

# Learning with Confident Examples: Rank Pruning for Robust Classification with Noisy Labels

Curtis G. Northcutt, \* Tailin Wu, \* Isaac L. Chuang

Massachusetts Institute of Technology  
Cambridge, MA 02139  
{cgn, tailin, ichuang}@mit.edu

## Abstract

$\tilde{P}\tilde{N}$  learning is the problem of binary classification when training examples may be mislabeled (flipped) uniformly with noise rate  $\rho_1$  for positive examples and  $\rho_0$  for negative examples. We propose Rank Pruning (RP) to solve  $\tilde{P}\tilde{N}$  learning and the open problem of estimating the noise rates, i.e. the fraction of wrong positive and negative labels. Unlike prior solutions, RP is time-efficient and general, requiring  $\mathcal{O}(T)$  for any unrestricted choice of probabilistic classifier with  $T$  fitting time. We prove RP has consistent noise estimation and equivalent expected risk as learning with uncorrupted labels in ideal conditions, and derive closed-form solutions when conditions are non-ideal. RP achieves state-of-the-art noise estimation and F1, error, and AUC-PR for both MNIST and CIFAR datasets, regardless of the amount of noise and performs similarly impressively when a large portion of training examples are noise drawn from a third distribution. To highlight, RP with a CNN classifier can predict if an MNIST digit is a *one* or *not* with only 0.25% error, and 0.46% error across all digits, even when 50% of positive examples are mislabeled and 50% of observed positive labels are mislabeled negative examples.

## 1 Introduction

Consider a student with no knowledge of animals tasked with learning to classify whether a picture contains a dog. A teacher shows the student example pictures of lone four-legged animals, stating whether the image contains a dog or not. Unfortunately, the teacher may often make mistakes, asymmetrically, with a significantly large false positive rate,  $\rho_1 \in [0, 1]$ , and significantly large false negative rate,  $\rho_0 \in [0, 1]$ . The teacher may also include “white noise” images with a uniformly random label. This information is unknown to the student, who only knows of the images and corrupted labels, but suspects that the teacher may make mistakes. Can the student (1) estimate the mistake rates,  $\rho_1$  and  $\rho_0$ , (2) learn to classify pictures with dogs accurately, and (3) do so efficiently (e.g. less than an hour for 50 images)? This allegory clarifies the challenges of  $\tilde{P}\tilde{N}$  learning for any classifier trained with corrupted labels, perhaps with intermixed noise examples. We elect the notation  $\tilde{P}\tilde{N}$  to emphasize that both the positive and negative sets may contain mislabeled examples, reserving  $P$  and  $N$  for uncorrupted sets.

This example illustrates a fundamental reliance of supervised learning on training labels (Michalski et al., 1986). Traditional learning performance degrades monotonically with label noise (Aha et al., 1991; Nettleton et al., 2010), necessitating semi-supervised approaches (Blanchard et al., 2010). Examples of noisy datasets are medical (Raviv & Intrator, 1996), human-labeled (Paolacci et al., 2010), and sensor (Lane et al., 2010) datasets. The problem of uncovering the same classifications as if the data was not mislabeled is our fundamental goal.

Towards this goal, we introduce Rank Pruning<sup>2</sup>, an algorithm for  $\tilde{P}\tilde{N}$  learning composed of two sequential parts: (1) estimation of the asymmetric noise rates  $\rho_1$  and  $\rho_0$  and (2) removal of mislabeled examples prior to training. The

\* Equal Contribution

<sup>2</sup> Rank Pruning is open-source and available at <https://github.com/cgnorthcutt/rankpruning>

**Table 1:** Variable definitions and descriptions for  $\tilde{P}\tilde{N}$  learning and PU learning. Related work contains a prominent author using each variable.  $\rho_1$  is also referred to as *contamination* in PU learning literature.

VARIABLE	CONDITIONAL	DESCRIPTION	DOMAIN	RELATED WORK
$\rho_0$	$P(s = 1 y = 0)$	FRACTION OF $N$ EXAMPLES MISLABELED AS POSITIVE	$\tilde{P}\tilde{N}$	LIU
$\rho_1$	$P(s = 0 y = 1)$	FRACTION OF $P$ EXAMPLES MISLABELED AS NEGATIVE	$\tilde{P}\tilde{N}$ , PU	LIU, CLAESEN
$\pi_0$	$P(y = 1 s = 0)$	FRACTION OF MISLABELED EXAMPLES IN $\tilde{N}$	$\tilde{P}\tilde{N}$	SCOTT
$\pi_1$	$P(y = 0 s = 1)$	FRACTION OF MISLABELED EXAMPLES IN $\tilde{P}$	$\tilde{P}\tilde{N}$	SCOTT
$c = 1 - \rho_1$	$P(s = 1 y = 1)$	FRACTION OF CORRECTLY LABELED $P$ IF $P(y = 1 s = 1) = 1$	PU	ELKAN

fundamental mantra of Rank Pruning is *learning with confident examples*, i.e. examples with a predicted probability of being positive *near* 1 when the label is positive or 0 when the label is negative. If we imagine non-confident examples as a noise class, separate from the confident positive and negative classes, then their removal should unveil a subset of the uncorrupted data.

An ancillary mantra of Rank Pruning is *removal by rank* which elegantly exploits ranking without sorting. Instead of pruning non-confident examples by predicted probability, we estimate the number of mislabeled examples in each class. We then remove the  $k^{th}$ -most or  $k^{th}$ -least examples, *ranked* by predicted probability, via the BFPRT algorithm (Blum et al., 1973) in  $\mathcal{O}(n)$  time, where  $n$  is the number of training examples. *Removal by rank* mitigates sensitivity to probability estimation and exploits the reduced complexity of learning to rank over probability estimation (Menon et al., 2012). Together, *learning with confident examples* and *removal by rank* enable robustness, i.e. invariance to erroneous input deviation.

Beyond prediction, confident examples help estimate  $\rho_1$  and  $\rho_0$ . Typical approaches require averaging predicted probabilities on a holdout set (Liu & Tao, 2016; Elkan & Noto, 2008) tying noise estimation to the accuracy of the predicted probabilities, which in practice may be confounded by added noise or poor model selection. Instead, we estimate  $\rho_1$  and  $\rho_0$  as a fraction of the predicted counts of confident examples in each class, encouraging robustness for variation in probability estimation. △ .

## 1.1 Related Work

Rank Pruning bridges framework, nomenclature, and application across  $PU$  and  $\tilde{P}\tilde{N}$  learning. In this section, we consider the contributions of Rank Pruning in both.

### 1.1.1 $PU$ Learning

Positive-unlabeled ( $PU$ ) learning is a binary classification task in which a subset of positive training examples are labeled, and the rest are unlabeled. For example, co-training (Blum & Mitchell, 1998; Nigam & Ghani, 2000) with labeled and unlabeled examples can be framed as a  $PU$  learning problem by assigning all unlabeled examples the label ‘0’.  $PU$  learning methods often assume corrupted negative labels for the unlabeled examples  $U$  such that  $PU$  learning is  $\tilde{P}\tilde{N}$  learning with no mislabeled examples in  $P$ , hence their naming conventions.

Early approaches to  $PU$  learning modified the loss functions via weighted logistic regression (Lee & Liu, 2003) and biased SVM (Liu et al., 2003) to penalize more when positive examples are predicted incorrectly. Bagging SVM (Mordet & Vert, 2014) and RESVM (Claesen et al., 2015) extended biased SVM to instead use an ensemble of classifiers trained by resampling  $U$  (and  $P$  for RESVM) to improve robustness (Breiman, 1996). RESVM claims state-of-the-art for  $PU$  learning, but is impractically inefficient for large datasets because it requires optimization of five parameters and suffers from the pitfalls of SVM model selection (Chapelle & Vapnik, 1999). Elkan & Noto (2008) introduce a formative time-efficient probabilistic approach (denoted *Elk08*) for  $PU$  learning that directly estimates  $1 - \rho_1$  by averaging predicted probabilities of a holdout set and dividing all predicted probabilities by  $1 - \rho_1$ . On the SwissProt database, *Elk08* was 621 times faster than biased SVM, which only requires two parameter optimization. However, *Elk08* noise rate estimation is sensitive to inexact probability estimation and both RESVM and *Elk08* assume  $P = \tilde{P}$  and do not generalize to  $\tilde{P}\tilde{N}$  learning. Rank Pruning leverages *Elk08* to initialize  $\rho_1$ , but then re-estimates  $\rho_1$  using confident examples for both robustness (RESVM) and efficiency (*Elk08*).

### 1.1.2 $\tilde{P}\tilde{N}$ Learning

Theoretical approaches for  $\tilde{P}\tilde{N}$  learning often have two steps: (1) estimate the noise rates,  $\rho_1$ ,  $\rho_0$ , and (2) use  $\rho_1$ ,  $\rho_0$  for prediction. To our knowledge, Rank Pruning is the only time-efficient solution for the open problem (Liu & Tao, 2016; Yang et al., 2012) of noise estimation.

We first consider relevant work in noise rate estimation. Scott et al. (2013) established a lower bound method for estimating the *inversed* noise rates  $\pi_1$  and  $\pi_0$  (defined in Table 1). However, the method can be intractable due to unbounded convergence and assumes that the positive and negative distributions are mutually irreducible. Under additional assumptions, Scott (2015) proposed a time-efficient method for noise rate estimation, but reported poor performance Liu & Tao (2016). Liu & Tao (2016) used the minimum predicted probabilities as the noise rates, which often yields futile estimates of  $\min = 0$ . Natarajan et al. (2013) provide no method for estimation and view the noise rates as parameters optimized with cross-validation, inducing a sacrificial accuracy, efficiency trade-off. In comparison, Rank Pruning noise rate estimation is time-efficient, consistent in ideal conditions, and robust to imperfect probability estimation.

Natarajan et al. (2013) developed two methods for prediction in the  $\tilde{P}\tilde{N}$  setting which modify the loss function. The first method constructs an unbiased estimator of the loss function for the true distribution from the noisy distribution, but the estimator may be non-convex even if the original loss function is convex. If the classifier’s loss function cannot be modified directly, this method requires splitting each example in two with class-conditional weights and ensuring split examples are in the same batch during optimization. For these reasons, we instead compare Rank Pruning with their second method (*Nat13*), which constructs a label-dependent loss function such that for 0-1 loss, the minimizers of *Nat13*’s risk and the risk for the true distribution are equivalent.

Liu & Tao (2016) generalized *Elk08* to the  $\tilde{P}\tilde{N}$  learning setting by modifying the loss function with per-example importance reweighting (*Liu16*), but reweighting terms are derived from predicted probabilities which may be sensitive to inexact estimation. To mitigate sensitivity, Liu & Tao (2016) examine the use of density ratio estimation (Sugiyama et al., 2012). Instead, Rank Pruning mitigates sensitivity by learning from confident examples selected by rank order, not predicted probability. For fairness of comparison across methods, we compare Rank Pruning with their probability-based approach.

Assuming perfect estimation of  $\rho_1$  and  $\rho_0$ , we, Natarajan et al. (2013), and Liu & Tao (2016) all prove that the expected risk for the modified loss function is equivalent to the expected risk for the perfectly labeled dataset. However, both Natarajan et al. (2013) and Liu & Tao (2016) effectively “flip” example labels in the construction of their loss function, providing no benefit for added random noise. In comparison, Rank Pruning will also remove added random noise because noise drawn from a third distribution is unlikely to appear confidently positive or negative. Table 2 summarizes our comparison of  $\tilde{P}\tilde{N}$  and *PU* learning methods.

Procedural efforts have improved robustness to mislabeling in the context of machine vision (Xiao et al., 2015), neural networks (Reed et al., 2015), and face recognition (Angelova et al., 2005). Though promising, these methods are restricted in theoretical justification and generality, motivating the need for Rank Pruning.

