compute the threshold  $LB_{y=1}$  and  $UB_{y=0}$  used by  $\tilde{P}_{y=1}$ ,  $\tilde{N}_{y=1}$ ,  $\tilde{P}_{y=0}$  and  $\tilde{N}_{y=0}$ . Since  $P$  and  $N$  have non-overlapping support, we have  $P(y = 1|x) = \mathbb{1}[[y = 1]]$ . Also using  $g(x) = g^*(x)$ , we have

$$\begin{aligned} LB_{y=1} &= E_{x \in \tilde{P}}[g(x)] = E_{x \in \tilde{P}}[P(s = 1|x)] \\ &= E_{x \in \tilde{P}}[P(s = 1|x, y = 1)P(y = 1|x) + P(s = 1|x, y = 0)P(y = 0|x)] \\ &= E_{x \in \tilde{P}}[P(s = 1|y = 1)P(y = 1|x) + P(s = 1|y = 0)P(y = 0|x)] \\ &= (1 - \rho_1)(1 - \pi_1) + \rho_0\pi_1 \end{aligned} \quad (3)$$

Similarly, we have

$$UB_{y=0} = (1 - \rho_1)\pi_0 + \rho_0(1 - \pi_0)$$

Since  $\pi_1 = P(y = 0|s = 1)$ , we have  $\pi_1 \in [0, 1]$ . Furthermore, we have the requirement that  $\rho_1 + \rho_0 < 1$ , then  $\pi_1 = 1$  will lead to  $\rho_1 = P(s = 0|y = 1) = 1 - P(s = 1|y = 1) = 1 - \frac{P(y=1|s=1)P(s=1)}{P(y=1)} = 1 - 0 = 1$

which violates the requirement of  $\rho_1 + \rho_0 < 1$ . Therefore,  $\pi_1 \in [0, 1)$ . Similarly, we can prove  $\pi_0 \in [0, 1)$ . Therefore, we see that both  $LB_{y=1}$  and  $UB_{y=0}$  are interpolations of  $(1 - \rho_1)$  and  $\rho_0$ :

$$\begin{aligned}\rho_0 &< LB_{y=1} \leq 1 - \rho_1 \\ \rho_0 &\leq UB_{y=0} < 1 - \rho_1\end{aligned}$$

The first equality holds iff  $\pi_1 = 0$  and the second equality holds iff  $\pi_0 = 0$ .

Using Lemma S1, we know that under the condition of  $g(x) = g^*(x)$  and non-overlapping support,  $g(x) = (1 - \rho_1) \cdot \mathbb{1}[[y = 1]] + \rho_0 \cdot \mathbb{1}[[y = 0]]$ . In other words,

$$\begin{aligned}g(x) &\geq LB_{y=1} \Leftrightarrow x \in P \\ g(x) &\leq UB_{y=0} \Leftrightarrow x \in N\end{aligned}$$

Since

$$\begin{cases} \tilde{P}_{y=1} = \{x \in \tilde{P} | g(x) \geq LB_{y=1}\} \\ \tilde{N}_{y=1} = \{x \in \tilde{N} | g(x) \geq LB_{y=1}\} \\ \tilde{P}_{y=0} = \{x \in \tilde{P} | g(x) \leq UB_{y=0}\} \\ \tilde{N}_{y=0} = \{x \in \tilde{N} | g(x) \leq UB_{y=0}\} \end{cases}$$

where  $\tilde{P} = \{x | s = 1\}$  and  $\tilde{N} = \{x | s = 0\}$ , we have

$$\begin{cases} \tilde{P}_{y=1} = \{x \in P | s = 1\}, \tilde{N}_{y=1} = \{x \in P | s = 0\} \\ \tilde{P}_{y=0} = \{x \in N | s = 1\}, \tilde{N}_{y=0} = \{x \in N | s = 0\} \end{cases}$$

## 1.2 PROOF OF LEMMA 2

We restate Theorem 2 here:

**Theorem 2** *When  $g$  is ideal, i.e.  $g(x) = g^*(x)$  and  $P$  and  $N$  have non-overlapping support, we have*

$$\hat{\rho}_1^{conf} = \rho_1, \hat{\rho}_0^{conf} = \rho_0$$

**Proof:** Using the definition of  $\hat{\rho}_1^{conf}$  in the main paper:

$$\hat{\rho}_1^{conf} = \frac{|\tilde{N}_{y=1}|}{|\tilde{N}_{y=1}| + |\tilde{P}_{y=1}|}, \hat{\rho}_0^{conf} = \frac{|\tilde{P}_{y=0}|}{|\tilde{P}_{y=0}| + |\tilde{N}_{y=0}|}$$

Since  $g(x) = g^*(x)$  and  $P$  and  $N$  have non-overlapping support, using Lemma 1, we know

$$\begin{cases} \tilde{P}_{y=1} = \{x \in P | s = 1\}, \tilde{N}_{y=1} = \{x \in P | s = 0\} \\ \tilde{P}_{y=0} = \{x \in N | s = 1\}, \tilde{N}_{y=0} = \{x \in N | s = 0\} \end{cases}$$

Since  $\rho_1 = P(s = 0 | y = 1)$  and  $\rho_0 = P(s = 1 | y = 0)$ , we immediately have

$$\hat{\rho}_1^{conf} = \frac{|\{x \in P | s = 0\}|}{|P|} = \rho_1, \hat{\rho}_0^{conf} = \frac{|\{x \in N | s = 1\}|}{|N|} = \rho_0$$

### 1.3 PROOF OF LEMMA 3

We rewrite Lemma 3 below:

**Lemma 3** *When  $g$  is unassuming, i.e.,  $\Delta g(x) := g(x) - g^*(x)$  can be nonzero, and  $P$  and  $N$  can have overlapping support, we have*

$$\begin{cases} LB_{y=1} = LB_{y=1}^* + E_{x \in \hat{P}}[\Delta g(x)] - \frac{(1-\rho_1-\rho_0)^2}{p_{s1}} \Delta p_o \\ UB_{y=0} = UB_{y=0}^* + E_{x \in \hat{N}}[\Delta g(x)] + \frac{(1-\rho_1-\rho_0)^2}{1-p_{s1}} \Delta p_o \\ \hat{\rho}_1^{conf} = \rho_1 + \frac{1-\rho_1-\rho_0}{|P| - |\Delta P_1| + |\Delta N_1|} |\Delta N_1| \\ \hat{\rho}_0^{conf} = \rho_0 + \frac{1-\rho_1-\rho_0}{|N| - |\Delta N_0| + |\Delta P_0|} |\Delta P_0| \end{cases} \quad (4)$$

where

$$\begin{cases} \Delta p_o := \frac{|P \cap N|}{|P \cup N|} \\ \Delta P_1 = \{x \in P | g(x) < LB_{y=1}\} \\ \Delta N_1 = \{x \in N | g(x) \geq LB_{y=1}\} \\ \Delta P_0 = \{x \in P | g(x) \leq UB_{y=0}\} \\ \Delta N_0 = \{x \in N | g(x) > UB_{y=0}\} \end{cases} \quad (5)$$

**Proof:** We first calculate  $LB_{y=1}$  and  $UB_{y=0}$  under unassuming condition, then calculate  $\hat{\rho}_i^{conf}$ ,  $i = 0, 1$  under unassuming condition.

Note that  $\Delta p_o$  can also be expressed as

$$\Delta p_o := \frac{|P \cap N|}{|P \cup N|} = P(\hat{y} = 1, y = 0) = P(\hat{y} = 0, y = 1)$$

Here  $P(\hat{y} = 1, y = 0) \equiv P(\hat{y} = 1 | y = 0)P(y = 0)$ , where  $P(\hat{y} = 1 | y = 0)$  means for a perfect classifier  $f^*(x) = P(y = 1 | x)$ , the expected probability that it will label a  $y = 0$  example as positive ( $\hat{y} = 1$ ).

