

Unified Risk Analysis for Weakly Supervised Learning

Chao-Kai Chiang
Masashi Sugiyama

*Department of Complexity Science and Engineering
Graduate School of Frontier Sciences
The University of Tokyo
5-1-5 Kashiwanoha, Kashiwa-shi, Chiba 277-8561, Japan*

*chaokai@k.u-tokyo.ac.jp
sugi@k.u-tokyo.ac.jp*

Reviewed on OpenReview: <https://openreview.net/forum?id=RGsdAwWuu6>

Abstract

Among the flourishing research of weakly supervised learning (WSL), we recognize the lack of a unified interpretation of the mechanism behind the weakly supervised scenarios, let alone a systematic treatment of the risk rewrite problem, a crucial step in the empirical risk minimization approach. In this paper, we introduce a framework providing a comprehensive understanding and a unified methodology for WSL. The formulation component of the framework, leveraging a contamination perspective, provides a unified interpretation of how weak supervision is formed and subsumes fifteen existing WSL settings. The induced reduction graphs offer comprehensive connections over WSLs. The analysis component of the framework, viewed as a decontamination process, provides a systematic method of conducting risk rewrite. In addition to the conventional inverse matrix approach, we devise a novel strategy called marginal chain aiming to decontaminate distributions. We justify the feasibility of the proposed framework by recovering existing rewrites reported in the literature.

1 Introduction

Accurate labels allow one to generalize to unseen data via empirical risk minimization (ERM) and analyze the generalization error in terms of the classification risk. In practice, there are various situations in which acquiring accurate labels is hard or even impossible. One obstacle preventing us from acquiring accurate labels is labeling restrictions, such as imperfect supervision due to imperceptibility, time constraints, annotation costs, and even data sensitivity. Another obstacle is the disruption by unavoidable noise from the environment.

To address the first obstacle of restrictions, various formulations have been studied under the notion of weakly supervised learning (WSL) (Zhou, 2018; Sugiyama et al., 2022). Based on various types of available label information, it evolves to thriving topics, including the conventional settings (Lu et al., 2019; 2020; 2021; Elkan & Noto, 2008; du Plessis et al., 2014; 2015; Niu et al., 2016; Kiryo et al., 2017; Sansone et al., 2019) that investigating the potential of unlabeled data, complementary-label learning (Ishida et al., 2017; 2019; Yu et al., 2018; Feng et al., 2020a; Katsura & Uchida, 2020; Chou et al., 2020), partial-label learning (Cour et al., 2011; Wang et al., 2019; Lv et al., 2020; Feng et al., 2020b; Wu et al., 2023), learning with confidence information (Ishida et al., 2018; Cao et al., 2021a;b; Berthon et al., 2021; Ishida et al., 2023), and learning with comparative information (Bao et al., 2018; Shimada et al., 2021; Feng et al., 2021; Cao et al., 2021b). Developing to resolve the second obstacle of noise, learning with noisy labels (LNL) can be categorized into two major formulations; one is called mutually contaminated distributions (MCD) (Scott et al., 2013; Menon et al., 2015; Katz-Samuels et al., 2019) in which class-conditional distributions contaminate each other, and the other is named class-conditional random label noise (CCN) (Natarajan et al., 2013; 2017) where a label is flipped by random noise.

Despite fruitful results and tremendous impact, we recognize a lack of global understanding and systematic treatment of WSL. From the perspective of *formulation*, there are only scattered links among WSLs. Lu et al. (2019) and Feng et al. (2021) showed that parameter substitution could reduce unlabeled-unlabeled to similar-unlabeled and positive-unlabeled settings. Figure 1 in Wu et al. (2023) showed relationships among four WSLs of partial- and complementary-labels. A similar observation can be found in the intersection of WSLs and LNLs. Several WSLs were shown to be special cases of the MCD model, and some other WSLs are special cases of the CCN model. For details, please refer to the discussions in Sections 8.2.3 and 9.2.4 of Sugiyama et al. (2022). These connections encourage us to consider the possibility that there exists a unique interpretation that explains the mechanism behind WSL. Luckily, from the *methodological* viewpoint, most of the existing WSL research adopted certain forms of the ERM approach. A crucial shared step is to perform the risk rewrite, a way of rephrasing the uncomputable risk to a computable one in terms of the data-generating distributions. A successful rewrite is the starting point of many downstream tasks, including but not limited to the following: Devising a practical or robust objective for training, comparing the strengths and properties of loss functions, proving the consistency, and analyzing generalization error bounds. However, many rewrite forms (summarized in Tables 4 and 5) look independent as if they are tailored to fit each problem’s unique form of supervision and are not adaptable to each other. These seemingly non-adaptable estimators post a practical challenge: When facing a new form of weak (or noisy) supervision, we do not have a guideline or general strategy to leverage developed methods to address the new situation.

These observations raise the following questions we aim to answer in this paper: What is the essence of WSL? From a formulation perspective, can a unique interpretation be found to explain the mechanism behind WSL? Does a methodology exist to address as many WSLs as possible?

This paper proposes a framework with the following contributions to answer the research questions.

1. To the best of our knowledge, the framework is the first systematic attempt to address how and why WSLs are related. The framework consists of a formulation component and an analysis component, and subsumes fifteen weakly supervised scenarios. Table 10 summarizes the results obtained from our framework. This paper brings forth the next two new insights.
2. The formulation component, modeling from a *contamination* perspective, provides a coherent interpretation of the weakly supervised data-generating processes. It produces a comprehensive relationship graph, shown in Figure 1, consisting of Tables 7, 8, and 9. Figure 1 summarizes the WSLs and reveals connections between scenarios that were previously unknown to the community. Figure 1 is our answer to the second research question. Figure 1 also unveils a distinctive confidence-based type of WSLs that do not belong to the prominent MCD or CCN categories.
3. The analysis component, leveraging the *decontamination* concept, establishes a generic methodology for conducting risk rewrites for all WSLs discussed in this paper. Thus, rewrite derivations that previously seemed irrelevant can now be systematically analyzed, and the final research question is answered. In addition, the analysis component also applies to new scenarios.
4. Regarding the technical contributions, the proposed framework distinguishes two approaches, the inversion method and the marginal chain method presented in Theorems 1 and 2, to implement the decontamination idea. The marginal chain method is a newly developed invertibility-free loss correction approach. We also illustrate the subtle adjustments to develop simplified and intuitive proofs for the existing risk rewrites. These alternative proofs have their own logic stemming from the proposed framework.

The idea of decontamination has been widely implemented and investigated. There are two major approaches, loss correction, and label correction, in LNL. Closest to the current paper, Cid-Sueiro (2012), van Rooyen & Williamson (2017), Katz-Samuels et al. (2019), Patrini et al. (2017), and van Rooyen & Williamson (2015) exploited the inverse matrix, sometimes known as the backward method (Patrini et al., 2017), to construct a corrected training loss to obtain an unbiased estimator. There were deep learning methods leveraging the contamination assumption, sometimes called the forward method (Patrini et al., 2017), to train a classifier (Patrini et al., 2017; Yu et al., 2018; Sukhbaatar & Fergus, 2015; Goldberger & Ben-Reuven, 2017; Berthon

et al., 2021). Besides modifying the loss function, one has two other strategies to manipulate the corrupted labels. The (iterative) pseudo-label method modified the labels for training (Ma et al., 2018; Tanaka et al., 2018; Reed et al., 2015). Filtering clean data points for training is the other option (Northcutt et al., 2017; 2021; Jiang et al., 2018; Han et al., 2018; Yu et al., 2019). Apart from classification, a different research branch studies conditions and methods for recovering the base distributions (Katz-Samuels et al., 2019; Blanchard & Scott, 2014; Blanchard et al., 2016).

The current work is close to the loss correction approach in LNL. Most previous loss correction methods exploited invertibility to construct the corrected losses. In contrast, the marginal chain approach we propose in this paper adopts the conditional probability formula to build the corrected losses. Many of the existing work targeted either the MCD or the CCN models. Scott & Zhang (2020), Berthon et al. (2021), Patrini et al. (2017), Goldberger & Ben-Reuven (2017), Sukhbaatar & Fergus (2015), Yu et al. (2018), Natarajan et al. (2013), Natarajan et al. (2017), Northcutt et al. (2017), and Northcutt et al. (2021) were based on the CCN model, and Katz-Samuels et al. (2019), Blanchard & Scott (2014), and Blanchard et al. (2016) were based on the MCD model. Menon et al. (2015), van Rooyen & Williamson (2017), and Katz-Samuels et al. (2019) studied multiple noise models at the same time. However, the current paper investigates the connections between MCD, CCN, and confidence-based settings simultaneously through the lens of matrix decontamination as broadly as possible to identify a generic methodology for WSLs. Different from the current paper aiming for risk minimization, research also studied various performance measures, such as the balanced error rate (Scott & Zhang, 2020; 2019; Menon et al., 2015; du Plessis et al., 2013), the area under the receiver operating characteristic curve (Charoenphakdee et al., 2019; Sakai et al., 2018; Menon et al., 2015), and cost-sensitive measures (Charoenphakdee et al., 2021; Natarajan et al., 2017). We choose the classification risk as the only measure due to the focus of this paper.

The remaining sections are organized as follows. Section 2 reviews ERM in supervised learning, the risk rewrite problem, and the existing results. Section 3 presents the proposed framework. We show that the proposed framework provides a unified way to formulate diverse weakly supervised scenarios in Section 4. Section 5 demonstrates how to instantiate the framework to conduct risk rewrite. We demonstrate the applicability of the proposed framework to new scenarios in Section 6. Finally, we conclude the paper and discuss outlooks in Section 7. The current organization of Sections 4 and 5 aims at connecting multiple WSLs under one framework. We note that this paper can serve multiple purposes for the study of WSLs. A summary of possible use cases of the paper is provided in Appendix A.

2 Preliminaries

Let (y, x) be a data example where the instance $x \in X$ and the label $y \in Y$. For binary classification, the label space Y is $\{p, n\}$, and for multiclass classification with K classes, $Y = \{1, 2, \dots, K\} := [K]$. The joint distribution is $P(Y, X)$, the class prior is $P(Y)$, the class-conditional distribution is $P(X|Y)$, and the class probability function is $P(Y|X)$. Given a space of hypotheses G , we denote the loss of a hypothesis $g \in G$ on predicting y of (y, x) as $\ell_{Y=y}(g(x))$. To accommodate concise expressions and readability for all WSLs considered in this paper simultaneously, we use alias notations when the context is unambiguous. Table 1 provides a set of common notations used in this paper.

We use (y, x) instead of the convention (x, y) to represent a data example because, in the current paper, we focus on discussing different types of supervision. Placing the label before the instance emphasizes the type of supervision under investigation in theorems and derivations.

2.1 Supervised Learning and the ERM Method

In supervised learning with K classes, the observed data is of the form

$$\{x_i^y\}_{i=1}^{n_y} \stackrel{\text{i.i.d.}}{P_{X|Y=y}}, \quad y \in [K].$$

Table 1: Alias of Common Notations.

Name of the notation	Expression	Aliases
Binary classes	$\{\mathfrak{p}, \mathfrak{n}\}$	
Multiple classes	$\{1, \dots, K\}$	$[K]$
Compound set of $[K]$	$2^{[K]} \setminus \{\cdot, [K]\}$	S
Joint distribution	$\mathbb{P}(Y = y, X = x)$	$P_{Y=y,x}, P_{Y=y,X},$ or $P_{Y,X}$
Hypothesis and its space	$g \in G$	
Loss of g	$\ell_{Y=y}(g(x))$	$\ell_y, \ell_y(X),$ or $\ell_Y(g(X))$
Classification risk	$\mathbb{E}_{Y,X} [\ell_Y(g(X))]$	$R(g)$
The j -th entry of vector V	$(V)_j$	V_j
Class prior	$\mathbb{P}(Y = y)$	π_y
Marginal	$\mathbb{P}(X)$	P_X
Class-conditional	$\mathbb{P}(X = x Y = y)$	$P_{X Y}, P_{X Y=y},$ or $P_{x Y=y}$
Confidence	$\mathbb{P}(Y = y X = x)$	$r_y(X), r_y(x),$ or $r(X)$ if $y = \mathfrak{p}$

Notation x_i^y denotes the shorthand of (y, x_i) . The goal of learning is to find a classifier $g \in G$ that minimizes the classification risk

$$R(g) := \mathbb{E}_{Y,X} [\ell_Y(g(X))] = \sum_{y=1}^K \int_{\mathcal{X}} P_{Y=y,x} \ell_{Y=y}(g(x)) dx. \quad (1)$$

To find such a classifier, ERM first constructs an empirical risk estimator with the data in hand:

$$\hat{R}(g) = \sum_{y=1}^K \frac{1}{n_y} \sum_{i=1}^{n_y} \pi_y \ell_{Y=y}(g(x_i^y)). \quad (2)$$

The estimator approximates $R(g)$ consistently since it can be shown that (2) approaches (1) as N (Tewari & Bartlett, 2014; Kiryo et al., 2017) and (Sugiyama et al., 2022, Chapter 3). Then, ERM takes $\hat{R}(g)$ as the training objective and optimizes it to find the optimal classifier

$$g^* = \arg \min_{g \in G} \hat{R}(g) \quad (3)$$

in the hypothesis space G as the output of ERM.

2.2 The Risk Rewrite Problem and Existing Results

Sections 2.2.1 to 2.2.17 review the learning scenarios including WSLs, MCD, and CCN that will be discussed in this paper. A knowledgeable reader may refer directly to summary Tables 2 through 6 and proceed to Section 3.

In every WSL scenario, the goal of learning is the same as supervised learning. However, the observed data is no longer as perfectly labeled as in supervised learning. That said, there are differences in the formulations of the observed data and the ways of estimating the classification risk. We begin with reviewing WSLs derived from binary classes. For $K = 2$, we denote $\mathcal{Y} := \{\mathfrak{p}, \mathfrak{n}\}$.

2.2.1 Positive-Unlabeled (PU) learning

The observed data in PU learning (du Plessis et al., 2015) is of the form

$$\begin{aligned} \{x_i^{\mathfrak{p}}\}_{i=1}^{n_{\mathfrak{p}}} &\stackrel{\text{i.i.d.}}{\sim} P_{\mathfrak{p}} := P_{X|Y=\mathfrak{p}}, \\ \{x_j^{\mathfrak{u}}\}_{j=1}^{n_{\mathfrak{u}}} &\stackrel{\text{i.i.d.}}{\sim} P_{\mathfrak{u}} := \pi_{\mathfrak{p}} P_{X|Y=\mathfrak{p}} + \pi_{\mathfrak{n}} P_{X|Y=\mathfrak{n}}, \end{aligned} \quad (4)$$

where x_j^u is viewed as the shorthand of (u, x_j) symbolizing the unlabeled data¹. The unlabeled data set $\{x_j^u\}_j$ consists of a mixture of samples from $P_{X|Y=p}$ and $P_{X|Y=n}$ with proportion π_p . Since the information of negatively sampled data is unavailable, (2) is uncomputable, causing directly optimizing (3) infeasibility. Therefore, to make ERM applicable, the *risk rewrite problem* (Sugiyama et al., 2022) asks:

Can one rephrase the classification risk $R(g)$ (1) in terms of the given data formulation?

du Plessis et al. (2015) rewrote the classification risk in terms of the data-generating distributions P_P and P_U as

$$R(g) = \mathbb{E}_P [\pi_p \ell_p - \pi_p \ell_n] + \mathbb{E}_U [\ell_n]. \quad (5)$$

2.2.2 Positive-confidence (Pconf) Learning Learning

The observed data in Pconf learning (Ishida et al., 2018) is of the form

$$\{x_i, r(x_i)\}_{i=1}^n,$$

where

$$\begin{aligned} x_i &\stackrel{\text{i.i.d.}}{P_P} := P_{X|Y=p}, \\ r(x_i) &:= P_{Y=p|X=x_i}. \end{aligned} \quad (6)$$

The function $r(x)$ represents how confident an example x would be positively labeled. Ishida et al. (2018) rewrote the classification risk as

$$R(g) = \pi_p \mathbb{E}_P \left[\ell_p + \frac{1 - r(X)}{r(X)} \ell_n \right]. \quad (7)$$

2.2.3 Unlabeled-Unlabeled (UU) learning

The observed data in UU learning (Lu et al., 2019) is of the form

$$\begin{aligned} \{x_i^{u_1}\}_{i=1}^{n_{u_1}} &\stackrel{\text{i.i.d.}}{P_{U_1}} := (1 - \gamma_1) P_{X|Y=p} + \gamma_1 P_{X|Y=n}, \\ \{x_j^{u_2}\}_{j=1}^{n_{u_2}} &\stackrel{\text{i.i.d.}}{P_{U_2}} := \gamma_2 P_{X|Y=p} + (1 - \gamma_2) P_{X|Y=n}, \end{aligned} \quad (8)$$

where $x_i^{u_1}$ (resp. $x_j^{u_2}$) being the shorthand of (u_1, x_i) (resp. (u_2, x_j)) represents x_i (resp. x_j) belonging to the unlabeled data whose mixture parameter is γ_1 (resp. γ_2). Notice a difference that the mixture proportion of the unlabeled data in PU learning is π_p . Lu et al. (2019) rewrote the classification risk in terms of the data-generating distributions P_{U_1} and P_{U_2} as follows: Assume $\gamma_1 + \gamma_2 = 1$. Then,

$$R(g) = \mathbb{E}_{U_1} \left[\frac{(1 - \gamma_2)\pi_p}{1 - \gamma_1 - \gamma_2} \ell_p + \frac{-\gamma_2\pi_n}{1 - \gamma_1 - \gamma_2} \ell_n \right] + \mathbb{E}_{U_2} \left[\frac{-\gamma_1\pi_p}{1 - \gamma_1 - \gamma_2} \ell_p + \frac{(1 - \gamma_1)\pi_n}{1 - \gamma_1 - \gamma_2} \ell_n \right]. \quad (9)$$

2.2.4 Similar-Unlabeled (SU) learning

The observed data in SU learning (Bao et al., 2018) is of the form

$$\begin{aligned} \left\{ \left(x_i^s, x_i^s \right) \right\}_{i=1}^{n_s} &\stackrel{\text{i.i.d.}}{P_S} := \frac{\pi_p^2 P_{X|Y=p} P_{X|Y=p} + \pi_n^2 P_{X|Y=n} P_{X|Y=n}}{\pi_p^2 + \pi_n^2}, \\ \{x_j^u\}_{j=1}^{n_u} &\stackrel{\text{i.i.d.}}{P_U} := \pi_p P_{X|Y=p} + \pi_n P_{X|Y=n}. \end{aligned} \quad (10)$$

The word ‘‘similar’’ means the examples in every (x^s, x^s) pair have the same label; either both are positive, or both are negative. Under the assumption $\pi_p = \pi_n$, Bao et al. (2018) rewrote the classification risk as

$$R(g) = (\pi_p^2 + \pi_n^2) \mathbb{E}_S \left[\frac{L(X) + L(X)}{2} \right] + \mathbb{E}_U [L_-(X)], \quad (11)$$

¹Seemingly being redundant, but it is helpful to use (u, x_j) to distinguish it from the positively labeled instance (p, x_i) .

where

$$\begin{aligned} L(X) &:= \frac{1}{\pi_p - \pi_n} \ell_p(X) - \frac{1}{\pi_p - \pi_n} \ell_n(X), \\ L_-(X) &:= -\frac{\pi_n}{\pi_p - \pi_n} \ell_p(X) + \frac{\pi_p}{\pi_p - \pi_n} \ell_n(X). \end{aligned}$$

2.2.5 Dissimilar-Unlabeled (DU) learning

The observed data in DU learning (Shimada et al., 2021) is of the form

$$\begin{aligned} \left\{ \left(x_i^d, x_i^d \right) \right\}_{i=1}^{n_d} \text{ i.i.d. } P_D &:= \frac{P_{X/Y=p} P_{X/Y=n} + P_{X/Y=n} P_{X/Y=p}}{2}, \\ \left\{ x_j^u \right\}_{j=1}^{n_u} \text{ i.i.d. } P_U &:= \pi_p P_{X/Y=p} + \pi_n P_{X/Y=n}. \end{aligned} \quad (12)$$

The word ‘‘dissimilar’’ means the examples in every (x^d, x^d) pair have distinct labels. Under the assumption $\pi_p = \pi_n$, Shimada et al. (2021) rewrote the classification risk as

$$R(g) = 2\pi_p\pi_n \mathbb{E}_D \left[-\frac{L(X) + L(X)}{2} \right] + \mathbb{E}_U [L_+(X)], \quad (13)$$

where

$$\begin{aligned} L(X) &= \frac{1}{\pi_p - \pi_n} \ell_p(X) - \frac{1}{\pi_p - \pi_n} \ell_n(X), \\ L_+(X) &:= \frac{\pi_p}{\pi_p - \pi_n} \ell_p(X) - \frac{\pi_n}{\pi_p - \pi_n} \ell_n(X). \end{aligned}$$

Note that $L(X)$ has been defined in the SU setting. We repeat it here for clarity.

2.2.6 Similar-Dissimilar (SD) learning

The observed data in SD learning (Shimada et al., 2021) is of the form

$$\begin{aligned} \left\{ \left(x_i^s, x_i^s \right) \right\}_{i=1}^{n_s} \text{ i.i.d. } P_S &:= \frac{\pi_p^2 P_{X/Y=p} P_{X/Y=p} + \pi_n^2 P_{X/Y=n} P_{X/Y=n}}{\pi_p^2 + \pi_n^2}, \\ \left\{ \left(x_i^d, x_i^d \right) \right\}_{i=1}^{n_d} \text{ i.i.d. } P_D &:= \frac{P_{X/Y=p} P_{X/Y=n} + P_{X/Y=n} P_{X/Y=p}}{2}. \end{aligned} \quad (14)$$

Under the assumption $\pi_p = \pi_n$, Shimada et al. (2021) rewrote the classification risk as

$$R(g) = (\pi_p^2 + \pi_n^2) \mathbb{E}_S \left[\frac{L_+(X) + L_+(X)}{2} \right] + 2\pi_p\pi_n \mathbb{E}_D \left[\frac{L_-(X) + L_-(X)}{2} \right], \quad (15)$$

where

$$\begin{aligned} L_+(X) &= \frac{\pi_p}{\pi_p - \pi_n} \ell_p(X) - \frac{\pi_n}{\pi_p - \pi_n} \ell_n(X), \\ L_-(X) &= -\frac{\pi_n}{\pi_p - \pi_n} \ell_p(X) + \frac{\pi_p}{\pi_p - \pi_n} \ell_n(X). \end{aligned}$$

Note that $L_+(X)$ and $L_-(X)$ have been defined in the DU and SU settings. We repeat them here for clarity.

2.2.7 Pairwise Comparison (Pcomp) Learning

The observed data in Pcomp learning (Feng et al., 2021) is of the form

$$\left\{ \left(x_i^{\text{pc}}, x_i^{\text{pc}} \right) \right\}_{i=1}^{n_{\text{pc}}} \text{ i.i.d. } P_{\text{PC}} := \frac{\pi_p^2 P_{X/Y=p} P_{X/Y=p} + \pi_p\pi_n P_{X/Y=p} P_{X/Y=n} + \pi_n^2 P_{X/Y=n} P_{X/Y=n}}{\pi_p^2 + \pi_p\pi_n + \pi_n^2}. \quad (16)$$

The pairwise comparison encodes a meaning that each x^{pc} “can not be more negative” than x^{pc} in the $(x^{\text{pc}}, x^{\text{pc}})$ pair. That is, the labels in $(x^{\text{pc}}, x^{\text{pc}})$ are of the form (p, p) , (p, n) , or (n, n) . Feng et al. (2021) rewrote the classification risk as

$$R(g) = \mathbb{E}_{\text{Sup}} [\ell_{\text{p}} - \pi_{\text{p}} \ell_{\text{n}}] + \mathbb{E}_{\text{Inf}} [-\pi_{\text{n}} \ell_{\text{p}} + \ell_{\text{n}}], \quad (17)$$

where the expectations are computed over the following distributions

$$\begin{aligned} P_{\text{Sup}} &:= \int_{\mathcal{X}} P_{\text{PC}} dx, \\ P_{\text{Inf}} &:= \int_{\mathcal{X}} P_{\text{PC}} dx. \end{aligned}$$

2.2.8 Similarity-Confidence Learning (Sconf) Learning

The observed data in Sconf learning (Cao et al., 2021b) is of the form

$$\left\{ x_i^{\text{sc}}, x_i^{\text{sc}}, r(x_i^{\text{sc}}, x_i^{\text{sc}}) \right\}_{i=1}^n,$$

where

$$\begin{aligned} x_i^{\text{sc}} \text{ i.i.d. } P_X &:= \pi_{\text{p}} P_{X|Y=\text{p}} + \pi_{\text{n}} P_{X|Y=\text{n}}, \\ x_i^{\text{sc}} \text{ i.i.d. } P_X &:= \pi_{\text{p}} P_{X|Y=\text{p}} + \pi_{\text{n}} P_{X|Y=\text{n}}, \\ r(x_i^{\text{sc}}, x_i^{\text{sc}}) &:= P_{Y=y_i^{\text{sc}}=Y=y_i^{\text{sc}} | X=x_i^{\text{sc}}, X=x_i^{\text{sc}}}. \end{aligned} \quad (18)$$

Cao et al. (2021b) rewrote the classification risk as

$$R(g) = \mathbb{E}_{X, X} \left[\frac{r(X, X) - \pi_{\text{n}}}{\pi_{\text{p}} - \pi_{\text{n}}} L_{\text{p}}(X, X) + \frac{\pi_{\text{p}} - r(X, X)}{\pi_{\text{p}} - \pi_{\text{n}}} L_{\text{n}}(X, X) \right], \quad (19)$$

where

$$\begin{aligned} L_{\text{p}}(X, X) &:= \frac{\ell_{\text{p}}(X) + \ell_{\text{p}}(X)}{2}, \\ L_{\text{n}}(X, X) &:= \frac{\ell_{\text{n}}(X) + \ell_{\text{n}}(X)}{2}. \end{aligned}$$

2.2.9 Complementary-Label (CL) Learning

One can also formulate weak supervision from multiclass classification. For K classes, we denote $\mathcal{Y} := [K]$.

The observed data in CL learning (Ishida et al., 2019) is of the form

$$\{(\bar{s}_i, x_i)\}_{i=1}^n \text{ i.i.d. } P_{\bar{S}, X} := \frac{1}{K-1} \sum_{Y=\bar{S}} P_{Y, X}. \quad (20)$$

As is named “complementary,” $\bar{s} \in [K]$ represents that the true label y of x cannot be \bar{s} . Ishida et al. (2019) rewrote the classification risk as

$$R(g) = \mathbb{E}_{\bar{S}, X} \left[\sum_{y=1}^K \ell_y - (K-1) \ell_{\bar{S}} \right]. \quad (21)$$

2.2.10 Multi-Complementary-Label (MCL) Learning

The observed data in MCL learning (Feng et al., 2020a) is of the form

$$\{(s_i, x_i)\}_{i=1}^n \text{ i.i.d. } P_{\bar{S}, X} := \begin{cases} \sum_{d=1}^{K-1} P_{|\bar{S}|=d} \cdot \frac{1}{\binom{K-1}{|\bar{S}|}} \sum_{Y/\bar{S}} P_{Y, X}, & \text{if } |\bar{S}| = d, \\ 0, & \text{otherwise.} \end{cases} \quad (22)$$

Generalized from CL, $\bar{s} \subseteq [K]$ in MCL is a set of classes of size $d \in [K-1]$, representing multiple exclusions. In other words, CL is the special case of MCL with $d = 1$. Feng et al. (2020a) rewrote the classification risk as

$$R(g) = \sum_{d=1}^{K-1} P_{|\bar{S}|=d} \mathbb{E}_{\bar{S}, X \mid |\bar{S}|=d} \left[\sum_{y/\bar{S}} \ell_y - \frac{K-1-|\bar{S}|}{|\bar{S}|} \sum_{\bar{s} \subseteq \bar{S}} \ell_{\bar{s}} \right]. \quad (23)$$

2.2.11 Provably Consistent Partial-Label (PCPL) Learning

The observed data in PCPL learning (Feng et al., 2020b) is of the form

$$\{(s_i, x_i)\}_{i=1}^n \text{ i.i.d. } P_{S, X} := \frac{1}{2^{K-1} - 1} \sum_{Y \subseteq S} P_{Y, X}. \quad (24)$$

A partial-label $s \subseteq [K]$ is a set of classes containing the true label y of x . Feng et al. (2020b) rewrote the classification risk as

$$R(g) = \frac{1}{2} \mathbb{E}_{S, X} \left[\sum_{y=1}^K \frac{P_{Y=y/X}}{\sum_{a \subseteq S} P_{Y=a/X}} \ell_y \right]. \quad (25)$$

2.2.12 Proper Partial-Label (PPL) Learning

The observed data in PPL learning (Wu et al., 2023) is of the form

$$\{(s_i, x_i)\}_{i=1}^n \text{ i.i.d. } P_{S, X} := C(S, X) \sum_{Y \subseteq S} P_{Y, X}. \quad (26)$$

The weight $\frac{1}{2^{K-1} - 1}$ in PCPL is generalized to $C(S, X)$, a function of the partial-label and the instance, allowing one to characterize the ‘‘properness’’ of a partial-label. Wu et al. (2023) rewrote the classification risk as

$$R(g) = \mathbb{E}_{S, X} \left[\sum_{y \subseteq S} \frac{P_{Y=y/X}}{\sum_{a \subseteq S} P_{Y=a/X}} \ell_y \right]. \quad (27)$$

2.2.13 Single-Class Confidence (SC-Conf) Learning

The observed data in SC-Conf learning (Cao et al., 2021a) is of the form

$$\{x_i, r_1(x_i), \dots, r_K(x_i)\}_{i=1}^n,$$

where

$$\begin{aligned} x_i & \text{ i.i.d. } P_{X|Y=y_s} \text{ with } y_s \subseteq [K], \\ r_k(x_i) & := P_{Y=k|X=x_i} \text{ for each } k \subseteq [K]. \end{aligned} \quad (28)$$

The constraint of SC-Conf is that the examples are sampled from a specific class y_s . The key to risk rewrite is the availability of confident information $r_k(x)$ about each class. Cao et al. (2021a) rewrote the classification risk as

$$R(g) = \pi_{y_s} \mathbb{E}_{X|Y=y_s} \left[\sum_{y=1}^K \frac{r_y(X)}{r_{y_s}(X)} \ell_y \right]. \quad (29)$$

2.2.14 Subset Confidence (Sub-Conf) Learning

The observed data in Sub-Conf learning (Cao et al., 2021a) is of the form

$$\{x_i, r_1(x_i), \dots, r_K(x_i)\}_{i=1}^n,$$

where

$$\begin{aligned} x_i & \text{ i.i.d. } P_{X/Y} \quad Y_s \in [K], \\ r_k(x_i) & := P_{Y=k|X=x_i} \text{ for each } k \in [K]. \end{aligned} \quad (30)$$

Sub-Conf is a relaxed setting of SC-Conf where the samples come from a set of classes Y_s . Cao et al. (2021a) rewrote the classification risk as

$$R(g) = \pi_{Y_s} \mathbb{E}_{X/Y} \left[\sum_{y=1}^K \frac{r_y(X)}{r_{Y_s}(X)} \ell_y \right], \quad (31)$$

where $\pi_{Y_s} := \sum_{j \in Y_s} \pi_j$, and $r_{Y_s}(X) := P_{Y \in Y_s|X} = \sum_{j \in Y_s} P_{Y=j|X}$.

2.2.15 Soft-Label Learning

Ishida et al. (2023) formulated soft-label learning under the binary setting, in which the observed data is of the form

$$\{x_i, r(x_i)\}_{i=1}^n,$$

where

$$\begin{aligned} x_i & \text{ i.i.d. } P_X := P_{Y=p,X} + P_{Y=n,X}, \\ r(x_i) & := P_{Y=p|X=x_i}. \end{aligned} \quad (32)$$

It is straightforward to obtain a corresponding formulation under the multiclass setting:

$$\{x_i, r_1(x_i), \dots, r_K(x_i)\}_{i=1}^n,$$

where

$$\begin{aligned} x_i & \text{ i.i.d. } P_X := \sum_{k=1}^K P_{Y=k,X}, \\ r_k(x_i) & := P_{Y=k|X=x_i} \text{ for each } k \in [K]. \end{aligned} \quad (33)$$

The difference between SC-Conf and multiclass soft-label (resp. the difference between Pconf and binary soft-label) is the sample distribution of x . We rewrote the classification risk as

$$R(g) = \mathbb{E}_X \left[\sum_{y=1}^K r_y(X) \ell_y \right]. \quad (34)$$

2.2.16 Summary of Existing WSL Formulations and Risk Rewrites

We summarize the weakly supervised scenarios discussed and their risk rewrite results. The formulations are divided into the binary classification settings in Table 2 and the multiclass classification settings in Table 3. We list the formulations in chronological order, according to their publication order. Tables 4 and 5 are the corresponding rewrites.

Table 2: Binary WSL formulations.

WSL	Formulation	Equation
PU	$\{x_i^p\}_{i=1}^{n_p}$ i.i.d. $P_P := P_{X Y=p}$, $\{x_j^u\}_{j=1}^{n_u}$ i.i.d. $P_U := \pi_p P_{X Y=p} + \pi_n P_{X Y=n}$.	(4)
Pconf	$\{x_i, r(x_i)\}_{i=1}^n$, where x_i i.i.d. $P_P := P_{X Y=p}$, $r(x_i) := P_{Y=p X=x_i}$.	(6)
UU	$\{x_i^{u_1}\}_{i=1}^{n_{u_1}}$ i.i.d. $P_{U_1} := (1 - \gamma_1) P_{X Y=p} + \gamma_1 P_{X Y=n}$, $\{x_j^{u_2}\}_{j=1}^{n_{u_2}}$ i.i.d. $P_{U_2} := \gamma_2 P_{X Y=p} + (1 - \gamma_2) P_{X Y=n}$.	(8)
SU	$\left\{ \left(x_i^s, x_i^s \right) \right\}_{i=1}^{n_s}$ i.i.d. $P_S := \frac{\pi_p^2 P_{X Y=p} P_{X Y=p} + \pi_n^2 P_{X Y=n} P_{X Y=n}}{\pi_p^2 + \pi_n^2}$, $\{x_j^u\}_{j=1}^{n_u}$ i.i.d. $P_U := \pi_p P_{X Y=p} + \pi_n P_{X Y=n}$.	(10)
DU	$\left\{ \left(x_i^d, x_i^d \right) \right\}_{i=1}^{n_d}$ i.i.d. $P_D := \frac{P_{X Y=p} P_{X Y=n} + P_{X Y=n} P_{X Y=p}}{2}$, $\{x_j^u\}_{j=1}^{n_u}$ i.i.d. $P_U := \pi_p P_{X Y=p} + \pi_n P_{X Y=n}$.	(12)
SD	$\left\{ \left(x_i^s, x_i^s \right) \right\}_{i=1}^{n_s}$ i.i.d. $P_S := \frac{\pi_p^2 P_{X Y=p} P_{X Y=p} + \pi_n^2 P_{X Y=n} P_{X Y=n}}{\pi_p^2 + \pi_n^2}$, $\left\{ \left(x_i^d, x_i^d \right) \right\}_{i=1}^{n_d}$ i.i.d. $P_D := \frac{P_{X Y=p} P_{X Y=n} + P_{X Y=n} P_{X Y=p}}{2}$.	(14)
Pcomp	$\left\{ \left(x_i^{pc}, x_i^{pc} \right) \right\}_{i=1}^{n_{pc}}$ i.i.d. P_{PC} $:= \frac{\pi_p^2 P_{X Y=p} P_{X Y=p} + \pi_p \pi_n P_{X Y=p} P_{X Y=n} + \pi_n^2 P_{X Y=n} P_{X Y=n}}{\pi_p^2 + \pi_p \pi_n + \pi_n^2}$.	(16)
Sconf	$\left\{ x_i^{sc}, x_i^{sc}, r \left(x_i^{sc}, x_i^{sc} \right) \right\}_{i=1}^{n_{sc}}$, where x_i^{sc} i.i.d. $P_X := \pi_p P_{X Y=p} + \pi_n P_{X Y=n}$, x_i^{sc} i.i.d. $P_X := \pi_p P_{X Y=p} + \pi_n P_{X Y=n}$, $r \left(x_i^{sc}, x_i^{sc} \right) := P_{Y=y_i^{sc}=Y=y_i^{sc} X=x_i^{sc}, X=x_i^{sc}}$.	(18)

Table 3: Multiclass WSL formulations.