**Table 2:** Summary of state-of-the-art and selected general solutions to  $\tilde{P}\tilde{N}$  and *PU* learning.

RELATED WORK	NOISE ESTIM.	$\tilde{P}\tilde{N}$	<i>PU</i>	ANY PROB. CLASSIFIER	PROB ESTIM. ROBUSTNESS	TIME EFFICIENT	THEORY SUPPORT	ADDED NOISE
ELKAN & NOTO (2008)	✓		✓	✓		✓	✓	
CLAESEN ET AL. (2015)			✓		✓			
SCOTT ET AL. (2013)	✓			✓	✓		✓	
NATARAJAN ET AL. (2013)		✓	✓	✓	✓	✓	✓	
LIU & TAO (2016)		✓	✓	✓		✓	✓	
<b>RANK PRUNING</b>	✓	✓	✓	✓	✓	✓	✓	✓

## 1.2 Contributions

In this paper, we describe the Rank Pruning algorithm for binary classification with imperfectly labeled training data. In particular, we:

- Develop a robust, time-efficient, general solution for both  $\tilde{P}\tilde{N}$  learning, i.e. binary classification with noisy labels, and estimation of the fraction of mislabeling in both the positive and negative training sets.
- Introduce the *learning with confident examples* mantra as a new way to think about robust classification and estimation with mislabeled training data.
- Prove that under assumptions, Rank Pruning achieves perfect noise estimation and equivalent expected risk as learning with correct labels. We provide closed-form solutions when those assumptions are relaxed.
- Demonstrate that Rank Pruning performance generalizes across the number of training examples, feature dimension, fraction of mislabeling, and fraction of added noise examples drawn from a third distribution.
- Improve the state-of-the-art of  $\tilde{P}\tilde{N}$  learning across F1 score, AUC-PR, and Error. In many cases, Rank Pruning achieves nearly the same F1 score as learning with correct labels when 50% of positive examples are mislabeled and 50% of observed positive labels are mislabeled negative examples.

## 2 Framing the $\tilde{P}\tilde{N}$ Learning Problem

In this section, we formalize the foundational definitions, assumptions, and goals of the  $\tilde{P}\tilde{N}$  learning problem illustrated by the student-teacher motivational example.

Given  $n$  observed training examples  $x \in \mathcal{R}^D$  with associated observed corrupted labels  $s \in \{0, 1\}$  and unobserved true labels  $y \in \{0, 1\}$ , we seek a binary classifier  $f$  that estimates the mapping  $x \rightarrow y$ . Unfortunately, if we fit the classifier using observed  $(x, s)$  pairs, we estimate the mapping  $x \rightarrow s$  and obtain  $g(x) = P(\hat{s} = 1|x)$ .

We define the observed noisy positive and negative sets as  $\tilde{P} = \{x|s = 1\}$ ,  $\tilde{N} = \{x|s = 0\}$  and the unobserved true positive and negative sets as  $P = \{x|y = 1\}$ ,  $N = \{x|y = 0\}$ . Define the hidden training data as  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ , drawn i.i.d. from some true distribution  $\mathcal{D}$ . We assume that a class-conditional Classification Noise Process (CNP) (Angluin & Laird, 1988) maps  $y$  true labels to  $s$  observed labels such that each label in  $P$  is flipped independently with probability  $\rho_1$  and each label in  $N$  is flipped independently with probability  $\rho_0$  ( $s \leftarrow \text{CNP}(y, \rho_1, \rho_0)$ ). The resulting observed, corrupted dataset is  $D_\rho = \{(x_1, s_1), (x_2, s_2), \dots, (x_n, s_n)\}$ . Therefore,  $(s \perp x)|y$  and  $P(s = s|y = y, x) = P(s = s|y = y)$ . In recent work, CNP is referred to as the random noise classification (RCN) noise model (Liu & Tao, 2016; Natarajan et al., 2013).

The noise rate  $\rho_1 = P(s = 0|y = 1)$  is the fraction of  $P$  examples mislabeled as negative and the noise rate  $\rho_0 = P(s = 1|y = 0)$  is the fraction of  $N$  examples mislabeled as positive. Note that  $\rho_1 + \rho_0 < 1$  is a necessary condition, otherwise more examples would be mislabeled than labeled correctly. Thus,  $\rho_0 < 1 - \rho_1$ . We elect a subscript of “0” to refer to the negative set and a subscript of “1” to refer to the positive set. Additionally, let  $p_{s1} = P(s = 1)$  be the fraction of corrupted labels that are positive and  $p_{y1} = P(y = 1)$  be the fraction of true labels that are positive. It follows that the inversed noise rates are  $\pi_1 = P(y = 0|s = 1) = \frac{\rho_0(1-p_{y1})}{p_{s1}}$  and  $\pi_0 = P(y = 1|s = 0) = \frac{\rho_1 p_{y1}}{(1-p_{s1})}$ . Combining these relations, given any pair in  $\{(\rho_0, \rho_1), (\rho_1, \pi_1), (\rho_0, \pi_0), (\pi_0, \pi_1)\}$ , the remaining two and  $p_{y1}$  are known.

We consider five levels of assumptions for  $P$ ,  $N$ , and  $g$ :

**Perfect Condition:**  $g$  is a “perfect” probability estimator iff  $g(x) = g^*(x)$  where  $g^*(x) = P(s = 1|x)$ . Equivalently, let  $g(x) = P(s = 1|x) + \Delta g(x)$ . Then  $g(x)$  is “perfect” when  $\Delta g(x) = 0$  and “imperfect” when  $\Delta g(x) \neq 0$ .  $g$  may be imperfect due to the method of estimation or due to added uniformly randomly labeled examples drawn from a third noise distribution.

**Non-overlapping Condition:**  $P$  and  $N$  have “non-overlapping support” if  $P(y = 1|x) = \mathbb{1}[[y = 1]]$ , where the indicator function  $\mathbb{1}[[a]]$  is 1 if the  $a$  is true, else 0.

**Ideal Condition<sup>1</sup>:**  $g$  is “ideal” when both perfect and non-overlapping conditions hold and  $(s \perp x)|y$  such that

<sup>1</sup> Eq. (1) is first derived in (Elkan & Noto, 2008).

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

**Range Separability Condition**  $g$  range separates  $P$  and  $N$  iff  $\forall x_1 \in P$  and  $\forall x_2 \in N$ , we have  $g(x_1) > g(x_2)$ .  
**Unassuming Condition:**  $g$  is “unassuming” when perfect and/or non-overlapping conditions may not be true.

Their relationship is: **Unassuming**  $\supset$  **Range Separability**  $\supset$  **Ideal** = **Perfect**  $\cap$  **Non-overlapping**.

We can now state the two goals of Rank Pruning for  $\tilde{P}\tilde{N}$  learning. **Goal 1** is to perfectly estimate  $\hat{\rho}_1 \triangleq \rho_1$  and  $\hat{\rho}_0 \triangleq \rho_0$  when  $g$  is ideal. When  $g$  is not ideal, to our knowledge perfect estimation of  $\rho_1$  and  $\rho_0$  is impossible and at best **Goal 1** is to provide exact expressions for  $\hat{\rho}_1$  and  $\hat{\rho}_0$  w.r.t.  $\rho_1$  and  $\rho_0$ . **Goal 2** is to use  $\hat{\rho}_1$  and  $\hat{\rho}_0$  to uncover the classifications of  $f$  from  $g$ . Both tasks must be accomplished given only observed  $(x, s)$  pairs.  $y, \rho_1, \rho_0, \pi_1$ , and  $\pi_0$  are hidden.

### 3 Rank Pruning

We develop the Rank Pruning algorithm to address our two goals. In Section 3.1, we propose a method for noise rate estimation and prove consistency when  $g$  is ideal. An estimator is “consistent” if it achieves perfect estimation in the expectation of infinite examples. In Section 3.2, we derive exact expressions for  $\hat{\rho}_1$  and  $\hat{\rho}_0$  when  $g$  is unassuming. In Section 3.3, we provide the entire algorithm, and in Section 3.5, prove that Rank Pruning has equivalent expected risk as learning with uncorrupted labels for both ideal  $g$  and non-ideal  $g$  with weaker assumptions. Throughout, we assume  $n \rightarrow \infty$  so that  $P$  and  $N$  are the hidden distributions, each with infinite examples. This is a necessary condition for Theorems. 2, 4 and Lemmas 1, 3.

#### 3.1 Deriving Noise Rate Estimators $\hat{\rho}_1^{conf}$ and $\hat{\rho}_0^{conf}$

We propose the *confident counts* estimators  $\hat{\rho}_1^{conf}$  and  $\hat{\rho}_0^{conf}$  to estimate  $\rho_1$  and  $\rho_0$  as a fraction of the predicted counts of confident examples in each class, encouraging robustness for variation in probability estimation. To estimate  $\rho_1 = P(s = 0|y = 1)$ , we count the number of examples with label  $s = 0$  that we are “confident” have label  $y = 1$  and divide it by the total number of examples that we are “confident” have label  $y = 1$ . More formally,

$$\hat{\rho}_1^{conf} := \frac{|\tilde{N}_{y=1}|}{|\tilde{N}_{y=1}| + |\tilde{P}_{y=1}|}, \quad \hat{\rho}_0^{conf} := \frac{|\tilde{P}_{y=0}|}{|\tilde{P}_{y=0}| + |\tilde{N}_{y=0}|} \tag{2}$$

such that

$$\begin{cases} \tilde{P}_{y=1} = \{x \in \tilde{P} \mid g(x) \geq LB_{y=1}\} \\ \tilde{N}_{y=1} = \{x \in \tilde{N} \mid g(x) \geq LB_{y=1}\} \\ \tilde{P}_{y=0} = \{x \in \tilde{P} \mid g(x) \leq UB_{y=0}\} \\ \tilde{N}_{y=0} = \{x \in \tilde{N} \mid g(x) \leq UB_{y=0}\} \end{cases} \tag{3}$$

where  $g$  is fit to the corrupted training set  $D_\rho$  to obtain  $g(x) = P(\hat{s} = 1|x)$ . The threshold  $LB_{y=1}$  is the predicted probability in  $g(x)$  above which we guess that an example  $x$  has hidden label  $y = 1$ , and similarly for upper bound  $UB_{y=0}$ .  $LB_{y=1}$  and  $UB_{y=0}$  partition  $\tilde{P}$  and  $\tilde{N}$  into four sets representing a *best guess* of a *subset* of examples having labels (1)  $s = 1, y = 0$ , (2)  $s = 1, y = 1$ , (3)  $s = 0, y = 0$ , (4)  $s = 0, y = 1$ . The threshold values are defined as

$$\begin{cases} LB_{y=1} := P(\hat{s} = 1 \mid s = 1) = E_{x \in \tilde{P}}[g(x)] \\ UB_{y=0} := P(\hat{s} = 1 \mid s = 0) = E_{x \in \tilde{N}}[g(x)] \end{cases}$$

where  $\hat{s}$  is the predicted label from a classifier fit to the observed data.  $|\tilde{P}_{y=1}|$  counts examples with label  $s = 1$  that are *most* likely to be correctly labeled ( $y = 1$ ) because  $LB_{y=1} = P(\hat{s} = 1|s = 1)$ . The three other terms in Eq. (3) follow similar reasoning. Importantly, the four terms do not sum to  $n$ , i.e.  $|\tilde{N}| + |\tilde{P}|$ , but  $\hat{\rho}_1^{conf}$  and  $\hat{\rho}_0^{conf}$  are valid estimates because mislabeling noise is assumed to be uniformly random. The choice of threshold values relies on the following two important equations:

$$\begin{aligned}
LB_{y=1} &= E_{x \in \tilde{P}}[g(x)] = E_{x \in \tilde{P}}[P(s=1|x)] \\
&= E_{x \in \tilde{P}}[P(s=1|x, y=1)P(y=1|x) + P(s=1|x, y=0)P(y=0|x)] \\
&= E_{x \in \tilde{P}}[P(s=1|y=1)P(y=1|x) + P(s=1|y=0)P(y=0|x)] \\
&= (1 - \rho_1)(1 - \pi_1) + \rho_0\pi_1
\end{aligned} \tag{4}$$

Similarly, we have

$$UB_{y=0} = (1 - \rho_1)\pi_0 + \rho_0(1 - \pi_0) \tag{5}$$

To our knowledge, although simple, this is the first time that the relationship in Eq. (4) (5) has been published, linking the work of [Elkan & Noto \(2008\)](#), [Liu & Tao \(2016\)](#), [Scott et al. \(2013\)](#) and [Natarajan et al. \(2013\)](#). From Eq. (4) (5), we observe that  $LB_{y=1}$  and  $UB_{y=0}$  are linear interpolations of  $1 - \rho_1$  and  $\rho_0$  and since  $\rho_0 < 1 - \rho_1$ , we have that  $\rho_0 < LB_{y=1} \leq 1 - \rho_1$  and  $\rho_0 \leq UB_{y=0} < 1 - \rho_1$ . When  $g$  is ideal we have that  $g(x) = (1 - \rho_1)$ , if  $x \in P$  and  $g(x) = \rho_0$ , if  $x \in N$ . Thus when  $g$  is ideal, the thresholds  $LB_{y=1}$  and  $UB_{y=0}$  in Eq. (3) will perfectly separate  $P$  and  $N$  examples within each of  $\tilde{P}$  and  $\tilde{N}$ . Lemma 1 immediately follows.

