
Learning from Label Proportions by Learning with Label Noise

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Learning from label proportions (LLP) is a weakly supervised classification problem
2 where data points are grouped into bags, and the label proportions within
3 each bag are observed instead of the instance-level labels. The task is to learn a
4 classifier to predict the labels of future individual instances. Prior work on LLP
5 for multi-class data has yet to develop a theoretically grounded algorithm. In this
6 work, we propose an approach to LLP based on a reduction to learning with label
7 noise, using the forward correction (FC) loss of Patrini et al. [30]. We establish an
8 excess risk bound and generalization error analysis for our approach, while also
9 extending the theory of the FC loss which may be of independent interest. Our
10 approach demonstrates improved empirical performance in deep learning scenarios
11 across multiple datasets and architectures, compared to the leading methods.

12 1 Introduction

13 In the weakly supervised problem of *learning from label proportions* (LLP), the learner is presented
14 with bags of instances, where each bag is annotated with the proportions of the different classes in
15 the bag. The learner’s objective is to produce a classifier that accurately assigns labels to individual
16 instances in the future. LLP arises in various applications including high energy physics [7], election
17 prediction [45], computer vision [4, 20], medical image analysis [2], remote sensing [8], activity
18 recognition [32], and reproductive medicine [12].

19 To date, most methods for LLP have addressed the setting of binary classification [50, 36, 39, 34,
20 41, 44, 24, 37, 38], although multiclass methods have also recently been investigated [9, 22, 46].
21 The dominant approach to LLP in the literature is “label proportion matching”: train a classifier to
22 accurately reproduce the observed label proportions on the training data, perhaps with additional
23 regularization. In the multiclass setting, the Kullback-Leibler (KL) divergence between the observed
24 and predicted label proportions is adopted by the leading approaches to assess proportion matching.
25 Unfortunately, while matching the observed label proportions is intuitive and can work well in some
26 settings, it has little theoretical basis [50, 38], especially in the multiclass setting, and there are natural
27 settings where it fails [50, 39].

28 Recently, Scott and Zhang [39] demonstrated a principled approach to LLP with performance
29 guarantees based on a reduction to learning with label noise (LLN) in the binary setting. Their basic
30 strategy was to pair bags, and view each pair of bags as an LLN problem, where the observed label
31 proportions are related to the “label flipping” or “noise transition” probabilities. Using an existing
32 technique for LLN based on loss correction, which allows the learner to train directly on the noisy
33 data, they formulated an overall objective based on a (weighted) sum of objectives for each pair of
34 bags. They established generalization error analysis and consistency for the method, and also showed
35 that in the context of kernel methods, their approach outperformed the leading kernel methods.

36 The objective of the present paper is to develop a theoretically grounded and practical approach to
37 multiclass LLP, drawing inspiration from Scott and Zhang [39]. The primary challenge stems from
38 the fact that Scott and Zhang [39] employed the so-called “backward correction” loss, which solves
39 LLN by scaling the *output* of a loss function of interest according to the noise transition probabilities
40 [28, 30, 35]. While this loss correction was demonstrated to work well for kernel methods in a binary
41 setting, Patrini et al. [30] showed that it is poorly suited for multiclass deep learning settings. To
42 remedy this, they proposed the “forward correction” loss, which scales the *inputs* to a loss function of
43 interest according to the noise transition probabilities. They established Fisher consistency of the
44 technique, and demonstrated its superiority empirically (see also [53]).

45 The present work is thus inspired by Scott and Zhang [39] but uses the forward correction (FC)
46 loss in a multiclass setting. This requires a number of technical modifications to the arguments of
47 Scott and Zhang [39]. Most notably, it now becomes necessary to demonstrate that the FC loss is
48 *calibrated* with respect to the 0-1 loss, a critical property needed for showing consistency. Such
49 analysis is inherently not needed when using the backward correction, due to the way in which
50 this correction yields an unbiased estimate of the uncorrupted risk. Furthermore, Scott and Zhang
51 [39] adopt a probabilistic setting known as a “mutual contamination model”, whereas we adopt the
52 more conventional setting based on a noise transition matrix. Finally, the multiclass setting involves
53 new estimation challenges not present in the binary case. These factors mean that our work is not a
54 straightforward extension of Scott and Zhang [39]. Indeed, the authors of a recent report acknowledge
55 that it is “difficult to extend [the method of Scott and Zhang [39]] to multiclass classification” [16].

56 **Additional related work:** Much work on LLP has focused on learning specific types of models,
57 including support vector machines [36, 50, 47, 33, 5, 19, 40], probabilistic models [18, 13, 45, 32,
58 12], random forests [41], neural networks [21, 1, 9, 22, 46], and clustering-based models, [3, 44].
59 Many of these works develop learning criteria that are specific to the model being learned.

60 On the theoretical front, Quadrianto et al. [34] and Patrini et al. [30] initiated the learning theoretic
61 study of LLP, introducing Rademacher style bounds for linear methods, but they do not address
62 consistency w.r.t. a classification performance measure. Yu et al. [51] provides support for label
63 proportion matching but only under the assumption that the bags are very pure. Saket [37] studies
64 learnability of linear threshold functions. Recently Saket et al. [38] introduced a condition under
65 which label proportion matching does provably well w.r.t. a squared error loss in the binary setting,
66 and developed an associated algorithm. This method does not scale easily to large datasets, and further
67 requires knowledge of how bags are grouped according to different bag-generating distributions.