WSL	Formulation	Equation
CL	$\{(\bar{s}_i, x_i)\}_{i=1}^n$ i.i.d. $P_{\bar{S},X} := \frac{1}{K-1} \sum_{Y=\bar{S}} P_{Y,X}$.	(20)
MCL	$\{(\bar{s}_i, x_i)\}_{i=1}^n$ i.i.d. $P_{\bar{S},X} := \begin{cases} \sum_{d=1}^{K-1} P_{ \bar{S} =d} \cdot \frac{1}{\binom{K-1}{ \bar{S} }} \sum_{Y/\bar{S}} P_{Y,X}, & \text{if } \bar{S} = d, \\ 0, & \text{otherwise.} \end{cases}$	(22)
PCPL	$\{(s_i, x_i)\}_{i=1}^n$ i.i.d. $P_{S,X} := \frac{1}{2^{K-1} - 1} \sum_{Y \subseteq S} P_{Y,X}$.	(24)
PPL	$\{(s_i, x_i)\}_{i=1}^n$ i.i.d. $P_{S,X} := C(S, X) \sum_{Y \subseteq S} P_{Y,X}$.	(26)
SC-Conf	$\{x_i, r_1(x_i), \dots, r_K(x_i)\}_{i=1}^n$, where x_i i.i.d. $P_{X/Y=y_s}$ with $y_s \in [K]$, $r_k(x_i) := P_{Y=k/X=x_i}$ for each $k \in [K]$.	(28)
Sub-Conf	$\{x_i, r_1(x_i), \dots, r_K(x_i)\}_{i=1}^n$, where x_i i.i.d. $P_{X/Y \in Y_s}$ with $Y_s \in [K]$, $r_k(x_i) := P_{Y=k/X=x_i}$ for each $k \in [K]$.	(30)
Soft-label	$\{x_i, r_1(x_i), \dots, r_K(x_i)\}_{i=1}^n$, where x_i i.i.d. P_X , $r_k(x_i) := P_{Y=k/X=x_i}$ for each $k \in [K]$.	(33)

Table 4: Risk rewrites for binary WSLs.

WSL	Risk rewrite for $R(g) = \mathbb{E}_{Y,X} [\ell_Y(g(X))]$ (1)	Equation
PU	$R(g) = \mathbb{E}_P [\pi_p \ell_p - \pi_n \ell_n] + \mathbb{E}_U [\ell_n]$.	(5)
Pconf	$R(g) = \pi_p \mathbb{E}_P \left[\ell_p + \frac{1 - r(X)}{r(X)} \ell_n \right]$.	(7)
UU	$R(g) = \mathbb{E}_{U_1} \left[\frac{(1 - \gamma_2)\pi_p}{1 - \gamma_1 - \gamma_2} \ell_p + \frac{-\gamma_2\pi_n}{1 - \gamma_1 - \gamma_2} \ell_n \right] + \mathbb{E}_{U_2} \left[\frac{-\gamma_1\pi_p}{1 - \gamma_1 - \gamma_2} \ell_p + \frac{(1 - \gamma_1)\pi_n}{1 - \gamma_1 - \gamma_2} \ell_n \right]$.	(9)
SU	$R(g) = (\pi_p^2 + \pi_n^2) \mathbb{E}_S \left[\frac{L(X) + L(X)}{2} \right] + \mathbb{E}_U [L_-(X)]$, where $L(X) := \frac{1}{\pi_p - \pi_n} \ell_p(X) - \frac{1}{\pi_p - \pi_n} \ell_n(X)$, $L_-(X) := -\frac{\pi_n}{\pi_p - \pi_n} \ell_p(X) + \frac{\pi_p}{\pi_p - \pi_n} \ell_n(X)$.	(11)
DU	$R(g) = 2\pi_p\pi_n \mathbb{E}_D \left[-\frac{L(X) + L(X)}{2} \right] + \mathbb{E}_U [L_+(X)]$, where $L(X)$ is defined in the SU setting, and $L_+(X) := \frac{\pi_p}{\pi_p - \pi_n} \ell_p(X) - \frac{\pi_n}{\pi_p - \pi_n} \ell_n(X)$.	(13)
SD	$R(g) = (\pi_p^2 + \pi_n^2) \mathbb{E}_S \left[\frac{L_+(X) + L_+(X)}{2} \right] + 2\pi_p\pi_n \mathbb{E}_D \left[\frac{L_-(X) + L_-(X)}{2} \right]$, where $L_+(X)$ and $L_-(X)$ are defined in the SU and DU settings.	(15)
Pcomp	$R(g) = \mathbb{E}_{\text{Sup}} [\ell_p - \pi_p \ell_n] + \mathbb{E}_{\text{Inf}} [-\pi_n \ell_p + \ell_n]$, where $P_{\text{Sup}} := \int_x \times P_{\text{PC}} dx$, $P_{\text{Inf}} := \int_x \times P_{\text{PC}} dx$.	(17)
Sconf	$R(g) = \mathbb{E}_{X,X} \left[\frac{r(X, X) - \pi_n}{\pi_p - \pi_n} \frac{\ell_p(X) + \ell_p(X)}{2} + \frac{\pi_p - r(X, X)}{\pi_p - \pi_n} \frac{\ell_n(X) + \ell_n(X)}{2} \right]$.	(19)

Table 5: Risk rewrites for multiclass WSLs.

WSL	Risk rewrite for $R(g) = \mathbb{E}_{Y,X} [\ell_Y(g(X))]$ (1)	Equation
CL	$R(g) = \mathbb{E}_{\bar{S},X} \left[\sum_{y=1}^K \ell_y - (K-1)\ell_{\bar{S}} \right].$	(21)
MCL	$R(g) = \sum_{d=1}^{K-1} P_{ \bar{S} =d} \mathbb{E}_{\bar{S},X \bar{S} =d} \left[\sum_{y/\bar{S}} \ell_y - \frac{K-1- \bar{S} }{ \bar{S} } \sum_{\bar{s}} \ell_{\bar{s}} \right].$	(23)
PCPL	$R(g) = \frac{1}{2} \mathbb{E}_{S,X} \left[\sum_{y=1}^K \frac{P_{Y=y X}}{\sum_a P_{Y=a X}} \ell_y \right].$	(25)
PPL	$R(g) = \mathbb{E}_{S,X} \left[\sum_y \frac{P_{Y=y X}}{\sum_a P_{Y=a X}} \ell_y \right].$	(27)
SC-Conf	$R(g) = \pi_{y_s} \mathbb{E}_{X Y=y_s} \left[\sum_{y=1}^K \frac{r_y(X)}{r_{y_s}(X)} \ell_y \right].$	(29)
Sub-Conf	$R(g) = \pi_{Y_s} \mathbb{E}_{X Y=Y_s} \left[\sum_{y=1}^K \frac{r_y(X)}{r_{Y_s}(X)} \ell_y \right].$	(31)
Soft-label	$R(g) = \mathbb{E}_X \left[\sum_{y=1}^K r_y(X) \ell_y \right].$	(34)

From the above tables, finding a way to reexpress the classification risk $R(g)$ (1) in terms of the data-generating distributions becomes the crux when applying ERM for most WSL studies. The rewrites also replace loss functions ℓ_Y defining (1) with various modified losses (shown inside the expectations). These modified loss functions are sometimes called corrected losses, which is why the approach is also called loss correction. Proposing a generic methodology that finds properly corrected losses to achieve risk rewrite in different scenarios is a main topic we would like to elaborate on in this paper.

2.2.17 Learning with Noisy Labels (LNL) Formulations

Next, we review two related formulations in LNL, the MCD and CCN settings, in Table 6. The observed instances in MCD and CCN are still labeled by $\{p, n\}$ but are polluted by certain noise models. We use \bar{Y} to represent a polluted label, compared to an unpolluted Y . In MCD, a small portion of the negatively labeled data $\gamma_p P_{X|Y=n}$ contaminates the positively labeled data $P_{X|Y=p}$. Likewise, a small portion of the positive data $\gamma_n P_{X|Y=p}$ contaminates the negatively labeled data $P_{X|Y=n}$ (Scott et al., 2013). In the CCN setting, a label Y is flipped to become \bar{Y} with probability $P_{\bar{Y}|Y,X}$ (Natarajan et al., 2013). Although they are formulated for the study of noisy labels, their formulations share similar structures with many WSLs above. In Section 4, we will use the similarities to categorize WSLs and provide a bird’s eye view to reveal connections among WSLs.

Table 6: MCD and CCN formulations.

Scenario	Formulation
MCD	$\{x_i^{\bar{p}}\}_{i=1}^{n_{\bar{p}}} \text{ i.i.d. } P_{X/\bar{Y}=\text{p}} := (1 - \gamma_{\text{p}}) P_{X/Y=\text{p}} + \gamma_{\text{p}} P_{X/Y=\text{n}}.$ $\{x_j^{\bar{n}}\}_{j=1}^{n_{\bar{n}}} \text{ i.i.d. } P_{X/\bar{Y}=\text{n}} := \gamma_{\text{n}} P_{X/Y=\text{p}} + (1 - \gamma_{\text{n}}) P_{X/Y=\text{n}}.$
CCN	$\{(\bar{y}_i, x_i)\}_{i=1}^n \text{ i.i.d. } P_{\bar{Y}=\bar{y}_i, X} := \sum_{k \in \{\text{p}, \text{n}\}} P_{\bar{Y}=\bar{y}_i, Y=k, X} P_{Y=k, X}, \bar{y}_i \in \{\text{p}, \text{n}\}.$

3 A Framework for Risk Rewrite

We illustrate the proposed framework in this section. Its job is to provide a unified treatment and understanding of WSL. It consists of a formulation component and an analysis component. The analysis component suggests a generic methodology to solve the risk rewrite problem. Moreover, diving into the formulation component’s logic, we can interpret multiple WSL formulations and the diverse risk rewrites from a single perspective.

Before introducing the framework, we first define several abstract notations that will be used throughout the paper. There are three main characters and one supporting character in the framework. The main characters are the vector of data-generating distributions \bar{P} , the vector of risk-defining distributions P , and the vector of base distributions B . The supporting character is the vector of loss functions L . The reason for using vectorized pseudonyms is that the proposed framework uses matrix multiplication as a basic mathematical operation. \bar{P} contains distributions that produce the observational data. For instance, $\bar{P} = \begin{pmatrix} P_{\text{P}} \\ P_{\text{U}} \end{pmatrix}$ in PU learning (Table 2) and $P_{\bar{S}=k, X}$ is the k -th entry of \bar{P} in CL learning (Table 3). We use classification risk (1) to illustrate our framework. So P consists of joint distributions $P_{X, Y}$. We can look at Tables 2 and 3 and see that there are basic elements that define a data-generating distribution. These are class-conditionals $P_{X/Y}$ in Table 2 and joint distributions $P_{X, Y}$ in Table 3. Since these basic elements do not necessarily coincide with the entries of P , we denote them as B . L consists of loss functions and its k -th entry is $\ell_{Y=k}(g(X))^2$.

3.1 The Formulation Component of the Framework

The construction of the formulation component is to study the connections among WSLs and provide a foundation for developing the generic methodology. We draw inspiration from Section 2.2. Each WSL formulation represents a type of weakened information of the joint distribution $P_{Y, X}$ in supervised learning. For instance, unlabeled data discards the label information (Lu et al., 2020; 2021), the complementary-label is a label that cannot be the ground truth (Ishida et al., 2017; Yu et al., 2018), and the similarity encodes a comparative relationship of two ground truth labels (Bao et al., 2018; Shimada et al., 2021; Cao et al., 2021b). Thus, we are motivated to search for a general way to link data-generating distributions with the joint distribution.

We start by linking the data-generating distributions \bar{P} and the base distributions B . This involves finding matrix correspondences to Tables 2 and 3. We assume that a matrix M_{corr} formalizes the link:

$$\bar{P} = M_{\text{corr}} B. \quad (35)$$

Taking PU learning (4) for example, M_{corr} aims to connect $\bar{P} = \begin{pmatrix} P_{\text{P}} \\ P_{\text{U}} \end{pmatrix}$ with $B = \begin{pmatrix} P_{X/Y=\text{p}} \\ P_{X/Y=\text{n}} \end{pmatrix}$. To keep the framework as abstract as possible, we would like to defer the discussion of all other \bar{P} and B until we realize their corresponding M_{corr} in Section 4.

The matrix formulation has two advantages. First, it provides a unified way to characterize a wide range of WSL settings. By studying the entries of a matrix, we can easily link one WSL scenario to another to form

²We reserve \bar{P} , P , and B for vectors of distributions and L and \bar{L} for vectors of loss functions. We address them as “the distributions” and “the losses” to avoid the verbose “the vector of distributions/losses.”

reduction graphs of WSLs. As the first main topic of this work, Section 4 shows, for a given WSL setting, how to find the corresponding matrix M_{corr} , and Tables 7 – 9 summarize fifteen WSL settings covered by our matrix formulation and depict a reduction graph rooted from M_{corr} . The following subsection illustrates the second advantage of aiding the construction of a generic methodology for conducting risk rewrite.

3.2 The Analysis Component of the Framework

Note that the conventional expression (1) can be simplified, by the inner product, to be $R(g) = \int_{x \times X} L \cdot P dx$. It is immediately possible to rewrite the risk under data-generating distributions \bar{P} by showing $L \cdot P = \bar{L} \cdot \bar{P}$, where \bar{L} is called the vector of corrected losses and its role will be clarified later. Therefore, it is imperative to establish the connection between \bar{P} and P . The goal can be achieved by linking B and P since we have assumed that $\bar{P} = M_{\text{corr}}B$ in the previous subsection. Recall from the beginning of this section that the base distributions B are of the forms $P_{X|Y}$ or $P_{X,Y}$, and the risk-defining distributions P are of the form $P_{X,Y}$. Given their label-relevant nature (i.e., they are either the joint distribution or the class-conditionals), we assume that there exists a transformation matrix M_{trsf} that satisfies $B = M_{\text{trsf}}P$. Thus, (35) becomes

$$\bar{P} = M_{\text{corr}}M_{\text{trsf}}P. \quad (36)$$

Having the freedom to choose M_{trsf} allows the framework to handle different base distributions B and to adapt to various performance measures that define P , as we will discuss in Sections 5.1.1, 5.2.1, and 6.1, respectively.

The reason why connecting P with \bar{P} (36) helps the construction of the corrected losses is that if we manage to find a way to compensate for the combined effect of M_{corr} and M_{trsf} , we can implement the compensation mechanism on the “corrected” losses \bar{L} . Specifically, suppose there exists a matrix M_{corr}^\dagger satisfying

$$P = M_{\text{corr}}^\dagger \bar{P}. \quad (37)$$

Then, the corrected losses defined by

$$\bar{L} := L \cdot M_{\text{corr}}^\dagger \quad (38)$$

allows us to rephrase the classification risk as

$$\begin{aligned} \int_{x \times X} \bar{L} \cdot \bar{P} dx &= \int_{x \times X} L \cdot M_{\text{corr}}^\dagger \bar{P} dx \\ &= \int_{x \times X} L \cdot P dx = R(g), \end{aligned} \quad (39)$$

providing a rewrite for $R(g)$ with respect to \bar{P} .

The above procedure describes a generic methodology for the risk rewrite problem. As the second main topic, we instantiate the framework by presenting the corresponding matrices M_{corr}^\dagger and M_{trsf} for each learning scenario in Section 5 to demonstrate its applicability.

3.3 Intuition of the Framework

The key equations discussed in Sections 3.1 and 3.2 are

$$\begin{aligned} \bar{P} &\stackrel{(35)}{=} M_{\text{corr}}B \stackrel{(36)}{=} M_{\text{corr}}M_{\text{trsf}}P, \\ \bar{L} \cdot \bar{P} &\stackrel{(38)}{=} L \cdot M_{\text{corr}}^\dagger M_{\text{corr}}M_{\text{trsf}}P \stackrel{(37)}{=} L \cdot P. \end{aligned}$$

The logic behind them is succinct and interpretive. First, from a formulation perspective, viewing matrix M_{corr} as a contamination matrix that corrupts the base B to become the contaminated \bar{P} (35), we interpret this *contamination mechanism* as sacrificing certain information in exchange for certain saved costs or privacy, reflecting the essence underlying WSL formulations. In addition to formulating the data-generating

mechanism, the link between B and the risk-defining distributions P (36) connects \bar{P} and P to motivate the methodology design. This novel viewpoint of connecting the data distributions via the explicit two-stage formulation facilitates the unification work in this paper.

Second, regarding the methodological design, it becomes easier to devise a countermeasure when the connection between \bar{P} and P is in good shape. The design of $\bar{L} = L M_{\text{corr}}^\dagger$ involves M_{corr}^\dagger , which captures a common idea behind risk rewrite: Restoration of the risk-defining distributions and the original loss functions is accomplished by the *decontamination* (37) provided by \bar{L} . Furthermore, the instantiations of $\bar{L} = L M_{\text{corr}}^\dagger$ (38) justify that the apparently different forms of corrected losses reported in the literature (i.e., referred papers that contribute to Tables 4 and 5, and those referred to as recoveries in Section 5) essentially stem from M_{corr}^\dagger . In summary, the proposed framework is abstract and flexible enough that we use it in the current paper to formulate the contamination mechanisms and provide a generic methodology for a wide range of WSLs.

3.4 Building Blocks: The Inversion and the Marginal Chain Approaches

We describe two building blocks, the inversion method and the marginal chain method, that will be used to devise M_{corr}^\dagger that satisfies (37) in each scenario we study later.

Proposition 1 (The inversion method). *Let P and \bar{P} be vectors. Suppose $\bar{P} = MP$ holds for an invertible matrix M . Then, choosing $M_{\text{corr}}^\dagger = M^{-1}$, we have $P = M_{\text{corr}}^\dagger \bar{P}$.*

Proof. For any invertible M , it is easy to see that, by assigning $M_{\text{corr}}^\dagger = M^{-1}$, one has

$$M_{\text{corr}}^\dagger \bar{P} = M^{-1} \bar{P} = M^{-1} MP = P.$$

□

We remark that this simple strategy was adopted in many LNL works. A handful of related papers are Cid-Sueiro (2012), Blanchard & Scott (2014), Menon et al. (2015), van Rooyen & Williamson (2015), Patrini et al. (2017), van Rooyen & Williamson (2017), and Katz-Samuels et al. (2019). Hence, it can be applied to WSLs that are special cases of certain LNL scenarios.

Proposition 2 (The marginal chain method). *Let $Y = k \in [K]$ be a class label, where $[K]$ is the set of classes associated with the classification risk. Let $S = \{s_1, s_2, \dots, s_{|S|}\} \subseteq 2^{[K]}$ be the set of class sets and S be the random variable of the observational outcome. Denote*

$$P = \begin{pmatrix} P_{Y=1,X} \\ \vdots \\ P_{Y=K,X} \end{pmatrix} \text{ and } \bar{P} = \begin{pmatrix} P_{S=s_1,X} \\ \vdots \\ P_{S=s_{|S|},X} \end{pmatrix}.$$

Then,

$$M = \begin{pmatrix} P_{S=s_1|Y=1,X} & P_{S=s_1|Y=2,X} & \cdots & P_{S=s_1|Y=K,X} \\ P_{S=s_2|Y=1,X} & P_{S=s_2|Y=2,X} & \cdots & P_{S=s_2|Y=K,X} \\ \vdots & \vdots & \ddots & \vdots \\ P_{S=s_{|S|}|Y=1,X} & P_{S=s_{|S|}|Y=2,X} & \cdots & P_{S=s_{|S|}|Y=K,X} \end{pmatrix} \quad (40)$$

satisfies $\bar{P} = MP$, and

$$M_{\text{corr}}^\dagger = \begin{pmatrix} P_{Y=1|S=s_1,X} & P_{Y=1|S=s_2,X} & \cdots & P_{Y=1|S=s_{|S|},X} \\ P_{Y=2|S=s_1,X} & P_{Y=2|S=s_2,X} & \cdots & P_{Y=2|S=s_{|S|},X} \\ \vdots & \vdots & \ddots & \vdots \\ P_{Y=K|S=s_1,X} & P_{Y=K|S=s_2,X} & \cdots & P_{Y=K|S=s_{|S|},X} \end{pmatrix} \quad (41)$$

satisfies $P = M_{\text{corr}}^\dagger \bar{P}$.

The role of S is to represent a weak supervision that encodes some combinatorial information about the unobservable true label Y . We will discuss this concept in detail in Sections 4.2 and 5.2.

Proof. It suffices to show $(MP)_j = \bar{P}_j$ for any $j \in [S]$. Taking the inner product of the j -th row of M and P , we have

$$\sum_{k=1}^K P_{S=s_j|Y=k,X} P_{Y=k,X} = \sum_{k=1}^K P_{S=s_j,Y=k,X} = P_{S=s_j,X}$$

that verifies (40).

Next, we prove $P = M_{\text{corr}}^\dagger \bar{P}$ by showing $(M_{\text{corr}}^\dagger \bar{P})_i = P_i$. For each $i \in [K]$,

$$\begin{aligned} (M_{\text{corr}}^\dagger \bar{P})_i &= (M_{\text{corr}}^\dagger MP)_i = \sum_{j=1}^{|S|} P_{Y=i|S=s_j,X} \sum_{k=1}^K P_{S=s_j|Y=k,X} P_{Y=k,X} \\ &\stackrel{(a)}{=} \sum_{j=1}^{|S|} P_{Y=i|S=s_j,X} P_{S=s_j,X} \\ &\stackrel{(b)}{=} P_{Y=i,X} = P_i. \end{aligned} \tag{42}$$

□

Besides finding the inverse matrix, we propose a new approach called the *marginal chain* to achieve (37). The development of this approach begins with the observation that $P_{S=s_j,X}$ in $\bar{P} = MP$ is a distribution where Y is marginalized out. It inspires an idea that one could perform another marginalization to restore the original distribution $P_{Y,X}$; specifically, by marginalizing out S . The design of M_{corr}^\dagger in (41) aims to carry out the idea. As shown by (a) and (b) in the proof, two consecutive marginalization steps on Y and then S give the name of the marginal chain.

Both the inversion and marginal chain methods have strengths and weaknesses. The inversion method only requires P as a real vector but needs the invertible assumption on the contamination matrix M . In contrast, the marginal chain method exploits that P , in fact, is a distributional vector, allowing it to find a decontamination matrix M_{corr}^\dagger even for a non-invertible M . A restriction of the marginal chain method is that the construction of M_{corr}^\dagger is regulated by probability equations.

We are ready to justify the proposed framework through the following two sections. Section 4 discusses weakly supervised scenarios that can be subsumed by the formulation component (35). Section 5 verifies the analysis component by instantiating (38) to conduct the risk rewrite for each scenario mentioned in Section 4. In both sections, we divide the scenarios into three categories. The first two are WSLs that can be viewed as special cases in either the prevalent MCD or CCN settings. The third category contains confidence-based scenarios. The notations listed in Table 1 will still be functional. For all notations and their abbreviations required in the coming sections, please refer to Appendix B.

4 Contamination as Weak Supervision

In this section, we instantiate the contamination matrix for each weakly supervised scenario listed in Table 2 and Table 3. Tables 7 – 9 summarize the contamination matrices developed in this section. Each table also represents a reduction graph of WSL settings. These reduction graphs cluster WSL settings into three main categories, providing a hierarchy of relationships. With this hierarchy, we can understand, compare with, and relate to different settings or even grow the hierarchy by adding new branches. Next are the notations

for reading the graphs. For two contamination mechanisms, U and V, we use $M_U \rightarrow M_V$ to denote “ M_U is reduced to M_V ” or “ M_U is realized as M_V ”, and $M_U \leftarrow M_V$ means “ M_U is generalized to M_V ”.

The proposed framework provides a generic strategy for formulating multiple weakly supervised scenarios. Thus, the proofs will have a certain degree of similarity. To avoid repeating similar proofs, we provide proofs that appear for the first time. For auxiliary lemmas and results whose proofs are similar to the previous ones, we refer to the omitted proofs in Appendix C. In particular, the omitted proofs in Section 4.1 can be found in Appendix C.1, and those in Section 4.2 can be found in Appendix C.2.

Table 7: Contamination matrices of MCD category in Section 4.1.

WSLs	Entry Parameter	Contamination Matrix	Reduction path	
MCD	γ_p, γ_n	M_{MCD} (45)	M_{corr}	M_{MCD}
UU	γ_1, γ_2	M_{UU} (49)	M_{corr}	$M_{\text{UU}} \rightarrow M_{\text{MCD}}$
PU	$\gamma_1 = 0, \gamma_2 = \pi_p$	M_{PU} (50)	M_{UU}	M_{PU}
SU	$\gamma_1 = \frac{\pi_n^2}{\pi_p^2 + \pi_n^2}, \gamma_2 = \pi_p$	M_{SU} (51)	M_{UU}	M_{SU}
Pcomp	$\gamma_1 = \frac{\pi_n^2}{\pi_p + \pi_n^2}, \gamma_2 = \frac{\pi_p^2}{\pi_p^2 + \pi_n}$	M_{Pcomp} (52)	M_{UU}	M_{Pcomp}
DU	$\gamma_1 = 1/2, \gamma_2 = \pi_p$	M_{DU} (53)	M_{UU}	M_{DU}
SD	$\gamma_1 = \frac{\pi_n^2}{\pi_p^2 + \pi_n^2}, \gamma_2 = 1/2$	M_{SD} (54)	M_{UU}	M_{SD}
Sconf	–	M_{Sconf} (55)	M_{corr}	M_{Sconf}

Table 8: Contamination matrices of CCN category in Section 4.2.

WSLs	Entry Parameter	Contamination Matrix	Reduction path	
CCN	$P_{\bar{Y} Y,X}$ (59)	M_{CCN} (60)	M_{corr}	M_{CCN}
Generalized CCN	$P_{S Y,X}$ (61)	M_{gCCN} (64)	M_{corr}	$M_{\text{CCN}} \rightarrow M_{\text{gCCN}}$
PPL	$C(S, X) [Y \quad S]$ (65)	M_{PPL} (66)	M_{gCCN}	M_{PPL}
PCPL	$\frac{1}{2^{K-1}-1} [Y \quad S]$	M_{PCPL} (68)	M_{gCCN}	$M_{\text{PPL}} \rightarrow M_{\text{PCPL}}$
MCL	$\frac{q_{ \bar{S} }}{\binom{K-1}{ \bar{S} }} [Y \quad \bar{S}]$ (75)	M_{MCL} (71)	M_{gCCN}	$M_{\text{PPL}} \rightarrow M_{\text{MCL}}$
CL	$ S = 1, \frac{1}{K-1} [Y \quad S]$	M_{CL} (76)	M_{gCCN} M_{CL}	$M_{\text{PPL}} \rightarrow M_{\text{MCL}}$

Table 9: Contamination matrices of confidence-based category in Section 4.3.

WSLs	Entry Parameter	Contamination Matrix	Reduction path	
Sub-Conf	$\frac{P_{Y=y_s X}}{P_{Y=k X}}$	M_{Sub} (80)	M_{corr}	M_{Sub}
SC	$Y_s = \{y_s\}$ in M_{Sub}	M_{SC} (81)	M_{Sub}	M_{SC}
Pconf	$K = 2, y_s = p$ in M_{SC}	M_{Pconf} (82)	M_{Sub}	$M_{\text{SC}} \rightarrow M_{\text{Pconf}}$
Soft	$\frac{1}{P_{Y=k X}}$	M_{Soft} (84)	M_{Sub}	M_{Soft}

4.1 MCD Scenarios

As listed in Table 6, in binary classification, the MCD model (Menon et al., 2015) corrupts the clean class-conditionals $P_{X/Y=p}$ and $P_{X/Y=n}$ via parameters γ_p and γ_n as follows:

$$\begin{aligned} P_{X/\bar{Y}=p} &:= (1 - \gamma_p) P_{X/Y=p} + \gamma_p P_{X/Y=n}, \\ P_{X/\bar{Y}=n} &:= \gamma_n P_{X/Y=p} + (1 - \gamma_n) P_{X/Y=n}, \end{aligned} \quad (43)$$

where $\gamma_p, \gamma_n \in [0, 1]$ and $\gamma_p + \gamma_n < 1$. Viewing the contamination targets $P_{X/Y=p}$ and $P_{X/Y=n}$ as the base distributions

$$B := \begin{pmatrix} P_{X/Y=p} \\ P_{X/Y=n} \end{pmatrix}$$

and denoting the vector of data-generating distributions as

$$\bar{P} := \begin{pmatrix} P_{X/\bar{Y}=p} \\ P_{X/\bar{Y}=n} \end{pmatrix},$$

we can express (43) in the following matrix form

$$\begin{pmatrix} P_{X/\bar{Y}=p} \\ P_{X/\bar{Y}=n} \end{pmatrix} = \begin{pmatrix} 1 - \gamma_p & \gamma_p \\ \gamma_n & 1 - \gamma_n \end{pmatrix} \begin{pmatrix} P_{X/Y=p} \\ P_{X/Y=n} \end{pmatrix}. \quad (44)$$

Comparing (44) with $\bar{P} = M_{\text{corr}} B$ (35), we find that the contamination matrix M_{corr} is realized as

$$M_{\text{MCD}} := \begin{pmatrix} 1 - \gamma_p & \gamma_p \\ \gamma_n & 1 - \gamma_n \end{pmatrix} \quad (45)$$

in the MCD setting.

4.1.1 Unlabeled-Unlabeled (UU) Learning (Lu et al., 2019)

Next, we show how to characterize UU learning by a contamination matrix. Naming

$$\pi_p P_{X/Y=p} + \pi_n P_{X/Y=n}$$

as P_U is feasible since $\pi_p P_{X/Y=p} + \pi_n P_{X/Y=n} = P_X$ generates data that statistically equals to data sampled from $P_{Y,X}$ with labels removed. Viewing π_p as the mixture rate of samples from $P_{X/Y=p}$ and $P_{X/Y=n}$, P_U is parameterized by π_p . Therefore, we can interpret (8), recalled as follows, as formulating two unlabeled data distributions w.r.t. mixture rates $(1 - \gamma_1)$ and γ_2 , respectively:

$$\begin{aligned} P_{U_1} &= (1 - \gamma_1) P_{X/Y=p} + \gamma_1 P_{X/Y=n}, \\ P_{U_2} &= \gamma_2 P_{X/Y=p} + (1 - \gamma_2) P_{X/Y=n}. \end{aligned}$$

Taking the class-conditionals as the base distributions

$$B := \begin{pmatrix} P_{X/Y=p} \\ P_{X/Y=n} \end{pmatrix} \quad (46)$$

and converting (8) to the matrix form, we express the data-generating distributions of UU learning

$$\bar{P} := \begin{pmatrix} P_{U_1} \\ P_{U_2} \end{pmatrix} \quad (47)$$

as

$$\begin{pmatrix} P_{U_1} \\ P_{U_2} \end{pmatrix} = \begin{pmatrix} 1 - \gamma_1 & \gamma_1 \\ \gamma_2 & 1 - \gamma_2 \end{pmatrix} \begin{pmatrix} P_{X|Y=p} \\ P_{X|Y=n} \end{pmatrix}, \quad (48)$$

and we arrive at the following lemma.

Lemma 3. *Let B (46) be the base distributions and \bar{P} (47) be the data-generating distributions. For $\gamma_1, \gamma_2 \in [0, 1]$ such that $\gamma_1 + \gamma_2 = 1$, the contamination matrix*

$$M_{UU} := \begin{pmatrix} 1 - \gamma_1 & \gamma_1 \\ \gamma_2 & 1 - \gamma_2 \end{pmatrix} \quad (49)$$

characterizes the data-generating process of UU learning (8) via (48).

Comparing (48) with the formulation framework $\bar{P} = M_{\text{corr}}B$ (35), we see that in UU learning, M_{corr} is realized as M_{UU} :

$$M_{\text{corr}} = M_{UU}.$$

Like MCD, we assume $\gamma_1 + \gamma_2 = 1$. Our assumption is equivalent to that of MCD since the case of swapping P_{corr} and Q_{corr} in Menon et al. (2015) corresponds to $\gamma_1 + \gamma_2 > 1$ in our case. For details, refer to the discussion in Section 2.2 of Menon et al. (2015). The need for $\gamma_1 + \gamma_2 = 1$ can be explained by examining the entries in M_{UU} . The constraint $\gamma_1 + \gamma_2 = 1$ guarantees distinct rows in M_{UU} , implying the observed data sets are sampled from two distinct distributions. On the contrary, allowing $\gamma_1 + \gamma_2 \neq 1$ ends up observing one unlabeled data set (i.e., $P_{U_1} = P_{U_2}$) since $1 - \gamma_1 = \gamma_2$. Lu et al. (2019) proved in Section 3 that it is impossible to conduct a risk rewrite if one only observes one unlabeled data set.

Assigning $\gamma_1 = \gamma_p$ and $\gamma_2 = \gamma_n$ implies that MCD and UU have essentially the same data-generating process from the contamination perspective, as (44) and (48) have the identical right-hand sides (i.e., the same contamination targets and the same contamination matrix). However, they bear different meanings in respective research topics (i.e., distinct notions on the left-hand sides of the equations): In MCD, one still observes data with labels, nonetheless noisy, while in the UU setting, one observes two distinct unlabeled data sets. We use “ ” to denote their relation in the UU row of Table 7.

Connecting UU learning with MCD, and later the generalized CCN with CCN in Section 4.2.1, allows us to categorize WSLs from the LNL perspective into Sections 4.1 and 4.2. In the rest of this subsection, we collect WSLs whose base distributions are class-conditionals and show M_{UU} instantiates their formulations via respective assignments of γ_1 and γ_2 .

4.1.2 Positive-Unlabeled (PU) Learning (Kiryo et al., 2017)

Recall from (4) that $P_P = P_{X|Y=p}$ and $P_U = P_X$. The following lemma describes the contamination matrix of PU learning.

Lemma 4. *Let B (46) be the base distributions and*

$$\bar{P} := \begin{pmatrix} P_P \\ P_U \end{pmatrix}$$

be the data-generating distributions. Define the contamination matrix

$$M_{PU} := \begin{pmatrix} 1 & 0 \\ \pi_p & \pi_n \end{pmatrix}. \quad (50)$$

Then, $\bar{P} = M_{PU}B$, and M_{PU} characterizes the data-generating process of PU learning (4).

Proof. We apply the same proof strategy as in Lemma 3. By definitions,

$$M_{\text{PU}}B = \begin{pmatrix} 1 & 0 \\ \pi_{\text{p}} & \pi_{\text{n}} \end{pmatrix} \begin{pmatrix} P_{X/Y=\text{p}} \\ P_{X/Y=\text{n}} \end{pmatrix} = \begin{pmatrix} 1 \cdot P_{X/Y=\text{p}} + 0 \cdot P_{X/Y=\text{n}} \\ \pi_{\text{p}} \cdot P_{X/Y=\text{p}} + \pi_{\text{n}} \cdot P_{X/Y=\text{n}} \end{pmatrix}.$$

Since $1 \cdot P_{X/Y=\text{p}} + 0 \cdot P_{X/Y=\text{n}} = P_{X/Y=\text{p}}$ and $\pi_{\text{p}} \cdot P_{X/Y=\text{p}} + \pi_{\text{n}} \cdot P_{X/Y=\text{n}} = P_X$, we obtain $M_{\text{PU}}B = \bar{P}$. \square

Comparing with (35), we see that the contamination matrix M_{corr} is instantiated as M_{PU} (50) in PU learning. Further, M_{PU} can be obtained by assigning $\gamma_1 = 0$ and $\gamma_2 = \pi_{\text{p}}$ in M_{UU} (49), and hence, we obtain the reduction path

$$M_{\text{corr}} \quad M_{\text{UU}} \quad M_{\text{PU}}.$$

4.1.3 Similar-Unlabeled (SU) Learning (Bao et al., 2018)

Recall P_{S} (10) generates the pair of data points (x, x) who share the same label. In addition to the pairwise distribution P_{S} , a pointwise distribution

$$P_{\bar{\text{S}}} = \frac{\pi_{\text{p}}^2 P_{X/Y=\text{p}} + \pi_{\text{n}}^2 P_{X/Y=\text{n}}}{\pi_{\text{p}}^2 + \pi_{\text{n}}^2}$$

is also defined for single data point x (Bao et al., 2018, Lemma 1). Therefore, we choose $P_{\bar{\text{S}}}$ as the data-generating distribution when constructing the contamination matrix in the following lemma.