**Lemma 1** *When  $g$  is ideal,*

$$\begin{aligned}
\tilde{P}_{y=1} &= \{x \in P \mid s = 1\}, \tilde{N}_{y=1} = \{x \in P \mid s = 0\}, \\
\tilde{P}_{y=0} &= \{x \in N \mid s = 1\}, \tilde{N}_{y=0} = \{x \in N \mid s = 0\}
\end{aligned} \tag{6}$$

Thus, when  $g$  is ideal, the thresholds in Eq. (3) partition the training set such that  $\tilde{P}_{y=1}$  and  $\tilde{N}_{y=0}$  contain the correctly labeled examples and  $\tilde{P}_{y=0}$  and  $\tilde{N}_{y=1}$  contain the mislabeled examples. Theorem 2 follows (for brevity, proofs of all theorems/lemmas are in Appendix A.1-A.5).

**Theorem 2** *When  $g$  is ideal,*

$$\hat{\rho}_1^{conf} = \rho_1, \hat{\rho}_0^{conf} = \rho_0 \tag{7}$$

Thus, when  $g$  is ideal, the *confident counts* estimators  $\hat{\rho}_1^{conf}$  and  $\hat{\rho}_0^{conf}$  are consistent estimators for  $\rho_1$  and  $\rho_0$  and we set  $\hat{\rho}_1 := \hat{\rho}_1^{conf}$ ,  $\hat{\rho}_0 := \hat{\rho}_0^{conf}$ . These steps comprise Rank Pruning noise rate estimation (see Alg. 1). There are two practical observations. First, for any  $g$  with  $T$  fitting time, computing  $\hat{\rho}_1^{conf}$  and  $\hat{\rho}_0^{conf}$  is  $\mathcal{O}(T)$ . Second,  $\hat{\rho}_1$  and  $\hat{\rho}_0$  should be estimated out-of-sample to avoid over-fitting, resulting in sample variations. In our experiments, we use 3-fold cross-validation, requiring at most  $2T = \mathcal{O}(T)$ .

### 3.2 Noise Estimation: Unassuming Case

Theorem 2 states that  $\hat{\rho}_i^{conf} = \rho_i$ ,  $\forall i \in \{0, 1\}$  when  $g$  is ideal. Though theoretically constructive, in practice this is unlikely. Next, we derive expressions for the estimators when  $g$  is unassuming, i.e.  $g$  may not be perfect and  $P$  and  $N$  may have overlapping support.

Define  $\Delta p_o := \frac{|P \cap N|}{|P \cup N|}$  as the fraction of overlapping examples in  $\mathcal{D}$  and remember that  $\Delta g(x) := g(x) - g^*(x)$ . Denote  $LB_{y=1}^* = (1 - \rho_1)(1 - \pi_1) + \rho_0\pi_1$ ,  $UB_{y=0}^* = (1 - \rho_1)\pi_0 + \rho_0(1 - \pi_0)$ . We have

**Lemma 3** *When  $g$  is unassuming, we have*

$$\begin{cases}
LB_{y=1} = LB_{y=1}^* + E_{x \in \tilde{P}}[\Delta g(x)] - \frac{(1 - \rho_1 - \rho_0)^2}{p_{s1}} \Delta p_o \\
UB_{y=0} = UB_{y=0}^* + E_{x \in \tilde{N}}[\Delta g(x)] + \frac{(1 - \rho_1 - \rho_0)^2}{1 - p_{s1}} \Delta p_o \\
\hat{\rho}_1^{conf} = \rho_1 + \frac{1 - \rho_1 - \rho_0}{|P| - |\Delta P_1| + |\Delta N_1|} |\Delta N_1| \\
\hat{\rho}_0^{conf} = \rho_0 + \frac{1 - \rho_1 - \rho_0}{|N| - |\Delta N_0| + |\Delta P_0|} |\Delta P_0|
\end{cases} \tag{8}$$

where



$$\begin{cases} \Delta P_1 = \{x \in P \mid g(x) < LB_{y=1}\} \\ \Delta N_1 = \{x \in N \mid g(x) \geq LB_{y=1}\} \\ \Delta P_0 = \{x \in P \mid g(x) \leq UB_{y=0}\} \\ \Delta N_0 = \{x \in N \mid g(x) > UB_{y=0}\} \end{cases}$$

The second term on the R.H.S. of the  $\hat{\rho}_i^{conf}$  expressions captures the deviation of  $\hat{\rho}_i^{conf}$  from  $\rho_i$ ,  $i = 0, 1$ . This term results from both imperfect  $g(x)$  and overlapping support. Because the term is non-negative,  $\hat{\rho}_i^{conf} \geq \rho_i$ ,  $i = 0, 1$  in the limit of infinite examples. In other words,  $\hat{\rho}_i^{conf}$  is an *upper bound* for the noise rates  $\rho_i$ ,  $i = 0, 1$ . From Lemma 3, it also follows:

**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 shows that  $\hat{\rho}_1^{conf}$  and  $\hat{\rho}_0^{conf}$  are robust to imperfect probability estimation. As long as  $\Delta g(x)$  does not exceed the distance between the threshold in Eq. (3) and the perfect  $g^*(x)$  value,  $\hat{\rho}_1^{conf}$  and  $\hat{\rho}_0^{conf}$  are consistent estimators for  $\rho_1$  and  $\rho_0$ . Our numerical experiments in Section 4 suggest this is reasonable for  $\Delta g(x)$ . The average  $|\Delta g(x)|$  for the MNIST training dataset across different  $(\rho_1, \pi_1)$  varies between 0.01 and 0.08 for a logistic regression classifier, 0.01~0.03 for a CNN classifier, and 0.05~0.10 for the CIFAR dataset with a CNN classifier. Thus, when  $LB_{y=1} - \rho_0$  and  $1 - \rho_1 - UB_{y=0}$  are above 0.1 for these datasets, from Theorem 4 we see that  $\hat{\rho}_i^{conf}$  still accurately estimates  $\rho_i$ .

### 3.3 The Rank Pruning Algorithm

Using  $\hat{\rho}_1$  and  $\hat{\rho}_0$ , we must uncover the classifications of  $f$  from  $g$ . In this section, we describe how Rank Pruning selects confident examples, removes the rest, and trains on the pruned set using a reweighted loss function.

First, we obtain the inverse noise rates  $\hat{\pi}_1, \hat{\pi}_0$  from  $\hat{\rho}_1, \hat{\rho}_0$ :

$$\hat{\pi}_1 = \frac{\hat{\rho}_0}{p_{s1}} \frac{1 - p_{s1} - \hat{\rho}_1}{1 - \hat{\rho}_1 - \hat{\rho}_0}, \quad \hat{\pi}_0 = \frac{\hat{\rho}_1}{1 - p_{s1}} \frac{p_{s1} - \hat{\rho}_0}{1 - \hat{\rho}_1 - \hat{\rho}_0} \quad (9)$$

Next, we prune the  $\hat{\pi}_1|\tilde{P}|$  examples in  $\tilde{P}$  with smallest  $g(x)$  and the  $\hat{\pi}_0|\tilde{N}|$  examples in  $\tilde{N}$  with highest  $g(x)$  and denote the pruned sets  $\tilde{P}_{conf}$  and  $\tilde{N}_{conf}$ . To prune, we define  $k_1$  as the  $(\hat{\pi}_1|\tilde{P}|)^{th}$  smallest  $g(x)$  for  $x \in \tilde{P}$  and  $k_0$  as the  $(\hat{\pi}_0|\tilde{N}|)^{th}$  largest  $g(x)$  for  $x \in \tilde{N}$ . BFPRT ( $\mathcal{O}(n)$ ) (Blum et al., 1973) is used to compute  $k_1$  and  $k_0$  and pruning is reduced to the following  $\mathcal{O}(n)$  filter:

$$\tilde{P}_{conf} := \{x \in \tilde{P} \mid g(x) \geq k_1\}, \quad \tilde{N}_{conf} := \{x \in \tilde{N} \mid g(x) \leq k_0\} \quad (10)$$

Lastly, we refit the classifier to  $X_{conf} = \tilde{P}_{conf} \cup \tilde{N}_{conf}$  by class-conditionally reweighting the loss function for examples in  $\tilde{P}_{conf}$  with weight  $\frac{1}{1-\hat{\rho}_1}$  and examples in  $\tilde{N}_{conf}$  with weight  $\frac{1}{1-\hat{\rho}_0}$  to recover the estimated balance of positive and negative examples. The entire Rank Pruning algorithm is presented in Alg. 1 and illustrated step-by-step on a synthetic dataset in Fig. 1.

We conclude this section with a formal discussion of the loss function and efficiency of Rank Pruning. Define  $\hat{y}_i$  as the predicted label of example  $i$  for the classifier fit to  $X_{conf}, s_{conf}$  and let  $l(\hat{y}_i, s_i)$  be the original loss function for  $x_i \in D_p$ . Then the loss function for Rank Pruning is simply the original loss function exerted on the pruned  $X_{conf}$ , with class-conditional weighting:

$$\tilde{l}(\hat{y}_i, s_i) = \frac{1}{1-\hat{\rho}_1} l(\hat{y}_i, s_i) \cdot \mathbb{1}[[x_i \in \tilde{P}_{conf}]] + \frac{1}{1-\hat{\rho}_0} l(\hat{y}_i, s_i) \cdot \mathbb{1}[[x_i \in \tilde{N}_{conf}]] \quad (11)$$

Effectively this loss function uses a zero-weight for pruned examples. Other than potentially fewer examples, the only difference in the loss function for Rank Pruning and the original loss function is the class-conditional weights. These

**Algorithm 1 Rank Pruning****Input:** Examples  $X$ , corrupted labels  $s$ , classifier  $\text{clf}$ **Part 1. Estimating Noise Rates:**(1.1)  $\text{clf.fit}(X, s)$ 

$$g(x) \leftarrow \text{clf.predict\_crossval\_probability}(\hat{s} = 1|x)$$

$$p_{s1} = \frac{\text{count}(s=1)}{\text{count}(s=0 \vee s=1)}$$

$$LB_{y=1} = E_{x \in \tilde{P}}[g(x)], UB_{y=0} = E_{x \in \tilde{N}}[g(x)]$$

$$(1.2) \hat{\rho}_1 = \hat{\rho}_1^{\text{conf}} = \frac{|\tilde{N}_{y=1}|}{|\tilde{N}_{y=1}| + |\tilde{P}_{y=1}|}, \hat{\rho}_0 = \hat{\rho}_0^{\text{conf}} = \frac{|\tilde{P}_{y=0}|}{|\tilde{P}_{y=0}| + |\tilde{N}_{y=0}|}$$

$$\hat{\pi}_1 = \frac{\hat{\rho}_0}{p_{s1}} \frac{1 - p_{s1} - \hat{\rho}_1}{1 - \hat{\rho}_1 - \hat{\rho}_0}, \hat{\pi}_0 = \frac{\hat{\rho}_1}{1 - p_{s1}} \frac{p_{s1} - \hat{\rho}_0}{1 - \hat{\rho}_1 - \hat{\rho}_0}$$

**Part 2. Prune Inconsistent Examples:**(2.1) Remove  $\hat{\pi}_1 |\tilde{P}|$  examples in  $\tilde{P}$  with least  $g(x)$ , Remove  $\hat{\pi}_0 |\tilde{N}|$  examples in  $\tilde{N}$  with greatest  $g(x)$ ,Denote the remaining training set  $(X_{\text{conf}}, s_{\text{conf}})$ (2.2)  $\text{clf.fit}(X_{\text{conf}}, s_{\text{conf}})$ , with sample weight  $w(x) = \frac{1}{1 - \hat{\rho}_1} \mathbb{1}[[s_{\text{conf}} = 1]] + \frac{1}{1 - \hat{\rho}_0} \mathbb{1}[[s_{\text{conf}} = 0]]$ **Output:**  $\text{clf}$ 

constant factors do not increase the complexity of the minimization of the original loss function. In other words, we can fairly report the running time of Rank Pruning in terms of the running time ( $\mathcal{O}(T)$ ) of the choice of probabilistic estimator. Combining noise estimation ( $\mathcal{O}(T)$ ), pruning ( $\mathcal{O}(n)$ ), and the final fitting ( $\mathcal{O}(T)$ ), Rank Pruning has a running time of  $\mathcal{O}(T) + \mathcal{O}(n)$ , which is  $\mathcal{O}(T)$  for typical classifiers.