#### (1) $LB_{y=1}$ and $UB_{y=0}$ under unassuming condition

Firstly, we calculate  $LB_{y=1}$  and  $UB_{y=0}$  with perfect probability estimation  $g^*(x)$ , but the support may overlap. Secondly, we allow the probability estimation to be imperfect, superimposed onto the overlapping support condition, and calculate  $LB_{y=1}$  and  $UB_{y=0}$ .

#### I. Calculating $LB_{y=1}$ and $UB_{y=0}$ when $g(x) = g^*(x)$ and support may overlap

With overlapping support, we no longer have  $P(y = 1 | x) = \mathbb{1}[[y = 1]]$ . Instead, we have

$$\begin{aligned} LB_{y=1} &= E_{x \in \hat{P}}[g^*(x)] = E_{x \in \hat{P}}[P(s = 1 | x)] \\ &= E_{x \in \hat{P}}[P(s = 1 | x, y = 1)P(y = 1 | x) + P(s = 1 | x, y = 0)P(y = 0 | x)] \\ &= E_{x \in \hat{P}}[P(s = 1 | y = 1)P(y = 1 | x) + P(s = 1 | y = 0)P(y = 0 | x)] \\ &= (1 - \rho_1) \cdot E_{x \in \hat{P}}[P(y = 1 | x)] + \rho_0 \cdot E_{x \in \hat{P}}[P(y = 0 | x)] \\ &= (1 - \rho_1) \cdot P(\hat{y} = 1 | s = 1) + \rho_0 \cdot P(\hat{y} = 0 | s = 1) \end{aligned}$$

Here  $P(\hat{y} = 1|s = 1)$  can be calculated using  $\Delta p_o$ :

$$\begin{aligned}
P(\hat{y} = 1|s = 1) &= \frac{P(\hat{y} = 1, s = 1)}{P(s = 1)} \\
&= \frac{P(\hat{y} = 1, y = 1, s = 1) + P(\hat{y} = 1, y = 0, s = 1)}{P(s = 1)} \\
&= \frac{P(s = 1|y = 1)P(\hat{y} = 1, y = 1) + P(s = 1|y = 0)P(\hat{y} = 1, y = 0)}{P(s = 1)} \\
&= \frac{(1 - \rho_1)(p_{y1} - \Delta p_o) + \rho_0 \Delta p_o}{p_{s1}} \\
&= (1 - \pi_1) - \frac{1 - \rho_1 - \rho_0}{p_{s1}} \Delta p_o
\end{aligned}$$

Hence,

$$P(\hat{y} = 0|s = 1) = 1 - P(\hat{y} = 1|s = 1) = \pi_1 + \frac{1 - \rho_1 - \rho_0}{p_{s1}} \Delta p_o$$

Therefore,

$$\begin{aligned}
LB_{y=1} &= (1 - \rho_1) \cdot P(\hat{y} = 1|s = 1) + \rho_0 \cdot P(\hat{y} = 0|s = 1) \\
&= (1 - \rho_1) \cdot \left( (1 - \pi_1) - \frac{1 - \rho_1 - \rho_0}{p_{s1}} \Delta p_o \right) + \rho_0 \cdot \left( \pi_1 + \frac{1 - \rho_1 - \rho_0}{p_{s1}} \Delta p_o \right) \\
&= LB_{y=1}^* - \frac{(1 - \rho_1 - \rho_0)^2}{p_{s1}} \Delta p_o
\end{aligned} \tag{6}$$

where  $LB_{y=1}^*$  is the  $LB_{y=1}$  value when  $g(x)$  is ideal. We see in Eq. (6) that the overlapping support introduces a non-positive correction to  $LB_{y=1}^*$  compared with the ideal condition.

Similarly, we have

$$UB_{y=0} = UB_{y=0}^* + \frac{(1 - \rho_1 - \rho_0)^2}{1 - p_{s1}} \Delta p_o \tag{7}$$

## II. Calculating $LB_{y=1}$ and $UB_{y=0}$ when $g$ is unassuming

Define  $\Delta g(x) := g(x) - g^*(x)$ . When the support may overlap, we have

$$\begin{aligned}
LB_{y=1} &= E_{x \in \tilde{P}}[g(x)] \\
&= E_{x \in \tilde{P}}[g^*(x)] + E_{x \in \tilde{P}}[\Delta g(x)] \\
&= LB_{y=1}^* - \frac{(1 - \rho_1 - \rho_0)^2}{p_{s1}} \Delta p_o + E_{x \in \tilde{P}}[\Delta g(x)]
\end{aligned} \tag{8}$$

Similarly, we have

$$\begin{aligned}
UB_{y=0} &= E_{x \in \tilde{N}}[g(x)] \\
&= E_{x \in \tilde{N}}[g^*(x)] + E_{x \in \tilde{N}}[\Delta g(x)] \\
&= UB_{y=0}^* + \frac{(1 - \rho_1 - \rho_0)^2}{1 - p_{s1}} \Delta p_o + E_{x \in \tilde{N}}[\Delta g(x)]
\end{aligned} \tag{9}$$

In summary, Eq. (8) (9) give the expressions for  $LB_{y=1}$  and  $UB_{y=0}$ , respectively, when  $g$  is unassuming.

**(2)  $\hat{\rho}_i^{conf}$  under unassuming condition**

Now let's calculate  $\hat{\rho}_i^{conf}$ ,  $i = 0, 1$ . For simplicity, define

$$\begin{cases} PP = \{x \in P | s = 1\} \\ PN = \{x \in P | s = 0\} \\ NP = \{x \in N | s = 1\} \\ NN = \{x \in N | s = 0\} \\ \Delta_{PP_1} = \{x \in PP | g(x) < LB_{y=1}\} \\ \Delta_{NP_1} = \{x \in NP | g(x) \geq LB_{y=1}\} \\ \Delta_{PN_1} = \{x \in PN | g(x) < LB_{y=1}\} \\ \Delta_{NN_1} = \{x \in NN | g(x) \geq LB_{y=1}\} \end{cases} \quad (10)$$

For  $\hat{\rho}_1^{conf}$ , we have:

$$\hat{\rho}_1^{conf} = \frac{|\tilde{N}_{y=1}|}{|\tilde{P}_{y=1}| + |\tilde{N}_{y=1}|}$$

Here

$$\begin{aligned} \tilde{P}_{y=1} &= \{x \in \tilde{P} | g(x) \geq LB_{y=1}\} \\ &= \{x \in PP | g(x) \geq LB_{y=1}\} \cup \{x \in NP | g(x) \geq LB_{y=1}\} \\ &= (PP \setminus \Delta_{PP_1}) \cup \Delta_{NP_1} \end{aligned}$$

Similarly, we have

$$\tilde{N}_{y=1} = (PN \setminus \Delta_{PN_1}) \cup \Delta_{NN_1}$$

Therefore

$$\begin{aligned} \hat{\rho}_1^{conf} &= \frac{|PN| - |\Delta_{PN_1}| + |\Delta_{NN_1}|}{[(|PP| - |\Delta_{PP_1}|) + (|PN| - |\Delta_{PN_1}|)] + (|\Delta_{NN_1}| + |\Delta_{NP_1}|)} \\ &= \frac{|PN| - |\Delta_{PN_1}| + |\Delta_{NN_1}|}{|P| - |\Delta_{P_1}| + |\Delta_{N_1}|} \end{aligned} \quad (11)$$

where in the second equality we have used the definition of  $\Delta_{P_1}$  and  $\Delta_{N_1}$  in Eq. (5).

Using the definition of  $\rho_1$ , we have

$$\begin{aligned} \frac{|PN| - |\Delta_{PN_1}|}{|P| - |\Delta_{P_1}|} &= \frac{|\{x \in PN | g(x) \geq LB_{y=1}\}|}{|\{x \in P | g(x) \geq LB_{y=1}\}|} \\ &= \frac{P(x \in PN, g(x) \geq LB_{y=1})}{P(x \in P, g(x) \geq LB_{y=1})} \\ &= \frac{P(x \in PN | x \in P, g(x) \geq LB_{y=1}) \cdot P(x \in P, g(x) \geq LB_{y=1})}{P(x \in P, g(x) \geq LB_{y=1})} \\ &= \frac{P(x \in PN | x \in P) \cdot P(x \in P, g(x) \geq LB_{y=1})}{P(x \in P, g(x) \geq LB_{y=1})} \\ &= \rho_1 \end{aligned}$$

Here we have used the property of CNP that  $(s \perp x)|y$ , leading to  $P(x \in PN|x \in P, g(x) \geq LB_{y=1}) = P(x \in PN|x \in P) = \rho_1$ .