Lemma 5. *Let B (46) be the base distributions and*

$$\bar{P} := \begin{pmatrix} P_{\bar{\text{S}}} \\ P_{\text{U}} \end{pmatrix}.$$

Then, the contamination matrix

$$M_{\text{SU}} := \begin{pmatrix} \frac{\pi_{\text{p}}^2}{\pi_{\text{p}}^2 + \pi_{\text{n}}^2} & \frac{\pi_{\text{n}}^2}{\pi_{\text{p}}^2 + \pi_{\text{n}}^2} \\ \pi_{\text{p}} & \pi_{\text{n}} \end{pmatrix}, \quad (51)$$

which satisfies $\bar{P} = M_{\text{SU}}B$, characterizes the data-generating distributions \bar{P} .

Further, M_{SU} can be obtained by assigning $\gamma_1 = \frac{\pi_{\text{n}}^2}{\pi_{\text{p}}^2 + \pi_{\text{n}}^2}$ and $\gamma_2 = \pi_{\text{p}}$ in M_{UU} (49), and hence, we obtain the reduction path

$$M_{\text{corr}} \quad M_{\text{UU}} \quad M_{\text{SU}}.$$

4.1.4 Pairwise Comparison (Pcomp) Learning (Feng et al., 2021)

In SU learning, we formulate the pointwise data-generating distributions $P_{\bar{\text{S}}}$ and P_{U} ; likewise, we use the following pointwise distributions of P_{PC} (16) to formulate Pcomp learning:

$$\begin{aligned} P_{\text{Sup}} &:= \int_x \times P_{\text{PC}} dx = \frac{\pi_{\text{p}} P_{X/Y=\text{p}} + \pi_{\text{n}}^2 P_{X/Y=\text{n}}}{\pi_{\text{p}} + \pi_{\text{n}}^2}, \\ P_{\text{Inf}} &:= \int_x \times P_{\text{PC}} dx = \frac{\pi_{\text{p}}^2 P_{X/Y=\text{p}} + \pi_{\text{n}} P_{X/Y=\text{n}}}{\pi_{\text{p}}^2 + \pi_{\text{n}}}. \end{aligned}$$

Lemma 6. *Let B (46) be the base distributions and*

$$\bar{P} := \begin{pmatrix} P_{\text{Sup}} \\ P_{\text{Inf}} \end{pmatrix}.$$

Then, the contamination matrix

$$M_{\text{Pcomp}} := \begin{pmatrix} \frac{\pi_p}{\pi_p + \pi_n^2} & \frac{\pi_n^2}{\pi_p + \pi_n^2} \\ \frac{\pi_p^2}{\pi_p^2 + \pi_n} & \frac{\pi_n}{\pi_p^2 + \pi_n} \end{pmatrix}, \quad (52)$$

which satisfies $\bar{P} = M_{\text{Pcomp}}B$, characterizes the data-generating distributions \bar{P} .

Further, M_{Pcomp} can be obtained by assigning $\gamma_1 = \frac{\pi_n^2}{\pi_p + \pi_n^2}$ and $\gamma_2 = \frac{\pi_p^2}{\pi_p^2 + \pi_n}$ in M_{UU} (49), and hence, we obtain the reduction path

$$M_{\text{corr}} \quad M_{\text{UU}} \quad M_{\text{Pcomp}}.$$

4.1.5 Similar-dissimilar-unlabeled (SDU) Learning (Shimada et al., 2021)

Dissimilar-unlabeled (DU) learning and similar-dissimilar (SD) learning are two critical components of SDU learning. Hence, we present the matrix formulations of M_{DU} and M_{SD} . Similar to the strategy taken in Sections 4.1.3 and 4.1.4, we use the following pointwise distribution

$$P_{\bar{\text{D}}} = \int_x P_{\text{D}} \, dx = \frac{P_{x|Y=p} + P_{x|Y=n}}{2}$$

(Shimada et al., 2021, (36) in Appendix A.1) in the following formulations.

We formulate the contamination matrix of DU learning via the following lemma.

Lemma 7. *Let B (46) be the base distributions and*

$$\bar{P} = \begin{pmatrix} P_{\bar{\text{D}}} \\ P_{\text{U}} \end{pmatrix}.$$

Then, the contamination matrix

$$M_{\text{DU}} = \begin{pmatrix} 1/2 & 1/2 \\ \pi_p & \pi_n \end{pmatrix}, \quad (53)$$

which satisfies $\bar{P} = M_{\text{DU}}B$, characterizes the data-generating distributions \bar{P} .

Furthermore, since M_{UU} (49) reduces to M_{DU} by assigning $\gamma_1 = 1/2$ and $\gamma_2 = \pi_p$, we have the reduction path

$$M_{\text{corr}} \quad M_{\text{UU}} \quad M_{\text{DU}}.$$

The next lemma formulates the contamination matrix of SD learning.

Lemma 8. *Let B (46) be the base distributions and*

$$\bar{P} = \begin{pmatrix} P_{\bar{\text{S}}} \\ P_{\bar{\text{D}}} \end{pmatrix}.$$

Then, the contamination matrix

$$M_{\text{SD}} = \begin{pmatrix} \frac{\pi_p^2}{\pi_p^2 + \pi_n^2} & \frac{\pi_n^2}{\pi_p^2 + \pi_n^2} \\ 1/2 & 1/2 \end{pmatrix} \quad (54)$$

which satisfies $\bar{P} = M_{\text{SD}}B$, characterizes the data-generating distributions \bar{P} .

Moreover, because M_{UU} (49) reduces to M_{SD} via $\gamma_1 = \frac{\pi_n^2}{\pi_p^2 + \pi_n^2}$ and $\gamma_2 = 1/2$, we obtain the reduction path

$$M_{\text{corr}} \quad M_{\text{UU}} \quad M_{\text{SD}}.$$

4.1.6 Similarity-Confidence (Sconf) Learning (Cao et al., 2021b)

Recall from the Sconf setting (18) that (x, x) is a pair of data sampled i.i.d. from $P_{X,X} := P_X P_X$. On seeing P_X , one might wonder if it is sufficient to express the data-generating distribution simply as $P_X = P_{Y=p,X} + P_{Y=n,X}$. This approach, however correct, does not consider all available information in the Sconf setting. Similar to M_{UU} that uses parameters γ_1 and γ_2 to characterize the data-generating process in UU learning, we use the following lemma that includes the confidence $r(x, x) := P_{y=y|x,x}$ to characterize Sconf learning. Let us abbreviate $r(X, X)$ as r , $P_{X|Y=p}$ as $P_{X/p}$, and $P_{X|Y=n}$ as $P_{X/n}$.

Lemma 9. Assume $\pi_p = 1/2$. Let B (46) be the base distributions and

$$\bar{P} := \begin{pmatrix} P_X P_X \\ P_X P_X \end{pmatrix}.$$

Then, the contamination matrix

$$M_{\text{Sconf}} := \begin{pmatrix} \frac{\pi_p(\pi_p^2 P_{X/p} - \pi_n^2 P_{X/n})}{r - \pi_n} & \frac{\pi_p(\pi_n^2 P_{X/n} - \pi_n^2 P_{X/p})}{r - \pi_n} \\ \frac{\pi_n(\pi_p^2 P_{X/n} - \pi_p^2 P_{X/p})}{\pi_p - r} & \frac{\pi_n(\pi_p^2 P_{X/p} - \pi_n^2 P_{X/n})}{\pi_p - r} \end{pmatrix} \quad (55)$$

which satisfies $\bar{P} = M_{\text{Sconf}}B$, characterizes the data-generating distributions \bar{P} .

Proof. Note that once

$$\left(\frac{r - \pi_n}{\pi_p}\right) P_X P_X = (\pi_p^2 P_{X/p} - \pi_n^2 P_{X/n}) P_{X/p} + (\pi_n^2 P_{X/n} - \pi_n^2 P_{X/p}) P_{X/n} \quad (56)$$

and

$$\left(\frac{\pi_p - r}{\pi_n}\right) P_X P_X = (\pi_p^2 P_{X/n} - \pi_p^2 P_{X/p}) P_{X/p} + (\pi_p^2 P_{X/p} - \pi_n^2 P_{X/n}) P_{X/n}, \quad (57)$$

is obtained, reorganizing the terms gives

$$\begin{pmatrix} P_X P_X \\ P_X P_X \end{pmatrix} = \begin{pmatrix} \frac{\pi_p(\pi_p^2 P_{X/p} - \pi_n^2 P_{X/n})}{r - \pi_n} & \frac{\pi_p(\pi_n^2 P_{X/n} - \pi_n^2 P_{X/p})}{r - \pi_n} \\ \frac{\pi_n(\pi_p^2 P_{X/n} - \pi_p^2 P_{X/p})}{\pi_p - r} & \frac{\pi_n(\pi_p^2 P_{X/p} - \pi_n^2 P_{X/n})}{\pi_p - r} \end{pmatrix} \begin{pmatrix} P_{X/p} \\ P_{X/n} \end{pmatrix} \quad (58)$$

and finishes the proof.

Therefore, we will focus on proving (56) and (57). According to (2) of Cao et al. (2021b), the confidence $r(X, X)$, measuring how likely X and X share the same label, is shown to be

$$r = r(X, X) = \frac{\pi_p^2 P_{X/p} P_{X/p} + \pi_n^2 P_{X/n} P_{X/n}}{P_X P_X}.$$

It implies

$$rP_X P_X = \pi_p^2 P_{X/p} P_{X/p} + \pi_n^2 P_{X/n} P_{X/n}$$

and

$$(1-r)P_X P_X = \pi_p \pi_n (P_{X/p} P_{X/n} + P_{X/n} P_{X/p}).$$

If $\pi_p = 1/2$, $\pi_p - r = 0$ and $r - \pi_n = 0$. As a result, (56) is achieved as follows

$$\begin{aligned} \left(\frac{r - \pi_n}{\pi_p}\right) P_X P_X &= \left(r - \frac{\pi_n}{\pi_p}(1-r)\right) P_X P_X \\ &= \pi_p^2 P_{X/p} P_{X/p} + \pi_n^2 P_{X/n} P_{X/n} - \frac{\pi_n}{\pi_p} \pi_p \pi_n (P_{X/p} P_{X/n} + P_{X/n} P_{X/p}) \\ &= \pi_p^2 P_{X/p} P_{X/p} - \pi_n^2 P_{X/p} P_{X/n} + \pi_n^2 P_{X/n} P_{X/n} - \pi_n^2 P_{X/n} P_{X/p}. \end{aligned}$$

Also, (57) is achieved by having

$$\begin{aligned} \left(\frac{\pi_p - r}{\pi_n}\right) P_X P_X &= \left(\frac{\pi_p}{\pi_n}(1-r) - r\right) P_X P_X \\ &= \frac{\pi_p}{\pi_n} \pi_p \pi_n (P_{X/p} P_{X/n} + P_{X/n} P_{X/p}) - \pi_p^2 P_{X/p} P_{X/p} - \pi_n^2 P_{X/n} P_{X/n} \\ &= \pi_p^2 P_{X/p} P_{X/n} - \pi_p^2 P_{X/p} P_{X/p} + \pi_p^2 P_{X/n} P_{X/p} - \pi_n^2 P_{X/n} P_{X/n}. \end{aligned}$$

□

The equality (58) implies that the inner product of the first row (resp. the second row) of M_{Sconf} and B represents a way (resp. another way) of obtaining $P_X P_X$. Although one might suspect that it is redundant to formulate $P_X P_X$ twice, we show in Section 5.1.6 this expression is crucial to rewrite the classification risk via the proposed framework. Furthermore, comparing $\bar{P} = M_{\text{Sconf}} B$ (58) with $\bar{P} = M_{\text{corr}} B$ (35), we have the reduction path

$$M_{\text{corr}} \quad M_{\text{Sconf}}.$$

Note that M_{Sconf} does not fit the intuition of mutual contamination perfectly; we list Sconf learning in this subsection as all settings share the same base distributions B (46).

4.2 CCN Scenarios

The formulation component (35) also applies to the CCN model. Unlike MCD contaminating conditionals (distributions of X), CCN corrupts class probability functions (labeling distributions). Next, we show how to formulate CCN via (35) and extend the formulation to characterize diverse weakly supervised settings.

In binary classification, CCN (Natarajan et al., 2013; 2017) corrupts the labels by flipping the positive (resp. negative) labels with probability $P_{\tilde{Y}=n/Y=p,X}$ (resp. $P_{\tilde{Y}=p/Y=n,X}$). Specifically,

$$\begin{aligned} P_{\tilde{Y}=p/X} &:= P_{\tilde{Y}=p/Y=p,X} P_{Y=p/X} + P_{\tilde{Y}=p/Y=n,X} P_{Y=n/X}, \\ P_{\tilde{Y}=n/X} &:= P_{\tilde{Y}=n/Y=p,X} P_{Y=p/X} + P_{\tilde{Y}=n/Y=n,X} P_{Y=n/X} \end{aligned} \tag{59}$$

define the contaminated class probability functions. Taking the contamination targets $P_{Y=p/X}$ and $P_{Y=n/X}$ as the base distributions

$$B := \begin{pmatrix} P_{Y=p/X} \\ P_{Y=n/X} \end{pmatrix}$$

and denoting the label-generating distributions $P_{\bar{Y}=p/X}$ and $P_{\bar{Y}=n/X}$ as

$$\bar{P} := \begin{pmatrix} P_{\bar{Y}=p/X} \\ P_{\bar{Y}=n/X} \end{pmatrix},$$

we express (59) in the matrix form as follows

$$\begin{pmatrix} P_{\bar{Y}=p/X} \\ P_{\bar{Y}=n/X} \end{pmatrix} = \begin{pmatrix} P_{\bar{Y}=p/Y=p,X} & P_{\bar{Y}=p/Y=n,X} \\ P_{\bar{Y}=n/Y=p,X} & P_{\bar{Y}=n/Y=n,X} \end{pmatrix} \begin{pmatrix} P_{Y=p/X} \\ P_{Y=n/X} \end{pmatrix}.$$

Comparing with the abstract form $\bar{P} = M_{\text{corr}}B$ (35), we see that

$$M_{\text{CCN}} := \begin{pmatrix} P_{\bar{Y}=p/Y=p,X} & P_{\bar{Y}=p/Y=n,X} \\ P_{\bar{Y}=n/Y=p,X} & P_{\bar{Y}=n/Y=n,X} \end{pmatrix} \quad (60)$$

instantiates the contamination matrix M_{corr} in the CCN setting.

4.2.1 Generalized CCN

The concept of contaminating a *single* label can be extended to generating a *compound* label in the multiclass classification setting. Let 2^Y be the power set of the label space $Y = [K]$. Define $S := 2^Y \setminus \{\emptyset, Y\}$ as the observable space of compound labels³. Since a compound label $S \subseteq S$ consists of an arbitrary number of class indices, one can view the probability of observing S for a given X is governed by several class probabilities $P_{Y=k/X}$. Therefore, generalizing the CCN formulation (59), we define the label-generating process of a compound label S as

$$P_{S/X} = \sum_{k=1}^K P_{S|Y=k,X} P_{Y=k/X},$$

where the role of $P_{S|Y,X}$ is the probability of converting a single label Y to a compound label $S \subseteq S$. Moreover, in CCN, the distribution P_X is not contaminated. Thus, by multiplying P_X on both sides, we obtain the data-generating distribution

$$P_{S,X} = \sum_{k=1}^K P_{S|Y=k,X} P_{Y=k,X}, \quad (61)$$

Viewing $P_{S|Y,X}$ as a contamination probability, we arrange $P_{S=s|Y=k,X}$ into a matrix in the following lemma to formulate the contamination matrix for our generalized CCN (gCCN) setting.

Lemma 10. *Let the elements in S be $\{s_1, s_2, \dots, s_{|S|}\}$. For the gCCN setting, denote the data-generating distributions as*

$$\bar{P} := \begin{pmatrix} P_{S=s_1,X} \\ \vdots \\ P_{S=s_{|S|},X} \end{pmatrix} \quad (62)$$

and the base distributions as

$$B := \begin{pmatrix} P_{Y=1,X} \\ \vdots \\ P_{Y=K,X} \end{pmatrix} = P. \quad (63)$$

³Removing \emptyset and Y is that the empty set does not contain any label information and Y is a trivial case.

Define

$$M_{\text{gCCN}} := \begin{pmatrix} P_{S=s_1|Y=1,X} & P_{S=s_1|Y=2,X} & \cdots & P_{S=s_1|Y=K,X} \\ P_{S=s_2|Y=1,X} & P_{S=s_2|Y=2,X} & \cdots & P_{S=s_2|Y=K,X} \\ \vdots & \vdots & \ddots & \vdots \\ P_{S=s_{|S|}|Y=1,X} & P_{S=s_{|S|}|Y=2,X} & \cdots & P_{S=s_{|S|}|Y=K,X} \end{pmatrix}. \quad (64)$$

Then, $\bar{P} = M_{\text{gCCN}}B$.

The lemma implies that M_{gCCN} is the contamination matrix characterizing \bar{P} of the gCCN setting. Also, note that $\bar{P} = M_{\text{gCCN}}B$ is essentially the matrix form of (61). Moreover, M_{gCCN} generalizes M_{CCN} (60) by extending the label spaces: Both the clean label Y and the contaminated label \bar{Y} belong to $\{p, n\}$ in CCN, while in the gCCN setting, the clean label $Y \in \{1, \dots, K\}$ and the compound label $S \in \{s_1, \dots, s_{|S|}\}$.

Proof. For each $j \in [|S|]$, we have

$$\left(M_{\text{gCCN}}B \right)_j = \sum_{k=1}^K P_{S=s_j|Y=k,X} P_{Y=k,X} = \sum_{k=1}^K P_{S=s_j, Y=k, X} = P_{S=s_j, X} = \bar{P}_j.$$

□

Comparing $\bar{P} = M_{\text{gCCN}}B$ with the formulation framework $\bar{P} = M_{\text{corr}}B$ (35), we have the reduction path

$$M_{\text{corr}} \quad M_{\text{CCN}} \quad M_{\text{gCCN}}.$$

Similar to M_{UU} (49), which induces multiple contamination matrices as special cases of the MCD model, M_{gCCN} also derives several contamination matrices formulating partial- or complementary-label settings, as we will show in the rest of this subsection.

4.2.2 Proper Partial-Label (PPL) Learning (Wu et al., 2023)

For a given example (y, x) and a compound label $s \in S$, we call s a partial-label of x if $y = s$. Statistically speaking, we assume $P_{Y=s|S,X} = 1$. Formally, according to Definition 1 of Wu et al. (2023), if the contamination probability can be defined as

$$P_{S|Y,X} := C(S, X) | [Y = S], \quad (65)$$

via a function $C : S \times \mathcal{X} \rightarrow \mathbb{R}$, we call such a partial-label scenario proper.

Since the discussion above only involves specifying $P_{S|Y,X}$, we replace the entries of M_{gCCN} (64) according to (65) to construct M_{PPL} :

$$\begin{pmatrix} C(s_1, X) | [Y = 1 \quad s_1] & C(s_1, X) | [Y = 2 \quad s_1] & \cdots & C(s_1, X) | [Y = K \quad s_1] \\ C(s_2, X) | [Y = 1 \quad s_2] & C(s_2, X) | [Y = 2 \quad s_2] & \cdots & C(s_2, X) | [Y = K \quad s_2] \\ \vdots & \vdots & \ddots & \vdots \\ C(s_{|S|}, X) | [Y = 1 \quad s_{|S|}] & C(s_{|S|}, X) | [Y = 2 \quad s_{|S|}] & \cdots & C(s_{|S|}, X) | [Y = K \quad s_{|S|}] \end{pmatrix}. \quad (66)$$

The following lemma justifies M_{PPL} as the contamination matrix for PPL learning.

Lemma 11. *Let the elements in S be $\{s_1, s_2, \dots, s_{|S|}\}$. For each $j \in [|S|]$, let the j -th entry of \bar{P} be*

$$\bar{P}_j = P_{S=s_j, X} := C(S = s_j, X) \sum_{k \neq s_j} P_{Y=k, X},$$

which denotes the data-generating distribution of (s_j, X) . Assume the base distributions B and the contamination matrix M_{PPL} are given by (63) and (66), respectively. Then, M_{PPL} satisfies $\bar{P} = M_{\text{PPL}}B$ and characterizes PPL learning (26).

The entry replacement that converts (64) to (66) through (65) gives the reduction path

$$M_{\text{corr}} \quad M_{\text{gCCN}} \quad M_{\text{PPL}}.$$

4.2.3 Provably Consistent Partial-Label (PCPL) Learning (Feng et al., 2020b)

In PCPL, the probability of each partial-label is assumed to be sampled uniformly from all feasible partial-labels. Since there are $2^{K-1} - 1$ feasible partial-labels for every y , the label-converting probability $P_{S=s|Y=y,X}$ is $\frac{1}{2^{K-1}-1}$ if $y = s^4$. It corresponds to assign $C(S, X) = \frac{1}{2^{K-1}-1}$ in (65). Hence, we obtain

$$C(S, X) | [Y \quad S] := \frac{1}{2^{K-1}-1} | [Y \quad S], \quad (67)$$

which reduces the label-converting process of PPL to that of PCPL and recovers (5) of Feng et al. (2020b).

Then, replacing entries in (66) via (67), we obtain the contamination matrix of PCPL learning

$$M_{\text{PCPL}} := \frac{1}{2^{K-1}-1} \begin{pmatrix} | [Y=1 \quad s_1] & | [Y=2 \quad s_1] & \cdots & | [Y=K \quad s_1] \\ | [Y=1 \quad s_2] & | [Y=2 \quad s_2] & \cdots & | [Y=K \quad s_2] \\ \vdots & \vdots & \ddots & \vdots \\ | [Y=1 \quad s_{|S|}] & | [Y=2 \quad s_{|S|}] & \cdots & | [Y=K \quad s_{|S|}] \end{pmatrix} \quad (68)$$

and the reduction path

$$M_{\text{corr}} \quad M_{\text{gCCN}} \quad M_{\text{PPL}} \quad M_{\text{PCPL}}.$$

M_{PCPL} characterizing the data-generating process of PCPL is justified by the following lemma, whose proof follows the same steps as that for Lemma 11.

Lemma 12. *Let the elements in S be $\{s_1, s_2, \dots, s_{|S|}\}$. For each $j \in [|S|]$, let the j -th entry of \bar{P} be*

$$\bar{P}_j = P_{S=s_j, X} := \frac{1}{2^{K-1}-1} \sum_{k \in s_j} P_{Y=k, X},$$

which denotes the data-generating distribution of (s_j, X) . Assume the base distributions B and the contamination matrix M_{PCPL} are given by (63) and (68), respectively. Then, M_{PCPL} satisfies $\bar{P} = M_{\text{PCPL}}B$ and characterizes PCPL learning (24).

4.2.4 Multi-Complementary-Label (MCL) Learning (Feng et al., 2020a)

Recall the discussions in Sections 2.2.9 and 2.2.10 that a complementary-label contains the exclusion information of a true label. That is, for a data example (y, x) , we call a set of class indices $\bar{s} \subseteq S = 2^Y \setminus \{y, Y\}$ an MCL of x if \bar{s} does not contain y .

Denote $S := \{\bar{s}_1, \bar{s}_2, \dots, \bar{s}_N\}$. The equivalence

$$\sum_{d=1}^{K-1} P_{|\bar{S}|=d} \cdot \frac{1}{\binom{K-1}{|\bar{S}|}} \sum_{Y/\bar{S}} P_{Y, X} | [|\bar{S}|=d] = \begin{cases} \sum_{d=1}^{K-1} P_{|\bar{S}|=d} \cdot \frac{1}{\binom{K-1}{|\bar{S}|}} \sum_{Y/\bar{S}} P_{Y, X}, & \text{if } |\bar{S}|=d, \\ 0, & \text{otherwise} \end{cases}$$

⁴There are $2^{Y \setminus \{y\}} \setminus \{\emptyset, Y\} = 2^{K-1} - 1$ combinations whose union with $\{y\}$ are partial-labels of y .

allows us to define the data-generating distribution of MCL (22) as

$$\bar{P} := \begin{pmatrix} P_{\bar{S}=\bar{s}_1, X} \\ \vdots \\ P_{\bar{S}=\bar{s}_N, X} \end{pmatrix}, \quad (69)$$

where for each $j \in [N]$,

$$P_{\bar{S}=\bar{s}_j, X} := \sum_{d=1}^{K-1} \frac{P_{|\bar{s}_j|=d}}{\binom{K-1}{|\bar{s}_j|}} \sum_{Y/\bar{s}_j} P_{Y, X} \mathbb{1}[|\bar{s}_j|=d]. \quad (70)$$

Let (y, x) be fixed. The data-generating process of MCL proposed by Feng et al. (2020a) is that one first samples a size d with probability $P_{|\bar{s}_j|=d}$, and then samples a \bar{s} uniformly at random from $\{\bar{s}_{d,1}, \bar{s}_{d,2}, \dots, \bar{s}_{d,N_d}\}$ S , where \bar{s}_d means a set of size d excluding y and N_d is the total number of those sets. Note that $N_d = \binom{K-1}{d}$ since we remove y from Y and then choose a set of size d to form a \bar{s}_d . Furthermore, we need a more complicated lower index system to distinguish $\{\bar{s}_{d,1}, \bar{s}_{d,2}, \dots, \bar{s}_{d,N_d}\}$ from $S = \{\bar{s}_1, \bar{s}_2, \dots, \bar{s}_N\}$ since d ranges from 1 to $K-1$ and $\sum_{d=1}^{K-1} N_d = \sum_{d=1}^{K-1} \binom{K-1}{d} = 2^K - 2 = |S|$. According to this mechanism, we construct M_{MCL} :

$$\begin{pmatrix} \frac{P_{|\bar{s}_1|=\bar{s}_1}|}{\binom{K-1}{|\bar{s}_1|}} \mathbb{1}[Y=1/\bar{s}_1] & \frac{P_{|\bar{s}_1|=\bar{s}_1}|}{\binom{K-1}{|\bar{s}_1|}} \mathbb{1}[Y=2/\bar{s}_1] & \cdots & \frac{P_{|\bar{s}_1|=\bar{s}_1}|}{\binom{K-1}{|\bar{s}_1|}} \mathbb{1}[Y=K/\bar{s}_1] \\ \frac{P_{|\bar{s}_1|=\bar{s}_2}|}{\binom{K-1}{|\bar{s}_2|}} \mathbb{1}[Y=1/\bar{s}_2] & \frac{P_{|\bar{s}_1|=\bar{s}_2}|}{\binom{K-1}{|\bar{s}_2|}} \mathbb{1}[Y=2/\bar{s}_2] & \cdots & \frac{P_{|\bar{s}_1|=\bar{s}_2}|}{\binom{K-1}{|\bar{s}_2|}} \mathbb{1}[Y=K/\bar{s}_2] \\ \vdots & \vdots & \ddots & \vdots \\ \frac{P_{|\bar{s}_1|=\bar{s}_N}|}{\binom{K-1}{|\bar{s}_N|}} \mathbb{1}[Y=1/\bar{s}_N] & \frac{P_{|\bar{s}_1|=\bar{s}_N}|}{\binom{K-1}{|\bar{s}_N|}} \mathbb{1}[Y=2/\bar{s}_N] & \cdots & \frac{P_{|\bar{s}_1|=\bar{s}_N}|}{\binom{K-1}{|\bar{s}_N|}} \mathbb{1}[Y=K/\bar{s}_N] \end{pmatrix}. \quad (71)$$

The following lemma justifies M_{MCL} as the contamination matrix for MCL learning.

Lemma 13. *Suppose the base distributions B , the contamination matrix M_{MCL} , and the data-generating distributions \bar{P} are given by (63), (71), and (70), respectively. Then, M_{MCL} satisfies $\bar{P} = M_{\text{MCL}}B$ and characterizes MCL (22).*

At the first sight, M_{MCL} (71) does not resemble M_{PCPL} (68) or M_{PPL} (66). The subtle connection can be established via a relation between partial-label and complementary-label. Recall

$$2^Y \setminus \{y, Y\} = S = \{\bar{s}_1, \bar{s}_2, \dots, \bar{s}_N\},$$

where \bar{s}_j is a MCL. From the partial-label perspective, we can establish the following set equality relationship:

$$2^Y \setminus \{y, Y\} = S = \{\bar{s}_1, \bar{s}_2, \dots, \bar{s}_N\} = \{s_1 := Y \setminus \bar{s}_1, s_2 := Y \setminus \bar{s}_2, \dots, s_N := Y \setminus \bar{s}_N\}. \quad (72)$$

This is because for every $\bar{s} \in S$, there is a $s \in S$ such that $s := Y \setminus \bar{s}$. The intuition behind (72) is that if \bar{s} is an MCL of x , then $s := Y \setminus \bar{s}$ must be a partial-label of x . Therefore, we can also use the set of partial-labels $\{s_1, s_2, \dots, s_N\}$ to denote S . The following lemma exploits this relation to show that M_{MCL} is indeed a special case of M_{PPL} .

Lemma 14. *Assign each (s, k) entry of M_{PPL} (66) with*

$$C(s, X) \mathbb{1}[Y=k \setminus s] := \frac{P_{|S|=|s|}}{\binom{K-1}{|s|-1}} \mathbb{1}[Y=k \setminus s]. \quad (73)$$

Then, the resulting matrix

$$M_{\text{MCL}} = \begin{pmatrix} \frac{P_{|S|=|s_1|}}{\binom{K-1}{|s_1|-1}} | [Y = 1 \quad s_1] & \frac{P_{|S|=|s_1|}}{\binom{K-1}{|s_1|-1}} | [Y = 2 \quad s_1] & \cdots & \frac{P_{|S|=|s_1|}}{\binom{K-1}{|s_1|-1}} | [Y = K \quad s_1] \\ \frac{P_{|S|=|s_2|}}{\binom{K-1}{|s_2|-1}} | [Y = 1 \quad s_2] & \frac{P_{|S|=|s_2|}}{\binom{K-1}{|s_2|-1}} | [Y = 2 \quad s_2] & \cdots & \frac{P_{|S|=|s_2|}}{\binom{K-1}{|s_2|-1}} | [Y = K \quad s_2] \\ \vdots & \vdots & \ddots & \vdots \\ \frac{P_{|S|=|s_N|}}{\binom{K-1}{|s_N|-1}} | [Y = 1 \quad s_N] & \frac{P_{|S|=|s_N|}}{\binom{K-1}{|s_N|-1}} | [Y = 2 \quad s_N] & \cdots & \frac{P_{|S|=|s_N|}}{\binom{K-1}{|s_N|-1}} | [Y = K \quad s_N] \end{pmatrix} \quad (74)$$

is equivalent to M_{MCL} (71) under the relationship (72).

Proof. Note that for every $j \in [N]$, $s_j = Y \setminus \bar{s}_j$. This implies $P_{|\bar{s}|=|\bar{s}_j|} = P_{|S|=|s_j|}$, $\binom{K-1}{|s_j|-1} = \binom{K-1}{|\bar{s}_j|}$, and $| [Y \quad s_j] = | [Y \quad \bar{s}_j]$ hold for every $j \in [N]$. Therefore, for each (j, k) entry in M_{MCL} and M_{MCL} (71), we have

$$\frac{P_{|S|=|s_j|}}{\binom{K-1}{|s_j|-1}} | [Y = k \quad s_j] = \frac{P_{|\bar{s}|=|\bar{s}_j|}}{\binom{K-1}{|\bar{s}_j|}} | [Y = k \quad \bar{s}_j]. \quad (75)$$

□

The assignment rule (73) implies the reduction path

$$M_{\text{corr}} \quad M_{\text{gCCN}} \quad M_{\text{PPL}} \quad M_{\text{MCL}}.$$

Comparing (73) with (67) of PCPL, we see that MCL and PCPL can be viewed as different ways of composing $P_{Y,X}$ to generate a partial-label, with weights $\frac{P_{|S|=|s|}}{\binom{K-1}{|s|-1}}$ and $\frac{1}{2^{K-1}-1}$, respectively.

4.2.5 Complementary-Label (CL) Learning (Ishida et al., 2019)

As a special case of MCL (Section 2.2.10), we can construct the contamination matrix M_{CL} from M_{MCL} . The set of all CLs is composed of MCL with size 1: $\{1\}, \dots, \{K\}$. Therefore, we assign values in M_{MCL} (71) as follows. For each $\bar{s} \in S$, $P_{|\bar{s}|=|\bar{s}_j|} = 1$ if $|\bar{s}| = 1$ and $P_{|\bar{s}|=|\bar{s}_j|} = 0$ if $|\bar{s}| > 1$. Dropping all-zero rows, we obtain from (71) the contamination matrix of CL learning

$$\begin{aligned} M_{\text{CL}} &:= \begin{pmatrix} \frac{1}{K-1} | [Y = 1 \quad \{1\}] & \frac{1}{K-1} | [Y = 2 \quad \{1\}] & \cdots & \frac{1}{K-1} | [Y = K \quad \{1\}] \\ \frac{1}{K-1} | [Y = 1 \quad \{2\}] & \frac{1}{K-1} | [Y = 2 \quad \{2\}] & \cdots & \frac{1}{K-1} | [Y = K \quad \{2\}] \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{K-1} | [Y = 1 \quad \{K\}] & \frac{1}{K-1} | [Y = 2 \quad \{K\}] & \cdots & \frac{1}{K-1} | [Y = K \quad \{K\}] \end{pmatrix} \\ &= \frac{1}{K-1} \begin{pmatrix} 0 & 1 & \cdots & 1 \\ 1 & 0 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 0 \end{pmatrix} \end{aligned} \quad (76)$$

and the reduction path

$$M_{\text{corr}} \quad M_{\text{gCCN}} \quad M_{\text{PPL}} \quad M_{\text{MCL}} \quad M_{\text{CL}}.$$

Furthermore, it is easy to verify that given B (63), for any $j \in [K]$,

$$(M_{\text{CL}}B)_j = \sum_{Y=j} \frac{1}{K-1} P_{Y,X} = P_{\bar{S}=j,X},$$

which corresponds to formulation (20). Hence, we have the following.

Lemma 15. M_{CL} (76) is the contamination matrix characterizing the data-generating distribution $P_{\bar{s},X}$ (20) of CL learning.

4.3 Confidence-based Scenarios

At first sight, there seems to be no connection between “contamination” and single-class classification (Cao et al., 2021a). However, the following derivation

$$\frac{P_{Y=y_s|X}}{P_{Y=j|X}} \cdot P_{Y=j,X} = \frac{P_{Y=y_s|X}}{P_{Y=j|X}} \cdot P_{Y=j|X} P_X = P_{Y=y_s,X} \quad (77)$$

reveals a way to *contaminate* a clean joint probability $P_{Y=j,X}$ to the joint probability $P_{Y=y_s,X}$ of a designated class y_s via confidence weighting $\frac{P_{Y=y_s|X}}{P_{Y=j|X}}$. As we will see in the rest of this subsection, the confidence weights are the key elements in formulating the contamination matrices for the confidence-based WSL settings.

4.3.1 Subset Confidence (Sub-Conf) Learning (Cao et al., 2021a)

Let $Y_s \subseteq [K]$ be a subset of classes. Viewing Y_s as a “superclass”, such that every instance x of (y, x) will be labeled Y_s if $y \in Y_s$, we can define its class prior as $P_{Y=Y_s} = \pi_{Y_s} := \sum_{y \in Y_s} \pi_y$ and its class probability function as $P_{Y=Y_s|X} := \sum_{y \in Y_s} P_{Y=y|X}$. Substituting the designated class y_s in (77) with the superclass Y_s ,

$$\frac{P_{Y=Y_s|X}}{P_{Y=j|X}} \cdot P_{Y=j,X} = \frac{P_{Y=Y_s|X}}{P_{Y=j|X}} \cdot P_{Y=j|X} P_X = P_{Y=Y_s,X} \quad (78)$$

shows that no matter what joint distribution $P_{Y=j,X}$ to begin with, the confidence weight $\frac{P_{Y=Y_s|X}}{P_{Y=j|X}}$ twists that joint distribution so that every observed data appears to be sampled from the same superclass distribution $P_{Y=Y_s,X}$. The following lemma leverages the observation to specify the contamination matrix M_{Sub} characterizing Sub-Conf learning.