**3.4 Rank Pruning: A simple summary**

Recognizing that formalization can create obfuscation, in this section we describe the entire algorithm in a few sentences. Rank Pruning takes as input training examples  $X$ , noisy labels  $s$ , and a probabilistic classifier  $\text{clf}$  and finds a subset of  $X, s$  that is likely to be correctly labeled, i.e. a subset of  $X, y$ . To do this, we first find two thresholds,  $LB_{y=1}$  and  $UB_{y=0}$ , to *confidently* guess the correctly and incorrectly labeled examples in each of  $\tilde{P}$  and  $\tilde{N}$ , forming four sets, then use the set sizes to estimate the noise rates  $\rho_1 = P(s = 0|y = 1)$  and  $\rho_0 = P(s = 1|y = 0)$ . We then use the noise rates to estimate the number of examples with observed label  $s = 1$  and hidden label  $y = 0$  and remove that number of examples from  $\tilde{P}$  by removing those with lowest predicted probability  $g(x)$ . We prune  $\tilde{N}$  similarly. Finally, the classifier is fit to the pruned set, which is intended to represent a subset of the correctly labeled data.

**3.5 Expected Risk Evaluation**

In this section, we prove Rank Pruning exactly uncovers the classifier  $f$  fit to hidden  $y$  labels when  $g$  range separates  $P$  and  $N$  and  $\rho_1$  and  $\rho_0$  are given.

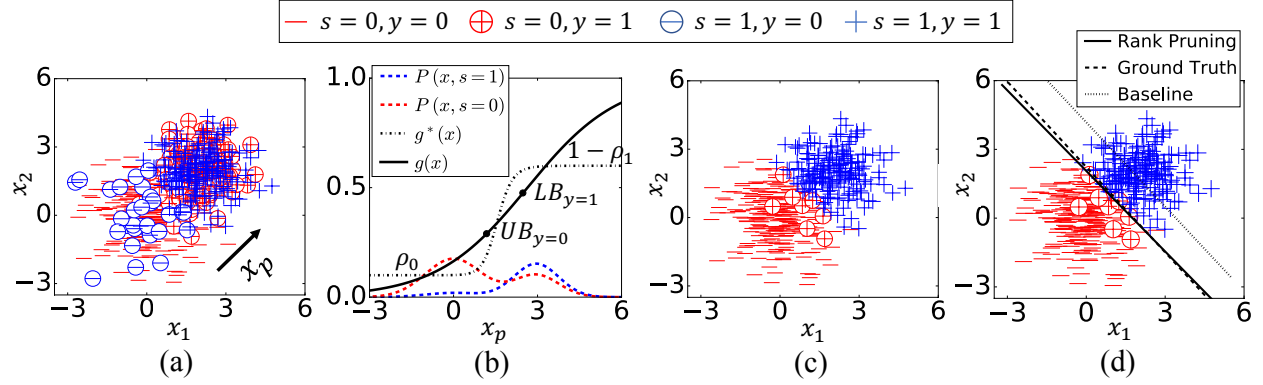
Denote  $f_\theta \in \mathcal{F} : x \rightarrow \hat{y}$  as a classifier's prediction function belonging to some function space  $\mathcal{F}$ , where  $\theta$  represents the classifier's parameters.  $f_\theta$  represents  $f$ , but without  $\theta$  necessarily fit to the training data.  $\hat{f}$  is the Rank Pruning estimate of  $f$ .

Denote the empirical risk of  $f_\theta$  w.r.t. the loss function  $\tilde{l}$  and corrupted data  $D_\rho$  as  $\hat{R}_{\tilde{l}, D_\rho}(f_\theta) = \frac{1}{n} \sum_{i=1}^n \tilde{l}(f_\theta(x_i), s_i)$ , and the expected risk of  $f_\theta$  w.r.t. the corrupted distribution  $\mathcal{D}_\rho$  as  $R_{\tilde{l}, \mathcal{D}_\rho}(f_\theta) = E_{(x,s) \sim \mathcal{D}_\rho}[\tilde{l}(f_\theta(x), s)]$ . Similarly, denote  $R_{l, \mathcal{D}}(f_\theta)$  as the expected risk of  $f_\theta$  w.r.t. the hidden distribution  $\mathcal{D}$  and loss function  $l$ . We show that using Rank Pruning, a classifier  $\hat{f}$  can be learned for the hidden data  $D$ , given the corrupted data  $D_\rho$ , by minimizing the empirical risk:

$$\hat{f} = \underset{f_\theta \in \mathcal{F}}{\text{argmin}} \hat{R}_{\tilde{l}, D_\rho}(f_\theta) = \underset{f_\theta \in \mathcal{F}}{\text{argmin}} \frac{1}{n} \sum_{i=1}^n \tilde{l}(f_\theta(x_i), s_i) \quad (12)$$

Under the *range separability* condition, we have





**Figure 1:** Illustration of Rank Pruning with a logistic regression classifier ( $\mathcal{LR}_\theta$ ). (a): The corrupted training set  $D_\rho$  with noise rates  $\rho_1 = 0.4$  and  $\rho_0 = 0.1$ . Corrupted colored labels ( $s = 1, s = 0$ ) are observed.  $y$  (+, -) is hidden. (b): The marginal distribution of  $D_\rho$  projected onto the  $x_p$  axis (indicated in (a)), and the  $\mathcal{LR}_\theta$ 's estimated  $g(x)$ , from which  $\hat{\rho}_1^{conf} = 0.4237$ ,  $\hat{\rho}_0^{conf} = 0.1144$  are estimated. (c): The pruned  $X_{conf}, s_{conf}$ . (d): The classification result by Rank Pruning ( $\hat{f} = \mathcal{LR}_\theta.\text{fit}(X_{conf}, s_{conf})$ ), ground truth classifier ( $f = \mathcal{LR}_\theta.\text{fit}(X, y)$ ), and baseline classifier ( $g = \mathcal{LR}_\theta.\text{fit}(X, s)$ ), with an accuracy of 94.16%, 94.16% and 78.83%, respectively.

**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

$$R_{\tilde{l}, D_\rho}(f_\theta) = R_{l, \mathcal{D}}(f_\theta) \quad (13)$$

where  $\tilde{l}(\hat{y}_i, s_i)$  is Rank Pruning's loss function (Eq. 11).

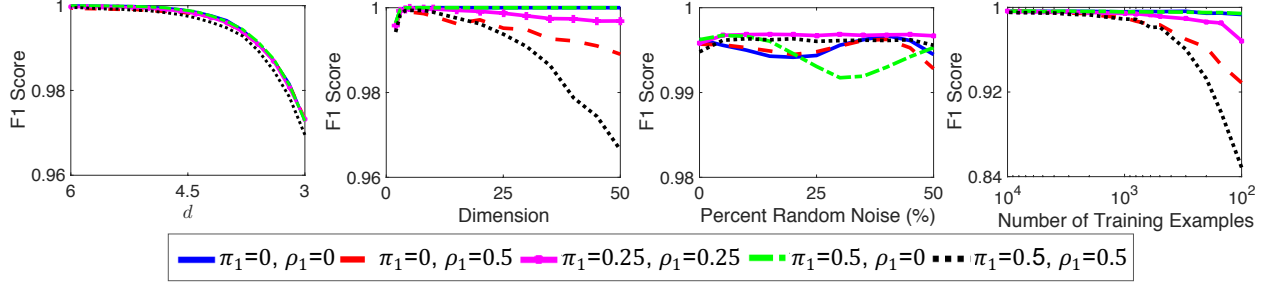
The proof of Theorem 5 is in Appendix A.5. Intuitively, Theorem 5 tells us that if  $g$  range separates  $P$  and  $N$ , then given exact noise rate estimates, Rank Pruning will exactly prune out the positive examples in  $\tilde{N}$  and negative examples in  $\tilde{P}$ , leading to the same expected risk as learning from uncorrupted labels. Thus, Rank Pruning can exactly uncover the classifications of  $f$  (with infinite examples) because the expected risk is equivalent for any  $f_\theta$ . Note Theorem 5 also holds when  $g$  is ideal, since *ideal*  $\subset$  *range separability*. In practice, *range separability* encompasses a wide range of imperfect  $g(x)$  scenarios, e.g.  $g(x)$  can have large fluctuation in both  $P$  and  $N$  or have systematic drift w.r.t. to  $g^*(x)$  due to underfitting.

## 4 Experimental Results

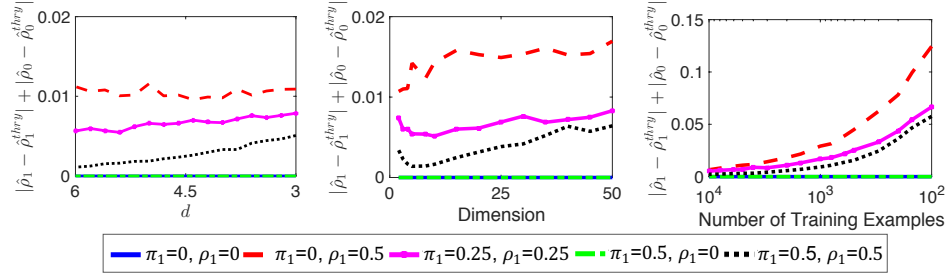
In Section 3, we developed a theoretical framework for Rank Pruning, proved exact noise estimation and equivalent expected risk when conditions are ideal, and derived closed-form solutions when conditions are non-ideal. Our theory suggests that, in practice, Rank Pruning should (1) accurately estimate  $\rho_1$  and  $\rho_0$ , (2) typically achieve as good or better F1, error and AUC-PR (Davis & Goadrich, 2006) as state-of-the-art methods, and (3) be robust to both mislabeling and added noise.

In this section, we support these claims with an evaluation of the comparative performance of Rank Pruning in non-ideal conditions across thousands of scenarios. These include less complex (MNIST) and more complex (CIFAR) datasets, simple (logistic regression) and complex (CNN) classifiers, the range of noise rates, added random noise, separability of  $P$  and  $N$ , input dimension, and number of training examples to ensure that Rank Pruning is a general, agnostic solution for  $\tilde{P}\tilde{N}$  learning.