Similarly, we have

$$\frac{|\Delta_{NN_1}|}{|\Delta N_1|} = 1 - \rho_0$$

Combining with Eq. (11), we have

$$\hat{\rho}_1^{conf} = \rho_1 + \frac{1 - \rho_1 - \rho_0}{|P| - |\Delta P_1| + |\Delta N_1|} |\Delta N_1| \quad (12)$$

Similarly, we have

$$\hat{\rho}_0^{conf} = \rho_0 + \frac{1 - \rho_1 - \rho_0}{|N| - |\Delta N_0| + |\Delta P_0|} |\Delta P_0| \quad (13)$$

From the two equations above, we see that

$$\hat{\rho}_1^{conf} \geq \rho_1, \hat{\rho}_0^{conf} \geq \rho_0 \quad (14)$$

In other words,  $\hat{\rho}_i^{conf}$  is an **upper bound** of  $\rho_i$ ,  $i = 0, 1$ . The equality for  $\hat{\rho}_1^{conf}$  holds if  $|\Delta N_1| = 0$ . The equality for  $\hat{\rho}_0^{conf}$  holds if  $|\Delta P_0| = 0$ .

## 1.4 PROOF OF LEMMA 4

Let's restate Theorem 4 below:

**Theorem 4** *Given non-overlapping support condition,*

*If  $\forall x \in N, \Delta g(x) < LB_{y=1} - \rho_0$ , then  $\hat{\rho}_1^{conf} = \rho_1$ .*

*If  $\forall x \in P, \Delta g(x) > -(1 - \rho_1 - UB_{y=0})$ , then  $\hat{\rho}_0^{conf} = \rho_0$ .*

Theorem 4 directly follows from Eq. (12) and (13). Assuming non-overlapping support, we have  $g^*(x) = P(s = 1|x) = (1 - \rho_1) \cdot \mathbb{1}[[y = 1]] + \rho_0 \cdot \mathbb{1}[[y = 0]]$ . In other words, the contribution of overlapping support to  $|\Delta N_1|$  and  $|\Delta P_0|$  is 0. The only source of deviation comes from imperfect  $g(x)$ .

For the first half of the theorem, since  $\forall x \in N, \Delta g(x) < LB_{y=1} - \rho_0$ , we have  $\forall x \in N, g(x) = \Delta g(x) + g^*(x) < (LB_{y=1} - \rho_0) + \rho_0 = LB_{y=1}$ , then  $|\Delta N_1| = |\{x \in N | g(x) \geq LB_{y=1}\}| = 0$ , so we have  $\hat{\rho}_1^{conf} = \rho_1$ .

Similarly, for the second half of the theorem, since  $\forall x \in P, \Delta g(x) > -(1 - \rho_1 - UB_{y=0})$ , then  $|\Delta P_0| = |\{x \in P | g(x) \leq UB_{y=0}\}| = 0$ , so we have  $\hat{\rho}_0^{conf} = \rho_0$ .

## 1.5 PROOF OF LEMMA 5

Theorem 5 reads as follows:

**Theorem 5** *If  $g$  range separates  $P$  and  $N$  and  $\hat{\rho}_i = \rho_i$ ,  $i = 0, 1$ , then for any classifier  $f_\theta$  and any bounded loss function  $l(\hat{y}_i, y_i)$ , we have*

$$R_{\tilde{\mathcal{D}}_\rho}(f_\theta) = R_{\mathcal{D}}(f_\theta) \quad (15)$$

where  $\tilde{l}(\hat{y}_i, s_i)$  is Rank Pruning's loss function given by

$$\tilde{l}(\hat{y}_i, s_i) = \frac{1}{1 - \hat{\rho}_1} l(\hat{y}_i, s_i) \cdot \mathbb{1}[[x_i \in \tilde{P}_{conf}]] + \frac{1}{1 - \hat{\rho}_0} l(\hat{y}_i, s_i) \cdot \mathbb{1}[[x_i \in \tilde{N}_{conf}]] \quad (16)$$

and  $\tilde{P}_{conf}$  and  $\tilde{N}_{conf}$  are given by

$$\tilde{P}_{conf} := \{x \in \tilde{P} \mid g(x) \geq k_1\}, \tilde{N}_{conf} := \{x \in \tilde{N} \mid g(x) \leq k_0\} \quad (17)$$

where  $k_1$  is the  $(\hat{\pi}_1|\tilde{P}|)^{th}$  smallest  $g(x)$  for  $x \in \tilde{P}$  and  $k_0$  is the  $(\hat{\pi}_0|\tilde{N}|)^{th}$  largest  $g(x)$  for  $x \in \tilde{N}$

**Proof:**

Since  $\tilde{P}$  and  $\tilde{N}$  are constructed from  $P$  and  $N$  with noise rates  $\pi_1$  and  $\pi_0$  using the class-conditional extension of the Classification Noise Process (Angluin & Laird, 1988), we have

$$\begin{cases} P = PP \cup PN \\ N = NP \cup NN \\ \tilde{P} = PP \cup NP \\ \tilde{N} = PN \cup NN \end{cases} \quad (18)$$

where

$$\begin{cases} PP = \{x \in P \mid s = 1\} \\ PN = \{x \in P \mid s = 0\} \\ NP = \{x \in N \mid s = 1\} \\ NN = \{x \in N \mid s = 0\} \end{cases} \quad (19)$$

satisfying

$$\begin{cases} PP \sim PN \sim P \\ NP \sim NN \sim N \\ \frac{|NP|}{|\tilde{P}|} = \pi_1, \frac{|PP|}{|\tilde{P}|} = 1 - \pi_1 \\ \frac{|PN|}{|\tilde{N}|} = \pi_0, \frac{|NN|}{|\tilde{N}|} = 1 - \pi_0 \\ \frac{|PN|}{|P|} = \rho_1, \frac{|PP|}{|P|} = 1 - \rho_1 \\ \frac{|NP|}{|N|} = \rho_0, \frac{|NN|}{|N|} = 1 - \rho_0 \end{cases} \quad (20)$$

Here the  $\sim$  means obeying the same distribution.

Since  $g$  range separates  $P$  and  $N$ , there exists a real number  $z$  such that  $\forall x_1 \in P$  and  $\forall x_0 \in N$ , we have  $g(x_1) > z > g(x_0)$ . Since  $P = PP \cup PN$ ,  $N = NP \cup NN$ , we have

$$\begin{aligned} \forall x \in PP, g(x) > z; \forall x \in PN, g(x) > z; \\ \forall x \in NP, g(x) < z; \forall x \in NN, g(x) < z \end{aligned} \quad (21)$$

Since  $\hat{\rho}_1 = \rho_1$  and  $\hat{\rho}_0 = \rho_0$ , we have

$$\begin{cases} \hat{\pi}_1 = \frac{\hat{\rho}_0}{p_{s1}} \frac{1-p_{s1}-\hat{\rho}_1}{1-\hat{\rho}_1-\hat{\rho}_0} = \frac{\rho_0}{p_{s1}} \frac{1-p_{s1}-\rho_1}{1-\rho_1-\rho_0} = \pi_1 \equiv \frac{\rho_0|N|}{|\tilde{P}|} \\ \hat{\pi}_0 = \frac{\hat{\rho}_1}{1-p_{s1}} \frac{p_{s1}-\hat{\rho}_0}{1-\hat{\rho}_1-\hat{\rho}_0} = \frac{\rho_1}{1-p_{s1}} \frac{p_{s1}-\rho_0}{1-\rho_1-\rho_0} = \pi_0 \equiv \frac{\rho_1|P|}{|\tilde{N}|} \end{cases} \quad (22)$$