Lemma 16. Denote the base distributions as

$$B := \begin{pmatrix} P_{Y=1,X} \\ \vdots \\ P_{Y=K,X} \end{pmatrix} = P \quad (79)$$

and the data-generating distributions as

$$\bar{P} := \begin{pmatrix} P_{Y=Y_s,X} \\ \vdots \\ P_{Y=Y_s,X} \end{pmatrix}.$$

Inserting the confidence weights into the identity matrix, we define the contamination matrix

$$M_{\text{Sub}} := \begin{pmatrix} \frac{P_{Y=Y_s|X}}{P_{Y=1|X}} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{P_{Y=Y_s|X}}{P_{Y=K|X}} \end{pmatrix}. \quad (80)$$

Then, $\bar{P} = M_{\text{Sub}}B$, and M_{Sub} characterizes the data-generating process of Sub-Conf learning (30).

Proof. For each $j \in [K]$, $(M_{\text{Sub}}B)_j = P_{Y=Y_s,X}$ follows from (78). Thus, $\bar{P} = M_{\text{Sub}}B$. It further implies all observed instances are labeled with the same superclass Y_s , meaning we can drop the observed labels, and the observed examples $\{x_i\}_{i=1}^n$ is equivalent to a set of i.i.d. samples from $P_{X|Y=Y_s}$ (30). \square

Comparing $\bar{P} = M_{\text{Sub}}B$ with the formulation framework $\bar{P} = M_{\text{corr}}B$ (35), we observe that in Sub-Conf learning, M_{corr} is realized as M_{Sub} :

$$M_{\text{corr}} = M_{\text{Sub}}.$$

4.3.2 Single-Class Confidence (SC-Conf) Learning (Cao et al., 2021a)

We compare the formulation of SC-Conf (28) with Sub-Conf (30) and observe that SC-Conf is a special case of Sub-Conf when $Y_s = \{y_s\}$ being a singleton. Thus, we straightforwardly obtain the matrix formulation of SC-Conf from Lemma 16 by replacing Y_s in (80) with y_s :

Lemma 17. *Let the base distributions B be defined by (79) and the data-generating distributions be defined by*

$$\bar{P} := \begin{pmatrix} P_{Y=y_s, X} \\ \vdots \\ P_{Y=y_s, X} \end{pmatrix}.$$

Define the contamination matrix

$$M_{\text{SC}} := \begin{pmatrix} \frac{P_{Y=y_s, X}}{P_{Y=1, X}} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{P_{Y=y_s, X}}{P_{Y=K, X}} \end{pmatrix} \quad (81)$$

by substituting Y_s in (80) with y_s . Then, $\bar{P} = M_{\text{SC}}B$ and M_{SC} characterizes the data-generating process of SC-Conf learning (28).

Since SC-Conf is a special case of Sub-Conf, we have the reduction path

$$M_{\text{corr}} = M_{\text{Sub}} = M_{\text{SC}}.$$

4.3.3 Positive-confidence (Pconf) Learning (Ishida et al., 2018)

Comparing (6) with (28), we see that Pconf is a special case of SC-Conf when $K = 2$ and $y_s = p$ since $r_n(X) = 1 - r_p(X)$. A further modification to (81) we obtain the contamination matrix M_{Pconf} characterizing Pconf learning.

Lemma 18. *Let $B := \begin{pmatrix} P_{Y=p, X} \\ P_{Y=n, X} \end{pmatrix} = P$ and $\bar{P} := \begin{pmatrix} P_{Y=p, X} \\ P_{Y=p, X} \end{pmatrix}$. Define*

$$M_{\text{Pconf}} := \begin{pmatrix} \frac{P_{Y=p, X}}{P_{Y=p, X}} & 0 \\ 0 & \frac{P_{Y=p, X}}{P_{Y=n, X}} \end{pmatrix}. \quad (82)$$

Then, $\bar{P} = M_{\text{Pconf}}B$, and M_{Pconf} characterizes the data-generating process of Pconf learning (6).

The entry replacement that converts (81) to (82) implies the reduction path

$$M_{\text{corr}} = M_{\text{Sub}} = M_{\text{SC}} = M_{\text{Pconf}}.$$

4.3.4 Soft-Label Learning (Ishida et al., 2023)

The difference between the soft-label and the previous confidence-based settings (Sub-Conf, SC-Conf, and Pconf) is how x is sampled. The sample distributions condition on the label information in the previous settings, while that in soft-label is P_X . Replacing the confidence weight $\frac{P_{Y=y_s|X}}{P_{Y=j|X}}$ in (77) with $\frac{1}{P_{Y=j|X}}$,

$$\frac{1}{P_{Y=j|X}} \cdot P_{Y=j,X} = P_X$$

explains how to convert $P_{Y=j,X}$ to P_X . Therefore, filling the j -th diagonal entry of the identity matrix with $\frac{1}{P_{Y=j|X}}$, we obtain the contamination matrix M_{Soft} for soft-label learning:

Lemma 19. *Let the base distributions B be defined by (79). Denote the data-generating distribution as*

$$\bar{P} := \begin{pmatrix} P_X \\ \vdots \\ P_X \end{pmatrix}. \quad (83)$$

Define the contamination matrix

$$M_{\text{Soft}} := \begin{pmatrix} \frac{1}{P_{Y=1|X}} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{P_{Y=K|X}} \end{pmatrix}. \quad (84)$$

Then, $\bar{P} = M_{\text{Soft}}B$, and M_{Soft} characterizes the data-generating process in (33).

Unlike SC-Conf and Pconf, which are special cases of Sub-Conf with Y_s taking only one label, the generation process of a soft-label can be viewed as assigning $Y_s := [K]$. Considering the entire label space results in $P_{Y=[K]|X} = 1$; it coincides with the meaning of P_X that samples x regardless of the labels. Although technically the soft-label setting is not a special case of Sub-Conf (recalling the $Y_s \in [K]$ assumption from Section 4.3.1), M_{Soft} (84) is reduced from M_{Sub} (80) by realizing $P_{Y=Y_s|X}$ as $P_{Y=[K]|X} = 1$. Therefore, we obtain the following reduction path

$$M_{\text{corr}} \quad M_{\text{Sub}} \quad M_{\text{Soft}}.$$

5 Risk Rewrite via Decontamination

We have demonstrated the capability of the proposed formulation component (35) in the last section. This section shows how the proposed framework provides a unified methodology for solving the risk rewrite problem. Specifically, given each contamination matrix described in Section 4, we show how to construct the corrected losses (38) to perform the risk rewrite via (39). We then recover each rewrite to the corresponding form reported in the literature to justify its feasibility. Because this paper focuses on a unified methodology for rewriting the classification risk instead of the designs of practical training objectives, we assume the required parameters are given or can be estimated accurately from the observed data.

Similar to the previous section, we only provide proofs that appear for the first time to avoid repeating similar proofs. For auxiliary lemmas and results whose proofs are similar to the previous ones, we refer to the omitted proofs in Appendix D. In particular, the omitted proofs in Section 5.1 can be found in Appendix D.1, those in Section 5.2 can be found in Appendix D.2, and those in Section 5.3 can be found in Appendix D.3.

5.1 MCD Scenarios

We apply the framework to conduct the risk rewrites for WSLs formulated in Section 4.1 and summarized in Table 7. A general approach is to show that the inversion method discussed in Proposition 1 provides the decontamination matrix M_{corr}^T required in (38).

5.1.1 Unlabeled-Unlabeled (UU) Learning

We justify the proposed framework for UU learning via the following steps.

Step 1: Corrected Loss Design and Risk Rewrite.

The three milestones in the proposed framework are (1) finding the contamination matrix M_{corr} that characterizes the data-generating process (35) of a weakly supervised scenario, (2) finding the decontamination matrix M_{corr}^\dagger that compensates for the contamination effect (37), which is then used in (3) the construction of corrected losses (38) for the risk rewrite (39).

Section 4.1.1 has reached the first milestone (35) as (48) of the form $\bar{P} = M_{\text{UU}}B$ finds

$$M_{\text{UU}} = \begin{pmatrix} 1 - \gamma_1 & \gamma_1 \\ \gamma_2 & 1 - \gamma_2 \end{pmatrix}$$

that connects the data-generating distributions $\bar{P} = \begin{pmatrix} P_{U_1} \\ P_{U_2} \end{pmatrix}$ and the base distributions $B = \begin{pmatrix} P_{X|Y=p} \\ P_{X|Y=n} \end{pmatrix}$.

Note that B is not the risk-defining distributions $P = \begin{pmatrix} P_{Y=p,X} \\ P_{Y=n,X} \end{pmatrix}$, we need an additional step before reaching the second milestone. To further link \bar{P} with P , we still need a M_{trsf} that satisfies $B = M_{\text{trsf}}P$. Introducing the prior matrix $\Pi = \begin{pmatrix} \pi_p & 0 \\ 0 & \pi_n \end{pmatrix}$, we see that choosing $M_{\text{trsf}} := \Pi^{-1}$ fulfills the need:

$$M_{\text{trsf}}P = \begin{pmatrix} \pi_p^{-1} & 0 \\ 0 & \pi_n^{-1} \end{pmatrix} \begin{pmatrix} P_{Y=p,X} \\ P_{Y=n,X} \end{pmatrix} = \begin{pmatrix} \frac{P_{Y=p,X}}{P_{Y=p}} \\ \frac{P_{Y=n,X}}{P_{Y=n}} \end{pmatrix} = \begin{pmatrix} P_{X|Y=p} \\ P_{X|Y=n} \end{pmatrix} = B.$$

Hence, we can instantiate $\bar{P} = M_{\text{corr}}M_{\text{trsf}}P$ (36) as

$$\begin{pmatrix} P_{U_1} \\ P_{U_2} \end{pmatrix} = M_{\text{UU}}\Pi^{-1} \begin{pmatrix} P_{Y=p,X} \\ P_{Y=n,X} \end{pmatrix} \quad (85)$$

in UU learning.

Next, we use Proposition 1 to derive the decontamination matrix M_{corr}^\dagger to reach the second milestone (37).

Corollary 20. Assume M_{UU} in (85) is invertible. Then, defining the decontamination matrix for UU learning as

$$M_{\text{UU}}^\dagger := \Pi M_{\text{UU}}^{-1}$$

gives rise to $M_{\text{UU}}^\dagger \bar{P} = P$.

Proof. Suggested by Proposition 1, the inverse matrix ΠM_{UU}^{-1} cancels out the contamination brought by $M_{\text{UU}}\Pi^{-1}$ in (85). Assigning $M_{\text{UU}}^\dagger = \Pi M_{\text{UU}}^{-1}$ and repeating the proof of Proposition 1, we have

$$M_{\text{UU}}^\dagger \bar{P} = \Pi M_{\text{UU}}^{-1} \bar{P} = \Pi M_{\text{UU}}^{-1} M_{\text{UU}} \Pi^{-1} P = P$$

that completes the proof. \square

Now we will move on to the third milestone. With M_{UU}^\dagger in hand, we devise the corrected losses \bar{L} to achieve the risk rewrite for UU learning. We denote the corrected loss at the \bar{k} -th entry of \bar{L} as $\bar{\ell}_{\bar{k}} := \ell_{\bar{Y}=\bar{k}}(g(X))$, where $\bar{k} \in \bar{Y}$ is a class of the observed data⁵. In UU learning, $\bar{Y} = \{U_1, U_2\}$. The following theorem proves rewrite (9) in Section 2.2.3.

⁵The definition of the corrected loss $\bar{\ell}_{\bar{k}}$ is in contrast to the original loss $\ell_k := \ell_{Y=k}(g(X))$.

Theorem 21. Let $\gamma_1, \gamma_2 > 0$ and $\gamma_1 + \gamma_2 = 1$. Then, M_{UU}^\dagger defined in Corollary 20 is feasible. Moreover, the vector of corrected losses suggested by (38)

$$\begin{pmatrix} \bar{\ell}_{\text{U}_1} & \bar{\ell}_{\text{U}_2} \end{pmatrix} = \bar{L} := L M_{\text{UU}}^\dagger$$

with

$$\begin{aligned} \bar{\ell}_{\text{U}_1} &= \frac{(1-\gamma_2)\pi_{\text{p}}}{1-\gamma_1-\gamma_2} \ell_{\text{p}} + \frac{-\gamma_2\pi_{\text{n}}}{1-\gamma_1-\gamma_2} \ell_{\text{n}}, \\ \bar{\ell}_{\text{U}_2} &= \frac{-\gamma_1\pi_{\text{p}}}{1-\gamma_1-\gamma_2} \ell_{\text{p}} + \frac{(1-\gamma_1)\pi_{\text{n}}}{1-\gamma_1-\gamma_2} \ell_{\text{n}} \end{aligned} \quad (86)$$

achieves the following risk rewrite:

$$R(g) = \mathbb{E}_{\text{U}_1} [\bar{\ell}_{\text{U}_1}] + \mathbb{E}_{\text{U}_2} [\bar{\ell}_{\text{U}_2}]. \quad (87)$$

Proof. Since $\gamma_1 + \gamma_2 = 1$,

$$M_{\text{UU}}^{-1} = \begin{pmatrix} 1-\gamma_1 & \gamma_1 \\ \gamma_2 & 1-\gamma_2 \end{pmatrix}^{-1} = \begin{pmatrix} \frac{1-\gamma_2}{1-\gamma_1-\gamma_2} & \frac{-\gamma_1}{1-\gamma_1-\gamma_2} \\ \frac{-\gamma_2}{1-\gamma_1-\gamma_2} & \frac{1-\gamma_1}{1-\gamma_1-\gamma_2} \end{pmatrix}$$

exists. Thus, it is feasible for us to define $M_{\text{UU}}^\dagger := \Pi M_{\text{UU}}^{-1}$ according to Corollary 20. Following (38), we construct

$$\bar{L} := L M_{\text{UU}}^\dagger$$

and obtain

$$\begin{aligned} \begin{pmatrix} \bar{\ell}_{\text{U}_1} & \bar{\ell}_{\text{U}_2} \end{pmatrix} &= L \Pi M_{\text{UU}}^{-1} \\ &= \begin{pmatrix} \ell_{\text{p}} & \ell_{\text{n}} \end{pmatrix} \begin{pmatrix} \pi_{\text{p}} & 0 \\ 0 & \pi_{\text{n}} \end{pmatrix} \begin{pmatrix} \frac{1-\gamma_2}{1-\gamma_1-\gamma_2} & \frac{-\gamma_1}{1-\gamma_1-\gamma_2} \\ \frac{-\gamma_2}{1-\gamma_1-\gamma_2} & \frac{1-\gamma_1}{1-\gamma_1-\gamma_2} \end{pmatrix} \\ &= \begin{pmatrix} \ell_{\text{p}} & \ell_{\text{n}} \end{pmatrix} \begin{pmatrix} \frac{(1-\gamma_2)\pi_{\text{p}}}{1-\gamma_1-\gamma_2} & \frac{-\gamma_1\pi_{\text{p}}}{1-\gamma_1-\gamma_2} \\ \frac{-\gamma_2\pi_{\text{n}}}{1-\gamma_1-\gamma_2} & \frac{(1-\gamma_1)\pi_{\text{n}}}{1-\gamma_1-\gamma_2} \end{pmatrix} \end{aligned} \quad (88)$$

that gives (86).

Next, with the critical component \bar{L} in hand, applying (39), we obtain

$$\begin{aligned} R(g) &= \int_{\mathcal{X}} \bar{L} \bar{P} \, dx \\ &= \int_{\mathcal{X}} (P_{\text{U}_1} \bar{\ell}_{\text{U}_1} + P_{\text{U}_2} \bar{\ell}_{\text{U}_2}) \, dx \\ &= \mathbb{E}_{\text{U}_1} [\bar{\ell}_{\text{U}_1}] + \mathbb{E}_{\text{U}_2} [\bar{\ell}_{\text{U}_2}], \end{aligned} \quad (89)$$

where the first equality holds since according to Corollary 20,

$$\bar{L} \bar{P} = L M_{\text{UU}}^\dagger \bar{P} = L P.$$

□

In (88), we do not need to specify the instance in ℓ_p and ℓ_n to be x^{u_1} or x^{u_2} since the equality holds for any instance x . We only need to distinguish x^{u_1} from x^{u_2} when the corrected losses multiply the data distributions. In particular, the detailed form of rewrite (89) using (8) is

$$\begin{aligned} R(g) &= \mathbb{E}_{U_1} [\bar{\ell}_{U_1}] + \mathbb{E}_{U_2} [\bar{\ell}_{U_2}] \\ &= \mathbb{E}_{x^{u_1}} \mathbb{P}_{U_1} \left[\frac{(1-\gamma_2)\pi_p}{1-\gamma_1-\gamma_2} \ell_p(X^{u_1}) + \frac{-\gamma_2\pi_n}{1-\gamma_1-\gamma_2} \ell_n(X^{u_1}) \right] \\ &\quad + \mathbb{E}_{x^{u_2}} \mathbb{P}_{U_2} \left[\frac{-\gamma_1\pi_p}{1-\gamma_1-\gamma_2} \ell_p(X^{u_2}) + \frac{(1-\gamma_1)\pi_n}{1-\gamma_1-\gamma_2} \ell_n(X^{u_2}) \right]. \end{aligned}$$

The freedom from specifying x in (88) eliminates the notational burden of distinguishing $\ell_Y(X^{u_1})$ from $\ell_Y(X^{u_2})$, allowing us to exploit the advantage of matrix multiplication while constructing the corrected losses. The freedom also enables separated treatments for the data distributions (e.g., formulating $\bar{P} = M_{UU}\Pi^{-1}P$) and the corrected losses (e.g., devising $\bar{L} = L - M_{UU}^T$).

Step 2: Recovering the previous result(s).

Lastly, we verify the feasibility of our rewrite by showing that our rewrite corresponds to an existing result. By parameter substitution, we replace γ_1 with $1 - \theta$, γ_2 with θ , π_n with $1 - \pi_p$, ℓ_p with $\ell(g(X))$, and ℓ_n with $\ell(-g(X))$. Then, (86) becomes

$$\begin{aligned} \frac{(1-\theta)\pi_p}{\theta-\theta} \ell(g(X)) + \frac{-\theta(1-\pi_p)}{\theta-\theta} \ell(-g(X)) &= \bar{\ell}_+(g(X)), \\ \frac{\theta(1-\pi_p)}{\theta-\theta} \ell(-g(X)) + \frac{-(1-\theta)\pi_p}{\theta-\theta} \ell(g(X)) &= \bar{\ell}_-(-g(X)), \end{aligned}$$

recovering the corrected loss functions (8) and the constants reported in Theorem 4 of Lu et al. (2019).

5.1.2 Positive-Unlabeled (PU) Learning

Recall that all WSLs discussed in Section 4.1 share the same base distributions B (46). Further, as shown in Table 7, the contamination matrix of every WSL scenario beneath UU learning except M_{Scnf} is a child of M_{UU} on the reduction graph. It means $\bar{P} = M_{UU}\Pi^{-1}P$ (85) is a general form for every child scenario in Table 7 (with different realizations of γ_1 and γ_2). Hence, we can reuse Theorem 21 to conduct the risk rewrite for every child scenario on the reduction graph. PU learning is the first of such examples.

Step 1: Corrected Loss Design and Risk Rewrite.

By the following corollary, we prove the rewrite (5) in Section 2.2.1.

Corollary 22. *For PU learning, the classification risk can be rewritten as*

$$R(g) = \mathbb{E}_P [\bar{\ell}_P] + \mathbb{E}_U [\bar{\ell}_U], \quad (90)$$

where

$$\begin{aligned} \bar{\ell}_P &= \pi_p \ell_p - \pi_p \ell_n, \\ \bar{\ell}_U &= \ell_n. \end{aligned} \quad (91)$$

Step 2: Recovering the previous result(s).

Since P_P is $P_{X|Y=p}$ and P_U is P_X , we swap the notations to obtain

$$\begin{aligned} R(g) &= \mathbb{E}_P [\bar{\ell}_P] + \mathbb{E}_U [\bar{\ell}_U] \\ &= \mathbb{E}_P [\pi_p \ell_p - \pi_p \ell_n] + \mathbb{E}_U [\ell_n] \\ &= \pi_p \mathbb{E}_{X|Y=p} [\ell_p] - \pi_p \mathbb{E}_{X|Y=p} [\ell_n] + \mathbb{E}_X [\ell_n] \end{aligned}$$

from (90), which corresponds to the risk estimators (2) in Kiryo et al. (2017) and (3) in du Plessis et al. (2015).

Moreover, with an additional symmetric assumption of $\ell_p + \ell_n = 1$, one further obtains

$$\begin{aligned}
R(g) &= \pi_p \mathbb{E}_{X|Y=p} [\ell_p] - \pi_p \mathbb{E}_{X|Y=p} [\ell_n] + \mathbb{E}_X [\ell_n] \\
&= \pi_p \mathbb{E}_{X|Y=p} [\ell_p] - \pi_p \mathbb{E}_{X|Y=p} [1 - \ell_p] + \mathbb{E}_X [\ell_n] \\
&= \pi_p \mathbb{E}_{X|Y=p} [\ell_p] - \pi_p \mathbb{E}_{X|Y=p} [1] + \pi_p \mathbb{E}_{X|Y=p} [\ell_p] + \mathbb{E}_X [\ell_n] \\
&= 2\pi_p \mathbb{E}_{X|Y=p} [\ell_p] - \pi_p + \mathbb{E}_X [\ell_n].
\end{aligned}$$

This expression recovers several risk rewrites such as (4) of Kiryo et al. (2017), (3) of Niu et al. (2016), (2) of du Plessis et al. (2015)⁶, and (3) of du Plessis et al. (2014).

5.1.3 Similar-Unlabeled (SU) Learning

According to Table 7, M_{SU} is a child of M_{UU} on the reduction graph. Thus, we can follow the same steps illustrated in Section 5.1.2 to justify the proposed framework.

Step 1: Corrected Loss Design and Risk Rewrite.

The following corollary combines (86) and (87) to conduct the risk rewrite.

Corollary 23. *Assume $\pi_p = 1/2$. For SU learning, the classification risk can be rewritten as*

$$R(g) = \mathbb{E}_{\bar{S}} [\bar{\ell}_{\bar{S}}] + \mathbb{E}_{\text{U}} [\bar{\ell}_{\text{U}}],$$

where

$$\begin{aligned}
\bar{\ell}_{\bar{S}} &= \frac{\pi_p^2 + \pi_n^2}{2\pi_p - 1} \ell_p - \frac{\pi_p^2 + \pi_n^2}{2\pi_p - 1} \ell_n, \\
\bar{\ell}_{\text{U}} &= -\frac{\pi_n}{2\pi_p - 1} \ell_p + \frac{\pi_p}{2\pi_p - 1} \ell_n.
\end{aligned} \tag{92}$$

Step 2: Recovering the previous result(s).

To recover Theorem 1 of Bao et al. (2018), we first need to restore $\mathbb{E}_{\text{S}}[\cdot]$ from $\mathbb{E}_{\bar{S}}[\cdot]$ in Corollary 23. The following lemma provides a means for us to do so.

Lemma 24. *Given B (46) and following the SU learning notations, we have*

$$M_{\text{SU}}B = \begin{pmatrix} P_{\bar{S}} \\ P_{\text{U}} \end{pmatrix} = \bar{P},$$

where

$$M_{\text{SU}} := \begin{pmatrix} \frac{\pi_p^2 \int_x \times P_{x|Y=p} dx}{\pi_p^2 + \pi_n^2} & \frac{\pi_n^2 \int_x \times P_{x|Y=n} dx}{\pi_p^2 + \pi_n^2} \\ \pi_p & \pi_n \end{pmatrix}.$$

Proof. Since $\int_x \times P_{x|Y=p} dx = 1$ and $\int_x \times P_{x|Y=n} dx = 1$, we have $M_{\text{SU}} = M_{\text{SU}}$, and hence $M_{\text{SU}}B = M_{\text{SU}}B = \begin{pmatrix} P_{\bar{S}} \\ P_{\text{U}} \end{pmatrix}$. The last equality follows from Lemma 5. \square

⁶As the 0-1 loss is symmetric.

Lemma 24 allows us to slightly revise the derivation (89) as follows:

$$\begin{aligned}
R(g) &= \int_{x \times x} \bar{L} \bar{P} dx = \int_{x \times x} \bar{L} M_{\text{SU}} B dx \\
&= \int_{x \times x} \begin{pmatrix} \bar{\ell}_{\bar{S}} & \bar{\ell}_{\bar{U}} \end{pmatrix} \begin{pmatrix} \frac{\pi_{\text{p}}^2 \int_x \frac{P_{x/Y=\text{p}} dx}{\pi_{\text{p}}^2 + \pi_{\text{n}}^2} & \frac{\pi_{\text{n}}^2 \int_x \frac{P_{x/Y=\text{n}} dx}{\pi_{\text{p}}^2 + \pi_{\text{n}}^2} \end{pmatrix} \begin{pmatrix} P_{x/Y=\text{p}} \\ P_{x/Y=\text{n}} \end{pmatrix} dx \\
&\stackrel{\text{(a)}}{=} \int_{x \times x} \int_{x \times x} P_{\text{S}} \bar{\ell}_{\bar{S}} dx dx + \int_{x \times x} P_{\text{U}} \bar{\ell}_{\bar{U}} dx \\
&= \mathbb{E}_{\text{S}} [\bar{\ell}_{\bar{S}}] + \mathbb{E}_{\text{U}} [\bar{\ell}_{\bar{U}}],
\end{aligned}$$

where equality (a) follows from the SU formulation (10).

Then, denoting

$$L(X) := \frac{1}{\pi_{\text{p}} - \pi_{\text{n}}} \ell_{\text{p}}(X) - \frac{1}{\pi_{\text{p}} - \pi_{\text{n}}} \ell_{\text{n}}(X), \quad (93)$$

$$L_{-}(X) := -\frac{\pi_{\text{n}}}{\pi_{\text{p}} - \pi_{\text{n}}} \ell_{\text{p}}(X) + \frac{\pi_{\text{p}}}{\pi_{\text{p}} - \pi_{\text{n}}} \ell_{\text{n}}(X) \quad (94)$$

and continuing with (92), we obtain

$$\begin{aligned}
\mathbb{E}_{\text{S}} [\bar{\ell}_{\bar{S}}] &= (\pi_{\text{p}}^2 + \pi_{\text{n}}^2) \mathbb{E}_{\text{S}} \left[\frac{1}{2\pi_{\text{p}} - 1} (\ell_{\text{p}} - \ell_{\text{n}}) \right] \\
&= (\pi_{\text{p}}^2 + \pi_{\text{n}}^2) \mathbb{E}_{\text{S}} [L(X)] \\
&\stackrel{\text{(b)}}{=} (\pi_{\text{p}}^2 + \pi_{\text{n}}^2) \mathbb{E}_{\text{S}} \left[\frac{L(X) + L_{-}(X)}{2} \right]
\end{aligned}$$

and

$$\begin{aligned}
\mathbb{E}_{\text{U}} [\bar{\ell}_{\bar{U}}] &= \mathbb{E}_{\text{U}} \left[-\frac{\pi_{\text{n}}}{2\pi_{\text{p}} - 1} \ell_{\text{p}} + \frac{\pi_{\text{p}}}{2\pi_{\text{p}} - 1} \ell_{\text{n}} \right] \\
&= \mathbb{E}_{\text{U}} [L_{-}(f(X))] \quad (95)
\end{aligned}$$

that prove rewrite (11) in Section 2.2.4 and recover Theorem 1 of Bao et al. (2018) by matching notations⁷. The following lemma justifies equality (b).

Lemma 25. *Let $(x, x) \in P_{\text{S}}$ defined by (10). Then, $\mathbb{E}_{\text{S}} \left[\frac{L(X)}{2} \right] = \mathbb{E}_{\text{S}} \left[\frac{L_{-}(X)}{2} \right]$.*

The derivation demonstrates the flexibility of the proposed framework in which a slight modification of M_{SU} recovers the pairwise distribution P_{S} required for $\mathbb{E}_{\text{S}} [\cdot]$. Moreover, the technique developed here significantly reduces the proof length in Appendix B of Bao et al. (2018). Later in Section 5.1.5, we apply the same trick to recover Theorem 1 of Shimada et al. (2021) for SDU learning.

We remark that the result recovered in this paper is merely Theorem 1 of Bao et al. (2018) but not the last expression in (5) of Bao et al. (2018), which later was implemented as the objective (10) for optimization. It is because, pointed out by Negishi (2023), the additional assumption $P_{\text{S}}(x, x) = P_{\bar{S}}(x)P_{\bar{S}}(x)$ required for achieving (5) of Bao et al. (2018) is impractical. We note that the remedy proposed by Negishi (2023) can be analyzed by the proposed framework, but we omit it due to the amount of overlap with the analyses in Sections 5.1.1 and 5.1.3.

⁷The matching to the notations of Bao et al. (2018) is as follows: π_{p} is π_{+} , π_{n} is π_{-} , $\pi_{\text{p}}^2 + \pi_{\text{n}}^2$ is π_{S} , P_{S} is p_{S} , P_{U} is p , ℓ_{p} is $\ell(f(X), +1)$, ℓ_{n} is $\ell(f(X), -1)$, $L(X)$ by definition is $\frac{1}{2\pi_{+}-1} (\ell(f(X), +1) - \ell(f(X), -1))$, and $L_{-}(f(X))$ by definition is $-\frac{\pi_{-}}{2\pi_{+}-1} \ell(f(X), +1) + \frac{\pi_{+}}{2\pi_{+}-1} \ell(f(X), -1)$.

5.1.4 Pairwise Comparison (Pcomp) Learning

We follow the steps illustrated in Section 5.1.2 to justify the proposed framework since, by Table 7, M_{Pcomp} is reduced from M_{UU} .

Step 1: Corrected Loss Design and Risk Rewrite.

The following corollary combines (86) and (87) to achieve rewrite (17) in Section 2.2.7.

Corollary 26. *For Pcomp learning, the classification risk can be rewritten as*

$$R(g) = \mathbb{E}_{\text{Sup}} [\bar{\ell}_{\text{Sup}}] + \mathbb{E}_{\text{Inf}} [\bar{\ell}_{\text{Inf}}],$$

where

$$\begin{aligned} \bar{\ell}_{\text{Sup}} &= \ell_{\text{p}} - \pi_{\text{p}} \ell_{\text{n}}, \\ \bar{\ell}_{\text{Inf}} &= -\pi_{\text{n}} \ell_{\text{p}} + \ell_{\text{n}}. \end{aligned} \tag{96}$$

Step 2: Recovering the previous result(s).

It is straightforward to recover Theorem 3 of Feng et al. (2021) by matching notations⁸. Since x is a variable and can be substituted by x , we express Corollary 26 as

$$R(g) = \mathbb{E}_x \mathcal{P}_{\text{Sup}} [\ell_{\text{p}}(x) - \pi_{\text{p}} \ell_{\text{n}}(x)] + \mathbb{E}_x \mathcal{P}_{\text{Inf}} [\ell_{\text{n}}(x) - \pi_{\text{n}} \ell_{\text{p}}(x)], \tag{97}$$

recovering (5) of Feng et al. (2021).

5.1.5 Similar-dissimilar-unlabeled (SDU) Learning

We justify the applicability of the proposed framework for DU and SD separately. Firstly, we start with DU learning, which is similar to SU learning in the sense that pairwise information is provided. From Lemmas 5 and 7, we see that the pairwise distributions are treated similarly. Thus, following the same steps in Section 5.1.3, we conduct the risk rewrite for DU learning.

Step 1: Corrected Loss Design and Risk Rewrite for DU Learning.

The following corollary is a variant of Corollary 23.

Corollary 27. *Assume $\pi_{\text{p}} = 1/2$. For DU learning, the classification risk can be rewritten as*

$$R(g) = \mathbb{E}_{\text{D}} [\bar{\ell}_{\text{D}}] + \mathbb{E}_{\text{U}} [\bar{\ell}_{\text{U}}],$$

where

$$\begin{aligned} \bar{\ell}_{\text{D}} &= 2\pi_{\text{p}}\pi_{\text{n}} \left(\frac{1}{\pi_{\text{n}} - \pi_{\text{p}}} \ell_{\text{p}} - \frac{1}{\pi_{\text{n}} - \pi_{\text{p}}} \ell_{\text{n}} \right), \\ \bar{\ell}_{\text{U}} &= -\frac{\pi_{\text{p}}}{\pi_{\text{n}} - \pi_{\text{p}}} \ell_{\text{p}} + \frac{\pi_{\text{n}}}{\pi_{\text{n}} - \pi_{\text{p}}} \ell_{\text{n}}. \end{aligned} \tag{98}$$

Step 2: Recovering the previous result(s) for DU Learning.

We reuse the trick in Lemma 24 for restoring the pairwise distribution \mathcal{P}_{S} to restore \mathcal{P}_{D} needed here, allowing us to recover the rewrite (15) in Theorem 1 of Shimada et al. (2021) and the first result in Theorem 7.3 of Sugiyama et al. (2022). The derivation resembles that of SU learning. We start with the next lemma, adapted from Lemma 24.

Lemma 28. *Given B (46) and following the DU learning notations, we have*

$$M_{\text{DU}} B = \begin{pmatrix} \mathcal{P}_{\text{D}} \\ \mathcal{P}_{\text{U}} \end{pmatrix} = \bar{P},$$

⁸The matching is as follows: \mathcal{P}_{Sup} is $\tilde{p}_+(x)$, \mathcal{P}_{Inf} is $\tilde{p}_-(x)$, ℓ_{p} is $\ell(f(x), +1)$, and ℓ_{n} is $\ell(f(x), -1)$.

where

$$M_{\text{DU}} := \begin{pmatrix} \frac{\int_{x \times X} P_{x/Y=n} dx}{2} & \frac{\int_{x \times X} P_{x/Y=p} dx}{2} \\ \pi_p & \pi_n \end{pmatrix}.$$

We apply Lemma 28 to slightly revise the derivation of (89) as follows:

$$\begin{aligned} R(g) &= \int_{x \times X} \bar{L} \bar{P} dx = \int_{x \times X} \bar{L} M_{\text{DU}} B dx \\ &= \int_{x \times X} \begin{pmatrix} \bar{\ell}_{\bar{D}} & \bar{\ell}_{\bar{U}} \end{pmatrix} \begin{pmatrix} \frac{\int_{x \times X} P_{x/Y=n} dx}{2} & \frac{\int_{x \times X} P_{x/Y=p} dx}{2} \\ \pi_p & \pi_n \end{pmatrix} \begin{pmatrix} P_{x/Y=p} \\ P_{x/Y=n} \end{pmatrix} dx \\ &= \int_{x \times X} \int_{x \times X} P_{\bar{D}} \bar{\ell}_{\bar{D}} dx dx + \int_{x \times X} P_{\bar{U}} \bar{\ell}_{\bar{U}} dx \\ &= \mathbb{E}_{\bar{D}} [\bar{\ell}_{\bar{D}}] + \mathbb{E}_{\bar{U}} [\bar{\ell}_{\bar{U}}], \end{aligned}$$

where the second to last equality follows from the DU formulation (12). Denoting

$$L_+(X) := \frac{\pi_p}{\pi_p - \pi_n} \ell_p(X) - \frac{\pi_n}{\pi_p - \pi_n} \ell_n(X), \quad (99)$$

recalling $L(X)$ from (93), and continuing with (98), we have

$$\begin{aligned} \mathbb{E}_{\bar{D}} [\bar{\ell}_{\bar{D}}] &= 2\pi_p \pi_n \mathbb{E}_{\bar{D}} \left[\frac{1}{\pi_n - \pi_p} \ell_p - \frac{1}{\pi_n - \pi_p} \ell_n \right] \\ &= 2\pi_p \pi_n \mathbb{E}_{\bar{D}} [-L(X)] \\ &\stackrel{(a)}{=} 2\pi_p \pi_n \mathbb{E}_{\bar{D}} \left[-\frac{L(X) + L(X)}{2} \right] \end{aligned}$$

and

$$\begin{aligned} \mathbb{E}_{\bar{U}} [\bar{\ell}_{\bar{U}}] &= \mathbb{E}_{\bar{U}} \left[-\frac{\pi_p}{\pi_n - \pi_p} \ell_p + \frac{\pi_n}{\pi_n - \pi_p} \ell_n \right] \\ &= \mathbb{E}_{\bar{U}} [L_+(X)] \end{aligned} \quad (100)$$

that prove rewrite (13) in Section 2.2.5. By matching notations, we recover (15) in Theorem 1 of Shimada et al. (2021)⁹. Equality (a) follows from the next lemma.