In our experiments, we adjust  $\pi_1$  instead of  $\rho_0$  because binary noisy classification problems (e.g. detection and recognition tasks) often have that  $|P| \ll |N|$ . This choice allows us to adjust both noise rates with respect to  $P$ , i.e. the fraction of true positive examples that are mislabeled as negative ( $\rho_1$ ) and the fraction of observed positive labels that are actually mislabeled negative examples ( $\pi_1$ ). The  $\tilde{P}\tilde{N}$  learning algorithms are trained with corrupted labels  $s$ , and tested on an unseen test set by comparing predictions  $\hat{y}$  with the true test labels  $y$  using F1 score, error, and



**Figure 2:** Comparison of Rank Pruning with different noise ratios  $(\pi_1, \rho_1)$  on a synthetic dataset for varying separability  $d$ , dimension, added random noise and number of training examples. Default settings for Fig. 2, 3 and 4:  $d = 4$ , 2-dimension, 0% random noise, and 5000 training examples with  $p_{y1} = 0.2$ . The lines are an average of 200 trials.



**Figure 3:** Sum of absolute difference between theoretically estimated  $\hat{\rho}_i^{theory}$  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. (8) do not address added noise examples.

AUC-PR metrics. We include all three to emphasize our apathy toward tuning results to any single metric. We provide F1 scores in this section with error and AUC-PR scores in Appendix C.

#### 4.1 Synthetic Dataset

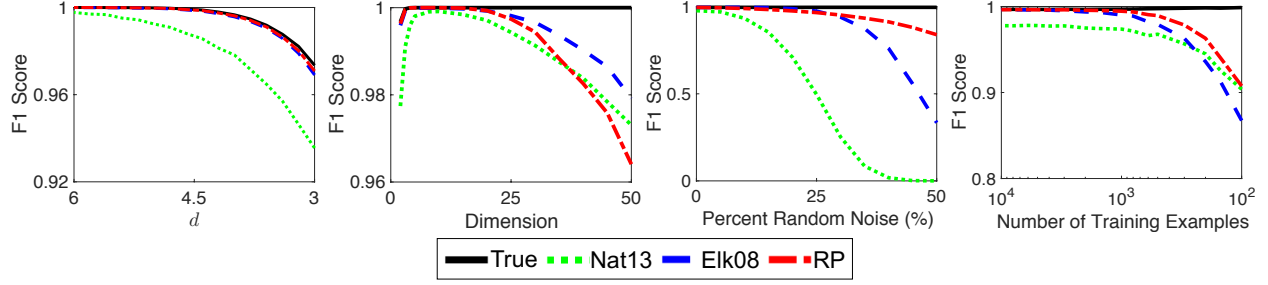
The synthetic dataset is comprised of a Gaussian positive class and a Gaussian negative classes such that negative examples ( $y = 0$ ) obey an  $m$ -dimensional Gaussian distribution  $N(\mathbf{0}, \mathbf{I})$  with unit variance  $\mathbf{I} = \text{diag}(1, 1, \dots, 1)$ , and positive examples obey  $N(d\mathbf{1}, 0.8\mathbf{I})$ , where  $d\mathbf{1} = (d, d, \dots, d)$  is an  $m$ -dimensional vector, and  $d$  measures the separability of the positive and negative set.

We test Rank Pruning by varying 4 different settings of the environment: separability  $d$ , dimension, number of training examples  $n$ , and percent (of  $n$ ) added random noise drawn from a uniform distribution  $U([-10, 10]^m)$ . In each scenario, we test 5 different  $(\pi_1, \rho_1)$  pairs:  $(\pi_1, \rho_1) \in \{(0, 0), (0, 0.5), (0.25, 0.25), (0.5, 0), (0.5, 0.5)\}$ . From Fig. 2, we observe that across these settings, the F1 score for Rank Pruning is fairly agnostic to magnitude of mislabeling (noise rates). As a validation step, in Fig. 3 we measure how closely our empirical estimates match our theoretical solutions in Eq. (8) and find near equivalence except when the number of training examples approaches zero.

For significant mislabeling ( $\rho_1 = 0.5$ ,  $\pi_1 = 0.5$ ), Rank Pruning often outperforms other methods (Fig. 4). In the scenario of different separability  $d$ , it achieves nearly the same F1 score as the ground truth classifier. Remarkably, from Fig. 2 and Fig. 4, we observe that when added random noise comprises 50% of total training examples, Rank Pruning still achieves  $F1 > 0.85$ , compared with  $F1 < 0.5$  for all other methods. This emphasizes a unique feature of Rank Pruning, it will also remove added random noise because noise drawn from a third distribution is unlikely to appear confidently positive or negative.

#### 4.2 MNIST and CIFAR Datasets

We consider the binary classification tasks of one-vs-rest for the MNIST (LeCun & Cortes, 2010) and CIFAR-10 (Krizhevsky et al.) datasets, e.g. the “car vs rest” task in CIFAR is to predict if an image is a “car” or “not”.  $\rho_1$  and  $\pi_1$  are given to all  $\tilde{P}\tilde{N}$  learning methods for fair comparison, except for  $RP_\rho$  which is Rank Pruning including noise rate



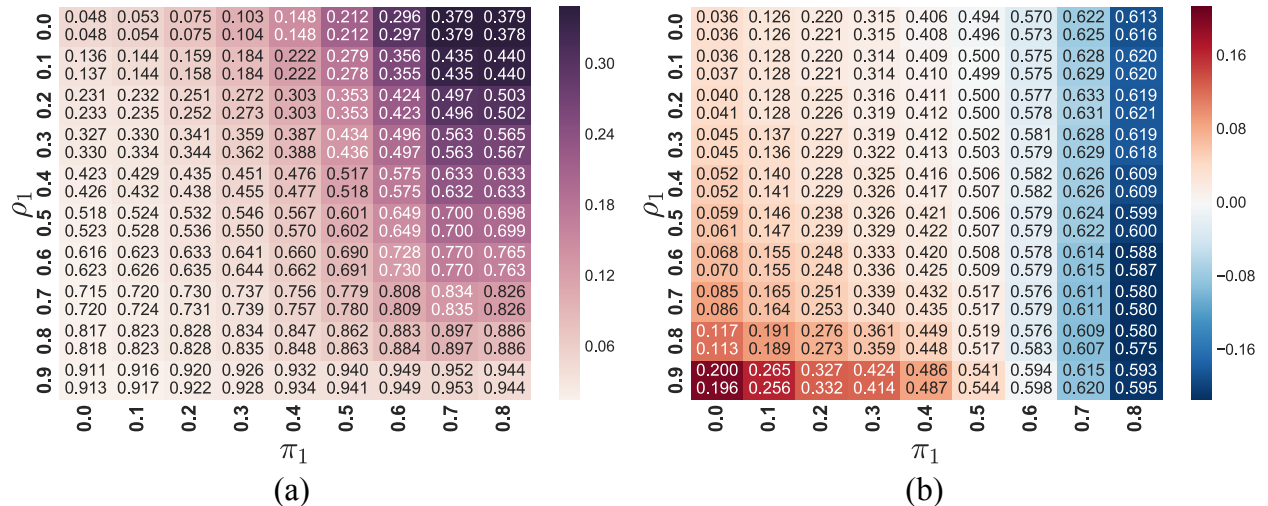
**Figure 4:** Comparison of  $\tilde{P}\tilde{N}$  methods for varying separability  $d$ , dimension, added random noise, and number of training examples for  $\pi_1 = 0.5$ ,  $\rho_1 = 0.5$  (given to all methods).

estimation.  $RP_\rho$  metrics measure our performance on the unadulterated  $\tilde{P}\tilde{N}$  learning problem.

As evidence that Rank Pruning is dataset and classifier agnostic, we demonstrate its superiority with both (1) a linear logistic regression model with unit L2 regularization and (2) an AlexNet CNN variant with max pooling and dropout, modified to have a two-class output. The CNN structure is adapted from [Chollet \(2016b\)](#) for MNIST and [Chollet \(2016a\)](#) for CIFAR. CNN training ends when a 10% holdout set shows no loss decrease for 10 epochs (max 50 for MNIST and 150 for CIFAR).

We consider noise rates  $\pi_1, \rho_1 \in \{(0, 0.5), (0.25, 0.25), (0.5, 0), (0.5, 0.5)\}$  for both MNIST and CIFAR, with additional settings for MNIST in Table 3 to emphasize Rank Pruning performance is noise rate agnostic. The  $\rho_1 = 0, \pi_1 = 0$  case is omitted because when given  $\rho_1, \pi_1$ , all methods have the same loss function as the ground truth classifier, resulting in nearly identical F1 scores. Note that in general, Rank Pruning does not require perfect probability estimation to achieve perfect F1-score. As an example, this occurs when  $P$  and  $N$  are range-separable, and the rank order of the sorted  $g(x)$  probabilities in  $P$  and  $N$  is consistent with the rank of the perfect probabilities, regardless of the actual values of  $g(x)$ .

For MNIST using logistic regression, we evaluate the consistency of our noise rate estimates with actual noise rates and theoretical estimates (Eq. 8) across  $\pi_1 \in [0, 0.8] \times \rho_1 \in [0, 0.9]$ . The computing time for one setting was  $\sim 10$  minutes on a single CPU core. The results for  $\hat{\rho}_1$  and  $\hat{\pi}_1$  (Fig. 5) are satisfyingly consistent, with mean absolute difference  $MD_{\hat{\rho}_1, \rho_1} = 0.105$  and  $MD_{\hat{\pi}_1, \pi_1} = 0.062$ , and validate our theoretical solutions ( $MD_{\hat{\rho}_1, \hat{\rho}_1^{thy}} = 0.0028$ ,  $MD_{\hat{\pi}_1, \hat{\pi}_1^{thy}} = 0.0058$ ). The deviation of the theoretical and empirical estimates reflects the assumption that we have



**Figure 5:** Rank Pruning  $\hat{\rho}_1$  and  $\hat{\pi}_1$  estimation consistency, averaged over all digits in MNIST. (a) Color depicts  $\hat{\rho}_1 - \rho_1$  with  $\hat{\rho}_1$  (upper) and theoretical  $\hat{\rho}_1^{thy}$  (lower) in each block. (b) Color depicts  $\hat{\pi}_1 - \pi_1$  with  $\hat{\pi}_1$  (upper) and  $\hat{\pi}_1^{thy}$  (lower) in each block.

**Table 3:** Comparison of F1 score for one-vs-rest MNIST and CIFAR-10 (averaged over all digits/images) using logistic regression. 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. Due to sensitivity to imperfect  $g(x)$ , *Liu16* often predicts the same label for all examples.

DATASET	CIFAR				MNIST															
	$\pi_1 =$				$\pi_1 = 0.0$				$\pi_1 = 0.25$				$\pi_1 = 0.5$				$\pi_1 = 0.75$			
	0.0	0.25	0.5	0.5	0.25	0.5	0.75	0.0	0.25	0.5	0.75	0.0	0.25	0.5	0.75	0.0	0.25	0.5	0.75	
MODEL, $\rho_1 =$	0.5	0.25	0.0	0.5	0.25	0.5	0.75	0.0	0.25	0.5	0.75	0.0	0.25	0.5	0.75	0.0	0.25	0.5	0.75	
TRUE	0.248	0.248	0.248	0.248	0.894	0.894	0.894	0.894	0.894	0.894	0.894	0.894	0.894	0.894	0.894	0.894	0.894	0.894	0.894	
RP <sub><math>\rho</math></sub>	<b>0.301</b>	<b>0.316</b>	<b>0.308</b>	<b>0.261</b>	<b>0.883</b>	<b>0.874</b>	<b>0.843</b>	<b>0.881</b>	<b>0.876</b>	<b>0.863</b>	<b>0.799</b>	0.823	0.831	<b>0.819</b>	<b>0.762</b>	0.583	0.603	0.587	0.532	
RP	<b>0.256</b>	<b>0.262</b>	<b>0.244</b>	0.209	<b>0.885</b>	<b>0.873</b>	<b>0.839</b>	<b>0.890</b>	<b>0.879</b>	<b>0.863</b>	<b>0.812</b>	<b>0.879</b>	<b>0.862</b>	<b>0.838</b>	<b>0.770</b>	<b>0.855</b>	<b>0.814</b>	<b>0.766</b>	0.617	
NAT13	0.226	0.219	0.194	0.195	0.860	0.830	0.774	0.865	0.836	0.802	0.748	0.839	0.810	0.777	0.721	0.809	0.776	0.736	<b>0.640</b>	
ELK08	0.221	0.226	0.228	<b>0.210</b>	0.862	0.830	0.771	0.864	0.847	0.819	0.762	0.843	0.835	0.814	0.736	0.674	0.669	0.599	0.473	
Liu16	0.182	0.182	0.000	0.182	0.021	0.000	0.000	0.000	0.147	0.147	0.073	0.000	0.164	0.163	0.163	0.047	0.158	0.145	0.164	

**Table 4:** F1 score comparison on MNIST and CIFAR-10 using a CNN. Except for  $RP_\rho$ ,  $\rho_1, \rho_0$  are given to all methods.