Therefore,  $\hat{\pi}_1|\tilde{P}| = \pi_1|\tilde{P}| = \rho_0|N|$ ,  $\hat{\pi}_0|\tilde{N}| = \pi_0|\tilde{N}| = \rho_1|P|$ . Using  $\tilde{P}_{conf}$  and  $\tilde{N}_{conf}$ 's definition in Eq. (17), and  $g(x)$ 's property in Eq. (21), we have

$$\tilde{P}_{conf} = PP \sim P, \tilde{N}_{conf} = NN \sim N \quad (23)$$

Hence  $P_{conf}$  and  $N_{conf}$  can be seen as a uniform downsampling of  $P$  and  $N$ , with a downsampling ratio of  $(1 - \rho_1)$  for  $P$  and  $(1 - \rho_0)$  for  $N$ . Then according to Eq. (16), the loss function  $\tilde{l}(\hat{y}_i, s_i)$  essentially sees a fraction of  $(1 - \rho_1)$  examples in  $P$  and a fraction of  $(1 - \rho_0)$  examples in  $N$ , with a final reweighting to restore the class balance. Then for any classifier  $f_\theta$  that maps  $x \rightarrow \hat{y}$  and any bounded loss function  $l(\hat{y}_i, y_i)$ , we have

$$\begin{aligned} R_{\tilde{l}, \mathcal{D}_\rho}(f_\theta) &= E_{(x,s) \sim \mathcal{D}_\rho}[\tilde{l}(f_\theta(x), s)] \\ &= \frac{1}{1-\hat{\rho}_1} \cdot E_{(x,s) \sim \mathcal{D}_\rho} [l(f_\theta(x), s) \cdot \mathbf{1}[[x \in \tilde{P}_{conf}]]] + \frac{1}{1-\hat{\rho}_0} \cdot E_{(x,s) \sim \mathcal{D}_\rho} [l(f_\theta(x), s) \cdot \mathbf{1}[[x \in \tilde{N}_{conf}]]] \\ &= \frac{1}{1-\rho_1} \cdot E_{(x,s) \sim \mathcal{D}_\rho} [l(f_\theta(x), s) \cdot \mathbf{1}[[x \in \tilde{P}_{conf}]]] + \frac{1}{1-\rho_0} \cdot E_{(x,s) \sim \mathcal{D}_\rho} [l(f_\theta(x), s) \cdot \mathbf{1}[[x \in \tilde{N}_{conf}]]] \\ &= \frac{1}{1-\rho_1} \cdot E_{(x,s) \sim \mathcal{D}_\rho} [l(f_\theta(x), s) \cdot \mathbf{1}[[x \in PP]]] + \frac{1}{1-\rho_0} \cdot E_{(x,s) \sim \mathcal{D}_\rho} [l(f_\theta(x), s) \cdot \mathbf{1}[[x \in NN]]] \\ &= \frac{1}{1-\rho_1} \cdot (1-\rho_1) \cdot E_{(x,y) \sim \mathcal{D}} [l(f_\theta(x), y) \cdot \mathbf{1}[[x \in P]]] + \\ &\quad \frac{1}{1-\rho_0} \cdot (1-\rho_0) \cdot E_{(x,y) \sim \mathcal{D}} [l(f_\theta(x), y) \cdot \mathbf{1}[[x \in N]]] \\ &= E_{(x,y) \sim \mathcal{D}} [l(f_\theta(x), y) \cdot \mathbf{1}[[x \in P]]] + l(f_\theta(x), y) \cdot \mathbf{1}[[x \in N]] \\ &= E_{(x,y) \sim \mathcal{D}} [l(f_\theta(x), y)] \\ &= R_{l, \mathcal{D}}(f_\theta) \end{aligned}$$

Therefore, we see that the expected risk for Rank Pruning with corrupted labels, is exactly the same as the expected risk for the true labels, for any bounded loss function  $l$  and classifier  $f_\theta$ . The reweighting ensures that after pruning, the two sets still remain unbiased w.r.t. to the true dataset.

Since the ideal condition is more strict than the range separability condition, we immediately have that when  $g$  is ideal and  $\hat{\rho}_i = \rho_i$ ,  $i = 0, 1$ ,  $R_{\tilde{l}, \mathcal{D}_\rho}(f_\theta) = R_{l, \mathcal{D}}(f_\theta)$  for any  $f_\theta$  and bounded loss function  $l$ .

## 2 ADDITIONAL FIGURES

Figure S1 shows the sum of absolute difference between theoretically estimated  $\hat{\rho}_i^{theory}$  (Eq. (8) in main paper) and empirical  $\hat{\rho}_i$ :  $|\hat{\rho}_1 - \hat{\rho}_1^{theory}| + |\hat{\rho}_0 - \hat{\rho}_0^{theory}|$ . The deviation of the theoretical and empirical estimates reflects the assumption that we have infinite examples, whereas empirically, the number of examples is finite.



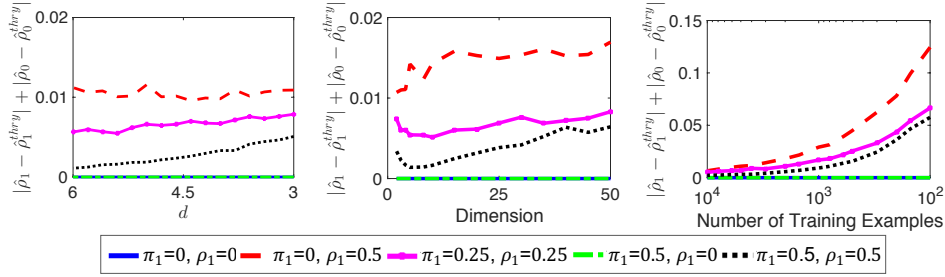


Figure S 1: Sum of absolute difference between theoretically estimated  $\hat{\rho}_i^{thry}$  and empirical  $\hat{\rho}_i$ ,  $i = 0, 1$ , with five different  $(\pi_1, \rho_1)$ , for varying separability  $d$ , dimension, and number of training examples. Note that no figure exists for percent random noise because the theoretical estimates in Eq. (4) do not address added noise examples. The default parameters are:  $d = 4$ , 2 dimensional input, 0% random noise, and 5000 training examples with a fraction of  $p_{y1} = 0.2$  examples as positive. The lines are an average of 200 trials.

Figure S2 shows the Rank Pruning’s noise rate estimation of  $\hat{\pi}_1$  for the MNIST dataset using a logistic regression classifier, for varying amount of  $(\hat{\pi}_1, \hat{\rho}_1)$ , averaging over 10 digits.

Figure S3 shows the average image for each digit for the binary classification problem “1” or “not 1” in MNIST with logistic regression and high noise ( $\rho_1 = 0.5, \pi_1 = 0.5$ ). The number on the bottom and on the right counts the total number of examples (images). From the figure we see that RP makes few mistakes, and when it does, the mistakes vary greatly in image from the typical digit.

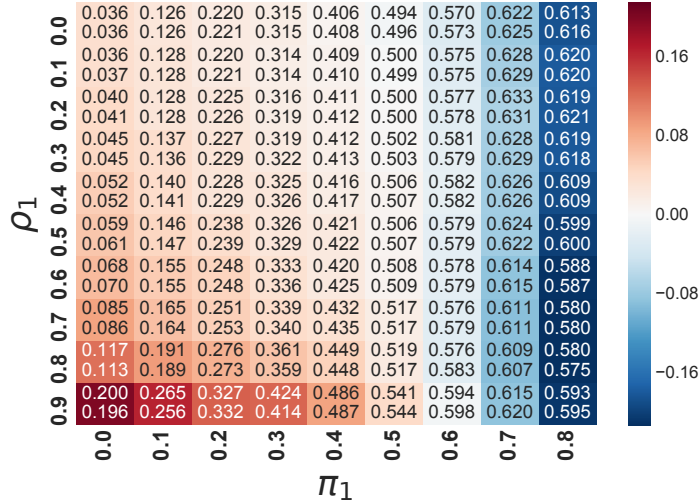


Figure S 2: Rank Pruning  $\hat{\pi}_1$  estimation consistency, averaged over all digits in MNIST. Color depicts  $\hat{\pi}_1 - \pi_1$  with  $\hat{\rho}_1$  (upper) and theoretical  $\hat{\pi}_1^{thry}$  (lower) in each block.