Lemma 29. *Let $(x, x) \in P_{\bar{D}}$ defined in (12). Then, $\mathbb{E}_{\bar{D}} \left[\frac{L(X)}{2} \right] = \mathbb{E}_{\bar{D}} \left[\frac{L(X)}{2} \right]$.*

Secondly, we consider the rewrite of SD learning. To do so, we apply the knowledge acquired from SU and DU learning (Corollaries 23 and 27).

Step 1: Corrected Loss Design and Risk Rewrite for SD Learning.

We provide another variant of Corollary 23 to conduct the risk rewrite.

Corollary 30. *Assume $\pi_p = 1/2$. For SD learning, the classification risk can be rewritten as*

$$R(g) = \mathbb{E}_{\bar{S}} [\bar{\ell}_{\bar{S}}] + \mathbb{E}_{\bar{D}} [\bar{\ell}_{\bar{D}}],$$

where

$$\begin{aligned} \bar{\ell}_{\bar{S}} &= (\pi_p^2 + \pi_n^2) \left(\frac{\pi_p}{\pi_p - \pi_n} \ell_p - \frac{\pi_n}{\pi_p - \pi_n} \ell_n \right), \\ \bar{\ell}_{\bar{D}} &= 2\pi_p \pi_n \left(-\frac{\pi_n}{\pi_p - \pi_n} \ell_p + \frac{\pi_p}{\pi_p - \pi_n} \ell_n \right). \end{aligned} \quad (101)$$

⁹The matching to the notations of Shimada et al. (2021) is as follows: π_p is π_+ , π_n is π_- , $\pi_p^2 + \pi_n^2$ is π_S , $2\pi_p \pi_n$ is π_D , P_S is $p_S(x, x)$, $P_{\bar{D}}$ is $p_{\bar{D}}(x, x)$, $P_{\bar{U}}$ is $p_{\bar{U}}(x)$, ℓ_p is $\ell(f(X), +1)$, ℓ_n is $\ell(f(X), -1)$, $L(X)$ is $\tilde{L}(f(X))$, $L_+(X)$ is $L(f(X), +1)$, and $L_-(X)$ is $L(f(X), -1)$.

Step 2: Recovering the previous result(s) for SD Learning.

We apply the same strategy as in Lemma 28 to obtain the needed P_S and P_D . We begin with the next lemma, adapted from Lemma 24, to recover (16) in Theorem 1 of Shimada et al. (2021) and the second result in Theorem 7.3 of Sugiyama et al. (2022).

Lemma 31. *Given B (46) and following the SD learning notations, we have*

$$M_{SD}B = \begin{pmatrix} P_{\tilde{S}} \\ P_{\tilde{D}} \end{pmatrix} = \bar{P},$$

where

$$M_{SD} := \begin{pmatrix} \frac{\pi_p^2 \int_x \times \frac{P_{x|Y=p} dx}{\pi_p^2 + \pi_n^2}}{\frac{\int_x \times \frac{P_{x|Y=n} dx}{2}}{2}} & \frac{\pi_n^2 \int_x \times \frac{P_{x|Y=n} dx}{\pi_p^2 + \pi_n^2}}{\frac{\int_x \times \frac{P_{x|Y=p} dx}{2}}{2}} \end{pmatrix}.$$

We apply Lemma 31 to slightly revise the derivation of (89) as follows:

$$\begin{aligned} R(g) &= \int_x \times \bar{L} \bar{P} dx = \int_x \times \bar{L} M_{SD}B dx \\ &= \int_x \times \begin{pmatrix} \bar{\ell}_{\tilde{S}} & \bar{\ell}_{\tilde{D}} \end{pmatrix} \begin{pmatrix} \frac{\pi_p^2 \int_x \times \frac{P_{x|Y=p} dx}{\pi_p^2 + \pi_n^2}}{\frac{\int_x \times \frac{P_{x|Y=n} dx}{2}}{2}} & \frac{\pi_n^2 \int_x \times \frac{P_{x|Y=n} dx}{\pi_p^2 + \pi_n^2}}{\frac{\int_x \times \frac{P_{x|Y=p} dx}{2}}{2}} \end{pmatrix} \begin{pmatrix} P_{x|Y=p} \\ P_{x|Y=n} \end{pmatrix} dx \\ &= \int_x \times \int_x \times P_S \bar{\ell}_{\tilde{S}} dx dx + \int_x \times \int_x \times P_D \bar{\ell}_{\tilde{D}} dx dx \\ &= E_S [\bar{\ell}_{\tilde{S}}] + E_D [\bar{\ell}_{\tilde{D}}], \end{aligned}$$

where the second to last equality follows from the SD formulation (14). Recalling $L_+(X)$ (99) and $L_-(X)$ (94) and continuing with (101),

$$\begin{aligned} E_S [\bar{\ell}_{\tilde{S}}] &= (\pi_p^2 + \pi_n^2) E_S \left[\frac{\pi_p}{\pi_p - \pi_n} \ell_p - \frac{\pi_n}{\pi_p - \pi_n} \ell_n \right] \\ &= (\pi_p^2 + \pi_n^2) E_S [L_+(X)] \\ &\stackrel{(b)}{=} (\pi_p^2 + \pi_n^2) E_S \left[\frac{L_+(X) + L_+(X)}{2} \right] \end{aligned}$$

and

$$\begin{aligned} E_D [\bar{\ell}_{\tilde{D}}] &= 2\pi_p \pi_n E_D \left[-\frac{\pi_n}{\pi_p - \pi_n} \ell_p + \frac{\pi_p}{\pi_p - \pi_n} \ell_n \right] \\ &= 2\pi_p \pi_n E_D [L_-(X)] \\ &\stackrel{(c)}{=} 2\pi_p \pi_n E_D \left[\frac{L_-(X) + L_-(X)}{2} \right] \end{aligned}$$

prove rewrite (15) in Section 2.2.6. We also recover (16) in Theorem 1 of Shimada et al. (2021) via matching notations. The required matches can be found in the paragraph before Lemma 29. The equality (b) holds by applying Lemma 25 with $L(X)$ replaced by $L_+(X)$, and (c) follows from Lemma 29 with $L(X)$ replaced by $L_-(X)$.

An intriguing observation worth mentioning is that the losses $L_+(X)$ and $L_-(X)$ applied to decontaminate the unlabeled data in SU and DU learning ((95) and (100)) are now used to decontaminate the similar and the dissimilar data in SD learning, respectively. One can also quickly draw the same conclusion from Table 4. Knowing the reason behind this observation would help to transfer one corrected loss developed in one scenario to another weakly supervised scenario.

5.1.6 Similarity-Confidence (Sconf) Learning

Since M_{Sconf} (55) is not a child of M_{UU} (49) on the reduction graph, a direct application of Theorem 21 is infeasible. Nevertheless, we demonstrate how our framework is applied to rewrite the classification risk for Sconf learning. We make a small adjustment to the framework that instead of showing $\bar{L} \bar{P} = L M^t \bar{P} = L P$, we show that for loss vector \bar{L} with a certain property,

$$\int_x \times \bar{L} \bar{P} dx = \bar{L} \tilde{M}_{\text{Sconf}} P. \quad (102)$$

The idea behind this approach is to accommodate x sampled from P_X (18). Suppose, informally, we have the equation above. Then, the right-hand side of (102) will produce $L P$ if we can compute a decontamination matrix $\tilde{M}_{\text{Sconf}}^t$ satisfying $\tilde{M}_{\text{Sconf}}^t \tilde{M}_{\text{Sconf}} = I$ and assign $\bar{L} := L \tilde{M}_{\text{Sconf}}^t$. Lastly, integrating over x on both sides, we obtain the key equation

$$\int_x \times \int_x \times \bar{L} \bar{P} dx dx = \int_x \times L P dx$$

for risk rewrite.

Step 1: Corrected Loss Design and Risk Rewrite.

Let us follow the notations in Section 4.1.6. We begin with two technical lemmas and leave their proofs to Appendix D.1. The first technical lemma shows how to achieve (102).

Lemma 32. *Assume the formulation $\bar{P} = M_{\text{Sconf}} B$ (58) is given. Suppose a vector of corrected losses \bar{L} of the form $(\tilde{\ell}_1(x) \tilde{\ell}_2(x))$ is independent of x . Then, we have*

$$\int_x \times \bar{L} \bar{P} dx = \bar{L} \tilde{M}_{\text{Sconf}} P, \quad (103)$$

where

$$\tilde{M}_{\text{Sconf}} = \begin{pmatrix} \int_x \frac{\pi_p^2 P_{x/p} - \pi_n^2 P_{x/n}}{r - \pi_n} dx & \int_x \frac{\pi_n^2 P_{x/n} - \pi_p^2 P_{x/p}}{r - \pi_n} dx \\ \int_x \frac{\pi_p^2 P_{x/n} - \pi_n^2 P_{x/p}}{\pi_p - r} dx & \int_x \frac{\pi_p^2 P_{x/p} - \pi_n^2 P_{x/n}}{\pi_p - r} dx \end{pmatrix}.$$

The second technical lemma computes the decontamination matrix.

Lemma 33. *Let*

$$\tilde{M}_{\text{Sconf}}^t := \begin{pmatrix} \frac{r - \pi_n}{\pi_p - \pi_n} & 0 \\ 0 & \frac{\pi_p - r}{\pi_p - \pi_n} \end{pmatrix}.$$

Then,

$$\tilde{M}_{\text{Sconf}}^t \tilde{M}_{\text{Sconf}} = I.$$

Next, we follow the sketch above to instantiate the corrected losses as

$$\bar{L} := L \tilde{M}_{\text{Sconf}}^t = \left(\frac{r - \pi_n}{\pi_p - \pi_n} \ell_p(X) \quad \frac{\pi_p - r}{\pi_p - \pi_n} \ell_n(X) \right).$$

Putting $\tilde{M}_{\text{Sconf}}^t$, \bar{L} , and (103) together, we have the following rewrite.

Theorem 34. *Assume $\pi_p = 1/2$. The classification risk of Sconf learning can be expressed by*

$$R(g) = \mathbb{E}_{X, X} \left[\frac{r - \pi_n}{\pi_p - \pi_n} \ell_p(X) + \frac{\pi_p - r}{\pi_p - \pi_n} \ell_n(X) \right]. \quad (104)$$

Proof. Integrating both sides of (103) over x and applying Lemma 33, we obtain

$$\begin{aligned} \int_x \times \int_x \times \bar{L} \bar{P} dx dx &= \int_x \times \bar{L} \tilde{M}_{\text{Sconf}} P dx \\ &= \int_x \times L \tilde{M}_{\text{Sconf}}^{-1} \tilde{M}_{\text{Sconf}} P dx = R(g). \end{aligned}$$

On the other hand, substituting \bar{L} with $\left(\frac{r-\pi_n}{\pi_p-\pi_n} \ell_p(X) \quad \frac{\pi_p-r}{\pi_p-\pi_n} \ell_n(X) \right)$ and \bar{P} with $\begin{pmatrix} P_X P_X \\ P_X P_X \end{pmatrix}$,

$$\begin{aligned} \int_x \times \int_x \times \bar{L} \bar{P} dx dx &= \int_x \times \int_x \times P_x P_x \left(\frac{r-\pi_n}{\pi_p-\pi_n} \ell_p(x) + \frac{\pi_p-r}{\pi_p-\pi_n} \ell_n(x) \right) dx dx \\ &= \mathbb{E}_{X,X} \left[\frac{r-\pi_n}{\pi_p-\pi_n} \ell_p(X) + \frac{\pi_p-r}{\pi_p-\pi_n} \ell_n(X) \right] \end{aligned}$$

completes the proof of the theorem. \square

Step 2: Recovering the previous result(s).

From the above derivation, we have achieved the first half of the rewrite in (19). Notice that (56) can be rephrased as

$$\left(\frac{r-\pi_n}{\pi_p} \right) P_X P_X = (\pi_p^2 P_{X/p} - \pi_n^2 P_{X/n}) P_{X/p} + (\pi_n^2 P_{X/n} - \pi_p^2 P_{X/p}) P_{X/n}$$

and that (57) can be rephrased as

$$\left(\frac{\pi_p-r}{\pi_n} \right) P_X P_X = (\pi_p^2 P_{X/n} - \pi_p^2 P_{X/p}) P_{X/p} + (\pi_p^2 P_{X/p} - \pi_n^2 P_{X/n}) P_{X/n}.$$

Thus, when $\pi_p = 1/2$, we can repeat the proof steps in Lemma 9 to rephrase (58) as

$$\begin{pmatrix} P_X P_X \\ P_X P_X \end{pmatrix} = \begin{pmatrix} \frac{\pi_p(\pi_p^2 P_{X/p} - \pi_n^2 P_{X/n})}{r-\pi_n} & \frac{\pi_p(\pi_n^2 P_{X/n} - \pi_n^2 P_{X/p})}{r-\pi_n} \\ \frac{\pi_n(\pi_p^2 P_{X/n} - \pi_p^2 P_{X/p})}{\pi_p-r} & \frac{\pi_n(\pi_p^2 P_{X/p} - \pi_n^2 P_{X/n})}{\pi_p-r} \end{pmatrix} \begin{pmatrix} P_X /p \\ P_X /n \end{pmatrix}.$$

Comparing the equation above with $\bar{P} = M_{\text{Sconf}} B$, we see that it is still feasible to formulate \bar{P} with X and X in M_{Sconf} and B of (58) swapped. Then, repeating the same argument in **Step 1** with x and x swapped, we obtain

$$R(g) = \mathbb{E}_{X,X} \left[\frac{r-\pi_n}{\pi_p-\pi_n} \ell_p(X) + \frac{\pi_p-r}{\pi_p-\pi_n} \ell_n(X) \right]. \quad (105)$$

Therefore, the following combines (104) and (105) to obtain

$$\begin{aligned} R(g) &= \frac{1}{2}(R(g) + R(g)) \\ &= \frac{1}{2} \mathbb{E}_{X,X} \left[\frac{r-\pi_n}{\pi_p-\pi_n} \ell_p(X) + \frac{\pi_p-r}{\pi_p-\pi_n} \ell_n(X) \right] + \frac{1}{2} \mathbb{E}_{X,X} \left[\frac{r-\pi_n}{\pi_p-\pi_n} \ell_p(X) + \frac{\pi_p-r}{\pi_p-\pi_n} \ell_n(X) \right] \\ &= \mathbb{E}_{X,X} \left[\frac{r-\pi_n}{\pi_p-\pi_n} \frac{\ell_p(X) + \ell_p(X)}{2} + \frac{\pi_p-r}{\pi_p-\pi_n} \frac{\ell_n(X) + \ell_n(X)}{2} \right] \end{aligned}$$

that recovers rewrite (19) in Section 2.2.8. By matching notations, we recover Theorem 3 of Cao et al. (2021b)¹⁰.

¹⁰The matching to the notations of Cao et al. (2021b) is as follows: π_p is π_+ , π_n is π_- , r is s , $\ell_p(X)$ is $\ell(g(X), +1)$, and $\ell_n(X)$ is $\ell(g(X), -1)$.

5.2 CCN Scenarios

The proposed framework is now applied to conduct the risk rewrites for WSLs discussed in Section 4.2 and summarized in Table 8. Counterintuitively, we demonstrate that finding an inverse matrix (e.g., Proposition 1) is not the only way to solve the risk rewrite problem. Introduced in Proposition 2, the new technique exploited in this subsection, marginal chain, calculates the decontamination matrix for (37) via applying the conditional probability formula twice during a chain of matrix multiplications.

5.2.1 Generalized CCN

We justify the proposed framework for generalized CCN learning via the following steps. Derived equations will be applied to solve the risk rewrite problem for WSLs discussed in Section 4.2.

Step 1: Corrected Loss Design.

Let us follow the notations in Section 4.2.1. Same as what we have illustrated in the beginning of Section 5.1.1, the proposed framework achieves three milestones to rewrite the risk. We apply Lemma 10 to achieve the first milestone, $\bar{P} = M_{\text{gCCN}} P$. This is done by noting that for generalized CCN, $\bar{P} = M_{\text{gCCN}} B$ and $B = P$ are given by Lemma 10 (i.e., do not need to handle M_{trsf} discussed in Section 3.2 since M_{trsf} is the identity matrix when $B = P$).

The second milestone is to find M_{gCCN}^\dagger to achieve $M_{\text{gCCN}}^\dagger \bar{P} = P$. Since M_{gCCN} (64) is identical to M (40), a direct application of Proposition 2 gives the decontamination matrix

$$M_{\text{gCCN}}^\dagger := \begin{pmatrix} P_{Y=1/S=s_1,X} & P_{Y=1/S=s_2,X} & \cdots & P_{Y=1/S=s_{|S|},X} \\ P_{Y=2/S=s_1,X} & P_{Y=2/S=s_2,X} & \cdots & P_{Y=2/S=s_{|S|},X} \\ \vdots & \vdots & \ddots & \vdots \\ P_{Y=K/S=s_1,X} & P_{Y=K/S=s_2,X} & \cdots & P_{Y=K/S=s_{|S|},X} \end{pmatrix} \quad (106)$$

that satisfies the $M_{\text{gCCN}}^\dagger \bar{P} = P$ requirement.

The final milestone is achieved by instantiating the corrected loss (38) as $\bar{L} := L M_{\text{gCCN}}^\dagger$. We denote the k -th entry of L is $\ell_{Y=k}$ with $k \in [K]$ and the j -th entry of \bar{L} is $\bar{\ell}_{S=s_j}$ with $j \in [|S|]$.

Despite Proposition 2's simplicity, the construction of M_{gCCN}^\dagger is somewhat surprising. M_{gCCN}^\dagger , to our best knowledge, contributes to a first loss correction result relaxing the invertibility constraint. Unlike M_{UU}^\dagger (Corollary 20), which needs to compute an inverse matrix, one can construct M_{gCCN}^\dagger by calculating each entry $P_{Y/S,X}$ in (106), to which, we point out a systematic way in Section 5.2.2.

Step 2: Classification Risk Rewrite.

With \bar{L} in hand, the following theorem provides an intermediate form of risk rewrite.

Theorem 35. *Let \bar{P} and P are given by (62) and (63), respectively. Denote $\bar{L} := L M_{\text{gCCN}}^\dagger$. Then, $\bar{L} \bar{P} = L P$ and*

$$R(g) = \int_X L P dx = \int_X \bar{L} \bar{P} dx. \quad (107)$$

Proof. Since M_{gCCN}^\dagger is given by Proposition 2, $M_{\text{gCCN}}^\dagger \bar{P} = P$. Thus, following the framework (39), we have $\bar{L} \bar{P} = L M_{\text{gCCN}}^\dagger \bar{P} = L P$ implying (107). \square

Theorem 35 will be applied to derive the respective rewrites for WSLs discussed in Section 4.2 in the rest of this subsection. In particular, we explain how to realize M_{gCCN}^\dagger (106) for a given CCN scenario. Then, the risk rewrite (107) automatically carries over for the scenario considered, and the respective \bar{L} specifies the corrected losses in the rewrite.

5.2.2 Proper Partial-Label (PPL) Learning

M_{gCCN}^\dagger provides an abstraction for us to construct the corrected losses \bar{L} . Next, we focus on deriving the actual form of $P_{Y/S,X}$ in M_{gCCN}^\dagger (106) to explicitly express $\bar{\ell}_S$ for PPL.

Step 1: Corrected Loss Design and Risk Rewrite.

Let us follow the notations in Section 4.2.2. The following lemma specifies the form of $P_{Y/S,X}$ to instantiate M_{gCCN}^\dagger .

Lemma 36. M_{PPL}^\dagger corresponds to realizing M_{gCCN}^\dagger (106) with

$$P_{Y=i/S=s_j,X} := \frac{P_{Y=i/X} [Y=i \quad s_j]}{\sum_a P_{Y=a/X}}. \quad (108)$$

Proof. Recall that the decontamination matrix of M_{gCCN} (64) is M_{gCCN}^\dagger (106) and M_{PPL} is a reduction of M_{gCCN} via $P_{S/Y,X} = C(S,X) | [Y \quad S]$ (65). Thus, to find out the (i,j) entry of M_{PPL}^\dagger , we need to find out the form of $P_{Y=i/S=s_j,X}$ subject to (65).

Applying Theorem 1 of Wu et al. (2023) directly gives

$$P_{Y=i/S=s_j,X} = \frac{P_{Y=i/X} [Y=i \quad s_j]}{\sum_a P_{Y=a/X}},$$

which completes the proof. For completeness, we provide a derivation as follows.

Note that $P_{S/Y,X} = C(S,X) | [Y \quad S]$ (65) implies

$$\begin{aligned} \sum_{b \in Y \setminus S} P_{S,Y=b/X} &= \sum_{b \in Y \setminus S} P_{S/Y=b,X} P_{Y=b/X} \\ &= \sum_{b \in Y \setminus S} C(S,X) | [b \quad S] P_{Y=b/X} \\ &= 0. \end{aligned}$$

Therefore, $P_{S/X} = \sum_a P_{S,Y=a/X} + \sum_{b \in Y \setminus S} P_{S,Y=b/X} = \sum_a P_{S,Y=a/X}$. Utilizing this fact, we obtain

$$\begin{aligned} P_{Y/S,X} &= \frac{P_{S,Y/X}}{P_{S/X}} = \frac{P_{S/Y,X} P_{Y/X}}{\sum_a P_{S/Y=a,X} P_{Y=a/X}} \\ &= \frac{C(S,X) | [Y \quad S] P_{Y/X}}{\sum_a C(S,X) | [Y=a \quad S] P_{Y=a/X}} \\ &= \frac{P_{Y/X} | [Y \quad S]}{\sum_a P_{Y=a/X}} \end{aligned}$$

that finishes the proof for Theorem 1 of Wu et al. (2023). \square

We have shown that M_{PPL}^\dagger is derived from M_{gCCN}^\dagger . Thus, we can follow Theorem 35 to construct the corrected losses using (108) and obtain the risk rewrite (27) for PPL in Section 2.2.12.

Corollary 37. Given M_{PPL}^\dagger defined by (108), we denote the corrected losses $\bar{L} := L M_{\text{PPL}}^\dagger$. Then, for PPL learning, the classification risk can be rewritten as

$$R(g) = \mathbb{E}_{S,X} [\bar{\ell}_S],$$

where

$$\bar{\ell}_S = \sum_i \frac{P_{Y=i/X}}{\sum_a P_{Y=a/X}} \ell_{Y=i}. \quad (109)$$

Proof. Given (108), the j -th entry of \bar{L} is of the form

$$\begin{aligned}\bar{\ell}_{S=s_j} &= \left(L \ M_{\text{PPL}}^\dagger \right)_j = \sum_{i=1}^K \frac{P_{Y=i|X} \mathbb{1}[Y=i \ s_j]}{\sum_{a \ s_j} P_{Y=a|X}} \ell_{Y=i} \\ &= \sum_{i \ s_j} \frac{P_{Y=i|X}}{\sum_{a \ s_j} P_{Y=a|X}} \ell_{Y=i}.\end{aligned}$$

Then, since M_{PPL}^\dagger is a realization of M_{gCCN}^\dagger according to Lemma 36, we continue (107) to express the risk as

$$R(g) = \int_{\mathcal{X}} \bar{L} \ \bar{P} dx = \int_{\mathcal{X}} \sum_{j=1}^{|\mathcal{S}|} P_{S=s_j, X} \bar{\ell}_{S=s_j} dx = \mathbb{E}_{S, X} [\bar{\ell}_S].$$

□

Step 2: Recovering the previous result(s).

We finish this part by pointing out Corollary 37 recovers Theorem 3 of Wu et al. (2023).

5.2.3 Provably Consistent Partial-Label (PCPL) Learning

It is fairly straightforward to apply the proposed framework to rewrite the classification risk. However, it is more involved in recovering the existing result.

Step 1: Corrected Loss Design and Risk Rewrite.

The argument for obtaining the risk rewrite for PCPL is similar to that of PPL. From Section 4.2.3 we know that PCPL is a special case of PPL that only differs in the choice of $C(S, X)$. Since $C(S, X)$ is independent of (108), M_{PCPL}^\dagger and M_{PPL}^\dagger are identical. Hence, following the notations in Section 4.2.3 and repeating the proof of Corollary 37, we obtain the risk rewrite for PCPL:

Corollary 38. *The decontamination matrix M_{PCPL}^\dagger for PCPL equals M_{PPL}^\dagger . If we define the corrected losses as $\bar{L} := L \ M_{\text{PCPL}}^\dagger$, the classification risk for PCPL learning can be rewritten as*

$$R(g) = \mathbb{E}_{S, X} [\bar{\ell}_S],$$

where

$$\bar{\ell}_S = \sum_{i \ S} \frac{P_{Y=i|X}}{\sum_{a \ S} P_{Y=a|X}} \ell_{Y=i}. \quad (110)$$

Step 2: Recovering the previous result(s).

In order to recover (8) of Feng et al. (2020b), we need to reorganize the sum in (110) by leveraging a unique property of a pair of partial-labels (s, \bar{s}) that complement each other. The following technical lemma states the required property, with proof deferred to Appendix D.2.

Lemma 39. *Let (s, \bar{s}) be a pair of partial-labels satisfying $s = Y \setminus \bar{s}$. Then,*

$$P_{S=s, X} \bar{\ell}_{S=s} + P_{S=\bar{s}, X} \bar{\ell}_{S=\bar{s}} = P_{S=s, X} \sum_{i=1}^K \frac{P_{Y=i|X} \ell_{Y=i}}{\sum_{a \ s} P_{Y=a|X}}.$$

Denote $s_j := Y \setminus s_j$ for every $s_j \in \mathcal{S}$. Then, Lemma 39 implies

$$\begin{aligned}\sum_{j=1}^{|\mathcal{S}|} 2P_{S=s_j, X} \bar{\ell}_{S=s_j} &= \sum_{j=1}^{|\mathcal{S}|} \left(P_{S=s_j, X} \bar{\ell}_{S=s_j} + P_{S=s_j, X} \bar{\ell}_{S=s_j} \right) \\ &= \sum_{j=1}^{|\mathcal{S}|} P_{S=s_j, X} \sum_{i=1}^K \frac{P_{Y=i|X} \ell_{Y=i}}{\sum_{a \ s_j} P_{Y=a|X}}.\end{aligned}$$

Hence, continuing from Corollary 38,

$$\begin{aligned}
\mathbb{E}_{S,X} [\bar{\ell}_S] &= \int_x \times \sum_{j=1}^{|S|} P_{S=s_j,x} \bar{\ell}_{S=s_j} dx \\
&= \frac{1}{2} \int_x \times \sum_{j=1}^{|S|} P_{S=s_j,x} \sum_{i=1}^K \frac{P_{Y=i/x} \ell_{Y=i}}{\sum_a P_{Y=a/x}} dx \\
&= \frac{1}{2} \mathbb{E}_{S,X} \left[\sum_{i=1}^K \frac{P_{Y=i/X}}{\sum_a P_{Y=a/X}} \ell_{Y=i} \right]
\end{aligned}$$

shows that the rewrite from the framework recovers (25) in Section 2.2.11. By matching notations, we also recover (8) of Feng et al. (2020b)¹¹.

5.2.4 Multi-Complementary-Label (MCL) Learning

Step 1: Corrected Loss Design and Risk Rewrite.

Let us follow the notations in Section 4.2.4. As discussed in Section 4.2.4, MCL is a special case of PPL. Thus, we can modify Lemma 36 based on the notations in Section 4.2.4 to construct the decontamination matrix M_{MCL}^\dagger for MCL. Then, following the same steps for proving Corollary 37, we instantiate \bar{L} to conduct the risk rewrite for MCL:

Corollary 40. *The (i, j) entry of the decontamination matrix M_{MCL}^\dagger is of the form*

$$P_{Y=i/\bar{s}=\bar{s}_j,X} = \frac{P_{Y=i/X} \mathbb{1}[Y=i/\bar{s}_j]}{\sum_{a/\bar{s}_j} P_{Y=a/X}}. \quad (111)$$

Define the corrected losses $\bar{L} := L M_{\text{MCL}}^\dagger$. Then, for MCL learning, the classification risk can be rewritten as

$$R(g) = \mathbb{E}_{\bar{S},X} [\bar{\ell}_{\bar{S}}],$$

where

$$\bar{\ell}_{\bar{S}} = \sum_{i/\bar{S}} \frac{P_{Y=i/X}}{\sum_{a/\bar{S}} P_{Y=a/X}} \ell_{Y=i}. \quad (112)$$

Step 2: Recovering the previous result(s).

Although legitimate, the risk rewrite (112) following the marginal chain approach appears different from Theorem 3 of Feng et al. (2020a), to which we resort to the inversion approach (Proposition 1) that finds another decontamination matrix, termed M_{MCL}^{-1} , to recover. As a preparation step, we denote N_d as the number of multi-complementary-labels with size d and group rows of M_{MCL} (71) by the size of labels as follows.

$$M_{\text{MCL}} = \begin{pmatrix} P_{|\bar{S}=1} M_1 \\ P_{|\bar{S}=2} M_2 \\ \vdots \\ P_{|\bar{S}=K-1} M_{K-1} \end{pmatrix}, \quad (113)$$

¹¹The matching to the notations of Feng et al. (2020b) is as follows: $P_{S,X}$ is $\tilde{p}(x, Y)$, $P_{Y=i/X}$ is $p(y = i/x)$, and $\ell_{Y=i}$ is $L(f(x), i)$.

where for $d \in [K - 1]$, each block is of the form¹²

$$M_d = \frac{1}{\binom{K-1}{d}} \begin{pmatrix} | [Y = 1 / \bar{s}_{d,1}] & | [Y = 2 / \bar{s}_{d,1}] & \cdots & | [Y = K / \bar{s}_{d,1}] \\ | [Y = 1 / \bar{s}_{d,2}] & | [Y = 2 / \bar{s}_{d,2}] & \cdots & | [Y = K / \bar{s}_{d,2}] \\ \vdots & \vdots & \ddots & \vdots \\ | [Y = 1 / \bar{s}_{d,N_d}] & | [Y = 2 / \bar{s}_{d,N_d}] & \cdots & | [Y = K / \bar{s}_{d,N_d}] \end{pmatrix}. \quad (114)$$

To maintain the equality $\bar{P} = M_{\text{MCL}}P$ established in Lemma 13, we also rearrange \bar{P} (69) as

$$\left(P_{\bar{S}=\bar{s}_{1,1},X} \quad \cdots \quad P_{\bar{S}=\bar{s}_{1,N_1},X} \quad \cdots \quad P_{\bar{S}=\bar{s}_{K-1,1},X} \quad \cdots \quad P_{\bar{S}=\bar{s}_{K-1,N_{K-1}},X} \right). \quad (115)$$

As a sanity check, we see that for any $d \in [K - 1]$ and $j \in [N_d]$,

$$\begin{aligned} \left(P_{|\bar{S}=d} M_d P \right)_j &= P_{|\bar{S}=d} \cdot \frac{1}{\binom{K-1}{d}} \sum_Y | [Y / \bar{s}_{d,j}] P_{Y,X} \\ &= \sum_{d=1}^{K-1} P_{\bar{S}_{d,j} \neq d} \cdot \frac{1}{\binom{K-1}{d}} \sum_{Y/\bar{s}_{d,j}} P_{Y,X} | [|\bar{s}_{d,j} = d] \\ &= P_{\bar{S}=\bar{s}_{d,j},X}. \end{aligned} \quad (116)$$

The next lemma is crucial for us to devise the decontamination matrix M_{MCL}^{-1} via the inversion approach. We defer its proof to the later part of this sub-subsection.

Lemma 41. *Let $i^* \in Y$ be fixed. Then, for every $d \in [K - 1]$,*

$$P_{Y=i^*,X} = \sum_{j=1}^{N_d} \left(1 - \frac{K-1}{d} | [Y = i^* \quad \bar{S} = \bar{s}_{d,j}] \right) P_{\bar{S}=\bar{s}_{d,j},X | |\bar{S}=d}.$$

Moreover, the inverse matrix M_d^{-1} of M_d (114) is of the form

$$\begin{pmatrix} 1 - \frac{K-1}{d} | [Y = 1 \quad \bar{s}_{d,1}] & 1 - \frac{K-1}{d} | [Y = 1 \quad \bar{s}_{d,2}] & \cdots & 1 - \frac{K-1}{d} | [Y = 1 \quad \bar{s}_{d,N_d}] \\ 1 - \frac{K-1}{d} | [Y = 2 \quad \bar{s}_{d,1}] & 1 - \frac{K-1}{d} | [Y = 2 \quad \bar{s}_{d,2}] & \cdots & 1 - \frac{K-1}{d} | [Y = 2 \quad \bar{s}_{d,N_d}] \\ \vdots & \vdots & \ddots & \vdots \\ 1 - \frac{K-1}{d} | [Y = K \quad \bar{s}_{d,1}] & 1 - \frac{K-1}{d} | [Y = K \quad \bar{s}_{d,2}] & \cdots & 1 - \frac{K-1}{d} | [Y = K \quad \bar{s}_{d,N_d}] \end{pmatrix}. \quad (117)$$

Applying the lemma, we construct

$$M_{\text{MCL}}^{-1} := \left(M_1^{-1} \quad M_2^{-1} \quad \cdots \quad M_{K-1}^{-1} \right) \quad (118)$$

and obtain $M_{\text{MCL}}^{-1} \bar{P} = P$ since $\bar{P} = M_{\text{MCL}}P$ (116) and

$$M_{\text{MCL}}^{-1} M_{\text{MCL}} = \sum_{d=1}^{K-1} M_d^{-1} P_{|\bar{S}=d} M_d = \sum_{d=1}^{K-1} P_{|\bar{S}=d} M_d^{-1} M_d = \sum_{d=1}^{K-1} P_{|\bar{S}=d} I = I.$$

We remark that M_{MCL}^{-1} plays the same role as M_{MCL}^T realized by (111), as they both are decontamination matrices (designed to convert \bar{P} back to P and used to construct the corrected losses \bar{L}). Distinct symbols

¹²Comparing to (71) where we use one index to denote a total of $|S|$ partial-labels, M_d uses a pair of indices d and j to denote the j -th partial-label with size d . It is easy to verify that $\sum_{d=1}^{K-1} N_d = \sum_{d=1}^{K-1} \binom{K-1}{d} = 2^K - 2 = |S|$.

are merely used to reflect the difference that M_{MCL}^\dagger results from the marginal chain method while M_{MCL}^{-1} comes from the inversion approach. Then, applying the framework (38), $\bar{L} := L M_{\text{MCL}}^{-1}$ leads to

$$\bar{L} \bar{P} = L M_{\text{MCL}}^{-1} \bar{P} = L P.$$

With the corrected losses \bar{L} in hand, the following theorem provides the risk rewrite (23) for MCL via the inversion approach and recovers Theorem 3 of Feng et al. (2020a)¹³.