MNIST/CIFAR IMAGE CLASS	TRUE	$\pi_1 = 0.0$ $\rho_1 = 0.5$					$\pi_1 = 0.25$ $\rho_1 = 0.25$					$\pi_1 = 0.5$									
		$RP_\rho$	RP	NAT13	ELK08	Liu16	$RP_\rho$	RP	NAT13	ELK08	Liu16	$RP_\rho$	RP	$\rho_1 = 0.0$ NAT13	ELK08	Liu16	$RP_\rho$	RP	$\rho_1 = 0.5$ NAT13	ELK08	Liu16
0	0.993	<b>0.991</b>	<b>0.988</b>	0.977	0.976	0.179	<b>0.991</b>	<b>0.992</b>	0.982	0.981	0.179	<b>0.991</b>	<b>0.992</b>	0.984	0.987	0.985	<b>0.989</b>	<b>0.989</b>	0.937	0.964	0.179
1	0.993	<b>0.990</b>	<b>0.991</b>	0.989	0.985	0.204	<b>0.992</b>	<b>0.992</b>	0.984	0.987	0.204	0.990	0.991	0.992	<b>0.993</b>	0.990	<b>0.989</b>	<b>0.989</b>	0.984	0.988	0.204
2	0.987	<b>0.973</b>	<b>0.976</b>	0.972	0.969	0.187	<b>0.984</b>	<b>0.983</b>	0.978	0.975	0.187	0.985	0.986	0.985	0.986	<b>0.988</b>	<b>0.971</b>	<b>0.975</b>	0.968	0.959	0.187
3	0.990	<b>0.984</b>	<b>0.984</b>	0.972	0.981	0.183	<b>0.986</b>	<b>0.986</b>	0.978	0.978	0.183	<b>0.990</b>	0.987	<b>0.989</b>	<b>0.989</b>	0.984	<b>0.981</b>	<b>0.979</b>	0.957	0.971	0.183
4	0.994	<b>0.981</b>	0.979	<b>0.981</b>	0.977	0.179	<b>0.985</b>	<b>0.987</b>	0.971	0.964	0.179	<b>0.987</b>	<b>0.990</b>	<b>0.990</b>	0.989	0.985	<b>0.977</b>	<b>0.982</b>	0.955	0.961	0.179
5	0.989	<b>0.982</b>	<b>0.980</b>	0.978	0.979	0.164	<b>0.985</b>	<b>0.982</b>	0.964	0.965	0.164	<b>0.988</b>	<b>0.987</b>	<b>0.987</b>	0.984	<b>0.987</b>	<b>0.965</b>	<b>0.968</b>	0.962	0.957	0.164
6	0.989	<b>0.986</b>	<b>0.985</b>	0.972	0.982	0.175	<b>0.985</b>	<b>0.987</b>	0.978	0.981	0.175	0.985	0.985	<b>0.988</b>	0.987	0.985	<b>0.983</b>	<b>0.982</b>	0.946	0.959	0.175
7	0.987	<b>0.981</b>	<b>0.980</b>	0.967	0.948	0.186	<b>0.976</b>	<b>0.975</b>	0.971	0.971	0.186	0.976	0.980	<b>0.985</b>	0.982	0.983	<b>0.973</b>	<b>0.968</b>	0.942	0.958	0.186
8	0.989	<b>0.975</b>	<b>0.978</b>	0.943	0.967	0.178	<b>0.982</b>	<b>0.981</b>	0.967	0.951	0.178	0.982	<b>0.984</b>	0.982	0.979	0.983	<b>0.977</b>	<b>0.975</b>	0.864	0.959	0.178
9	0.982	<b>0.966</b>	<b>0.974</b>	0.972	0.935	0.183	<b>0.976</b>	<b>0.974</b>	0.967	0.967	0.183	0.976	0.975	0.974	<b>0.978</b>	0.970	<b>0.959</b>	0.940	0.931	<b>0.942</b>	0.183
AVG <sub>MN</sub>	0.989	<b>0.981</b>	<b>0.981</b>	0.972	0.970	0.182	<b>0.984</b>	<b>0.984</b>	0.974	0.972	0.182	0.985	<b>0.986</b>	<b>0.986</b>	0.985	0.984	<b>0.976</b>	<b>0.975</b>	0.945	0.962	0.182
PLANE	0.755	<b>0.689</b>	<b>0.634</b>	0.619	0.585	0.182	<b>0.695</b>	<b>0.702</b>	0.671	0.640	0.182	<b>0.757</b>	<b>0.746</b>	0.716	0.735	0.000	<b>0.628</b>	<b>0.635</b>	0.459	0.598	0.182
AUTO	0.891	<b>0.791</b>	<b>0.785</b>	0.761	0.768	0.000	<b>0.832</b>	<b>0.824</b>	0.771	0.783	0.182	0.862	0.866	<b>0.869</b>	0.865	0.000	<b>0.749</b>	<b>0.720</b>	0.582	0.501	0.182
BIRD	0.669	<b>0.504</b>	<b>0.483</b>	0.445	0.389	0.182	<b>0.543</b>	<b>0.515</b>	0.469	0.426	0.182	<b>0.577</b>	<b>0.619</b>	0.543	0.551	0.000	<b>0.447</b>	<b>0.409</b>	0.366	0.387	0.182
CAT	0.487	<b>0.350</b>	0.279	0.310	<b>0.313</b>	0.000	<b>0.426</b>	0.317	<b>0.350</b>	0.345	0.182	<b>0.489</b>	<b>0.433</b>	0.426	0.347	0.000	<b>0.394</b>	0.282	0.240	<b>0.313</b>	0.182
DEER	0.726	<b>0.593</b>	<b>0.540</b>	0.455	0.522	0.182	<b>0.585</b>	0.554	0.480	<b>0.569</b>	0.182	0.614	0.630	<b>0.643</b>	0.633	0.000	<b>0.458</b>	0.375	0.310	<b>0.383</b>	0.182
DOG	0.569	<b>0.544</b>	<b>0.577</b>	0.429	0.456	0.000	<b>0.579</b>	0.559	0.569	<b>0.576</b>	0.182	0.647	0.637	<b>0.667</b>	0.630	0.000	<b>0.516</b>	0.461	0.412	<b>0.465</b>	0.182
FROG	0.815	<b>0.746</b>	0.727	<b>0.733</b>	0.718	0.000	<b>0.729</b>	<b>0.750</b>	0.630	0.584	0.182	0.767	<b>0.782</b>	0.777	0.770	0.000	<b>0.635</b>	<b>0.615</b>	0.589	0.524	0.182
HORSE	0.805	<b>0.690</b>	0.670	0.624	<b>0.672</b>	0.182	<b>0.710</b>	0.669	<b>0.683</b>	0.627	0.182	0.761	<b>0.776</b>	0.769	0.753	0.000	<b>0.672</b>	<b>0.569</b>	0.551	0.461	0.182
SHIP	0.851	<b>0.791</b>	<b>0.783</b>	0.719	0.758	0.182	<b>0.810</b>	<b>0.801</b>	0.758	0.723	0.182	0.816	0.822	0.830	<b>0.831</b>	0.000	<b>0.715</b>	<b>0.738</b>	0.569	0.632	0.182
TRUCK	0.861	<b>0.744</b>	<b>0.722</b>	0.655	0.665	0.182	<b>0.814</b>	<b>0.826</b>	0.798	0.774	0.182	0.812	0.830	<b>0.826</b>	0.824	0.000	<b>0.654</b>	0.543	0.575	<b>0.584</b>	0.182
AVG <sub>CF</sub>	0.743	<b>0.644</b>	<b>0.620</b>	0.575	0.585	0.109	<b>0.672</b>	<b>0.652</b>	0.618	0.605	0.182	<b>0.710</b>	<b>0.714</b>	0.707	0.694	0.000	<b>0.587</b>	<b>0.535</b>	0.465	0.485	0.182

infinite examples, whereas empirically, the number of examples is finite.

We emphasize two observations from our analysis on CIFAR and MNIST. First, Rank Pruning performs well in nearly every scenario and boasts the most dramatic improvement over prior state-of-the-art in the presence of extreme noise ( $\pi_1 = 0.5, \rho_1 = 0.5$ ). This is easily observed in the right-most quadrant of Table 4. The  $\pi_1 = 0.5, \rho_1 = 0$  quadrant is nearest to  $\pi_1 = 0, \rho_1 = 0$  and mostly captures CNN prediction variation because  $|\tilde{P}| \ll |\tilde{N}|$ .

Second,  $RP_\rho$  often achieves equivalent (MNIST in Table 4) or significantly higher (CIFAR in Tables 3 and 4) F1 score than Rank Pruning when  $\rho_1$  and  $\pi_1$  are provided, particularly when noise rates are large. This effect is exacerbated for harder problems (lower F1 score for the ground truth classifier) like the “cat” in CIFAR or the “9” digit in MNIST likely because these problems are more complex, resulting in less confident predictions, and therefore more pruning.

Remember that  $\rho_1^{conf}$  and  $\rho_0^{conf}$  are upper bounds when  $g$  is unassuming. Noise rate overestimation accounts for the complexity of harder problems. As a downside, Rank Pruning may remove correctly labeled examples that “confuse” the classifier, instead fitting only the confident examples in each class. We observe this on CIFAR in Table 3 where logistic regression severely underfits so that  $RP_\rho$  has significantly higher F1 score than the ground truth classifier. Although Rank Pruning with noisy labels seemingly outperforms the ground truth model, if we lower the classification threshold to 0.3 instead of 0.5, the performance difference goes away by accounting for the lower probability predictions.

## 5 Discussion

To our knowledge, Rank Pruning is the first time-efficient algorithm, w.r.t. classifier fitting time, for  $\tilde{P}\tilde{N}$  learning that achieves similar or better F1, error, and AUC-PR than current state-of-the-art methods across practical scenarios for synthetic, MNIST, and CIFAR datasets, with logistic regression and CNN classifiers, across all noise rates,  $\rho_1, \rho_0$ , for varying added noise, dimension, separability, and number of training examples. By *learning with confident examples*, we discover provably consistent estimators for noise rates,  $\rho_1, \rho_0$ , derive theoretical solutions when  $g$  is unassuming,

and accurately uncover the classifications of  $f$  fit to hidden labels, perfectly when  $g$  range separates  $P$  and  $N$ .

We recognize that disambiguating whether we are in the unassuming or range separability condition may be desirable. Although knowing  $g^*(x)$  and thus  $\Delta g(x)$  is impossible, if we assume randomly uniform noise, and toggling the  $LB_{y=1}$  threshold does not change  $\rho_1^{conf}$ , then  $g$  range separates  $P$  and  $N$ . When  $g$  is unassuming, Rank Pruning is still robust to imperfect  $g(x)$  within a range separable subset of  $P$  and  $N$  by training with confident examples even when noise rate estimates are inexact.

An important contribution of Rank Pruning is generality, both in classifier and implementation. The use of logistic regression and a generic CNN in our experiments emphasizes that our findings are not dependent on model complexity. We evaluate thousands of scenarios to avoid findings that are an artifact of problem setup. A key point of Rank Pruning is that we only report the simplest, non-parametric version. For example, we use 3-fold cross-validation to compute  $g(x)$  even though we achieved improved performance with larger folds. We tried many variants of pruning and achieved significant higher F1 for MNIST and CIFAR, but to maintain generality, we present only the basic model.