### 3 ADDITIONAL TABLES

Here we provide additional tables for the comparison of error, Precision-Recall AUC (AUC-PR, Davis & Goadrich (2006)), and F1 score for the algorithms *RP*, *Nat13*, *Elk08*, *Liu16* with  $\rho_1, \rho_0$  given to all methods

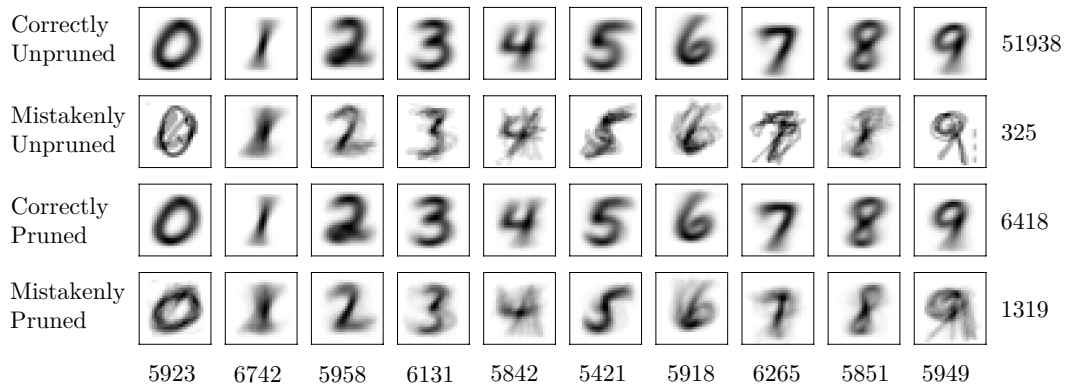


Figure S 3: Average image for each digit for the binary classification problem “1” or “not 1” in MNIST with logistic regression and significant mislabeling ( $\rho_1 = 0.5, \pi_1 = 0.5$ ). The right and bottom numbers count the total number of example images averaged in the corresponding row or column.

for fair comparison. Additionally, we provide the performance of the ground truth classifier (*true*) trained with uncorrupted labels  $(X, y)$ , as well as the complete Rank Pruning algorithm ( $RP_\rho$ ) trained using the noise rates estimated by Rank Pruning. The top model scores are in bold with  $RP_\rho$  in red if its performance is better than non-RP models. The  $\pi_1 = 0$  quadrant in each table represents the “PU learning” case of  $\tilde{P}\tilde{N}$  learning.

Whenever  $g(x) = P(\hat{s} = 1|x)$  is estimated for any algorithm, we use a 3-fold cross-validation to estimate the probability  $g(x)$ . For improved performance, a higher fold may be used.

For the logistic regression classifier, we use scikit-learn’s LogisticRegression class (scikit learn (2016)) with default settings (L2 regularization with inverse strength  $C = 1$ ).

For the convolutional neural networks (CNN), for MNIST we use the structure in Chollet (2016b) and for CIFAR-10, we use the structure in Chollet (2016a). A 10% holdout set is used to monitor the weighted validation loss (using the sample weight given by each algorithm) and ends training when there is no decrease for 10 epochs, with a maximum of 50 epochs for MNIST and 150 epochs for CIFAR-10.

The following list comprises the MNIST and CIFAR-10 experimental result tables for error, AUC-PR and F1 score metrics:

Table S1: Error for MNIST with logistic regression as classifier.

Table S2: AUC-PR for MNIST with logistic regression as classifier.

Table S3: Error for MNIST with CNN as classifier.

Table S4: AUC-PR for MNIST with CNN as classifier.

Table S5: F1 score for CIFAR-10 with logistic regression as classifier.

Table S6: Error for CIFAR-10 with logistic regression as classifier.

Table S7: AUC-PR for CIFAR-10 with logistic regression as classifier.

Table S8: Error for CIFAR-10 with CNN as classifier.

Table S9: AUC-PR for CIFAR-10 with CNN as classifier.

Due to its sensitivity to imperfect probability estimation, here *Liu16* always predicts all labels to be positive or negative, resulting in the same metric score for every digit/image in each scenario. Since  $p_{y1} \simeq 0.1$ , when

predicting all labels as positive, *Liu16* has an F1 score of 0.182, error of 0.90, and AUC-PR of 0.55; when predicting all labels as negative, *Liu16* has an F1 score of 0.0, error of 0.1, and AUC-PR of 0.55.

## 4 ADDITIONAL RELATED WORK

In this section we include tangentially related work which was unable to make it into the final manuscript.

### 4.1 ONE-CLASS CLASSIFICATION

One-class classification (Moya et al., 1993) is distinguished from binary classification by a training set containing examples from only one class, making it useful for outlier and novelty detection (Hempstalk et al., 2008). This can be framed as  $\tilde{P}\tilde{N}$  learning when outliers take the form of mislabeled examples. The predominant approach, one-class SVM, fits a hyper-boundary around the training class (Platt et al., 1999), but often performs poorly due to boundary over-sensitivity (Manevitz & Yousef, 2002) and fails when the training class contains mislabeled examples.

### 4.2 $\tilde{P}\tilde{N}$ LEARNING FOR IMAGE RECOGNITION AND DEEP LEARNING

Variations of  $\tilde{P}\tilde{N}$  learning have been used in the context of machine vision to improve robustness to mislabeling (Xiao et al., 2015). In a face recognition task with 90% of non-faces mislabeled as faces, a bagging model combined with consistency voting was used to remove images with poor voting consistency (Angelova et al., 2005). However, no theoretical justification was provided. In the context of deep learning, consistency of predictions for inputs with mislabeling enforces can be enforced by combining a typical cross-entropy loss with an auto-encoder loss (Reed et al., 2015). This method enforces label consistency by constraining the network to uncover the input examples given the output prediction, but is restricted in architecture and generality.

Table S 1: Comparison of **error** for one-vs-rest MNIST (averaged over all digits) using a **logistic regression** classifier. Except for  $RP_\rho$ ,  $\rho_1$ ,  $\rho_0$  are given to all methods. Top model scores are in bold with  $RP_\rho$  in red if better (smaller) than non-RP models.

MODEL, $\rho_1 =$	$\pi_1 = 0$			$\pi_1 = 0.25$				$\pi_1 = 0.5$				$\pi_1 = 0.75$			
	0.25	0.50	0.75	0.00	0.25	0.50	0.75	0.00	0.25	0.50	0.75	0.00	0.25	0.50	0.75
TRUE	0.020	0.020	0.020	0.020	0.020	0.020	0.020	0.020	0.020	0.020	0.020	0.020	0.020	0.020	0.020
RP $_\rho$	<b>0.023</b>	<b>0.025</b>	<b>0.031</b>	<b>0.024</b>	<b>0.025</b>	<b>0.027</b>	<b>0.038</b>	0.040	0.037	0.039	0.049	0.140	0.128	0.133	0.151
RP	<b>0.022</b>	<b>0.025</b>	<b>0.031</b>	<b>0.021</b>	<b>0.024</b>	<b>0.027</b>	<b>0.035</b>	<b>0.023</b>	<b>0.027</b>	<b>0.031</b>	<b>0.043</b>	<b>0.028</b>	<b>0.036</b>	<b>0.045</b>	0.069
NAT13	0.025	0.030	0.038	0.025	0.029	0.034	0.042	0.030	0.033	0.038	0.047	0.035	0.039	0.046	<b>0.067</b>
ELK08	0.025	0.030	0.038	0.026	0.028	0.032	0.042	0.030	0.031	0.035	0.051	0.092	0.093	0.123	0.189
LIU16	0.187	0.098	0.100	0.100	0.738	0.738	0.419	0.100	0.820	0.821	0.821	0.098	0.760	0.741	0.820

Table S 2: Comparison of **AUC-PR** for one-vs-rest MNIST (averaged over all digits) using a **logistic regression** classifier. Except for  $RP_\rho$ ,  $\rho_1$ ,  $\rho_0$  are given to all methods. Top model scores are in bold with  $RP_\rho$  in red if greater than non-RP models.