Theorem 42. *For MCL learning, the classification risk can be expressed as follows.*

$$R(g) = \mathbb{E}_{\bar{S}, X} [\bar{\ell}_{\bar{S}}] = \sum_{d=1}^{K-1} P_{|\bar{S}|=d} \mathbb{E}_{\bar{S}, X || \bar{S}|=d} [\bar{\ell}_{\bar{S}}],$$

where

$$\bar{\ell}_{\bar{S}} = \sum_{i/\bar{S}} \ell_{Y=i} - \frac{K-1-|\bar{S}|}{|\bar{S}|} \sum_{\bar{s} \in \bar{S}} \ell_{Y=\bar{s}}.$$

Proof. We first establish

$$R(g) = \int_X \bar{L} \bar{P} dx = \mathbb{E}_{\bar{S}, X} [\bar{\ell}_{\bar{S}}]$$

since $\bar{L} \bar{P} = L P$, where \bar{P} is specified in (115) and $\bar{L} = L M_{\text{MCL}}^{-1}$ with the \bar{S} -th entry being $\bar{\ell}_{\bar{S}}$. Also, recall that $P_{\bar{S}, X} = \sum_{d=1}^{K-1} P_{|\bar{S}|=d} P_{\bar{S}, X || \bar{S}|=d}$ in Section 4.2.4. Thus, decomposing the probability by the size of \bar{S} , we have

$$\mathbb{E}_{\bar{S}, X} [\bar{\ell}_{\bar{S}}] = \sum_{d=1}^{K-1} P_{|\bar{S}|=d} \mathbb{E}_{\bar{S}, X || \bar{S}|=d} [\bar{\ell}_{\bar{S}}].$$

Lastly, M_{MCL}^{-1} (118) and M_d^{-1} (117) imply, when $\bar{S} = \bar{s}_{d,j}$,

$$\bar{\ell}_{\bar{S}=\bar{s}_{d,j}} = \left(L M_d^{-1} \right)_j = \sum_{i=1}^K \ell_{Y=i} \left(1 - \frac{K-1}{d} \mathbb{1}[Y=i \mid \bar{s}_{d,j}] \right) = \sum_{i=1}^K \ell_{Y=i} - \frac{K-1}{d} \sum_{i \in \bar{s}_{d,j}} \ell_{Y=i}.$$

A simple reorganization and substituting d with $|\bar{S}|$ shows

$$\bar{\ell}_{\bar{S}} = \sum_{i/\bar{S}} \ell_{Y=i} + \sum_{i \in \bar{S}} \ell_{Y=i} - \frac{K-1}{|\bar{S}|} \sum_{i \in \bar{S}} \ell_{Y=i} = \sum_{i/\bar{S}} \ell_{Y=i} - \frac{K-1-|\bar{S}|}{|\bar{S}|} \sum_{i \in \bar{S}} \ell_{Y=i}.$$

□

Now we return to the postponed proof.

Proof. of Lemma 41. We start with identifying M_d^{-1} . Denote $\{\bar{s}_{d,1}, \dots, \bar{s}_{d,N_d}\}$, the set of multi-complementary-labels of size d , as \bar{S}_d . Let us focus on the sized- d data-generating distribution

$$\bar{P}_d = \begin{pmatrix} P_{\bar{S}=\bar{s}_{d,1}, X || \bar{S}|=d} \\ \vdots \\ P_{\bar{S}=\bar{s}_{d,N_d}, X || \bar{S}|=d} \end{pmatrix}.$$

¹³The matching to the notations of Feng et al. (2020a) is as follows: $P_{\bar{S}, X || \bar{S}|=d}$ is $\bar{p}(x, \bar{Y}/s = d)$, $P_{|\bar{S}|=d}$ is $p(s = d)$, and $\bar{\ell}_{\bar{S}}$ is $\bar{L}_d(f(x), \bar{Y})$.

Note that \bar{P}_d corresponds to extracting the entries from (115) that generate sized- d data and then dividing them by $P_{|\bar{S}|=d}$. Thus, $\bar{P} = M_{\text{MCL}}P$ in Lemma 13 implies $\bar{P}_d = M_d P$ and its j -th entry is expressed as

$$P_{\bar{S}=\bar{s}_{d,j},X||\bar{S}|=d} = \frac{1}{\binom{K-1}{d}} \sum_{i=1}^K \mathbb{1}[Y = i / \bar{s}_{d,j}] P_{Y=i,X}. \quad (119)$$

The equality hints to us that if one manages to collect certain multi-complementary-labels \bar{s} to form an equation resembling $\sum_{\bar{s}} P_{\bar{S}=\bar{s},X||\bar{S}|=d} = c_3 \cdot P_{Y=i,X}$ for some constant c_3 , then a reciprocal operation $\frac{1}{c_3}$ recovers $P_{Y=i,X}$ we need (recall we want to find M_d^{-1} achieving $M_d^{-1}\bar{P}_d = P$). To achieve such a goal, we fix on class i^* and collect elements in \bar{S}_d that do not contain i^* to form $E_d^{i^*} := \{\bar{s}_{d,j} / \bar{s}_{d,j} \mid \bar{S}_d, i^* / \bar{s}_{d,j}\}$ to connect $P_{\bar{S},X||\bar{S}|=d}$ with $P_{Y=i^*,X}$ as follows. Summing (119) over all elements in $E_d^{i^*}$, we obtain

$$\begin{aligned} \sum_{\bar{s} \in E_d^{i^*}} P_{\bar{S}=\bar{s},X||\bar{S}|=d} &= \sum_{\bar{s} \in E_d^{i^*}} \frac{1}{\binom{K-1}{d}} \sum_{i=1}^K \mathbb{1}[Y = i / \bar{s}] P_{Y=i,X} \\ &= \frac{1}{\binom{K-1}{d}} \left[\binom{K-2}{d} \sum_{i=1}^K P_{Y=i,X} + \binom{K-1}{d} P_{Y=i^*,X} \right]. \end{aligned}$$

The last equality holds since there are $\binom{K-2}{d}$ multi-complementary-labels $\bar{s} \in \bar{S}_d$ such that $i = i^*$ and neither of them is in \bar{s} (i.e., i / \bar{s} and i^* / \bar{s}), and there are $\binom{K-1}{d}$ multi-complementary-labels $\bar{s} \in \bar{S}_d$ such that $i = i^*$ and i is not in \bar{s} . Then, we regroup the sums by pulling $\binom{K-2}{d} P_{Y=i^*,X}$ out of $\binom{K-1}{d} P_{Y=i^*,X}$ to combine with $\binom{K-2}{d} \sum_{i=1}^K P_{Y=i,X}$. It leads to

$$\begin{aligned} \sum_{\bar{s} \in E_d^{i^*}} P_{\bar{S}=\bar{s},X||\bar{S}|=d} &= \frac{1}{\binom{K-1}{d}} \left[\binom{K-2}{d} \sum_{i=1}^K P_{Y=i,X} + \binom{K-2}{d-1} P_{Y=i^*,X} \right] \\ &= \frac{K-1-d}{K-1} P_X + \frac{d}{K-1} P_{Y=i^*,X}. \end{aligned} \quad (120)$$

Denoting $\bar{S}_d \setminus E_d^{i^*} = \{\bar{s}_{d,j} / \bar{s}_{d,j} \mid \bar{S}_d, i^* / \bar{s}_{d,j}\}$ as $I_d^{i^*}$ and rearranging terms in the above equation according to the reciprocal idea illustrated above, we have

$$\begin{aligned} P_{Y=i^*,X} &= \frac{K-1}{d} \left(\sum_{\bar{s} \in E_d^{i^*}} P_{\bar{S}=\bar{s},X||\bar{S}|=d} - \frac{K-1-d}{K-1} P_X \right) \\ &\stackrel{(a)}{=} \frac{K-1}{d} \left(P_X - \sum_{\bar{s} \in I_d^{i^*}} P_{\bar{S}=\bar{s},X||\bar{S}|=d} - \frac{K-1-d}{K-1} P_X \right) \\ &= P_X - \frac{K-1}{d} \sum_{\bar{s} \in I_d^{i^*}} P_{\bar{S}=\bar{s},X||\bar{S}|=d}. \end{aligned} \quad (121)$$

Equality (a) holds since $|\bar{S}|$ and X are independent (Feng et al., 2020a), which implies

$$P_X = P_{X||\bar{S}|=d} = \sum_{\bar{s} \in \bar{S}_d} P_{\bar{S}=\bar{s},X||\bar{S}|=d} = \sum_{\bar{s} \in E_d^{i^*}} P_{\bar{S}=\bar{s},X||\bar{S}|=d} + \sum_{\bar{s} \in I_d^{i^*}} P_{\bar{S}=\bar{s},X||\bar{S}|=d}.$$

Continuing the derivation, we have

$$\begin{aligned} P_{Y=i^*,X} &= \sum_{j=1}^{N_d} P_{\bar{S}=\bar{s}_{d,j},X||\bar{S}|=d} - \sum_{j=1}^{N_d} \frac{K-1}{d} \mathbb{1}[Y = i^* / \bar{s}_{d,j}] P_{\bar{S}=\bar{s}_{d,j},X||\bar{S}|=d} \\ &= \sum_{j=1}^{N_d} \left(1 - \frac{K-1}{d} \mathbb{1}[Y = i^* / \bar{s}_{d,j}] \right) P_{\bar{S}=\bar{s}_{d,j},X||\bar{S}|=d}, \end{aligned} \quad (122)$$

proving the first part of the lemma.

The derivation of turning (120) to (121) is a reciprocal action. Thus, if we view $1 - \frac{K-1}{d} \mathbb{1}[Y = i^* \mid \bar{s}_{d,j}]$ as the (i^*, j) entry of some matrix M , (122) can be interpreted as $P_{i^*} = \left(M \bar{P}_d \right)_{i^*}$, suggesting $M M_d = I$ since $\bar{P}_d = M_d P$. We formalize this intuition in the next lemma.

Lemma 43. *Let M be of the form (117), and recall M_d is defined by (114). Then, $M M_d = I$, meaning $M = M_d^{-1}$.*

The above lemma finishes the proof of Lemma 41. \square

Proof. of Lemma 43. Let d be fixed. Denoted by $A_{i,k}$, the (i, k) entry of $M M_d$, is the inner product of i -th row of M (117) and the k -th column of M_d (114)

$$A_{i,k} = \sum_{j=1}^{N_d} \left(1 - \frac{K-1}{d} \mathbb{1}[Y = i \mid \bar{s}_{d,j}] \right) \left(\frac{1}{\binom{K-1}{d}} \mathbb{1}[Y = k \mid \bar{s}_{d,j}] \right) = \sum_{j=1}^{N_d} c_{i,k}.$$

In the following, we will show that the calculation results in the identity matrix

$$A_{i,k} = \begin{cases} 1, & \text{if } i = k, \\ 0, & \text{if } i \neq k, \end{cases}$$

to complete the proof.

When $i = k$, we have 4 possible cases: (i) Both i and k are in $\bar{s}_{d,j}$, (ii) Both of them are not in $\bar{s}_{d,j}$, (iii) $i \in \bar{s}_{d,j}$ and $k \notin \bar{s}_{d,j}$, and (iv) $i \notin \bar{s}_{d,j}$ and $k \in \bar{s}_{d,j}$. For cases (i) and (iv), the coefficients $c_{i,k}$ are 0 since $\mathbb{1}[k \in \bar{s}_{d,j}] = 0$ if $k \notin \bar{s}_{d,j}$. For case (ii), the coefficient $c_{i,k}$ is $\frac{1}{\binom{K-1}{d}}$. The number of such $\bar{s}_{d,j}$ is $\binom{K-2}{d}$ since we are counting the ways of forming a set of size d from $K-2$ elements. For case (iii), the coefficient $c_{i,k}$ is $\left(1 - \frac{K-1}{d}\right) \frac{1}{\binom{K-1}{d}}$. The number of such $\bar{s}_{d,j}$ is $\binom{K-2}{d-1}$ since we are counting the ways of forming a set of size $d-1$ from $K-2$ elements. Thus, if $i = k$,

$$\begin{aligned} A_{i,k} &= \frac{1}{\binom{K-1}{d}} \binom{K-2}{d} + \left(1 - \frac{K-1}{d}\right) \frac{1}{\binom{K-1}{d}} \binom{K-2}{d-1} \\ &= \frac{\binom{K-2}{d}}{\binom{K-1}{d}} + \frac{\binom{K-2}{d-1}}{\binom{K-1}{d}} - \frac{K-1}{d} \frac{\binom{K-2}{d-1}}{\binom{K-1}{d}} = 0 \end{aligned}$$

since

$$\binom{K-2}{d} + \binom{K-2}{d-1} = \binom{K-1}{d} = \frac{K-1}{d} \binom{K-2}{d-1}.$$

When $i \neq k$, we have 2 possible cases: (i) Both i and k are in $\bar{s}_{d,j}$, (ii) Both are not in $\bar{s}_{d,j}$. For case (i), the coefficient $c_{i,k}$ is 0. For case (ii), the coefficient $c_{i,k}$ is $\frac{1}{\binom{K-1}{d}}$, and the number of such $\bar{s}_{d,j}$ is $\binom{K-1}{d}$, as we want to form a set of size d from $K-1$ candidates. Therefore, if $i \neq k$,

$$A_{i,k} = \frac{1}{\binom{K-1}{d}} \binom{K-1}{d} = 1.$$

\square

We want to elaborate more on the role of Theorem 1 of Wu et al. (2023) in the analyses discussed in Section 5.2. Firstly, as shown in the proof of Lemma 41, it aids the execution of the inversion approach (Proposition 1). The properness $C(S, X) \mathbb{1}[Y = S]$ (65) can be instantiated to define the entries of M_d (114),

which in turn establishes the key equation (120) enabling us to identify the entries of M_d^{-1} (122). Composing M_d^{-1} , we obtain M_{MCL}^{-1} , a crucial element for applying our framework (38).

Secondly, Theorem 1 of Wu et al. (2023) contributes to the marginal chain approach (Proposition 2) as well. The key equations (108) and (111) realised from Theorem 1 of Wu et al. (2023) provide the entries of M_{PPL}^\dagger (Lemma 36, Section 5.2.2), M_{PCPL}^\dagger (Section 5.2.3), and M_{MCL}^\dagger (Section 5.2.4) when applying (38). Therefore, the combined advantage of our framework and Theorem 1 of Wu et al. (2023) provides CCN scenarios unified analyses whose key steps can also be rationally interpreted. Moreover, as will be shown later, we compare the marginal chain and the inversion approaches via a CL example in Section 5.2.5. A CL example is the simplest way to convey the differences between the two methods without burying the essence in complicated derivations.

5.2.5 Complementary-Label (CL) Learning

Step 1: Corrected Loss Design and Risk Rewrite.

Note that the parameters chosen for the construction of M_{CL} (76) in Section 4.2.5 reduces M_{MCL} (113) to be M_1 of (114). That is, assigning $P_{|\bar{S}|=d} = 1$ for $d = 1$, $P_{|\bar{S}|=d} = 0$ for $d > 1$, and $\bar{s}_{1,j} = \{j\}$ for all $j \in [K]$ in (113), we have

$$M_{\text{MCL}} \quad M_1 = \frac{1}{K-1} \begin{pmatrix} 0 & 1 & \cdots & 1 \\ 1 & 0 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 0 \end{pmatrix} = M_{\text{CL}}.$$

Hence, the proof steps of Theorem 42 carry over to CL learning. With a simple rearranging on

$$\begin{aligned} \bar{\ell}_{\bar{S}} &= \sum_{i/\bar{S}} \ell_{Y=i} - \frac{K-1-|\bar{S}|}{|\bar{S}|} \sum_{\bar{s} \in \bar{S}} \ell_{Y=\bar{s}} \\ &= \sum_{i=1}^K \ell_{Y=i} - \frac{K-1}{|\bar{S}|} \sum_{\bar{s} \in \bar{S}} \ell_{Y=\bar{s}} \end{aligned}$$

and assigning $|\bar{S}| = 1$, we arrive at (21):

Corollary 44. *For CL learning, the classification risk can be expressed as*

$$R(g) = \mathbb{E}_{\bar{S}, X} [\bar{\ell}_{\bar{S}}] = \mathbb{E}_{\bar{S}, X} \left[\sum_{i=1}^K \ell_{Y=i} - (K-1)\ell_{\bar{S}} \right].$$

Step 2: Recovering the previous result(s).

The rewrite above recovers Theorem 1 of Ishida et al. (2019) if we substitute \bar{S} with \bar{Y} and $\ell_{\bar{S}}$ with $\ell(\bar{Y}, g(X))$. Moreover, if we choose $d = 1$ and $\bar{s}_{1,j} = \{j\}$ for all $j \in [K]$, the decontamination matrix provided by (117) becomes

$$M_1^{-1} = \begin{pmatrix} -(K-2) & 1 & \cdots & 1 \\ 1 & -(K-2) & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & -(K-2) \end{pmatrix}, \quad (123)$$

which translates the corrected losses $\bar{L} = L M_1^{-1}$ as

$$L \left(-(K-2)\mathbf{I}_K + \mathbf{1}\mathbf{1}^\top \right),$$

recovering (9) of Ishida et al. (2019).

Comparing inversion with marginal chain via an example.

We use a simple CL example to demonstrate the differences between the inversion (Proposition 1) and the marginal chain (Proposition 2) approaches and explain how the intuition of decontamination is implemented. Here, we focus on comparing how a decontamination matrix M_{corr}^\dagger achieves $M_{\text{corr}}^\dagger \bar{P} = P$ (37) since when the equality is established, the downstream construction of the corrected losses and the risk rewrite follow the framework. For this example, let us choose $K = 4$ and simplify $\mathcal{P}_{Y=k, X}$ as p_k . Applying (76), the contamination process defining the data-generating distributions is expressed as

$$\bar{P} = M_{\text{CL}} P = \frac{1}{3} \begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix} \begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{pmatrix} = \begin{pmatrix} \frac{p_2+p_3+p_4}{3} \\ \frac{p_1+p_3+p_4}{3} \\ \frac{p_1+p_2+p_4}{3} \\ \frac{p_1+p_2+p_3}{3} \end{pmatrix}.$$

Equation (123), simplified from (117), provides the decontamination matrix from the inversion approach:

$$M_{\text{CL}}^{-1} = \begin{pmatrix} -2 & 1 & 1 & 1 \\ 1 & -2 & 1 & 1 \\ 1 & 1 & -2 & 1 \\ 1 & 1 & 1 & -2 \end{pmatrix}.$$

Then, the inversion approach (Proposition 1) achieves the decontamination (37) by showing

$$\begin{aligned} M_{\text{CL}}^{-1} \bar{P} &= \frac{1}{3} \begin{pmatrix} -2 & 1 & 1 & 1 \\ 1 & -2 & 1 & 1 \\ 1 & 1 & -2 & 1 \\ 1 & 1 & 1 & -2 \end{pmatrix} \begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix} \begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{pmatrix} \\ &= \frac{1}{3} \begin{pmatrix} 3 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 3 \end{pmatrix} \begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{pmatrix} = \begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{pmatrix} = P. \end{aligned} \tag{124}$$

On the other hand, equation (111) produces the decontamination matrix from the marginal chain approach:

$$M_{\text{CL}}^\dagger = \begin{pmatrix} \frac{0 \cdot p_1}{p_2+p_3+p_4} & \frac{p_1}{p_1+p_3+p_4} & \frac{p_1}{p_1+p_2+p_4} & \frac{p_1}{p_1+p_2+p_3} \\ \frac{p_2}{p_2+p_3+p_4} & \frac{0 \cdot p_2}{p_1+p_3+p_4} & \frac{p_2}{p_1+p_2+p_4} & \frac{p_2}{p_1+p_2+p_3} \\ \frac{p_3}{p_2+p_3+p_4} & \frac{p_3}{p_1+p_3+p_4} & \frac{0 \cdot p_3}{p_1+p_2+p_4} & \frac{p_3}{p_1+p_2+p_3} \\ \frac{p_4}{p_2+p_3+p_4} & \frac{p_4}{p_1+p_3+p_4} & \frac{p_4}{p_1+p_2+p_4} & \frac{0 \cdot p_4}{p_1+p_2+p_3} \end{pmatrix}.$$

Then, the marginal chain approach (Proposition 2) achieves the decontamination (37) by showing

$$\begin{aligned}
M_{\text{CL}}^\dagger \bar{P} &= \begin{pmatrix} \frac{0 \cdot p_1}{p_2+p_3+p_4} & \frac{p_1}{p_1+p_3+p_4} & \frac{p_1}{p_1+p_2+p_4} & \frac{p_1}{p_1+p_2+p_3} \\ \frac{p_2}{p_2+p_3+p_4} & \frac{0 \cdot p_2}{p_1+p_3+p_4} & \frac{p_2}{p_1+p_2+p_4} & \frac{p_2}{p_1+p_2+p_3} \\ \frac{p_3}{p_2+p_3+p_4} & \frac{p_3}{p_1+p_3+p_4} & \frac{0 \cdot p_3}{p_1+p_2+p_4} & \frac{p_3}{p_1+p_2+p_3} \\ \frac{p_4}{p_2+p_3+p_4} & \frac{p_4}{p_1+p_3+p_4} & \frac{p_4}{p_1+p_2+p_4} & \frac{0 \cdot p_4}{p_1+p_2+p_3} \end{pmatrix} \begin{pmatrix} \frac{p_2+p_3+p_4}{3} \\ \frac{p_1+p_3+p_4}{3} \\ \frac{p_1+p_2+p_4}{3} \\ \frac{p_1+p_2+p_3}{3} \end{pmatrix} \\
&= \begin{pmatrix} \frac{p_1+p_1+p_1}{3} \\ \frac{p_2+p_2+p_2}{3} \\ \frac{p_3+p_3+p_3}{3} \\ \frac{p_4+p_4+p_4}{3} \end{pmatrix} = \begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \end{pmatrix} = P.
\end{aligned} \tag{125}$$

Comparing (124) and (125), we see that the intuition of decontamination is realized differently. The inversion approach (124) directly cancels out the effect of M_{corr} without relying on any property of P . In contrast, the marginal chain method (125) leverages the fact that P is a probability vector and carries out a procedure similar to importance reweighting to resolve the contamination. Both methods have respective merits, and we hope the comparison will inspire new thoughts leveraging certain properties of P for the corrected loss design and the study of decontamination.

5.3 Confidence-based Scenarios

The proposed framework is now applied to conduct the risk rewrites for WSLs discussed in Section 4.3 and summarized in Table 9.

5.3.1 Subset Confidence (Sub-Conf) Learning

Step 1: Corrected Loss Design and Risk Rewrite.

Let us follow the notations in Section 4.3.1. Recall that Lemma 16 has reached the first milestone (36) by showing $\bar{P} = M_{\text{Sub}}P$. To reach the second milestone (37), we apply Proposition 1 to construct the decontamination matrix M_{Sub}^\dagger to cancel out the contamination caused by M_{Sub} (80) as follows.

Lemma 45. *Assume $P_{Y \setminus \gamma_s/X} > 0$ for all possible outcomes of X . Choosing*

$$M_{\text{Sub}}^\dagger := M_{\text{Sub}}^{-1} = \begin{pmatrix} \frac{P_{Y=1/X}}{P_{Y \setminus \gamma_s/X}} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \frac{P_{Y=K/X}}{P_{Y \setminus \gamma_s/X}} \end{pmatrix}, \tag{126}$$

we have $M_{\text{Sub}}^\dagger \bar{P} = P$, where \bar{P} and $M_{\text{Sub}}P$ are given by Lemma 16.

Proof. The assumption $P_{Y \setminus \gamma_s/X} > 0$ implies M_{Sub} is invertible. As suggested by Proposition 1, we define $M_{\text{Sub}}^\dagger := M_{\text{Sub}}^{-1}$. Then,

$$M_{\text{Sub}}^\dagger M_{\text{Sub}} = \begin{pmatrix} \frac{P_{Y=1/X}}{P_{Y \setminus \gamma_s/X}} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \frac{P_{Y=K/X}}{P_{Y \setminus \gamma_s/X}} \end{pmatrix} \begin{pmatrix} \frac{P_{Y \setminus \gamma_s/X}}{P_{Y=1/X}} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \frac{P_{Y \setminus \gamma_s/X}}{P_{Y=K/X}} \end{pmatrix} = I$$

implies $M_{\text{Sub}}^\dagger \bar{P} = M_{\text{Sub}}^{-1} M_{\text{Sub}}P = P$ that proves the lemma. \square

With M_{Sub}^\dagger in hand, the next theorem defines the corrected losses \bar{L} and achieves the risk rewrite (31).

Theorem 46. For Sub-Conf learning, the classification risk can be written as

$$R(g) = \pi_{Y_s} \mathbb{E}_{X/Y_s} \left[\sum_{i=1}^K \frac{r_i(X)}{r_{Y_s}(X)} \ell_i \right].$$

Proof. Given Lemma 45, we can define $\bar{L} := L M_{\text{Sub}}^\dagger$ so that

$$\bar{L}_i = \left(L M_{\text{Sub}}^\dagger \right)_i = \frac{P_{Y=i/X}}{P_{Y_s/X}} \ell_i$$

for each $i \in [K]$ and

$$\bar{L} \bar{P} = L M_{\text{Sub}}^\dagger \bar{P} = L P.$$

Therefore, we can apply (39) to obtain (31) as follows.

$$\begin{aligned} R(g) &= \int_x \int_X L P dx = \int_x \int_X \bar{L} \bar{P} dx = \int_x \int_X \sum_{i=1}^K \frac{P_{Y=i/X}}{P_{Y_s/X}} \ell_i \cdot P_{Y_s} P_{X/Y_s} dx \\ &= P_{Y_s} \mathbb{E}_{X/Y_s} \left[\sum_{i=1}^K \frac{P_{Y=i/X}}{P_{Y_s/X}} \ell_i \right] \\ &= \pi_{Y_s} \mathbb{E}_{X/Y_s} \left[\sum_{i=1}^K \frac{r_i(X)}{r_{Y_s}(X)} \ell_i \right]. \end{aligned}$$

The last equality follows the notations in Section 2.2.14. \square

Step 2: Recovering the previous result(s).

Notation matching gives

$$R(g) = \pi_{Y_s} \mathbb{E}_{p(x/y_s)} \left[\sum_{y=1}^K \frac{r^y(x)}{r_{Y_s}(x)} \ell(g(x), y) \right],$$

recovering Theorem 6 of Cao et al. (2021a)¹⁴.

5.3.2 Single-Class Confidence (SC-Conf) Learning

Step 1: Corrected Loss Design and Risk Rewrite.

The SC-Conf derivation resembles that in Section 5.3.1 since M_{SC} is a child of M_{Sub} on the reduction graph. Thus, following the notations in Section 4.3.2, assuming $P_{Y=y_s/X} > 0$ for all possible outcomes of X , and replacing the set Y_s in M_{Sub}^\dagger (126) with a singleton y_s , we have

$$M_{\text{SC}}^\dagger := \begin{pmatrix} \frac{P_{Y=1/X}}{P_{Y=y_s/X}} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{P_{Y=K/X}}{P_{Y=y_s/X}} \end{pmatrix}$$

satisfying $M_{\text{SC}}^\dagger \bar{P} = P$. We also obtain $\bar{L} = L M_{\text{SC}}^\dagger$ and $\bar{L} \bar{P} = L P$ by inheriting the proof of Lemma 45.

Then, a modification to Theorem 46 by replacing $\bar{L}_i = \frac{P_{Y=i/X}}{P_{Y_s/X}} \ell_i$ with

$$\bar{L}_i = \left(L M_{\text{SC}}^\dagger \right)_i = \frac{P_{Y=i/X}}{P_{Y=y_s/X}} \ell_i = \frac{r_i(X)}{r_{y_s}(X)} \ell_i$$

proves the risk rewrite (29) for SC-Conf learning:

¹⁴The matching is as follows: P_{X/Y_s} is $p(x/y_s)$, $r_i(X)$ is $r^i(X)$, $r_{Y_s}(X)$ is $r^{y_s}(X)$, and ℓ_i is $\ell(g(X), i)$.

Corollary 47. For SC-Conf learning, the classification risk can be written as

$$R(g) = \pi_{y_s} \mathbb{E}_{X|Y=y_s} \left[\sum_{i=1}^K \frac{r_i(X)}{r_{y_s}(X)} \ell_i \right].$$

Step 2: Recovering the previous result(s).

By matching notations, we obtain

$$R(g) = \pi_{y_s} \mathbb{E}_{p(x/y_s)} \left[\sum_{y=1}^K \frac{r^y(x)}{r^{y_s}(x)} \ell(g(x), y) \right],$$

recovering Theorem 1 of Cao et al. (2021a)¹⁵.

5.3.3 Positive-confidence (Pconf) Learning

Step 1: Corrected Loss Design and Risk Rewrite.

Let us follow the notations in Section 4.3.3. Recall that M_{Pconf} is a child of M_{SC} on the reduction graph with $K = 2$ and $y_s = p$. Thus, assuming $P_{Y=p/X} > 0$ for all possible outcomes of X and replacing K and y_s in Section 5.3.2 accordingly, we obtain the decontamination matrix

$$M_{\text{Pconf}}^\dagger := \begin{pmatrix} \frac{P_{Y=p/X}}{P_{Y=p/X}} & 0 \\ 0 & \frac{P_{Y=n/X}}{P_{Y=p/X}} \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & \frac{1-r(X)}{r(X)} \end{pmatrix}$$

and the rewrite (7) reviewed in Section 2.2.2.

Corollary 48. For Pconf learning, the classification risk can be written as

$$R(g) = \pi_p \mathbb{E}_p \left[\ell_p + \frac{1-r(X)}{r(X)} \ell_n \right].$$

Step 2: Recovering the previous result(s).

By matching notations, we obtain

$$R(g) = \pi_+ \mathbb{E}_+ \left[\ell(g(x)) + \frac{1-r(x)}{r(x)} \ell(-g(x)) \right],$$

recovering Theorem 1 of Ishida et al. (2018)¹⁶.

5.3.4 Soft-Label Learning

Step 1: Corrected Loss Design and Risk Rewrite.

We follow the notations in Section 4.3.4. As discussed in Section 4.3.4, M_{Soft} is a special case of M_{Sub} when $\mathcal{Y}_s := [K]$. Thus, reducing M_{Sub}^\dagger (126) by assigning $P_{Y=y_s/X} = P_{Y=[K]/X} = 1$, we obtain the the decontamination matrix

$$M_{\text{Soft}}^\dagger := \begin{pmatrix} P_{Y=1/X} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & P_{Y=K/X} \end{pmatrix}.$$

Then, we follow the same argument in Theorem 46 to achieve the rewrite (34) by the next corollary.

Corollary 49. For soft-label learning, the classification risk can be written as

$$R(g) = \mathbb{E}_X \left[\sum_{i=1}^K P_{Y=i/X} \ell_i \right] = \mathbb{E}_X \left[\sum_{i=1}^K r_i(X) \ell_i \right].$$

¹⁵The matching is as follows: $P_{X|Y=y_s}$ is $p(x/y_s)$, $r_i(X)$ is $r^i(X)$, $r_{y_s}(X)$ is $r^{y_s}(X)$, and ℓ_i is $\ell(g(X), i)$.

¹⁶The matching is as follows: π_p is π_+ , $P_{X|Y=p}$ is $p(x/y = +1)$, $P_{Y=p/X}$ is $r(x)$, $P_{Y=n/X}$ is $1 - r(x)$, ℓ_p is $\ell(g(x))$, and ℓ_n is $\ell(-g(x))$.

Step 2: Recovering the previous result(s).

Ishida et al. (2023) did not focus on the classification risk rewrite problem. We can modify Corollary 49 to provide a risk rewrite for binary soft-label learning mentioned by Ishida et al. (2023). Taking $K = 2$, we have

$$R(g) = \mathbb{E}_X [P_{Y=p/X} \ell_p + P_{Y=n/X} \ell_n] = \mathbb{E}_X [r(X) \ell_p + (1 - r(X)) \ell_n].$$

6 New Risk Rewrites

The proposed framework can be applied to derive risk rewrites for new scenarios. In this section, we show how the proposed framework accommodates a different performance metric and adapts to a noisy environment.

6.1 The Balanced Error Rate

In the previous sections, we chose the most common metric, the classification risk, to elaborate our framework. Here, we demonstrate that the framework adapts to another common but different performance metric. The balanced error rate (BER) is defined as

$$R_{\text{BER}}(g) := \frac{\mathbb{E}_X [P_{X/Y=p} [\ell_p(g(X))] + \mathbb{E}_X [P_{X/Y=n} [\ell_n(g(X))]]}{2} \quad (127)$$

(Scott & Zhang, 2020). Next, we show how to apply the framework to obtain the risk rewrite of UU learning under BER.

Recall that in Section 5.1.1, for classification risk, $M_{\text{trsf}} = \begin{pmatrix} 1/\pi_p & 0 \\ 0 & 1/\pi_n \end{pmatrix}$ was chosen to link the base

distributions $B = \begin{pmatrix} P_{X/Y=p} \\ P_{X/Y=n} \end{pmatrix}$ and the risk-defining distributions $P = \begin{pmatrix} P_{Y=p,X} \\ P_{Y=n,X} \end{pmatrix}$:

$$\begin{pmatrix} P_{X/Y=p} \\ P_{X/Y=n} \end{pmatrix} = \begin{pmatrix} 1/\pi_p & 0 \\ 0 & 1/\pi_n \end{pmatrix} \begin{pmatrix} P_{Y=p,X} \\ P_{Y=n,X} \end{pmatrix}.$$

By definition (127), the risk-defining distributions become $\begin{pmatrix} P_{X/Y=p}/2 \\ P_{X/Y=n}/2 \end{pmatrix}$. Since the base distributions B

remain unchanged, we must redefine $M_{\text{trsf}} := \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$, so that

$$\begin{pmatrix} P_{X/Y=p} \\ P_{X/Y=n} \end{pmatrix} = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} P_{X/Y=p}/2 \\ P_{X/Y=n}/2 \end{pmatrix}$$

satisfies $B = M_{\text{trsf}} P$ for the BER setting.

With the newly defined $M_{\text{trsf}} = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$ in hand, we can follow Step 1 in Section 5.1.1 to obtain the data-generating process

$$\begin{pmatrix} P_{U_1} \\ P_{U_2} \end{pmatrix} = M_{\text{UU}} \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} P_{X/Y=p}/2 \\ P_{X/Y=n}/2 \end{pmatrix}.$$

Then, by assigning $M_{\text{UU-BER}}^\dagger := \begin{pmatrix} 1/2 & 0 \\ 0 & 1/2 \end{pmatrix} M_{\text{UU}}^{-1}$ and following the similar derivation in the proof of Theorem 21, we achieve the risk rewrite for UU learning under BER:

$$R_{\text{BER}}(g) = \mathbb{E}_{U_1} [\bar{\ell}_{U_1}] + \mathbb{E}_{U_2} [\bar{\ell}_{U_2}],$$

where

$$\begin{aligned}\bar{\ell}_{U_1} &= \frac{1 - \gamma_2}{2(1 - \gamma_1 - \gamma_2)} \ell_p + \frac{-\gamma_2}{2(1 - \gamma_1 - \gamma_2)} \ell_n, \\ \bar{\ell}_{U_2} &= \frac{-\gamma_1}{2(1 - \gamma_1 - \gamma_2)} \ell_p + \frac{1 - \gamma_1}{2(1 - \gamma_1 - \gamma_2)} \ell_n.\end{aligned}$$

6.2 Learning with Label Noise

In this section, we show that the framework can easily handle the risk rewrite under label noise. We will take PU learning as an example. We use P to denote the clean distributions and Q for the noisy ones. The scenario is formulated as follows. According to Lemma 4, the data-generating process of the noisy PU learning is of the form

$$\begin{pmatrix} Q_P \\ Q_U \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ \pi_p & \pi_n \end{pmatrix} \begin{pmatrix} Q_{X/Y=p} \\ Q_{X/Y=n} \end{pmatrix},$$

where $\pi_p = Q_{Y=p}$ and $\pi_n = Q_{Y=n}$. We choose the MCD setting (Section 2.2.17) to formulate the label noise:

$$\begin{pmatrix} Q_{X/Y=p} \\ Q_{X/Y=n} \end{pmatrix} = \begin{pmatrix} 1 - \alpha_p & \alpha_p \\ \alpha_n & 1 - \alpha_n \end{pmatrix} \begin{pmatrix} P_{X/Y=p} \\ P_{X/Y=n} \end{pmatrix},$$

where $\alpha_p = 0$ and $\alpha_n = 0$ are parameters describing the degree of noise perturbation. As in the previous section, we also choose BER as the performance metric.