At its core, Rank Pruning is a simple, robust, and general solution for noisy binary classification by *learning with confident examples*, but it also challenges how we think about training data. For example, SVM showed how a decision boundary can be recovered from only support vectors. Yet, when training data contains significant mislabeling, confident examples, many of which are far from the boundary, are informative for uncovering the true relationship  $P(y = 1|x)$ . Although modern affordances of “big data” emphasize the value of *more* examples for training, through Rank Pruning we instead encourage a rethinking of learning with *confident* examples.

## References

- Aha, D. W., Kibler, D., and Albert, M. K. Instance-based learning algorithms. *Mach. Learn.*, 6(1):37–66, 1991.
- Angelova, A., Abu-Mostafam, Y., and Perona, P. Pruning training sets for learning of object categories. In *CVPR*, volume 1, pp. 494–501. IEEE, 2005.
- Angluin, D. and Laird, P. Learning from noisy examples. *Machine Learning*, 2(4):343–370, 1988.
- Blanchard, G., Lee, G., and Scott, C. Semi-supervised novelty detection. *J. Mach. Learn. Res.*, 11:2973–3009, December 2010. ISSN 1532-4435.
- Blum, A. and Mitchell, T. Combining labeled and unlabeled data with co-training. In *11th Conf. on COLT*, pp. 92–100, New York, NY, USA, 1998. ACM.
- Blum, M., Floyd, R. W., Pratt, V., Rivest, R. L., and Tarjan, R. E. Time bounds for selection. *J. Comput. Syst. Sci.*, 7(4):448–461, August 1973. ISSN 0022-0000.
- Breiman, L. Bagging predictors. *Machine Learning*, 24(2):123–140, August 1996. ISSN 0885-6125.
- Chapelle, O. and Vapnik, V. Model selection for support vector machines. In *Proc. of 12th NIPS*, pp. 230–236, Cambridge, MA, USA, 1999.
- Chollet, F. *Keras CIFAR CNN*, 2016a. [bit.ly/2mVKR3d](https://bit.ly/2mVKR3d).
- Chollet, F. *Keras MNIST CNN*, 2016b. [bit.ly/2nKiqJv](https://bit.ly/2nKiqJv).
- Claesen, M., Smet, F. D., Suykens, J. A., and Moor, B. D. A robust ensemble approach to learn from positive and unlabeled data using {SVM} base models. *Neurocomputing*, 160:73 – 84, 2015. ISSN 0925-2312.
- Davis, J. and Goadrich, M. The relationship between precision-recall and roc curves. In *Proc. of 23rd ICML*, pp. 233–240, NYC, NY, USA, 2006. ACM.
- Elkan, C. and Noto, K. Learning classifiers from only positive and unlabeled data. In *Proc. of 14th KDD*, pp. 213–220, NYC, NY, USA, 2008. ACM.
- Hempstalk, K., Frank, E., and Witten, I. H. One-class classification by combining density and class probability estimation. In *Proc. of ECML-PKDD*, pp. 505–519, Berlin, Heidelberg, 2008. Springer-Verlag.
- Krizhevsky, A., Nair, V., and Hinton, G. Cifar-10 (canadian institute for advanced research).
- Lane, N. D., Miluzzo, E., Lu, H., Peebles, D., Choudhury, T., and Campbell, A. T. A survey of mobile phone sensing. *IEEE Communications*, 48(9), 2010.
- LeCun, Y. and Cortes, C. MNIST handwritten digit database. 2010.
- Lee, W. and Liu, B. Learning with positive and unlabeled examples using weighted logistic regression. In *Proc. of 20th ICML*, volume 1, pp. 448–455, 12 2003.
- Liu, B., Dai, Y., Li, X., Lee, W. S., and Yu, P. S. Building text classifiers using positive and unlabeled examples. In *Proc. of 3rd ICDM*, pp. 179–, Washington, DC, USA, 2003. IEEE Computer Society.
- Liu, T. and Tao, D. Classification with noisy labels by importance reweighting. *IEEE Trans. Pattern Anal. Mach. Intell.*, 38(3):447–461, March 2016.
- Manevitz, L. M. and Yousef, M. One-class svms for document classification. *JMLR*, 2:139–154, March 2002.
- Menon, A. K., Jiang, X., Vembu, S., Elkan, C., and Ohno-Machado, L. Predicting accurate probabilities with a ranking loss. *CoRR*, abs/1206.4661, 2012.
- Michalski, S. R., Carbonell, G. J., and Mitchell, M. T. *ML an AI Approach*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1986.
- Mordelet, F. and Vert, J. P. A bagging svm to learn from positive and unlabeled examples. *Pattern Recogn. Lett.*, 37: 201–209, February 2014. ISSN 0167-8655.
- 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.
- Natarajan, N., Dhillon, I. S., Ravikumar, P. K., and Tewari, A. Learning with noisy labels. In *Adv. in NIPS 26*, pp. 1196–1204. Curran Associates, Inc., 2013.

- Nettleton, D. F., Orriols-Puig, A., and Fornells, A. A study of the effect of different types of noise on the precision of supervised learning techniques. *Artificial Intelligence Review*, 33(4):275–306, 2010.
- Nigam, K. and Ghani, R. Understanding the behavior of co-training. In *KDD Workshop*, 2000.
- Paolacci, G., Chandler, J., and Ipeirotis, P. G. Running experiments on amazon mechanical turk. *Judgment and Decision Making*, 5(5):411–419, 2010.
- Platt, J., Schölkopf, B., Shawe-Taylor, J., Smola, A. J., and Williamson, R. C. Estimating support of a high dimensional distribution. Technical report, MSR, 1999.
- Raviv, Y. and Intrator, N. Bootstrapping with noise: An effective regularization technique. *Connection Science*, 8(3-4):355–372, 1996.
- Reed, S. E., Lee, H., Anguelov, D., Szegedy, C., Erhan, D., and Rabinovich, A. Training deep neural networks on noisy labels with bootstrapping. In *ICLR*, 2015.
- scikit learn. *LogisticRegression Class at scikit-learn*, 2016.
- Scott, C. A rate of convergence for mixture proportion estimation, with application to learning from noisy labels. *JMLR*, 38:838–846, 2015. ISSN 1532-4435.
- Scott, C., Blanchard, G., and Handy, G. Classification with asymmetric label noise: Consistency and maximal denoising. In *COLT*, pp. 489–511, 2013.
- Sugiyama, M., Suzuki, T., and Kanamori, T. *Density Ratio Estimation in ML*. Cambridge University Press, New York, NY, USA, 1st edition, 2012.
- Xiao, T., Xia, T., Yang, Y., Huang, C., and Wang, X. Learning from massive noisy labeled data for image classification. In *CVPR*, 2015.
- Yang, T., Mahdavi, M., Jin, R., Zhang, L., and Zhou, Y. Multiple kernel learning from noisy labels by stochastic programming. In *Proc. of 29th ICML*, pp. 233–240, New York, NY, USA, 2012. ACM.



# Appendix

## A 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, other 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 A1** When  $g$  is ideal, i.e.  $g(x) = g^*(x)$  and  $P$  and  $N$  have non-overlapping support, we have

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

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

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

### A.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

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

**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

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

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 A1, 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

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

## A.2 Proof of Theorem 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

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

## A.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

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

where

$$\begin{cases} LB_{y=1}^* = (1 - \rho_1)(1 - \pi_1) + \rho_0\pi_1 \\ UB_{y=0}^* = (1 - \rho_1)\pi_0 + \rho_0(1 - \pi_0) \\ \Delta p_o := \frac{|P \cap N|}{|P \cup N|} \\ \Delta P_1 = \{x \in P | g(x) < LB_{y=1}\} \\ \Delta N_1 = \{x \in N | g(x) \geq LB_{y=1}\} \\ \Delta P_0 = \{x \in P | g(x) \leq UB_{y=0}\} \\ \Delta N_0 = \{x \in N | g(x) > UB_{y=0}\} \end{cases} \quad (\text{A.5})$$

**Proof:** We first calculate  $LB_{y=1}$  and  $UB_{y=0}$  under unassuming conditions, 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

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

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,

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

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

Similarly, we have

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

## 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

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

Similarly, we have

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

In summary, Eq. (A.8) (A.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

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

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

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

where in the second equality we have used the definition of  $\Delta P_1$  and  $\Delta N_1$  in Eq. (A.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  $(s \perp\!\!\!\perp x) | y$ , 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. (A.11), we have

$$\hat{\rho}_1^{conf} = \rho_1 + \frac{1 - \rho_1 - \rho_0}{|P| - |\Delta P_1| + |\Delta N_1|} |\Delta N_1| \tag{A.12}$$

Similarly, we have

$$\hat{\rho}_0^{conf} = \rho_0 + \frac{1 - \rho_1 - \rho_0}{|N| - |\Delta N_0| + |\Delta P_0|} |\Delta P_0| \tag{A.13}$$

From the two equations above, we see that

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

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  $|\Delta P_0| = 0$ .

#### A.4 Proof of Theorem 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. (A.12) and (A.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  $\hat{\rho}_0^{conf} = \rho_0$ .

#### A.5 Proof of Theorem 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*

$$R_{\tilde{l}, \mathcal{D}_\rho}(f_\theta) = R_{l, \mathcal{D}}(f_\theta) \quad (\text{A.15})$$

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

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

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

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

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

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

where

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

satisfying

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

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

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

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

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

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. (A.17), and  $g(x)$ 's property in Eq. (A.21), we have

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

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. (A.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

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

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  $l$ .

## B Additional Figures

Figure B1 shows the average image for each digit for the 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 Rank Pruning makes few mistakes, and when it does, the mistakes vary greatly in image from the typical digit.

## C 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 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 experimental result tables for error, AUC-PR and F1 score metrics:

Table C1: Error for MNIST with logisitic regression as classifier.

Table C2: AUC-PR for MNIST with logisitic regression as classifier.

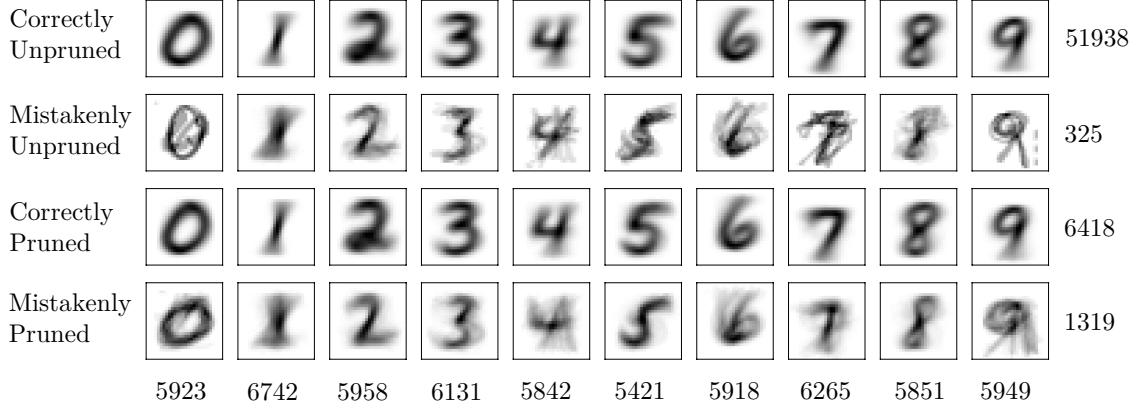


Figure B 1: Average image for each digit for the binary classification problem “1” or “not 1” in MNIST with logistic regression and significant mislabeling ( $\rho_1 = 0.5, \pi_1 = 0.5$ ). The right and bottom numbers count the total number of example images averaged in the corresponding row or column.

Table C3: Error for MNIST with CNN as classifier.

Table C4: AUC-PR for MNIST with CNN as classifier.

Table C5: F1 score for CIFAR-10 with logistic regression as classifier.

Table C6: Error for CIFAR-10 with logistic regression as classifier.