MODEL, $\rho_1 =$	$\pi_1 = 0$			$\pi_1 = 0.25$				$\pi_1 = 0.5$				$\pi_1 = 0.75$			
	0.25	0.50	0.75	0.00	0.25	0.50	0.75	0.00	0.25	0.50	0.75	0.00	0.25	0.50	0.75
TRUE	0.935	0.935	0.935	0.935	0.935	0.935	0.935	0.935	0.935	0.935	0.935	0.935	0.935	0.935	0.935
$RP_\rho$	0.921	<b>0.913</b>	<b>0.882</b>	<b>0.928</b>	<b>0.920</b>	<b>0.906</b>	<b>0.853</b>	<b>0.903</b>	<b>0.902</b>	<b>0.879</b>	<b>0.803</b>	0.851	0.835	<b>0.788</b>	0.640
RP'	<b>0.922</b>	<b>0.913</b>	<b>0.882</b>	<b>0.930</b>	<b>0.921</b>	<b>0.906</b>	<b>0.858</b>	<b>0.922</b>	<b>0.903</b>	<b>0.883</b>	<b>0.811</b>	<b>0.893</b>	<b>0.841</b>	<b>0.799</b>	0.621
NAT13	<b>0.922</b>	0.908	0.878	0.918	0.909	0.890	0.839	0.899	0.892	0.862	0.794	0.863	0.837	0.784	<b>0.645</b>
ELK08	0.921	0.903	0.864	0.917	0.908	0.884	0.821	0.898	0.892	0.861	0.763	0.852	0.837	0.772	0.579
LIU16	0.498	0.549	0.550	0.550	0.500	0.550	0.505	0.550	0.550	0.550	0.549	0.503	0.512	0.550	0.550

Table S 3: Comparison of **error** for one-vs-rest MNIST (averaged over all digits) using a **CNN** classifier. Except for  $RP_\rho$ ,  $\rho_1$ ,  $\rho_0$  are given to all methods. Top model scores are in bold with  $RP_\rho$  in red if better (smaller) than non-RP models.

IMAGE TRUE	$\pi_1 = 0$ $\rho_1 = 0.5$				$\pi_1 = 0.25$ $\rho_1 = 0.25$				$\pi_1 = 0.5$ $\rho_1 = 0$				$\pi_1 = 0.5$ $\rho_1 = 0.5$			
	$RP_\rho$	RP	NAT13	ELK08 LIU16	$RP_\rho$	RP	NAT13	ELK08 LIU16	$RP_\rho$	RP	NAT13	ELK08 LIU16	$RP_\rho$	RP	NAT13	ELK08 LIU16
0	0.0013	<b>0.0018</b>	<b>0.0023</b>	0.0045 0.0047 0.9020	<b>0.0017</b>	<b>0.0016</b>	0.0034 0.0036 0.9020	<b>0.0017</b>	<b>0.0016</b>	0.0031 0.0026 0.0029	<b>0.0021</b>	<b>0.0022</b>	0.0116 0.0069 0.9020			
1	0.0015	<b>0.0022</b>	<b>0.0020</b>	0.0025 0.0034 0.8865	<b>0.0019</b>	<b>0.0019</b>	0.0035 0.0030 0.8865	0.0023	0.0020	0.0018	<b>0.0016</b>	0.0023	<b>0.0025</b>	<b>0.0025</b>	0.0036 0.0027 0.8865	
2	0.0027	<b>0.0054</b>	<b>0.0049</b>	0.0057 0.0062 0.8968	<b>0.0032</b>	<b>0.0035</b>	0.0045 0.0051 0.8968	0.0030	0.0029	0.0031	0.0029	<b>0.0024</b>	<b>0.0059</b>	<b>0.0050</b>	0.0066 0.0083 0.8968	
3	0.0020	<b>0.0032</b>	<b>0.0032</b>	0.0055 0.0038 0.8990	<b>0.0029</b>	<b>0.0029</b>	0.0043 0.0043 0.8990	<b>0.0021</b>	0.0027	<b>0.0023</b>	<b>0.0023</b>	0.0032	<b>0.0038</b>	<b>0.0042</b>	0.0084 0.0057 0.8990	
4	0.0012	<b>0.0037</b>	0.0040	<b>0.0038</b>	0.0044 0.9018	<b>0.0029</b>	<b>0.0025</b>	0.0055 0.0069 0.9018	0.0026	0.0020	<b>0.0019</b>	0.0021	0.0030	<b>0.0044</b>	<b>0.0035</b>	0.0086 0.0077 0.9018
5	0.0019	<b>0.0032</b>	<b>0.0035</b>	0.0039 0.0038 0.9108	<b>0.0027</b>	<b>0.0031</b>	0.0062 0.0060 0.9108	<b>0.0021</b>	0.0024	0.0024	0.0028	<b>0.0023</b>	<b>0.0061</b>	<b>0.0056</b>	0.0066 0.0074 0.9108	
6	0.0021	<b>0.0027</b>	<b>0.0028</b>	0.0053 0.0035 0.9042	<b>0.0028</b>	<b>0.0025</b>	0.0042 0.0036 0.9042	0.0029	0.0029	<b>0.0022</b>	0.0024	0.0028	<b>0.0032</b>	<b>0.0035</b>	0.0098 0.0075 0.9042	
7	0.0026	<b>0.0039</b>	<b>0.0041</b>	0.0066 0.0103 0.8972	<b>0.0050</b>	<b>0.0052</b>	0.0058 0.0058 0.8972	0.0049	0.0040	<b>0.0030</b>	0.0037	0.0035	<b>0.0054</b>	<b>0.0064</b>	0.0113 0.0085 0.8972	
8	0.0022	<b>0.0047</b>	<b>0.0043</b>	0.0106 0.0063 0.9026	<b>0.0034</b>	<b>0.0036</b>	0.0062 0.0091 0.9026	0.0036	<b>0.0030</b>	0.0035	0.0041	0.0032	<b>0.0044</b>	<b>0.0048</b>	0.0234 0.0077 0.9026	
9	0.0036	0.0067	<b>0.0052</b>	0.0056 0.0124 0.8991	<b>0.0048</b>	<b>0.0051</b>	0.0065 0.0064 0.8991	0.0048	0.0050	0.0051	<b>0.0043</b>	0.0059	<b>0.0081</b>	0.0114 0.0131	<b>0.0112</b>	0.8991
AVG	0.0021	<b>0.0038</b>	<b>0.0036</b>	0.0054 0.0059 0.9000	<b>0.0031</b>	<b>0.0032</b>	0.0050 0.0054 0.9000	0.0030	<b>0.0028</b>	<b>0.0028</b>	0.0029 0.0032	<b>0.0046</b>	<b>0.0049</b>	0.0103 0.0074	0.9000	

Table S 4: Comparison of **AUC-PR** for one-vs-rest MNIST (averaged over all digits) using a **CNN** classifier. Except for  $RP_\rho$ ,  $\rho_1$ ,  $\rho_0$  are given to all methods. Top model scores are in bold with  $RP_\rho$  in red if greater than non-RP models.