By cascading matrices, the framework effortlessly links the data-generating distributions and the risk-defining distributions:

$$\begin{pmatrix} Q_P \\ Q_U \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ \pi_p & \pi_n \end{pmatrix} \begin{pmatrix} 1 - \alpha_p & \alpha_p \\ \alpha_n & 1 - \alpha_n \end{pmatrix} \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} P_{X/Y=p}/2 \\ P_{X/Y=n}/2 \end{pmatrix}.$$

Knowing the contamination process, we apply the inversion method (Proposition 1) to construct the decontamination matrix

$$\begin{aligned}M_{\text{PU-BER}}^\dagger &:= \begin{pmatrix} 1/2 & 0 \\ 0 & 1/2 \end{pmatrix} \frac{1}{1 - \alpha_p - \alpha_n} \begin{pmatrix} 1 - \alpha_n & -\alpha_p \\ -\alpha_n & 1 - \alpha_p \end{pmatrix} \begin{pmatrix} 1 & 0 \\ -\pi_p/\pi_n & 1/\pi_n \end{pmatrix} \\ &= \begin{pmatrix} \frac{(1 - \alpha_n)\pi_n + \alpha_p\pi_p}{2(1 - \alpha_p - \alpha_n)\pi_n} & \frac{-\alpha_p}{2(1 - \alpha_p - \alpha_n)\pi_n} \\ \frac{-\alpha_n\pi_n - (1 - \alpha_p)\pi_p}{2(1 - \alpha_p - \alpha_n)\pi_n} & \frac{1 - \alpha_p}{2(1 - \alpha_p - \alpha_n)\pi_n} \end{pmatrix} := \begin{pmatrix} c_1 & c_2 \\ c_3 & c_4 \end{pmatrix}.\end{aligned}$$

Then, applying equations (38) and (39), we can rewrite the BER for PU learning as

$$R_{\text{BER}}(g) = \mathbb{E}_{X \sim Q_P} [\bar{\ell}_P] + \mathbb{E}_{X \sim Q_U} [\bar{\ell}_U],$$

where

$$\begin{aligned}\bar{\ell}_P &= c_1 \ell_p + c_3 \ell_n, \\ \bar{\ell}_U &= c_2 \ell_p + c_4 \ell_n.\end{aligned}$$

In the two subsections above, we have shown that the proposed framework can address risk rewrite under a different performance metric and a complex system. In the outlook part of the next section, we will further discuss the potential of the framework.

7 Conclusion and Outlook

We set out with the questions wishing to determine if there is a common way to interpret the formation of weak supervision and search for a generic treatment to solve WSL, to understand the essence of WSL. In response, we proposed a framework that unifies the formulations and analyses of a set of WSL scenarios to provide a common ground to connect, compare, and understand various weakly-supervised signals. The formulation component of the proposed framework, viewing WSL from a contamination perspective, associates a WSL data-generating process with a base distribution vector multiplied by a contamination matrix. By instantiating the contamination matrices of WSLs, we revealed a comprehensive reduction graph, Figure 1, connecting existing WSLs. Each vertex contains a contamination matrix and the section index of the WSL scenario which the matrix characterizes. Each edge represents the reduction relation of two WSLs. We can see three major branches from the abstract M_{corr} , corresponding to Tables 7, 8, and 9 we discussed in Section 4. The analysis component of the proposed framework, tackling the problem from a decontamination viewpoint, working with the technical building blocks Theorems 1 and 2 constitute a generic treatment to solve the risk rewrite problem. Section 5 discussed in depth how the analysis component conducts risk rewrite and recovers existing results for WSLs.

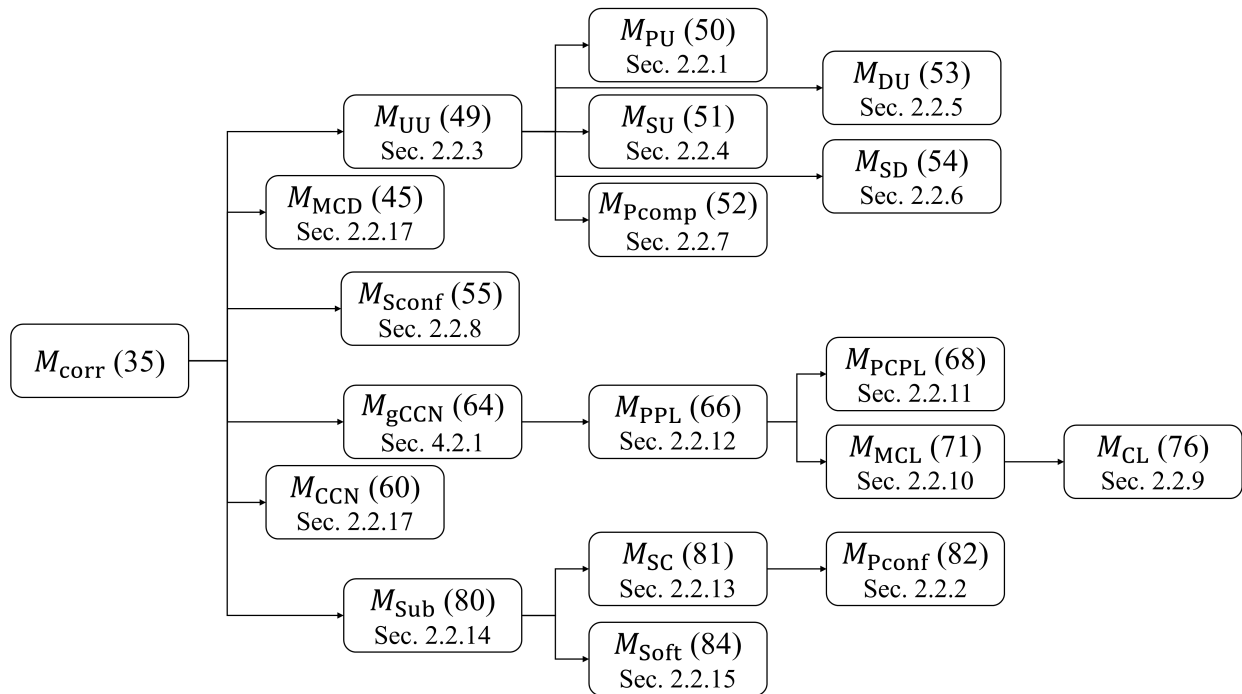


Figure 1: Depicting the reduction map from Tables 7, 8, and 9.

The application of the proposed framework results in a set of theorems. We summarize them in Table 10. The Formulation column consists of the results of the formulation component (35). The Decontamination and the Corrected losses columns correspond to the results of the analysis component ((37), (38), and (39)). The Recovery column justifies the framework by recovering results from the literature. Crucial results are marked red. Since the analyses of different scenarios are subsumed under a single framework, we now have a basis for transferring a technique developed for one scenario to another. In addition, these alternative proofs provide different ways of dissecting the risk, which in turn could aid in the development of a training objective by examining multiple risk decomposition approaches.

Table 10: Theorem Structure.

Model	Formulation (Find M s.t. $\bar{P} = MB.$)	Decontamination (Find M^\dagger s.t. $P = M^\dagger \bar{P}.$)	Corrected losses (Rewrite via $\bar{L} = L M^\dagger$ and $\bar{P}.$)	Recovery
Abstract model	(35)	(37) Proposition 1 and Proposition 2	(38) and (39)	
MCD				
UU	Lemma 3	Corollary 20	Theorem 21	(Notation swap.)
PU	Lemma 4	(Immediate reduction.)	Corollary 22	(Notation swap.)
SU	Lemma 5	(Immediate reduction.)	Corollary 23	Lemmas 24 and 25
Pcomp	Lemma 6	(Immediate reduction.)	Corollary 26	(Notation swap.)
DU	Lemma 7	(Immediate reduction.)	Corollary 27	Lemmas 28 and 29
SD	Lemma 8	(Immediate reduction.)	Corollary 30	Lemma 31
Sconf	Lemma 9	Lemmas 32 and 33	Theorem 34	(Notation swap.)
CCN				
gCCN	Lemma 10	(106) and Proposition 2	Theorem 35	(Notation swap.)
PPL	Lemma 11	Lemma 36	Corollary 37	(Notation swap.)
PCPL	Lemma 12	Corollary 38	Corollary 38	Lemma 39
MCL	Lemma 13	Corollary 40	Corollary 40	Theorem 42, Lemmas 41 and 43
CL	Lemma 15	(Immediate reduction.)	Corollary 44	(Notation swap.)
Sub-Conf	Lemma 16	Lemma 45	Theorem 46	(Notation swap.)
SC-Conf	Lemma 17	(Immediate reduction.)	Corollary 47	(Notation swap.)
Pconf	Lemma 18	(Immediate reduction.)	Corollary 48	(Notation swap.)
Soft	Lemma 19	(Immediate reduction.)	Corollary 49	(N/A.)

The proposed framework is abstract and flexible; hence, we would like to discuss its potential from the following aspects. Firstly, the performance measure focused on in this paper is the classification risk. With proper choices of P and L , our framework can be extended to other performance metrics, such as the balanced error rate, one-versus-rest risk, and cost-sensitive measures (Rifkin & Klautau, 2004; Zhang, 2004; Brodersen et al., 2010; du Plessis et al., 2014; Menon et al., 2015; Blanchard et al., 2016; Natarajan et al., 2017; Scott & Zhang, 2020). We have demonstrated the applicability of the proposed framework for the balanced error rate in Section 6. Secondly, we can explore the formulation capability by exploiting the power of matrix operations. Cascading matrices allow us to formulate complex scenarios, such as data containing preference relations collected in a noisy environment. Matrix addition allows us to categorize different contamination mechanisms into cases to capture the structural properties of a problem. A complicated scenario could undergo a sophisticated formulation procedure, but once we have the resulting contamination matrix, the problem boils down to calculating the corresponding decontamination matrix. Thirdly, the MCD scenarios discussed in this paper (Sections 4.1 and 5.1) belong to binary classification. A way of extending an MCD formulation to multiclass classification is to extend M_{MCD} (45) from a 2×2 matrix to a $K \times K$ one, in which $K^2 - K$ mixture rates are used to characterize the extended M_{gMCD} : the (i, j) entry is $\gamma_{i,j}$ if $i = j$ and is $1 - \sum_{j=i} \gamma_{i,j}$ for the i -th entry on the diagonal. Fourthly, the label-flipping probabilities $P_{\bar{Y}|Y}$ in Natarajan et al. (2017) and Feng et al. (2020b) assume that the contaminated label \bar{Y} is independent of X condition on the true label Y . The formulation matrices, M_{CCN} (60) and M_{gCCN} (64), in contrast, take X into consideration. This formulation enables us to tackle the instance-dependent problem (Berthon et al., 2021) and the effect of sampling strategies such as SAR and SCAR (Elkan & Noto, 2008; Coudray et al., 2023) in the future. Fifthly, we hope that the analysis technique developed in Section 5.2, which combines marginal chain and properness, opens up a new possibility to search for invertibility-free methods for the risk rewrite problem. We also project its potential in research regarding the broader sense of contamination

and decontamination. Sixthly, the properness of Wu et al. (2023) provides an efficient technique to compute $P_{Y|S,X}$ needed in M_{gCCN} (106). It would be intriguing to know if there are any other alternatives. Finally but not least, the proposed framework operating under matrix multiplication belongs to a broader question of under what circumstances does a function f^{\dagger} exist with $P = f^{\dagger}(\bar{P})$ if $\bar{P} = f(P)$.

Acknowledgments

The authors were supported by the Institute for AI and Beyond, UTokyo. The first author would like to thank Professor Takashi Ishida (UTokyo) for valuable insights and discussions in extending the coverage of the framework, and the colleagues Xin-Qiang Cai, Masahiro Negishi, Wei Wang, and Yivan Zhang (in alphabetical order) for comments in improving the manuscript.

References

- Han Bao, Gang Niu, and Masashi Sugiyama. Classification from pairwise similarity and unlabeled data. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018*, pp. 461–470, Stockholmsmässan, Stockholm, Sweden, 2018.
- Antonin Berthon, Bo Han, Gang Niu, Tongliang Liu, and Masashi Sugiyama. Confidence scores make instance-dependent label-noise learning possible. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, Virtual Event*, pp. 825–836, 2021.
- Gilles Blanchard and Clayton Scott. Decontamination of mutually contaminated models. In *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics, AISTATS 2014*, pp. 1–9, Reykjavik, Iceland, 2014.
- Gilles Blanchard, Marek Flaska, Gregory Handy, Sara Pozzi, and Clayton Scott. Classification with asymmetric label noise: Consistency and maximal denoising. *Electronic Journal of Statistics*, 10(2):2780–2824, 2016.
- Kay Henning Brodersen, Cheng Soon Ong, Klaas Enno Stephan, and Joachim M. Buhmann. The balanced accuracy and its posterior distribution. In *20th International Conference on Pattern Recognition, ICPR 2010*, pp. 3121–3124, Istanbul, Turkey, 2010.
- Yuzhou Cao, Lei Feng, Senlin Shu, Yitian Xu, Bo An, Gang Niu, and Masashi Sugiyama. Multi-class classification from single-class data with confidences. *CoRR*, abs/2106.08864, 2021a. URL <https://arxiv.org/abs/2106.08864>.
- Yuzhou Cao, Lei Feng, Yitian Xu, Bo An, Gang Niu, and Masashi Sugiyama. Learning from similarity-confidence data. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, Virtual Event*, pp. 1272–1282, 2021b.
- Nontawat Charoenphakdee, Jongyeong Lee, and Masashi Sugiyama. On symmetric losses for learning from corrupted labels. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, Long Beach, California, USA*, pp. 961–970, 2019.
- Nontawat Charoenphakdee, Zhenghang Cui, Yivan Zhang, and Masashi Sugiyama. Classification with rejection based on cost-sensitive classification. In *Proceedings of 38th International Conference on Machine Learning, ICML 2021, Virtual Event*, pp. 1507–1517, 2021.
- Yu-Ting Chou, Gang Niu, Hsuan-Tien Lin, and Masashi Sugiyama. Unbiased risk estimators can mislead: A case study of learning with complementary labels. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, Virtual Event*, pp. 1929–1938, 2020.
- Jesús Cid-Sueiro. Proper losses for learning from partial labels. In *Advances in Neural Information Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012*, pp. 1574–1582, Lake Tahoe, Nevada, USA, 2012.

- Olivier Coudray, Christine Keribin, Pascal Massart, and Patrick Pamphile. Risk bounds for positive-unlabeled learning under the selected at random assumption. *Journal of Machine Learning Research*, 24:107:1–107:31, 2023. URL <https://jmlr.org/papers/v24/22-067.html>.
- Timothee Cour, Ben Sapp, and Ben Taskar. Learning from partial labels. *Journal of Machine Learning Research*, 12:1501–1536, 2011. URL <https://dl.acm.org/doi/10.5555/1953048.2021049>.
- Marthinus Christoffel du Plessis, Gang Niu, and Masashi Sugiyama. Clustering unclustered data: Unsupervised binary labeling of two datasets having different class balances. In *Proceedings of Conference on Technologies and Applications of Artificial Intelligence, TAAI 2013*, pp. 1–6, Taipei, Taiwan, 2013.
- Marthinus Christoffel du Plessis, Gang Niu, and Masashi Sugiyama. Analysis of learning from positive and unlabeled data. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014*, pp. 703–711, Montréal, Québec, Canada, 2014.
- Marthinus Christoffel du Plessis, Gang Niu, and Masashi Sugiyama. Convex formulation for learning from positive and unlabeled data. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015*, pp. 1386–1394, Lille, France, 2015.
- Charles Elkan and Keith Noto. Learning classifiers from only positive and unlabeled data. In *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 213–220, Las Vegas, Nevada, USA, 2008.
- Lei Feng, Takuo Kaneko, Bo Han, Gang Niu, Bo An, and Masashi Sugiyama. Learning with multiple complementary labels. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, Virtual Event*, pp. 3072–3081, 2020a.
- Lei Feng, Jiaqi Lv, Bo Han, Miao Xu, Gang Niu, Xin Geng, Bo An, and Masashi Sugiyama. Provably consistent partial-label learning. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, Virtual Event, 2020b*. URL <https://proceedings.neurips.cc/paper/2020/hash/7bd28f15a49d5e5848d6ec70e584e625-Abstract.html>.
- Lei Feng, Senlin Shu, Nan Lu, Bo Han, Miao Xu, Gang Niu, Bo An, and Masashi Sugiyama. Pointwise binary classification with pairwise confidence comparisons. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, Virtual Event*, pp. 3252–3262, 2021.
- Jacob Goldberger and Ehud Ben-Reuven. Training deep neural-networks using a noise adaptation layer. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, Conference Track Proceedings*, Toulon, France, 2017. URL <https://openreview.net/forum?id=H12GRgxcg>.
- Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor W. Tsang, and Masashi Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018*, pp. 8536–8546, Montréal, Québec, Canada, 2018.
- Takashi Ishida, Gang Niu, Weihua Hu, and Masashi Sugiyama. Learning from complementary labels. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017*, pp. 5639–5649, Long Beach, California, USA, 2017.
- Takashi Ishida, Gang Niu, and Masashi Sugiyama. Binary classification from positive-confidence data. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018*, pp. 5921–5932, Montréal, Québec, Canada, 2018.
- Takashi Ishida, Gang Niu, Aditya Krishna Menon, and Masashi Sugiyama. Complementary-label learning for arbitrary losses and models. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019*, pp. 2971–2980, Long Beach, California, USA, 2019.

- Takashi Ishida, Ikko Yamane, Nontawat Charoenphakdee, Gang Niu, and Masashi Sugiyama. Is the performance of my deep network too good to be true? A direct approach to estimating the bayes error in binary classification. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, 2023*. URL <https://openreview.net/forum?id=FZdJQgy05rz>.
- Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018*, pp. 2309–2318, Stockholmsmässan, Stockholm, Sweden, 2018.
- Yasuhiro Katsura and Masato Uchida. Bridging ordinary-label learning and complementary-label learning. In *Proceedings of The 12th Asian Conference on Machine Learning, ACML 2020*, pp. 161–176, Bangkok, Thailand, 2020.
- Julian Katz-Samuels, Gilles Blanchard, and Clayton Scott. Decontamination of mutual contamination models. *Journal of Machine Learning Research*, 20:41:1–41:57, 2019. URL <http://jmlr.org/papers/v20/17-576.html>.
- Ryuichi Kiryo, Gang Niu, Marthinus Christoffel du Plessis, and Masashi Sugiyama. Positive-unlabeled learning with non-negative risk estimator. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017*, pp. 1675–1685, Long Beach, California, USA, 2017.
- Nan Lu, Gang Niu, Aditya Krishna Menon, and Masashi Sugiyama. On the minimal supervision for training any binary classifier from only unlabeled data. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, Louisiana, USA, New Orleans, Louisiana, USA, 2019*. URL <https://openreview.net/forum?id=B1xWcjOqYm>.
- Nan Lu, Tianyi Zhang, Gang Niu, and Masashi Sugiyama. Mitigating overfitting in supervised classification from two unlabeled datasets: A consistent risk correction approach. In *The 23rd International Conference on Artificial Intelligence and Statistics, AISTATS 2020, Online Event*, pp. 1115–1125, 2020. URL <http://proceedings.mlr.press/v108/lu20c.html>.
- Nan Lu, Shida Lei, Gang Niu, Issei Sato, and Masashi Sugiyama. Binary classification from multiple unlabeled datasets via surrogate set classification. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, Virtual Event*, pp. 7134–7144, 2021. URL <http://proceedings.mlr.press/v139/lu21c.html>.
- Jiaqi Lv, Miao Xu, Lei Feng, Gang Niu, Xin Geng, and Masashi Sugiyama. Progressive identification of true labels for partial-label learning. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, Virtual Event*, pp. 6500–6510, 2020.
- Xingjun Ma, Yisen Wang, Michael E. Houle, Shuo Zhou, Sarah M. Erfani, Shu-Tao Xia, Sudanthi N. R. Wijewickrema, and James Bailey. Dimensionality-driven learning with noisy labels. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018*, pp. 3361–3370, Stockholmsmässan, Stockholm, Sweden, 2018.
- Aditya Krishna Menon, Brendan van Rooyen, Cheng Soon Ong, and Bob Williamson. Learning from corrupted binary labels via class-probability estimation. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015*, pp. 125–134, Lille, France, 2015.
- Nagarajan Natarajan, Inderjit S. Dhillon, Pradeep Ravikumar, and Ambuj Tewari. Learning with noisy labels. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013*, pp. 1196–1204, Lake Tahoe, Nevada, USA, 2013.
- Nagarajan Natarajan, Inderjit S. Dhillon, Pradeep Ravikumar, and Ambuj Tewari. Cost-sensitive learning with noisy labels. *Journal of Machine Learning Research*, 18:155:1–155:33, 2017. URL <http://jmlr.org/papers/v18/15-226.html>.

- Masahiro Negishi. Pairwise-constraint classification in weakly supervised machine learning: Risk-consistent approach and classifier-consistent approach. Senior thesis, 2023.
- Gang Niu, Marthinus Christoffel du Plessis, Tomoya Sakai, Yao Ma, and Masashi Sugiyama. Theoretical comparisons of positive-unlabeled learning against positive-negative learning. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016*, pp. 1199–1207, Barcelona, Spain, 2016.
- Curtis G. Northcutt, Tailin Wu, and Isaac L. Chuang. Learning with confident examples: Rank pruning for robust classification with noisy labels. In *Proceedings of the Thirty-Third Conference on Uncertainty in Artificial Intelligence, UAI 2017*, Sydney, Australia, 2017. URL <http://auai.org/uai2017/proceedings/papers/35.pdf>.
- Curtis G. Northcutt, Lu Jiang, and Isaac L. Chuang. Confident learning: Estimating uncertainty in dataset labels. *Journal of Artificial Intelligence Research*, 70:1373–1411, 2021.
- Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu. Making deep neural networks robust to label noise: A loss correction approach. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, pp. 2233–2241, Honolulu, Hawaii, USA, 2017.
- Scott E. Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, and Andrew Rabinovich. Training deep neural networks on noisy labels with bootstrapping. In *3rd International Conference on Learning Representations, ICLR 2015, Workshop Track Proceedings*, San Diego, California, USA, 2015. URL <http://arxiv.org/abs/1412.6596>.
- Ryan M. Rifkin and Aldebaro Klautau. In defense of one-vs-all classification. *Journal of Machine Learning Research*, 5:101–141, 2004. URL <https://jmlr.org/papers/v5/rifkin04a.html>.
- Tomoya Sakai, Gang Niu, and Masashi Sugiyama. Semi-supervised AUC optimization based on positive-unlabeled learning. *Machine Learning*, 107(4):767–794, 2018.
- Emanuele Sansone, Francesco G. B. De Natale, and Zhi-Hua Zhou. Efficient training for positive unlabeled learning. *IEEE transactions on pattern analysis and machine intelligence*, 41(11):2584–2598, 2019.
- Clayton Scott and Jianxin Zhang. Learning from multiple corrupted sources, with application to learning from label proportions. *CoRR*, abs/1910.04665, 2019. URL <http://arxiv.org/abs/1910.04665>.
- Clayton Scott and Jianxin Zhang. Learning from label proportions: A mutual contamination framework. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, Virtual Event*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/fcde14913c766cf307c75059e0e89af5-Abstract.html>.
- Clayton Scott, Gilles Blanchard, and Gregory Handy. Classification with asymmetric label noise: Consistency and maximal denoising. In *COLT 2013 - The 26th Annual Conference on Learning Theory*, pp. 489–511, Princeton University, New Jersey, USA, 2013.
- Takuya Shimada, Han Bao, Issei Sato, and Masashi Sugiyama. Classification from pairwise similarities/dissimilarities and unlabeled data via empirical risk minimization. *Neural Computation*, 33(5):1234–1268, 2021.
- Masashi Sugiyama, Han Bao, Takashi Ishida, Nan Lu, Tomoya Sakai, and Gang Niu. *Machine Learning from Weak Supervision: An Empirical Risk Minimization Approach*. Adaptive Computation and Machine Learning series. MIT Press, 2022. ISBN 9780262047074.
- Sainbayar Sukhbaatar and Rob Fergus. Learning from noisy labels with deep neural networks. In *3rd International Conference on Learning Representations, ICLR 2015, Workshop Track Proceedings*, San Diego, California, USA, 2015. URL <http://arxiv.org/abs/1406.2080>.

- Daiki Tanaka, Daiki Ikami, Toshihiko Yamasaki, and Kiyoharu Aizawa. Joint optimization framework for learning with noisy labels. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018*, pp. 5552–5560, Salt Lake City, Utah, USA, 2018.
- Ambuj Tewari and Peter L. Bartlett. Chapter 14 - learning theory. In *Academic Press Library in Signal Processing: Volume 1*, pp. 775–816. Elsevier, 2014. doi: <https://doi.org/10.1016/B978-0-12-396502-8.00014-0>.
- Brendan van Rooyen and Robert C. Williamson. Learning in the presence of corruption. *CoRR*, abs/1504.00091, 2015. URL <http://arxiv.org/abs/1504.00091>.
- Brendan van Rooyen and Robert C. Williamson. A theory of learning with corrupted labels. *Journal of Machine Learning Research*, 18:228:1–228:50, 2017. URL <http://jmlr.org/papers/v18/16-315.html>.
- Qian-Wei Wang, Yu-Feng Li, and Zhi-Hua Zhou. Partial label learning with unlabeled data. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019*, pp. 3755–3761, Macao, China, 2019.
- Zhenguo Wu, Jiaqi Lv, and Masashi Sugiyama. Learning with proper partial labels. *Neural Computation*, 35(1):58–81, 2023. URL https://doi.org/10.1162/neco_a_01554.
- Xingrui Yu, Bo Han, Jiangchao Yao, Gang Niu, Ivor W. Tsang, and Masashi Sugiyama. How does disagreement help generalization against label corruption? In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019*, pp. 7164–7173, Long Beach, California, USA, 2019.
- Xiyu Yu, Tongliang Liu, Mingming Gong, and Dacheng Tao. Learning with biased complementary labels. In *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, Proceedings, Part I*, pp. 69–85, 2018.
- Tong Zhang. Statistical analysis of some multi-category large margin classification methods. *Journal of Machine Learning Research*, 5:1225–1251, 2004. URL <https://jmlr.org/papers/volume5/zhang04b/zhang04b.pdf>.
- Zhi-Hua Zhou. A brief introduction to weakly supervised learning. *National science review*, 5(1):44–53, 2018.

A Use Cases for This Paper

There are several ways to use this paper in the study of WSLs. We provide some common use cases as follows:

Use Case 1: If you want a quick overview of WSL. Reading Section 2 would provide a quick catch-up on the formulations and results of various WSL scenarios.

Use Case 2: If you want to get the high-level idea of this paper. One could start with Section 3, which provides the abstract form of the proposed framework. Then, Sections 4.1.1 and 5.1.1 give the practical application of the framework to UU learning. If you are more familiar with PU learning, then you are welcome to proceed to Sections 4.1.2 and 5.1.1 to see how our framework interprets PU learning. In addition, these reading steps will reveal the connection between UU and PU.

Use Case 3: If you want to know how to apply the framework to a specific scenario. One could read Sections 4.a.b and 5.a.b at the same time, and come back to the rest of Sections 4 and 5 only as needed. Note that “a.b” represents the index of a subsection.

Use Case 4: If you want to know the connections discovered and the analysis techniques developed in this paper. Section 4 provides detailed discussions of the formulations, and Figure 1 summarizes the relationship. Section 5 provides detailed explanations of how our framework is applied to rewrite the classification risk. The extensive analyses from Section 5 are summarized in Table 10.

B Notations

Table 11: Notations and Aliases.

Name of the notation	Expression	Aliases	Convention
Example	(y, x)		(x, y)
Binary classes	$\{p, n\}$		$\{+1, -1\}$
Multiple classes	$\{1, \dots, K\}$	$[K]$	
Compound classes of $[K]$	$2^{[K]} \setminus \{\emptyset, [K]\}$	S	
A subset of classes	$Y_s \subseteq [K]$		
Joint distribution	$\Pr(Y = y, X = x)$	$P_{Y=y,x}, P_{Y=y,X},$ or $P_{Y,X}$	$\Pr(x, y)$
Class prior	$\Pr(Y = y)$	π_y	
Marginal	$\Pr(X)$	P_X	
Class-conditional	$\Pr(X = x Y = y)$	$P_{x y}, P_{X y}, P_{x Y=y},$ or $P_{X Y=y}$	
Class probability	$\Pr(Y = y X = x)$	$P_{Y=y x}, P_{Y=y X},$ or $P_{Y X}$	$\eta(x)$
Confidence	$\Pr(Y = y X = x)$	$r_y(X), r_y(x),$ or $r(X)$ if $y = p$	$r^y(x)$ or $r(x)$
Sample size probability	$\Pr(S = d)$	$P_{ S =d}$ or $q_{ S }$	
Hypothesis and its space	$g \in G$		
Loss of g	$\ell_{Y=y}(g(x))$	$\ell_y, \ell_y(X),$ or $\ell_Y(g(X))$	$\ell(g(X), Y)$
Classification risk	$\mathbb{E}_{Y,X} [\ell_Y(g(X))]$	$R(g)$	$\mathbb{E}_{X,Y} [\ell(g(X), Y)]$
The j -th entry of vector V	$(V)_j$	V_j	
Indicator function of E	$\mathbb{1}[E]$		
Complement of set s	$Y \setminus s$	\bar{s}	
Identity matrix	I		
MCD parameters	γ_p and γ_n		
UU parameters	γ_1 and γ_2		$1 - \theta$ and θ
CCN parameters	$P_{\bar{Y} Y,X}$	$P_{S Y,X}$ or $P_{\bar{S} Y,X}$	ρ_+ and ρ_-

C Omitted Proofs in Section 4

Omitted proofs in Section 4 are provided in this appendix. We first restate a claim in the main body of the paper, and then provide the corresponding proof.

C.1 Omitted Proofs in Section 4.1

Proof of Lemma 5

Lemma 5. *Let B (46) be the base distributions and*

$$\bar{P} := \begin{pmatrix} P_{\bar{S}} \\ P_U \end{pmatrix}.$$

Then, the contamination matrix

$$M_{\text{SU}} := \begin{pmatrix} \frac{\pi_p^2}{\pi_p^2 + \pi_n^2} & \frac{\pi_n^2}{\pi_p^2 + \pi_n^2} \\ \pi_p & \pi_n \end{pmatrix}, \quad (51)$$

which satisfies $\bar{P} = M_{\text{SU}}B$, characterizes the data-generating distributions \bar{P} .

Proof. The proof steps follow that of Lemma 4. By definitions,

$$M_{\text{SU}}B = \begin{pmatrix} \frac{\pi_p^2}{\pi_p^2 + \pi_n^2} & \frac{\pi_n^2}{\pi_p^2 + \pi_n^2} \\ \pi_p & \pi_n \end{pmatrix} \begin{pmatrix} P_{X|Y=p} \\ P_{X|Y=n} \end{pmatrix}.$$

Since

$$\frac{\pi_p^2}{\pi_p^2 + \pi_n^2} P_{X|Y=p} + \frac{\pi_n^2}{\pi_p^2 + \pi_n^2} P_{X|Y=n} = \frac{\pi_p^2 P_{X|Y=p} + \pi_n^2 P_{X|Y=n}}{\pi_p^2 + \pi_n^2} = P_{\bar{S}}$$

and

$$\pi_p P_{X|Y=p} + \pi_n P_{X|Y=n} = P_X = P_U,$$

the first entry of the resulting vector equals $P_{\bar{S}}$ and the second entry equals P_U , we achieve $\bar{P} = M_{\text{SU}}B$. \square

Proof of Lemma 6

Lemma 6. *Let B (46) be the base distributions and*

$$\bar{P} := \begin{pmatrix} P_{\text{Sup}} \\ P_{\text{Inf}} \end{pmatrix}.$$

Then, the contamination matrix

$$M_{\text{Pcomp}} := \begin{pmatrix} \frac{\pi_p}{\pi_p + \pi_n} & \frac{\pi_n}{\pi_p + \pi_n} \\ \frac{\pi_p^2}{\pi_p^2 + \pi_n} & \frac{\pi_n}{\pi_p^2 + \pi_n} \end{pmatrix}, \quad (52)$$

which satisfies $\bar{P} = M_{\text{Pcomp}}B$, characterizes the data-generating distributions \bar{P} .

Proof. By definitions,

$$M_{\text{Pcomp}}B = \begin{pmatrix} \frac{\pi_p}{\pi_p + \pi_n^2} & \frac{\pi_n^2}{\pi_p + \pi_n^2} \\ \frac{\pi_p^2}{\pi_p^2 + \pi_n} & \frac{\pi_n}{\pi_p^2 + \pi_n} \end{pmatrix} \begin{pmatrix} P_{X/Y=p} \\ P_{X/Y=n} \end{pmatrix}.$$

Since

$$\frac{\pi_p}{\pi_p + \pi_n^2} P_{X/Y=p} + \frac{\pi_n^2}{\pi_p + \pi_n^2} P_{X/Y=n} = \frac{\pi_p P_{X/Y=p} + \pi_n^2 P_{X/Y=n}}{\pi_p + \pi_n^2} = P_{\text{Sup}}$$

and

$$\frac{\pi_p^2}{\pi_p^2 + \pi_n} P_{X/Y=p} + \frac{\pi_n}{\pi_p^2 + \pi_n} P_{X/Y=n} = \frac{\pi_p^2 P_{X/Y=p} + \pi_n P_{X/Y=n}}{\pi_p^2 + \pi_n} = P_{\text{Inf}},$$

the first entry of the resulting vector equals P_{Sup} and the second entry of the resulting vector equals P_{Inf} , which establishes $M_{\text{Pcomp}}B = \bar{P}$. \square

Proof of Lemma 7

Lemma 7. Let B (46) be the base distributions and

$$\bar{P} = \begin{pmatrix} P_{\bar{D}} \\ P_{\bar{U}} \end{pmatrix}.$$

Then, the contamination matrix

$$M_{\text{DU}} = \begin{pmatrix} 1/2 & 1/2 \\ \pi_p & \pi_n \end{pmatrix}, \quad (53)$$

which satisfies $\bar{P} = M_{\text{DU}}B$, characterizes the data-generating distributions \bar{P} .

Proof. Similar to the proofs of Lemmas 5 and 6, we begin with

$$M_{\text{DU}}B = \begin{pmatrix} 1/2 & 1/2 \\ \pi_p & \pi_n \end{pmatrix} \begin{pmatrix} P_{X/Y=p} \\ P_{X/Y=n} \end{pmatrix}.$$

Since $(P_{X/Y=p} + P_{X/Y=n})/2 = P_{\bar{D}}$ and $\pi_p P_{X/Y=p} + \pi_n P_{X/Y=n} = P_X = P_{\bar{U}}$, we have $M_{\text{DU}}B = \bar{P}$. \square

Proof of Lemma 8

Lemma 8. Let B (46) be the base distributions and

$$\bar{P} = \begin{pmatrix} P_{\bar{S}} \\ P_{\bar{D}} \end{pmatrix}.$$

Then, the contamination matrix

$$M_{\text{SD}} = \begin{pmatrix} \frac{\pi_p^2}{\pi_p^2 + \pi_n^2} & \frac{\pi_n^2}{\pi_p^2 + \pi_n^2} \\ 1/2 & 1/2 \end{pmatrix} \quad (54)$$

which satisfies $\bar{P} = M_{\text{SD}}B$, characterizes the data-generating distributions \bar{P} .

Proof. First, we begin with

$$M_{\text{SD}}B = \begin{pmatrix} \frac{\pi_p^2}{\pi_p^2 + \pi_n^2} & \frac{\pi_n^2}{\pi_p^2 + \pi_n^2} \\ 1/2 & 1/2 \end{pmatrix} \begin{pmatrix} P_{X/Y=p} \\ P_{X/Y=n} \end{pmatrix}.$$

Then, we have the lemma by reusing the calculations in the proofs of Lemmas 5 and 7. \square

C.2 Omitted Proofs in Section 4.2

Proof of Lemma 11

Lemma 11. Let the elements in S be $\{s_1, s_2, \dots, s_{|S|}\}$. For each $j \in [|S|]$, let the j -th entry of \bar{P} be

$$\bar{P}_j = P_{S=s_j, X} := C(S = s_j, X) \sum_{k \in \mathcal{S}_j} P_{Y=k, X},$$

which denotes the data-generating distribution of (s_j, X) . Assume the base distributions B and the contamination matrix M_{PPL} are given by (63) and (66), respectively. Then, M_{PPL} satisfies $\bar{P} = M_{\text{PPL}}B$ and characterizes PPL learning (26).