Table C7: AUC-PR for CIFAR-10 with logistic regression as classifier.

Table C8: Error for CIFAR-10 with CNN as classifier.

Table C9: AUC-PR for CIFAR-10 with CNN as classifier.

Due to 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.

## D Additional Related Work

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

### D.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.

### D.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 C 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 C 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 C 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	$\pi_1 = 0$					$\pi_1 = 0.25$					$\pi_1 = 0.5$				
	$RP_\rho$	RP	$\rho_1 = 0.5$	NAT13	ELK08	LIU16	$RP_\rho$	RP	$\rho_1 = 0.25$	NAT13	ELK08	LIU16	$RP_\rho$	RP	$\rho_1 = 0.5$
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
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>
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
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>
4	0.0012	<b>0.0037</b>	<b>0.0040</b>	<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
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	<b>0.0023</b>
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
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
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
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>
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

Table C 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	$\pi_1 = 0$ $\rho_1 = 0.5$					$\pi_1 = 0.25$ $\rho_1 = 0.25$					$\pi_1 = 0.5$										
	$RP_\rho$	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	<b>0.9982</b>	<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>	0.9964	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 C 5: Comparison of **F1 score** for one-vs-rest CIFAR-10 (averaged over all images) 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.

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	0.272	<b>0.311</b>	<b>0.252</b>	0.217	0.220	0.182	<b>0.329</b>	<b>0.275</b>	0.222	0.224	0.182	<b>0.330</b>	<b>0.265</b>	0.231	0.259	0.0	<b>0.266</b>	<b>0.188</b>	0.183	0.187	0.182
AUTO	0.374	<b>0.389</b>	<b>0.355</b>	0.318	0.320	0.182	<b>0.388</b>	<b>0.368</b>	0.321	0.328	0.182	<b>0.372</b>	<b>0.355</b>	0.308	0.341	0.0	<b>0.307</b>	0.287	0.287	<b>0.297</b>	0.182
BIRD	0.136	<b>0.241</b>	0.167	0.143	0.136	<b>0.182</b>	<b>0.248</b>	<b>0.185</b>	0.137	0.137	0.182	<b>0.258</b>	<b>0.147</b>	0.100	0.126	0.0	<b>0.206</b>	0.153	0.132	0.150	<b>0.182</b>
CAT	0.122	<b>0.246</b>	0.170	0.141	0.150	<b>0.182</b>	<b>0.232</b>	0.163	0.112	0.127	<b>0.182</b>	<b>0.241</b>	<b>0.125</b>	0.068	0.103	0.0	<b>0.209</b>	0.148	0.119	0.157	<b>0.182</b>
DEER	0.166	<b>0.250</b>	<b>0.184</b>	0.153	0.164	<b>0.182</b>	<b>0.259</b>	0.175	0.146	0.163	<b>0.182</b>	<b>0.259</b>	<b>0.177</b>	0.126	0.164	0.0	<b>0.222</b>	0.162	0.132	0.164	<b>0.182</b>
DOG	0.139	<b>0.245</b>	0.174	0.146	0.148	<b>0.182</b>	<b>0.262</b>	0.171	0.115	0.126	<b>0.182</b>	<b>0.254</b>	<b>0.152</b>	0.075	0.120	0.0	<b>0.203</b>	0.151	0.128	0.137	<b>0.182</b>
FROG	0.317	<b>0.322</b>	<b>0.315</b>	0.289	0.281	0.182	<b>0.350</b>	<b>0.319</b>	0.283	0.299	0.182	<b>0.346</b>	<b>0.305</b>	0.239	0.279	0.0	<b>0.308</b>	0.252	0.244	<b>0.269</b>	0.182
HORSE	0.300	<b>0.300</b>	<b>0.299</b>	0.283	0.263	0.182	<b>0.334</b>	<b>0.313</b>	0.272	0.281	0.182	<b>0.322</b>	<b>0.310</b>	0.260	0.292	0.0	<b>0.275</b>	<b>0.258</b>	0.240	0.245	0.182
SHIP	0.322	<b>0.343</b>	<b>0.322</b>	0.297	0.272	0.182	<b>0.385</b>	<b>0.319</b>	0.287	0.289	0.182	<b>0.350</b>	<b>0.303</b>	0.250	0.293	0.0	<b>0.304</b>	<b>0.248</b>	0.230	0.237	0.182
TRUCK	0.330	<b>0.359</b>	<b>0.323</b>	0.273	0.261	0.182	<b>0.369</b>	<b>0.327</b>	0.293	0.290	0.182	<b>0.343</b>	<b>0.302</b>	0.278	0.299	0.0	<b>0.313</b>	0.246	0.252	<b>0.262</b>	0.182
AVG	0.248	<b>0.301</b>	<b>0.256</b>	0.226	0.221	0.182	<b>0.316</b>	<b>0.262</b>	0.219	0.226	0.182	<b>0.308</b>	<b>0.244</b>	0.194	0.228	0.000	<b>0.261</b>	0.209	0.195	<b>0.210</b>	0.182

Table C 6: Comparison of **error** for one-vs-rest CIFAR-10 (averaged over all images) 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. Here the logistic regression classifier severely underfits CIFAR, resulting in Rank Pruning pruning out some correctly labeled examples that “confuse” the classifier, hence in this scenario, RP and  $RP_\rho$  generally have slightly smaller precision, much higher recall, and hence larger F1 scores than other models and even the ground truth classifier (Table C5). Due to the class imbalance ( $p_{y1} = 0.1$ ) and their larger recall, RP and  $RP_\rho$  here have larger error than the other 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	$RP_\rho$	RP	NAT13	ELK08	LIU16
PLANE	0.107	0.287	0.133	0.123	<b>0.122</b>	0.900	0.177	0.128	<b>0.119</b>	0.123	0.900	0.248	0.124	0.110	0.118	<b>0.100</b>	0.202	0.147	<b>0.142</b>	0.160	0.900					
AUTO	0.099	0.184	0.120	<b>0.110</b>	<b>0.110</b>	0.900	0.132	0.114	<b>0.105</b>	0.109	0.900	0.189	0.110	0.105	0.110	<b>0.100</b>	0.159	0.129	<b>0.125</b>	0.139	0.900					
BIRD	0.117	0.354	0.148	<b>0.133</b>	<b>0.131</b>	0.900	0.217	0.135	<b>0.120</b>	0.125	0.900	0.277	0.135	0.115	0.123	<b>0.100</b>	0.226	0.147	<b>0.139</b>	0.158	0.900					
CAT	0.114	0.351	0.138	<b>0.129</b>	<b>0.129</b>	0.900	0.208	0.139	<b>0.122</b>	0.125	0.900	0.303	0.132	0.114	0.122	<b>0.100</b>	0.225	0.151	<b>0.141</b>	0.158	0.900					
DEER	0.112	0.336	0.143	<b>0.128</b>	0.130	0.900	0.194	0.135	<b>0.120</b>	0.122	0.900	0.271	0.133	0.118	0.126	<b>0.100</b>	0.209	0.150	<b>0.147</b>	0.161	0.900					
DOG	0.119	0.370	0.150	<b>0.136</b>	0.138	0.900	0.205	0.142	<b>0.129</b>	0.132	0.900	0.288	0.135	0.120	0.128	<b>0.100</b>	0.229	0.154	<b>0.147</b>	0.168	0.900					
FROG	0.107	0.228	0.128	<b>0.117</b>	<b>0.117</b>	0.900	0.155	0.124	<b>0.113</b>	0.115	0.900	0.228	0.118	0.110	0.116	<b>0.100</b>	0.167	0.137	<b>0.130</b>	0.142	0.900					
HORSE	0.104	0.251	0.127	<b>0.114</b>	0.116	0.900	0.153	0.123	<b>0.110</b>	0.112	0.900	0.224	0.116	0.108	0.113	<b>0.100</b>	0.178	0.134	<b>0.129</b>	0.144	0.900					
SHIP	0.112	0.239	0.134	<b>0.121</b>	0.126	0.900	0.160	0.131	<b>0.119</b>	0.123	0.900	0.236	0.122	0.113	0.120	<b>0.100</b>	0.193	0.145	<b>0.139</b>	0.159	0.900					
TRUCK	0.106	0.210	0.130	<b>0.121</b>	0.122	0.900	0.145	0.125	<b>0.113</b>	0.117	0.900	0.213	0.121	0.108	0.117	<b>0.100</b>	0.165	0.142	<b>0.134</b>	0.150	0.900					
AVG	0.110	0.281	0.135	<b>0.123</b>	0.124	0.900	0.175	0.130	<b>0.117</b>	0.120	0.900	0.248	0.125	0.112	0.119	<b>0.100</b>	0.195	0.144	<b>0.137</b>	0.154	0.900					

Table C 7: Comparison of **AUC-PR** for one-vs-rest CIFAR-10 (averaged over all images) 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. Since  $p_{y1} = 0.1$ , here *Liu16* always predicts all labels as positive or negative, resulting in a constant AUC-PR of 0.550.

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	0.288	0.225	0.224	0.225	0.207	<b>0.550</b>	0.261	0.235	0.225	0.217	<b>0.550</b>	0.285	0.251	0.245	0.248	<b>0.550</b>	0.196	0.171	0.171	0.159	<b>0.550</b>
AUTO	0.384	0.350	0.317	0.312	0.316	<b>0.550</b>	0.342	0.335	0.331	0.331	<b>0.550</b>	0.328	0.348	0.334	0.333	<b>0.550</b>	0.256	0.257	0.259	0.261	<b>0.550</b>
BIRD	0.198	0.160	0.169	0.166	0.161	<b>0.550</b>	0.188	0.185	0.179	0.177	<b>0.550</b>	0.186	0.173	0.174	0.175	<b>0.550</b>	0.150	0.154	0.150	0.147	<b>0.550</b>
CAT	0.188	0.164	0.175	0.174	0.175	<b>0.550</b>	0.163	0.169	0.168	0.170	<b>0.550</b>	0.148	0.156	0.154	0.152	<b>0.550</b>	0.145	0.143	0.140	0.145	<b>0.550</b>
DEER	0.215	0.161	0.177	0.180	0.183	<b>0.550</b>	0.194	0.180	0.180	0.182	<b>0.550</b>	0.174	0.175	0.176	0.175	<b>0.550</b>	0.151	0.152	0.146	0.151	<b>0.550</b>
DOG	0.188	0.162	0.161	0.165	0.155	<b>0.550</b>	0.175	0.160	0.161	0.158	<b>0.550</b>	0.173	0.169	0.162	0.164	<b>0.550</b>	0.145	0.142	0.139	0.133	<b>0.550</b>
FROG	0.318	0.246	0.264	0.262	0.258	<b>0.550</b>	0.292	0.277	0.272	0.273	<b>0.550</b>	0.276	0.274	0.277	0.277	<b>0.550</b>	0.239	0.212	0.206	0.212	<b>0.550</b>
HORSE	0.319	0.242	0.267	0.269	0.260	<b>0.550</b>	0.283	0.264	0.264	0.263	<b>0.550</b>	0.288	0.282	0.279	0.278	<b>0.550</b>	0.223	0.218	0.208	0.207	<b>0.550</b>
SHIP	0.317	0.257	0.267	0.271	0.248	<b>0.550</b>	0.296	0.266	0.267	0.259	<b>0.550</b>	0.279	0.268	0.259	0.262	<b>0.550</b>	0.220	0.212	0.207	0.191	<b>0.550</b>
TRUCK	0.329	0.288	0.261	0.271	0.263	<b>0.550</b>	0.298	0.275	0.286	0.284	<b>0.550</b>	0.289	0.272	0.276	0.277	<b>0.550</b>	0.241	0.213	0.208	0.204	<b>0.550</b>
AVG	0.274	0.226	0.228	0.229	0.223	<b>0.550</b>	0.249	0.235	0.233	0.231	<b>0.550</b>	0.243	0.237	0.234	0.234	<b>0.550</b>	0.197	0.187	0.183	0.181	<b>0.550</b>

Table C 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 C 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>	<b>0.483</b>	<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