IMAGE TRUE	$\pi_1 = 0$ $\rho_1 = 0.5$				$\pi_1 = 0.25$ $\rho_1 = 0.25$				$\pi_1 = 0.5$ $\rho_1 = 0$				$\pi_1 = 0.5$ $\rho_1 = 0.5$			
	$RP_\rho$	RP	NAT13	ELK08 LIU16	$RP_\rho$	RP	NAT13	ELK08 LIU16	$RP_\rho$	RP	NAT13	ELK08 LIU16	$RP_\rho$	RP	NAT13	ELK08 LIU16
0	0.9998	<b>0.9992</b>	<b>0.9990</b>	0.9986 0.9982 0.5490	<b>0.9996</b>	<b>0.9996</b>	0.9986 0.9979 0.5490	<b>0.9989</b>	<b>0.9995</b>	0.9976 0.9979 0.9956	<b>0.9984</b>	<b>0.9982</b>	0.9963 0.9928 0.5490			
1	0.9999	<b>0.9995</b>	<b>0.9995</b>	0.9976 0.9974 0.5568	<b>0.9996</b>	0.9993	<b>0.9995</b>	<b>0.9995</b>	0.5568	<b>0.9995</b>	<b>0.9998</b>	0.9982 0.9972 0.9965	<b>0.9995</b>	<b>0.9994</b>	0.9978 0.9985 0.5568	
2	0.9994	<b>0.9971</b>	<b>0.9969</b>	0.9917 0.9942 0.5516	<b>0.9980</b>	<b>0.9977</b>	0.9971 0.9945 0.5516	<b>0.9988</b>	<b>0.9992</b>	0.9958 0.9934 0.9940	<b>0.9938</b>	<b>0.9947</b>	0.9847 0.9873 0.5516			
3	0.9996	<b>0.9986</b>	<b>0.9987</b>	0.9983 0.9984 0.5505	<b>0.9991</b>	<b>0.9989</b>	0.9982 0.9980 0.5505	<b>0.9993</b>	<b>0.9994</b>	0.9991 0.9971 0.9974	<b>0.9969</b>	<b>0.9959</b>	0.9951 <b>0.9959</b>	0.5505		
4	0.9997	0.9982	<b>0.9989</b>	0.9939 0.9988 0.0891	<b>0.9992</b>	<b>0.9991</b>	0.9976 0.9965 0.5491	<b>0.9994</b>	<b>0.9996</b>	0.9985 0.9978 0.9986	<b>0.9983</b>	<b>0.9977</b>	0.9961 0.9919 0.5491			
5	0.9993	<b>0.9982</b>	<b>0.9976</b>	0.9969 0.9956 0.5446	<b>0.9986</b>	<b>0.9987</b>	0.9983 0.9979 0.5446	<b>0.9984</b>	<b>0.9982</b>	0.9971 0.9963 0.9929	<b>0.9958</b>	<b>0.9965</b>	0.9946 0.9934 0.5446			
6	0.9987	<b>0.9976</b>	<b>0.9970</b>	0.9928 0.9931 0.5479	<b>0.9974</b>	<b>0.9980</b>	0.9956 0.9959 0.5479	<b>0.9968</b>	<b>0.9983</b>	0.9933 0.9950 0.9905	<b>0.9964</b>	0.9957	0.9942 <b>0.9961</b>	0.5479		
7	0.9989	<b>0.9973</b>	<b>0.9972</b>	0.9965 0.9944 0.0721	0.9968	0.9973	0.9966 <b>0.9979</b>	0.5514	0.9969	<b>0.9983</b>	0.9961 0.9958 0.9974	<b>0.9933</b>	<b>0.9937</b>	0.9896 0.9886 0.5514		
8	0.9996	<b>0.9974</b>	<b>0.9964</b>	<b>0.9964</b>	0.9946 0.5487	<b>0.9981</b>	<b>0.9981</b>	0.9973 0.9971 0.5487	<b>0.9983</b>	0.9988	0.9984 0.9976 <b>0.9989</b>	<b>0.9976</b>	<b>0.9975</b>	0.9873 0.9893 0.5487		
9	0.9979	<b>0.9931</b>	<b>0.9951</b>	0.9901 0.9922 0.5504	<b>0.9935</b>	<b>0.9951</b>	0.9933 0.9920 0.5504	<b>0.9961</b>	<b>0.9951</b>	0.9924 0.9922 0.9912	<b>0.9877</b>	<b>0.9876</b>	0.9819 0.9828 0.5504			
AVG	0.9993	<b>0.9976</b>	<b>0.9976</b>	0.9953 0.9957 0.4561	<b>0.9980</b>	<b>0.9982</b>	0.9972 0.9967 0.5500	<b>0.9983</b>	<b>0.9986</b>	0.9966 0.9960 0.9953	<b>0.9958</b>	<b>0.9957</b>	0.9918 0.9917 0.5500			



Table S 8: Comparison of **error** for one-vs-rest CIFAR-10 (averaged over all images) using a CNN classifier. Except for  $RP_\rho$ ,  $\rho_1$ ,  $\rho_0$  are given to all methods. Top model scores are in bold with  $RP_\rho$  in red if better (smaller) than non-RP models.

IMAGE TRUE	$\pi_1 = 0$ $\rho_1 = 0.5$					$\pi_1 = 0.25$ $\rho_1 = 0.25$					$\pi_1 = 0.5$ $\rho_1 = 0$					$\pi_1 = 0.5$ $\rho_1 = 0.5$					
	$RP_\rho$	RP	NAT13	ELK08	LIU16	$RP_\rho$	RP	NAT13	ELK08	LIU16	$RP_\rho$	RP	NAT13	ELK08	LIU16	$RP_\rho$	RP	NAT13	ELK08	LIU16	
PLANE	0.044	<b>0.054</b>	<b>0.057</b>	0.059	0.063	0.900	<b>0.050</b>	<b>0.051</b>	0.054	0.057	0.900	<b>0.048</b>	<b>0.045</b>	0.049	0.048	0.100	<b>0.063</b>	<b>0.061</b>	0.074	0.065	0.900
AUTO	0.021	<b>0.040</b>	<b>0.037</b>	0.041	0.043	0.100	<b>0.032</b>	<b>0.034</b>	0.040	0.039	0.900	0.028	<b>0.026</b>	<b>0.026</b>	<b>0.026</b>	0.100	<b>0.047</b>	<b>0.049</b>	0.062	0.070	0.900
BIRD	0.055	0.083	<b>0.078</b>	0.080	0.082	0.900	<b>0.074</b>	<b>0.074</b>	0.077	0.078	0.900	0.072	<b>0.066</b>	0.072	0.070	0.100	0.124	<b>0.084</b>	0.089	0.093	0.900
CAT	0.077	0.108	<b>0.091</b>	0.092	0.095	0.100	0.111	0.090	<b>0.086</b>	0.089	0.900	0.113	<b>0.084</b>	0.086	0.088	0.100	0.117	0.098	<b>0.094</b>	0.100	0.900
DEER	0.049	0.081	<b>0.078</b>	<b>0.078</b>	0.079	0.900	0.080	<b>0.069</b>	0.075	0.070	0.900	0.076	0.062	<b>0.061</b>	0.062	0.100	0.106	<b>0.086</b>	0.091	0.093	0.900
DOG	0.062	<b>0.075</b>	<b>0.071</b>	0.079	0.080	0.100	0.071	0.069	0.070	<b>0.067</b>	0.900	0.069	0.061	<b>0.057</b>	0.076	0.100	0.103	<b>0.081</b>	0.084	0.086	0.900
FROG	0.038	0.050	<b>0.048</b>	<b>0.048</b>	0.054	0.100	<b>0.047</b>	<b>0.052</b>	0.056	0.062	0.900	0.045	<b>0.040</b>	0.042	0.043	0.100	<b>0.058</b>	<b>0.062</b>	0.066	0.071	0.900
HORSE	0.035	<b>0.050</b>	<b>0.052</b>	0.057	0.054	0.900	<b>0.048</b>	<b>0.051</b>	0.052	0.057	0.900	0.045	<b>0.040</b>	0.042	0.046	0.100	<b>0.065</b>	<b>0.063</b>	0.066	0.075	0.900
SHIP	0.028	<b>0.042</b>	<b>0.042</b>	0.046	<b>0.042</b>	0.900	<b>0.037</b>	<b>0.036</b>	0.042	0.047	0.900	0.035	0.033	<b>0.031</b>	0.033	0.100	<b>0.051</b>	<b>0.049</b>	0.064	0.058	0.900
TRUCK	0.027	<b>0.044</b>	<b>0.046</b>	0.054	0.056	0.900	<b>0.034</b>	<b>0.032</b>	0.038	0.043	0.900	<b>0.034</b>	<b>0.031</b>	0.034	0.034	0.100	<b>0.060</b>	0.066	0.067	<b>0.065</b>	0.900
AVG	0.043	<b>0.063</b>	<b>0.060</b>	0.064	0.065	0.580	<b>0.059</b>	<b>0.056</b>	0.059	0.061	0.900	0.056	<b>0.049</b>	0.050	0.053	0.100	0.080	<b>0.070</b>	0.076	0.077	0.900

Table S 9: Comparison of **AUC-PR** for one-vs-rest CIFAR-10 (averaged over all images) using a CNN classifier. Except for  $RP_\rho$ ,  $\rho_1$ ,  $\rho_0$  are given to all methods. Top model scores are in bold with  $RP_\rho$  in red if greater than non-RP models.