Proof. For each $j \in [|S|]$,

$$\left(M_{\text{PPL}}B\right)_j = \sum_{k=1}^K C(s_j, X) \mathbb{1}[Y = k \in \mathcal{S}_j] P_{Y=k, X} = C(s_j, X) \sum_{k \in \mathcal{S}_j} P_{Y=k, X} = \bar{P}_j.$$

Note that $C(s_j, X) \sum_{k \in \mathcal{S}_j} P_{Y=k, X}$ corresponds to (26) when the observed partial-label is s_j . \square

Proof of Lemma 12

Lemma 12. Let the elements in S be $\{s_1, s_2, \dots, s_{|S|}\}$. For each $j \in [|S|]$, let the j -th entry of \bar{P} be

$$\bar{P}_j = P_{S=s_j, X} := \frac{1}{2^{K-1} - 1} \sum_{k \in \mathcal{S}_j} P_{Y=k, X},$$

which denotes the data-generating distribution of (s_j, X) . Assume the base distributions B and the contamination matrix M_{PCPL} are given by (63) and (68), respectively. Then, M_{PCPL} satisfies $\bar{P} = M_{\text{PCPL}}B$ and characterizes PCPL learning (24).

Proof. For each $j \in [|S|]$,

$$\left(M_{\text{PCPL}}B\right)_j = \sum_{k=1}^K \frac{1}{2^{K-1} - 1} \mathbb{1}[Y = k \in \mathcal{S}_j] P_{Y=k, X} = \frac{1}{2^{K-1} - 1} \sum_{k \in \mathcal{S}_j} P_{Y=k, X} = \bar{P}_j.$$

Note that $\frac{1}{2^{K-1} - 1} \sum_{k \in \mathcal{S}_j} P_{Y=k, X}$ corresponds to (24) when the observed partial-label is s_j . \square

Proof of lemma 13

Lemma 13. Suppose the base distributions B , the contamination matrix M_{MCL} , and the data-generating distributions \bar{P} are given by (63), (71), and (70), respectively. Then, M_{MCL} satisfies $\bar{P} = M_{\text{MCL}}B$ and characterizes MCL (22).

Proof. For each $j \in [N]$, we have

$$\begin{aligned}
(M_{\text{MCLB}})_j &= \sum_Y \frac{P_{|\bar{s}_j|=|\bar{s}_j|} | [Y / \bar{s}_j] P_{Y,X}}{\binom{K-1}{|\bar{s}_j|}} \\
&= \sum_Y \frac{\sum_{d=1}^{K-1} | [\bar{s}_j|=d] P_{|\bar{s}_j|=d} | [Y / \bar{s}_j] P_{Y,X}}{\binom{K-1}{|\bar{s}_j|}} \\
&= \sum_{d=1}^{K-1} \frac{P_{|\bar{s}_j|=d}}{\binom{K-1}{|\bar{s}_j|}} \sum_Y | [Y / \bar{s}_j] P_{Y,X} | [\bar{s}_j|=d] \\
&= \sum_{d=1}^{K-1} \frac{P_{|\bar{s}_j|=d}}{\binom{K-1}{|\bar{s}_j|}} \sum_{Y/\bar{s}_j} P_{Y,X} | [\bar{s}_j|=d] \\
&= P_{\bar{S}=\bar{s}_j, X}.
\end{aligned}$$

□

D Omitted Proofs in Section 5

Omitted proofs in Section 5 are provided in this appendix. We first restate a claim in the main body of the paper, and then provide the corresponding proof.

D.1 Omitted Proofs in Section 5.1

Proof of Corollary 22

Corollary 22. *For PU learning, the classification risk can be rewritten as*

$$R(g) = \mathbb{E}_{\mathbb{P}} [\bar{\ell}_{\mathbb{P}}] + \mathbb{E}_{\mathbb{U}} [\bar{\ell}_{\mathbb{U}}], \quad (90)$$

where

$$\begin{aligned}
\bar{\ell}_{\mathbb{P}} &= \pi_{\mathbb{P}} \ell_{\mathbb{P}} - \pi_{\mathbb{P}} \ell_{\mathbb{N}}, \\
\bar{\ell}_{\mathbb{U}} &= \ell_{\mathbb{N}}.
\end{aligned} \quad (91)$$

Proof. According to Table 7, $M_{\mathbb{P}\mathbb{U}}$ is a child of $M_{\mathbb{U}\mathbb{U}}$ on the reduction graph. Thus, replacing the subscripts $\{\mathbb{U}_1, \mathbb{U}_2\}$ of \bar{P} and \bar{L} with $\{\mathbb{P}, \mathbb{U}\}$ and assigning $\gamma_1 = 0$ and $\gamma_2 = \pi_{\mathbb{P}}$ as what we choose in Section 4.1.2, we follow the proof of Theorem 21 to conduct the risk rewrite: We first obtain the corrected losses (91) by plugging the assigned values into (86). Then, repeating the steps in (89), we achieve (90). □

Proof of Corollary 23

Corollary 23. *Assume $\pi_{\mathbb{P}} = 1/2$. For SU learning, the classification risk can be rewritten as*

$$R(g) = \mathbb{E}_{\bar{\mathbb{S}}} [\bar{\ell}_{\bar{\mathbb{S}}}] + \mathbb{E}_{\mathbb{U}} [\bar{\ell}_{\mathbb{U}}],$$

where

$$\begin{aligned}
\bar{\ell}_{\bar{\mathbb{S}}} &= \frac{\pi_{\mathbb{P}}^2 + \pi_{\mathbb{N}}^2}{2\pi_{\mathbb{P}} - 1} \ell_{\mathbb{P}} - \frac{\pi_{\mathbb{P}}^2 + \pi_{\mathbb{N}}^2}{2\pi_{\mathbb{P}} - 1} \ell_{\mathbb{N}}, \\
\bar{\ell}_{\mathbb{U}} &= -\frac{\pi_{\mathbb{N}}}{2\pi_{\mathbb{P}} - 1} \ell_{\mathbb{P}} + \frac{\pi_{\mathbb{P}}}{2\pi_{\mathbb{P}} - 1} \ell_{\mathbb{N}}.
\end{aligned} \quad (92)$$

Proof. By Table 7, $M_{\bar{\mathbb{S}}\mathbb{U}}$ is a child of $M_{\mathbb{U}\mathbb{U}}$. Substituting the subscripts $\{\mathbb{U}_1, \mathbb{U}_2\}$ with subscripts $\{\bar{\mathbb{S}}, \mathbb{U}\}$ and choosing $\gamma_1 = \frac{\pi_{\mathbb{N}}^2}{\pi_{\mathbb{P}}^2 + \pi_{\mathbb{N}}^2}$ and $\gamma_2 = \pi_{\mathbb{P}}$ as we did in Section 4.1.3, we obtain the corrected losses (92) by plugging the assigned values into (86). We note that $\pi_{\mathbb{P}} = 1/2$ ensures the choices of γ_1 and γ_2 above satisfy the $\gamma_1 + \gamma_2 = 1$ assumption discussed in Section 4.1.1. Then, we achieve the rewrite $R(g) = \mathbb{E}_{\bar{\mathbb{S}}} [\bar{\ell}_{\bar{\mathbb{S}}}] + \mathbb{E}_{\mathbb{U}} [\bar{\ell}_{\mathbb{U}}]$ by repeating the derivation for (89). □

Proof of Lemma 25

Lemma 25. Let $(x, x) \quad P_S$ defined by (10). Then, $E_S \left[\frac{L(X)}{2} \right] = E_S \left[\frac{L(X)}{2} \right]$.

Proof. For clarity, we simplify P_S as $c_1 P_{X|Y=p} P_{X|Y=p} + c_2 P_{X|Y=n} P_{X|Y=n}$, with $c_1 = \frac{\pi_p^2}{\pi_p^2 + \pi_n^2}$ and $c_2 = \frac{\pi_n^2}{\pi_p^2 + \pi_n^2}$. The lemma follows from

$$\begin{aligned} & E_S [L(X)] \\ &= \int_x \times \int_x \times P_S L(x) \, dx \, dx \\ &= \int_x \times \int_x \times (c_1 P_{x|Y=p} P_{x|Y=p} + c_2 P_{x|Y=n} P_{x|Y=n}) L(x) \, dx \, dx \\ &= c_1 \int_x \times P_{x|Y=p} L(x) \, dx \int_x \times P_{x|Y=p} \, dx + c_2 \int_x \times P_{x|Y=n} L(x) \, dx \int_x \times P_{x|Y=n} \, dx \\ &= c_1 \int_x \times P_{x|Y=p} L(x) \, dx + c_2 \int_x \times P_{x|Y=n} L(x) \, dx, \end{aligned}$$

and similarly,

$$E_S [L(X)] = c_1 \int_x \times P_{x|Y=p} L(x) \, dx + c_2 \int_x \times P_{x|Y=n} L(x) \, dx .$$

□

Proof of Corollary 26

Corollary 26. For P_{comp} learning, the classification risk can be rewritten as

$$R(g) = E_{\text{Sup}} [\bar{\ell}_{\text{Sup}}] + E_{\text{Inf}} [\bar{\ell}_{\text{Inf}}] ,$$

where

$$\begin{aligned} \bar{\ell}_{\text{Sup}} &= \ell_p - \pi_p \ell_n, \\ \bar{\ell}_{\text{Inf}} &= -\pi_n \ell_p + \ell_n. \end{aligned} \tag{96}$$

Proof. Since M_{UU} reduces to M_{Pcomp} with $\gamma_1 = \frac{\pi_n^2}{\pi_p + \pi_n^2}$ and $\gamma_2 = \frac{\pi_p^2}{\pi_p^2 + \pi_n}$ according to Table 7, we replace the subscripts $\{U_1, U_2\}$ with $\{\text{Sup}, \text{Inf}\}$ and instantiate (86) with the assigned values to obtain the corrected losses $\bar{\ell}_{\text{Sup}}$ and $\bar{\ell}_{\text{Inf}}$ (96). Then, repeating the same steps in (89), we have the corollary. □

Proof of Corollary 27

Corollary 27. Assume $\pi_p = 1/2$. For DU learning, the classification risk can be rewritten as

$$R(g) = E_{\tilde{D}} [\bar{\ell}_{\tilde{D}}] + E_U [\bar{\ell}_U] ,$$

where

$$\begin{aligned} \bar{\ell}_{\tilde{D}} &= 2\pi_p \pi_n \left(\frac{1}{\pi_n - \pi_p} \ell_p - \frac{1}{\pi_n - \pi_p} \ell_n \right), \\ \bar{\ell}_U &= -\frac{\pi_p}{\pi_n - \pi_p} \ell_p + \frac{\pi_n}{\pi_n - \pi_p} \ell_n. \end{aligned} \tag{98}$$

Proof. By Table 7, M_{DU} is reduced from M_{UU} . Thus, replacing $\{U_1, U_2\}$ with $\{\tilde{D}, U\}$, and assigning $\gamma_1 = 1/2$ and $\gamma_2 = \pi_p$, we obtain the corrected losses (98) by plugging the assigned values into (86). Note that $\pi_p = 1/2$ implies that the γ_1 and γ_2 assignments are feasible. Then, repeating the steps in (89), we have the corollary. □

Proof of Lemma 28

Lemma 28. Given B (46) and following the DU learning notations, we have

$$M_{\text{DU}}B = \begin{pmatrix} P_{\tilde{\text{D}}} \\ P_{\text{U}} \end{pmatrix} = \bar{P},$$

where

$$M_{\text{DU}} := \begin{pmatrix} \frac{\int_{x \times X} P_{x|Y=n} dx}{2} & \frac{\int_{x \times X} P_{x|Y=p} dx}{2} \\ \pi_p & \pi_n \end{pmatrix}.$$

Proof. Since $\int_{x \times X} P_{x|Y=n} dx = 1$ and $\int_{x \times X} P_{x|Y=p} dx = 1$, we have $M_{\text{DU}} = M_{\text{DU}}$ and hence $M_{\text{DU}}B = M_{\text{DU}}B = \begin{pmatrix} P_{\tilde{\text{D}}} \\ P_{\text{U}} \end{pmatrix}$. The last equality follows from Lemma 7. \square

Proof of Lemma 29

Lemma 29. Let $(x, x) \quad P_{\text{D}}$ defined in (12). Then, $E_{\text{D}} \left[\frac{L(X)}{2} \right] = E_{\text{D}} \left[\frac{L(X)}{2} \right]$.

Proof. Recall $P_{\text{D}} = \frac{1}{2}(P_{x|Y=p}P_{x|Y=n} + P_{x|Y=n}P_{x|Y=p})$. Following the similar argument in Lemma 25,

$$\begin{aligned} E_{\text{D}} \left[\frac{L(X)}{2} \right] &= \int_{x \times X} \int_{x \times X} (P_{x|Y=p}P_{x|Y=n} + P_{x|Y=n}P_{x|Y=p}) \frac{L(x)}{4} dx dx \\ &= \int_{x \times X} (P_{x|Y=p} + P_{x|Y=n}) \frac{L(x)}{4} dx \end{aligned}$$

and

$$E_{\text{D}} \left[\frac{L(X)}{2} \right] = \int_{x \times X} (P_{x|Y=n} + P_{x|Y=p}) \frac{L(x)}{4} dx$$

prove the lemma. \square

Proof of Lemma 30

Corollary 30. Assume $\pi_p = 1/2$. For SD learning, the classification risk can be rewritten as

$$R(g) = E_{\tilde{\text{S}}} [\bar{\ell}_{\tilde{\text{S}}}] + E_{\tilde{\text{D}}} [\bar{\ell}_{\tilde{\text{D}}}],$$

where

$$\begin{aligned} \bar{\ell}_{\tilde{\text{S}}} &= (\pi_p^2 + \pi_n^2) \left(\frac{\pi_p}{\pi_p - \pi_n} \ell_p - \frac{\pi_n}{\pi_p - \pi_n} \ell_n \right), \\ \bar{\ell}_{\tilde{\text{D}}} &= 2\pi_p\pi_n \left(-\frac{\pi_n}{\pi_p - \pi_n} \ell_p + \frac{\pi_p}{\pi_p - \pi_n} \ell_n \right). \end{aligned} \quad (101)$$

Proof. By Table 7, M_{SD} is reduced from M_{UU} . Thus, replacing $\{\text{U}_1, \text{U}_2\}$ with $\{\tilde{\text{S}}, \tilde{\text{D}}\}$, and assigning $\gamma_1 = \frac{\pi_n^2}{\pi_p^2 + \pi_n^2}$ and $\gamma_2 = 1/2$, we obtain the corrected losses (101) by plugging the assigned values into (86). Note that $\pi_p = 1/2$ implies that the γ_1 and γ_2 assignments are feasible. Then, repeating the steps in (89), we have the corollary. \square

Proof of Lemma 31

Lemma 31. Given B (46) and following the SD learning notations, we have

$$M_{\text{SD}}B = \begin{pmatrix} P_{\bar{S}} \\ P_{\bar{D}} \end{pmatrix} = \bar{P},$$

where

$$M_{\text{SD}} := \begin{pmatrix} \frac{\pi_p^2 \int_x \chi_{P_x/Y=p} dx}{\pi_p^2 + \pi_n^2} & \frac{\pi_n^2 \int_x \chi_{P_x/Y=n} dx}{\pi_p^2 + \pi_n^2} \\ \frac{\int_x \chi_{P_x/Y=n} dx}{2} & \frac{\int_x \chi_{P_x/Y=p} dx}{2} \end{pmatrix}.$$

Proof. Since $\int_x \chi_{P_x/Y=p} dx = 1$ and $\int_x \chi_{P_x/Y=n} dx = 1$, we have $M_{\text{SD}} = M_{\text{SD}}$ and hence $M_{\text{SD}}B = M_{\text{SD}}B = \begin{pmatrix} P_{\bar{S}} \\ P_{\bar{D}} \end{pmatrix}$. The last equality follows from Lemma 8. \square

Proof of Lemma 32

Lemma 32. Assume the formulation $\bar{P} = M_{\text{Sconf}}B$ (58) is given. Suppose a vector of corrected losses \bar{L} of the form $(\tilde{\ell}_1(x) \ \tilde{\ell}_2(x))$ is independent of x . Then, we have

$$\int_x \chi_{\bar{L} \ \bar{P}} dx = \bar{L} \ \tilde{M}_{\text{Sconf}}P, \quad (103)$$

where

$$\tilde{M}_{\text{Sconf}} = \begin{pmatrix} \int_x \frac{\pi_p^2 P_{x/p} - \pi_n^2 P_{x/n}}{r - \pi_n} dx & \int_x \frac{\pi_n^2 P_{x/n} - \pi_p^2 P_{x/p}}{r - \pi_n} dx \\ \int_x \frac{\pi_p^2 P_{x/n} - \pi_n^2 P_{x/p}}{\pi_p - r} dx & \int_x \frac{\pi_n^2 P_{x/p} - \pi_p^2 P_{x/n}}{\pi_p - r} dx \end{pmatrix}.$$

Proof. of Lemma 32. We replace \bar{P} using $\bar{P} = M_{\text{Sconf}}B$ (58). Since \bar{L} is independent of x , we can move the integral over x into M_{Sconf} to obtain

$$\begin{aligned} \int_x \chi_{\bar{L} \ \bar{P}} dx &= \int_x \bar{L} \begin{pmatrix} \frac{\pi_p(\pi_p^2 P_{x/p} - \pi_n^2 P_{x/n})}{r - \pi_n} & \frac{\pi_n(\pi_n^2 P_{x/n} - \pi_p^2 P_{x/p})}{r - \pi_n} \\ \frac{\pi_n(\pi_p^2 P_{x/n} - \pi_p^2 P_{x/p})}{\pi_p - r} & \frac{\pi_p(\pi_n^2 P_{x/p} - \pi_n^2 P_{x/n})}{\pi_p - r} \end{pmatrix} \begin{pmatrix} P_{X/p} \\ P_{X/n} \end{pmatrix} dx \\ &= \bar{L} \begin{pmatrix} \int_x \frac{\pi_p^2 P_{x/p} - \pi_n^2 P_{x/n}}{r - \pi_n} dx & \int_x \frac{\pi_n^2 P_{x/n} - \pi_p^2 P_{x/p}}{r - \pi_n} dx \\ \int_x \frac{\pi_p^2 P_{x/n} - \pi_n^2 P_{x/p}}{\pi_p - r} dx & \int_x \frac{\pi_n^2 P_{x/p} - \pi_p^2 P_{x/n}}{\pi_p - r} dx \end{pmatrix} \begin{pmatrix} \pi_p P_{X/p} \\ \pi_n P_{X/n} \end{pmatrix}. \end{aligned}$$

Since $\pi_p P_{X/p} = P_{X,Y=p}$ and $\pi_n P_{X/n} = P_{X,Y=n}$, $\begin{pmatrix} \pi_p P_{X/p} \\ \pi_n P_{X/n} \end{pmatrix} = P$. Furthermore, comparing the equality in the above derivation with (103), we have

$$\tilde{M}_{\text{Sconf}} = \begin{pmatrix} \int_x \frac{\pi_p^2 P_{x/p} - \pi_n^2 P_{x/n}}{r - \pi_n} dx & \int_x \frac{\pi_n^2 P_{x/n} - \pi_p^2 P_{x/p}}{r - \pi_n} dx \\ \int_x \frac{\pi_p^2 P_{x/n} - \pi_n^2 P_{x/p}}{\pi_p - r} dx & \int_x \frac{\pi_n^2 P_{x/p} - \pi_p^2 P_{x/n}}{\pi_p - r} dx \end{pmatrix}$$

that completes the proof. \square

Proof of Lemma 33

Lemma 33. *Let*

$$\tilde{M}_{\text{Sconf}}^\dagger := \begin{pmatrix} \frac{r - \pi_n}{\pi_p - \pi_n} & 0 \\ 0 & \frac{\pi_p - r}{\pi_p - \pi_n} \end{pmatrix}.$$

Then,

$$\tilde{M}_{\text{Sconf}}^\dagger \tilde{M}_{\text{Sconf}} = I.$$

Proof. We prove the lemma by examining each entry of $\tilde{M}_{\text{Sconf}}^\dagger \tilde{M}_{\text{Sconf}}$. The value of (1, 1) entry is

$$\begin{aligned} \frac{r - \pi_n}{\pi_p - \pi_n} \int_x \frac{\pi_p^2 P_{x/p} - \pi_n^2 P_{x/n}}{r - \pi_n} dx &= \frac{1}{\pi_p - \pi_n} \left(\pi_p^2 \int_x P_{x/p} dx - \pi_n^2 \int_x P_{x/n} dx \right) \\ &= \frac{\pi_p^2 - \pi_n^2}{\pi_p - \pi_n} = 1. \end{aligned}$$

The (2, 2) entry has value

$$\begin{aligned} \frac{\pi_p - r}{\pi_p - \pi_n} \int_x \frac{\pi_p^2 P_{x/p} - \pi_n^2 P_{x/n}}{\pi_p - r} dx &= \frac{1}{\pi_p - \pi_n} \left(\pi_p^2 \int_x P_{x/p} dx - \pi_n^2 \int_x P_{x/n} dx \right) \\ &= \frac{\pi_p^2 - \pi_n^2}{\pi_p - \pi_n} = 1. \end{aligned}$$

The (1, 2) entry is zero since $\int_x (\pi_n^2 P_{x/n} - \pi_p^2 P_{x/p}) dx = 0$. Similarly, since $\int_x (\pi_p^2 P_{x/n} - \pi_n^2 P_{x/p}) dx = 0$, the (2, 1) entry is also zero. \square

D.2 Omitted Proofs in Section 5.2**Proof of Corollary 38**

Corollary 38. *The decontamination matrix M_{PCPL}^\dagger for PCPL equals M_{PPL}^\dagger . If we define the corrected losses as $\bar{L} := L M_{\text{PCPL}}^\dagger$, the classification risk for PCPL learning can be rewritten as*

$$R(g) = \mathbb{E}_{S, X} [\bar{\ell}_S],$$

where

$$\bar{\ell}_S = \sum_i \frac{P_{Y=i|X}}{\sum_a P_{Y=a|X}} \ell_{Y=i}. \quad (110)$$

Proof. The proof follows the standard argument: First, find out M_{PCPL}^\dagger , then construct the corrected losses to rewrite the risk.

Since M_{PCPL} is reduced from M_{PPL} , we can exploit Lemma 36. Note that the only difference $C(S, X)$ between the formulations of PCPL and PPL cancels itself out in the derivation of $P_{Y=i|S=s_j, X}$ (please refer to the proof of Lemma 36 for a detailed derivation), the (i, j) entry of M_{PCPL}^\dagger coincides with that of M_{PPL}^\dagger for all i and j , proving the first assertion.

Since M_{PPL}^\dagger and M_{PCPL}^\dagger are identical, $L M_{\text{PPL}}^\dagger = L M_{\text{PCPL}}^\dagger$ gives (110):

$$\begin{aligned} \bar{\ell}_{S=s_j} &= \left(L M_{\text{PCPL}}^\dagger \right)_j = \sum_{i=1}^K \frac{P_{Y=i|X} \mathbb{1}[Y=i \mid S=s_j]}{\sum_a P_{Y=a|X}} \ell_{Y=i} \\ &= \sum_{i \mid s_j} \frac{P_{Y=i|X}}{\sum_a P_{Y=a|X}} \ell_{Y=i}. \end{aligned}$$

Being identical to M_{PPL}^\dagger also means that M_{PCPL}^\dagger is derived from M_{gCCN}^\dagger . Thus, we can continue (107) to rewrite the risk by repeating the proof of Corollary 37:

$$R(g) = \int_{\mathcal{X}} \int_{\mathcal{X}} \bar{L} \bar{P} dx = \int_{\mathcal{X}} \sum_{j=1}^{|\mathcal{S}|} P_{S=s_j, X} \bar{\ell}_{S=s_j} dx = \mathbb{E}_{S, X} [\bar{\ell}_S].$$

□

Proof of Lemma 39

Lemma 39. *Let (s, s) be a pair of partial-labels satisfying $s = Y \setminus s$. Then,*

$$P_{S=s, X} \bar{\ell}_{S=s} + P_{S=s, X} \bar{\ell}_{S=s} = P_{S=s, X} \sum_{i=1}^K \frac{P_{Y=i|X} \ell_{Y=i}}{\sum_{a \in s} P_{Y=a|X}}.$$

Proof. Given M_{PCPL} (68), we apply $\bar{P} = M_{\text{PCPL}} B$ to obtain

$$P_{S=s, X} = \frac{\sum_{k=1}^K \mathbb{1}[Y = k \setminus s] P_{Y=k, X}}{2^{K-1} - 1}.$$

We also have

$$\bar{\ell}_{S=s} = \frac{\sum_{i \in s} P_{Y=i|X} \ell_{Y=i}}{\sum_{a \in s} P_{Y=a|X}}$$

according to (110). Since

$$\frac{\sum_{k=1}^K \mathbb{1}[Y = k \setminus s] P_{Y=k, X}}{\sum_{a \in s} P_{Y=a|X}} = P_X = \frac{\sum_{k=1}^K \mathbb{1}[Y = k \setminus s] P_{Y=k, X}}{\sum_{a \in s} P_{Y=a|X}},$$

$$\begin{aligned} P_{S=s, X} \bar{\ell}_{S=s} &= \frac{\sum_{k=1}^K \mathbb{1}[Y = k \setminus s] P_{Y=k, X}}{2^{K-1} - 1} \frac{\sum_{i \in s} P_{Y=i|X} \ell_{Y=i}}{\sum_{a \in s} P_{Y=a|X}} \\ &= \frac{\sum_{k=1}^K \mathbb{1}[Y = k \setminus s] P_{Y=k, X}}{2^{K-1} - 1} \frac{\sum_{i \in s} P_{Y=i|X} \ell_{Y=i}}{\sum_{a \in s} P_{Y=a|X}}. \end{aligned}$$

Thus,

$$\begin{aligned} P_{S=s, X} \bar{\ell}_{S=s} + P_{S=s, X} \bar{\ell}_{S=s} &= \frac{\sum_{k=1}^K \mathbb{1}[Y = k \setminus s] P_{Y=k, X}}{2^{K-1} - 1} \frac{\sum_{i \in s} P_{Y=i|X} \ell_{Y=i}}{\sum_{a \in s} P_{Y=a|X}} \\ &\quad + \frac{\sum_{k=1}^K \mathbb{1}[Y = k \setminus s] P_{Y=k, X}}{2^{K-1} - 1} \frac{\sum_{i \in s} P_{Y=i|X} \ell_{Y=i}}{\sum_{a \in s} P_{Y=a|X}} \\ &= P_{S=s, X} \frac{\sum_{i \in s} P_{Y=i|X} \ell_{Y=i} + \sum_{i \in s} P_{Y=i|X} \ell_{Y=i}}{\sum_{a \in s} P_{Y=a|X}} \\ &= P_{S=s, X} \sum_{i=1}^K \frac{P_{Y=i|X} \ell_{Y=i}}{\sum_{a \in s} P_{Y=a|X}} \end{aligned}$$

proves the lemma. □

Proof of Corollary 40

Corollary 40. *The (i, j) entry of the decontamination matrix M_{MCL}^\dagger is of the form*

$$P_{Y=i|\bar{S}=\bar{s}_j, X} = \frac{P_{Y=i|X} \mathbb{1}[Y=i/\bar{s}_j]}{\sum_{a/\bar{s}_j} P_{Y=a|X}}. \quad (111)$$

Define the corrected losses $\bar{L} := L M_{\text{MCL}}^\dagger$. Then, for MCL learning, the classification risk can be rewritten as

$$R(g) = \mathbb{E}_{\bar{S}, X} [\bar{\ell}_{\bar{S}}],$$

where

$$\bar{\ell}_{\bar{S}} = \sum_{i/\bar{S}} \frac{P_{Y=i|X}}{\sum_{a/\bar{S}} P_{Y=a|X}} \ell_{Y=i}. \quad (112)$$

Proof. The proof follows the standard strategy in Section 5.2: We will first find out M_{MCL}^\dagger , and then construct the corrected losses \bar{L} to rewrite the risk.

Based on the notion in (71), we denote the (j, i) entry of M_{MCL} as

$$\frac{P_{|\bar{S}|=|\bar{s}_j|}}{\binom{K-1}{|\bar{s}_j|}} \mathbb{1}[Y=i/\bar{s}_j] = C(\bar{s}_j, X) \mathbb{1}[Y=i/\bar{s}_j] = P_{\bar{S}=\bar{s}_j|Y=i, X}. \quad (128)$$

Expressing M_{MCL} via (128) allows us to apply the argument for (106) to show that the (i, j) entry of M_{MCL}^\dagger is of the form $P_{Y=i|\bar{S}=\bar{s}_j, X}$. Specifically, assigning $(M)_{j,i}$ in (40) as $P_{\bar{S}=\bar{s}_j|Y=i, X}$ and applying marginal chain (i.e., Proposition 2), we obtain

$$\begin{aligned} \left(M_{\text{MCL}}^\dagger M_{\text{MCL}} P \right)_i &= \sum_{j=1}^{|\bar{S}|} \left(M_{\text{MCL}}^\dagger \right)_{i,j} \sum_{k=1}^K (M_{\text{MCL}})_{j,k} P_k \\ &= \sum_{j=1}^{|\bar{S}|} P_{Y=i|\bar{S}=\bar{s}_j, X} \sum_{k=1}^K P_{\bar{S}=\bar{s}_j|Y=k, X} P_{Y=k, X} \\ &= P_{Y=i, X} = P_i. \end{aligned}$$

Then, we follow the same argument in Lemma 36 to calculate $P_{Y=i|\bar{S}=\bar{s}_j, X}$ subject to (128). Note that $P_{\bar{S}|Y, X} = C(\bar{S}, X) \mathbb{1}[Y/\bar{S}]$ in (128) implies

$$\begin{aligned} \sum_b P_{\bar{S}, Y=b|X} &= \sum_b P_{\bar{S}|Y=b, X} P_{Y=b|X} \\ &= \sum_b C(\bar{S}, X) \mathbb{1}[Y=b/\bar{S}] P_{Y=b|X} \\ &= 0. \end{aligned}$$

Thus, $P_{\bar{S}|X} = \sum_b P_{\bar{S}, Y=b|X} + \sum_{a/\bar{S}} P_{\bar{S}, Y=a|X} = \sum_{a/\bar{S}} P_{\bar{S}, Y=a|X}$. The fact further implies

$$\begin{aligned} P_{Y|\bar{S}, X} &= \frac{P_{\bar{S}, Y|X}}{P_{\bar{S}|X}} = \frac{P_{\bar{S}|Y, X} P_{Y|X}}{\sum_{a/\bar{S}} P_{\bar{S}|Y=a, X} P_{Y=a|X}} \\ &= \frac{C(\bar{S}, X) \mathbb{1}[Y/\bar{S}] P_{Y|X}}{\sum_{a/\bar{S}} C(\bar{S}, X) \mathbb{1}[Y=a/\bar{S}] P_{Y=a|X}} \\ &= \frac{P_{Y|X} \mathbb{1}[Y/\bar{S}]}{\sum_{a/\bar{S}} P_{Y=a|X}}. \end{aligned}$$

Therefore, for $Y = i$ and $\bar{S} = \bar{s}_j$, we achieve

$$P_{Y=i/\bar{S}=\bar{s}_j, X} = \frac{P_{Y=i/X} \mathbb{1}[Y = i / \bar{s}_j]}{\sum_{a/\bar{s}_j} P_{Y=a/X}}$$

that proves (111).

With M_{MCL}^\dagger in hand, we repeat the same argument in Corollary 37 to obtain

$$\begin{aligned} \bar{\ell}_{\bar{S}=\bar{s}_j} &= \left(L \ M_{\text{MCL}}^\dagger \right)_j = \sum_{i=1}^K \frac{P_{Y=i/X} \mathbb{1}[Y = i / \bar{s}_j]}{\sum_{a/\bar{s}_j} P_{Y=a/X}} \ell_{Y=i} \\ &= \sum_{i/\bar{s}_j} \frac{P_{Y=i/X}}{\sum_{a/\bar{s}_j} P_{Y=a/X}} \ell_{Y=i} \end{aligned}$$

and

$$\begin{aligned} R(g) &= \int_{x \times X} L \ P \ dx = \int_{x \times X} L \ M_{\text{MCL}}^\dagger M_{\text{MCL}} P \ dx \\ &= \int_{x \times X} \bar{L} \ \bar{P} \ dx = \int_{x \times X} \sum_{j=1}^{|\mathcal{S}|} P_{\bar{S}=\bar{s}_j, X} \bar{\ell}_{\bar{S}=\bar{s}_j} \ dx = \mathbb{E}_{\bar{S}, X} [\bar{\ell}_{\bar{S}}] \end{aligned}$$

to complete the risk rewrite of MCL. \square

D.3 Omitted Proofs in Section 5.3

Proof of Corollary 47

Corollary 47. *For SC-Conf learning, the classification risk can be written as*

$$R(g) = \pi_{y_s} \mathbb{E}_{X/Y=y_s} \left[\sum_{i=1}^K \frac{r_i(X)}{r_{y_s}(X)} \ell_i \right].$$

Proof. The corollary follows from notation substitution and the same argument for Theorem 46. Specifically, we replace M_{Sub} with M_{SC} , Y_s with y_s , and $\frac{P_{Y=i/X}}{P_{Y=y_s/X}} \ell_i$ with $\frac{P_{Y=i/X}}{P_{Y=y_s/X}} \ell_i$. \square

Proof of Corollary 48

Corollary 48. *For Pconf learning, the classification risk can be written as*

$$R(g) = \pi_p \mathbb{E}_P \left[\ell_p + \frac{1-r(X)}{r(X)} \ell_n \right].$$

Proof. Since

$$M_{\text{Pconf}}^\dagger M_{\text{Pconf}} = \begin{pmatrix} \frac{P_{Y=p/X}}{P_{Y=p/X}} & 0 \\ 0 & \frac{P_{Y=p/X}}{P_{Y=n/X}} \end{pmatrix} \begin{pmatrix} \frac{P_{Y=p/X}}{P_{Y=p/X}} & 0 \\ 0 & \frac{P_{Y=n/X}}{P_{Y=p/X}} \end{pmatrix} = I,$$

we define $\bar{L} := L \ M_{\text{Pconf}}^\dagger$ and apply (39) to rewrite the risk as follows

$$\begin{aligned} R(g) &= \int_{x \times X} L \ P \ dx = \int_{x \times X} \bar{L} \ \bar{P} \ dx = \int_{x \times X} \left(\ell_p \ \frac{1-r(X)}{r(X)} \ell_n \right) \begin{pmatrix} P_{Y=p, X} \\ P_{Y=p, X} \end{pmatrix} \ dx \\ &= P_{Y=p} \mathbb{E}_{X/Y=p} \left[\ell_p + \frac{1-r(X)}{r(X)} \ell_n \right] \\ &= \pi_p \mathbb{E}_P \left[\ell_p + \frac{1-r(X)}{r(X)} \ell_n \right]. \end{aligned}$$

□

Proof of Corollary 49

Corollary 49. *For soft-label learning, the classification risk can be written as*

$$R(g) = \mathbb{E}_X \left[\sum_{i=1}^K P_{Y=i|X} \ell_i \right] = \mathbb{E}_X \left[\sum_{i=1}^K r_i(X) \ell_i \right].$$

Proof. Defining $\bar{L} := L - M_{\text{Soft}}^\dagger$ and recalling \bar{P} in (83), we apply (39) to obtain

$$\begin{aligned} R(g) &= \int_{x \times \mathcal{X}} L \cdot P dx = \int_{x \times \mathcal{X}} \bar{L} \cdot \bar{P} dx = \int_{x \times \mathcal{X}} \sum_{i=1}^K P_{Y=i|X} \ell_i \cdot P_X dx \\ &= \mathbb{E}_X \left[\sum_{i=1}^K P_{Y=i|X} \ell_i \right]. \end{aligned}$$

□