IMAGE TRUE	$\pi_1 = 0$ $\rho_1 = 0.5$					$\pi_1 = 0.25$ $\rho_1 = 0.25$					$\pi_1 = 0.5$ $\rho_1 = 0$					$\pi_1 = 0.5$ $\rho_1 = 0.5$					
	$RP_\rho$	RP	NAT13	ELK08	LIU16	$RP_\rho$	RP	NAT13	ELK08	LIU16	$RP_\rho$	RP	NAT13	ELK08	LIU16	$RP_\rho$	RP	NAT13	ELK08	LIU16	
PLANE	0.856	0.779	0.780	<b>0.784</b>	0.756	0.550	<b>0.808</b>	<b>0.797</b>	0.770	0.742	0.550	<b>0.813</b>	<b>0.824</b>	0.792	0.794	0.550	<b>0.710</b>	<b>0.722</b>	0.662	0.682	0.550
AUTO	0.954	0.874	<b>0.889</b>	0.878	0.833	0.550	<b>0.905</b>	<b>0.900</b>	0.871	0.866	0.550	<b>0.931</b>	<b>0.927</b>	0.924	0.910	0.550	<b>0.824</b>	<b>0.814</b>	0.756	0.702	0.550
BIRD	0.761	0.559	0.566	<b>0.569</b>	0.568	0.550	<b>0.619</b>	<b>0.618</b>	0.584	0.597	0.550	<b>0.623</b>	<b>0.679</b>	0.613	0.619	0.115	0.465	0.492	0.436	0.434	<b>0.550</b>
CAT	0.601	0.387	0.447	<b>0.463</b>	0.433	0.550	0.423	0.454	0.487	0.480	<b>0.550</b>	0.483	<b>0.512</b>	0.493	0.473	0.050	0.373	0.375	0.382	0.371	<b>0.550</b>
DEER	0.820	<b>0.620</b>	0.600	<b>0.615</b>	0.573	0.550	0.646	<b>0.660</b>	0.610	0.657	0.550	0.658	<b>0.707</b>	0.700	0.703	0.550	0.434	0.487	0.414	0.435	<b>0.550</b>
DOG	0.758	<b>0.629</b>	<b>0.662</b>	0.617	0.573	0.550	<b>0.673</b>	<b>0.667</b>	0.658	0.660	0.550	0.705	0.722	<b>0.741</b>	0.705	0.550	0.541	0.545	0.496	0.519	<b>0.550</b>
FROG	0.891	0.812	<b>0.815</b>	0.812	0.776	0.550	<b>0.821</b>	<b>0.827</b>	0.808	0.749	0.550	<b>0.841</b>	<b>0.851</b>	0.828	0.831	0.550	<b>0.753</b>	<b>0.710</b>	0.691	0.620	0.550
HORSE	0.897	<b>0.810</b>	<b>0.817</b>	0.799	0.779	0.550	<b>0.824</b>	<b>0.809</b>	0.801	0.772	0.550	<b>0.826</b>	<b>0.844</b>	0.818	0.819	0.550	<b>0.736</b>	<b>0.699</b>	<b>0.699</b>	0.600	0.550
SHIP	0.922	<b>0.870</b>	0.862	<b>0.864</b>	0.853	0.550	<b>0.889</b>	<b>0.885</b>	0.843	0.848	0.550	0.889	<b>0.897</b>	0.891	0.887	0.550	<b>0.800</b>	<b>0.808</b>	0.767	0.741	0.550
TRUCK	0.929	<b>0.845</b>	<b>0.848</b>	0.824	0.787	0.550	<b>0.887</b>	<b>0.894</b>	0.873	0.853	0.550	<b>0.904</b>	<b>0.902</b>	0.898	0.883	0.550	<b>0.740</b>	<b>0.709</b>	0.695	0.690	0.550
AVG	0.839	0.719	<b>0.729</b>	0.722	0.693	0.550	<b>0.750</b>	<b>0.751</b>	0.730	0.722	0.550	0.767	<b>0.787</b>	0.770	0.762	0.457	<b>0.637</b>	<b>0.636</b>	0.600	0.579	0.550

## References

- Angelova, Anelia, Abu-Mostafam, Yaser, and Perona, Pietro. Pruning training sets for learning of object categories. In *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, volume 1, pp. 494–501. IEEE, 2005.
- Angluin, Dana and Laird, Philip. Learning from noisy examples. *Machine Learning*, 2(4):343–370, 1988.
- Chollet, Francois. *Keras CIFAR CNN*, 2016a. URL <http://bit.ly/2mVKR3d>.
- Chollet, Francois. *Keras MNIST CNN*, 2016b. URL <http://bit.ly/2nKiqJv>.
- Davis, Jesse and Goadrich, Mark. The relationship between precision-recall and roc curves. In *Proceedings of the 23rd International Conference on Machine Learning, ICML '06*, pp. 233–240, New York, NY, USA, 2006. ACM. ISBN 1-59593-383-2. doi: 10.1145/1143844.1143874. URL <http://doi.acm.org/10.1145/1143844.1143874>.
- Hempstalk, Kathryn, Frank, Eibe, and Witten, Ian H. One-class classification by combining density and class probability estimation. In *Proceedings of the 2008 European Conference on Machine Learning and Knowledge Discovery in Databases - Part I, ECML PKDD '08*, pp. 505–519, Berlin, Heidelberg, 2008. Springer-Verlag. ISBN 978-3-540-87478-2. doi: 10.1007/978-3-540-87479-9\_51. URL [http://dx.doi.org/10.1007/978-3-540-87479-9\\_51](http://dx.doi.org/10.1007/978-3-540-87479-9_51).
- Manevitz, Larry M. and Yousef, Malik. One-class svms for document classification. *Journal of Machine Learning Research*, 2:139–154, March 2002. ISSN 1532-4435. URL <http://dl.acm.org/citation.cfm?id=944790.944808>.
- Moya, M. M., Koch, M. W., and Hostetler, L. D. One-class classifier networks for target recognition applications. *NASA STI/Recon Technical Report N*, 93, 1993.
- Platt, John, Schlkopf, Bernhard, Shawe-Taylor, John, Smola, Alex J., and Williamson, Robert C. Estimating the support of a high-dimensional distribution. Technical report, Microsoft Research, November 1999. URL <https://www.microsoft.com/en-us/research/publication/estimating-the-support-of-a-high-dimensional-distribution/>.
- Reed, Scott E., Lee, Honglak, Anguelov, Dragomir, Szegedy, Christian, Erhan, Dumitru, and Rabinovich, Andrew. Training deep neural networks on noisy labels with bootstrapping. In *ICLR 2015*, 2015. URL <http://arxiv.org/abs/1412.6596>.
- scikit learn. *LogisticRegression Class at scikit-learn*, 2016. URL <http://bit.ly/2o3y6r5>.
- Xiao, Tong, Xia, Tian, Yang, Yi, Huang, Chang, and Wang, Xiaogang. Learning from massive noisy labeled data for image classification. In *CVPR*, 2015.