68 A handful of recent papers have studied multiclass LLP in deep learning scenarios. Dulac-Arnold et al.
69 [9] study the KL loss for label proportion matching, and a variant based on optimal transport. Liu et al.
70 [22, 23] examine an approach based on generative adversarial models. Tsai and Lin [46] study the use
71 of a regularizer derived from semi-supervised learning. One challenge common to these approaches is
72 that their implementations employ mini-batches of bags, which becomes computationally prohibitive
73 for deep architectures when the batch size is still very small, e.g., 2 or 3. In contrast, our approach
74 avoids this issue. Finally, a recent technical report presents a risk analysis for multiclass LLP under
75 the assumption of fixed bag size, which we do not require [16]. Their method is not tractable for
76 large bag sizes in which case they approximate their objective “using the bag-level loss proposed in
77 the existing research.”

78 **Contributions and Outline:** Our contributions and the paper structure are summarized as follows.
79 In Section 2, we review the FC loss as a solution to LLN. In Section 3, we extend the theory of the
80 FC loss for LLN. In particular, we show that the FC loss is “uniformly calibrated” with respect to
81 the 0-1 loss using the framework of Steinwart [43], establish an excess risk bound, and determine an
82 explicit lower bound on the calibration function in terms of the noise transition matrix. In Section 4,
83 we extend the results of Section 3 to the setting with multiple noise transition matrices, which form
84 the basis of our approach to LLP. In particular, we establish an excess risk bound and generalization
85 error analysis for learning with multiple noise transition matrices, which in turn enables proofs of
86 consistency. In Section 5, we state the probabilistic model for reducing LLP to LLN with multiple
87 different noise transition matrices and present the LLPFC algorithms. Experiments with deep neural
88 networks are presented in Section 6, where we observe that our approach outperforms competing
89 methods by a substantial margin. Proofs appear in the supplemental material.

90 2 Learning with Label Noise and the Forward Correction Loss

91 This section sets notation and introduces the FC loss as a solution to learning with label noise. Let \mathcal{X}
 92 be the feature space and $\mathcal{Y} = \{1, 2, \dots, C\}$ be the label space, $C \in \mathbb{N}$. We define the C -simplex as
 93 $\Delta^C = \{p \in \mathbb{R}^C : p_i \geq 0, \forall i = 1, 2, \dots, C, \sum_{i=1}^C p_i = 1\}$ and denote its interior by $\overset{\circ}{\Delta}^C$. Let P be a
 94 probability measure on the space $\mathcal{X} \times \mathcal{Y}$.

95 Viewing P as the “clean” probability measure, a noisy probability measure with label-dependent
 96 label noise can be constructed from P in terms of a $C \times C$ column-stochastic matrix T , referred
 97 to as the *noise transition matrix*. Formally, we define a measure \bar{P}_T on $\mathcal{X} \times \mathcal{Y} \times \mathcal{Y}$ by requiring
 98 \forall events $\mathcal{A} \subset \mathcal{X}, \bar{P}_T(\mathcal{A} \times \{i\} \times \{j\}) = P(\mathcal{A} \times \{i\})t_{j,i}$ where $t_{j,i}$ is the element at the j -th row
 99 and i -th column of T . Let (X, Y, \tilde{Y}) have joint distribution \bar{P}_T where X is the feature vector, Y
 100 is the “clean” label, and \tilde{Y} is the “noisy” label. Thus the element of T at row i and column j is
 101 $t_{i,j} = \bar{P}_T(\tilde{Y} = i | Y = j)$. In addition, P is the marginal distribution of (X, Y) . Define P_T to be the
 102 marginal distribution of (X, \tilde{Y}) . Let \mathcal{F} be the collection of all measurable functions from \mathcal{X} to Δ^C .

103 The existence of a regular conditional distribution is guaranteed by the Disintegration Theorem (*e.g.*
 104 Theorem 6.4 in Kallenberg [14]) under suitable properties (*e.g.* when \mathcal{X} is a Radon space). While
 105 the existence of regular conditional probability is beyond the scope of this paper, we assume fixed
 106 regular conditional distributions for Y and \tilde{Y} given X exist, denoted by $P(\cdot | \cdot) : \mathcal{Y} \times \mathcal{X} \rightarrow [0, 1]$
 107 and $P_T(\cdot | \cdot) : \mathcal{Y} \times \mathcal{X} \rightarrow [0, 1]$, respectively. Given $x \in \mathcal{X}$, we define the probability vectors
 108 $\eta(x) = [P(1 | x), \dots, P(C | x)]^{tr}$ and $\eta_T(x) = [P_T(1 | x), \dots, P_T(C | x)]^{tr}$ where we use *tr* to
 109 denote transposition. It directly follows that $\eta_T(x) = T\eta(x)$.

110 We use \mathbb{R}_+ to denote the positive real numbers. The goal of LLN is to learn a classifier that optimizes
 111 a performance measure defined *w.r.t.* P , given access to corrupted training data $(X_i, \tilde{Y}_i) \stackrel{i.i.d.}{\sim} P_T$. In
 112 this work we assume T is known or can be estimated, as is the case when we apply LLN techniques
 113 to LLP (see Section 5). A more formal formulation of LLP is given in Section 5.

114 When attempting to minimize the risk associated to the 0-1 loss and the clean distribution P , it
 115 is common to employ a smooth or convex surrogate loss. For LLN problems, the idea of a *loss*
 116 *correction* is to modify the surrogate loss so that when optimized using the *noisy* data, it still achieves
 117 the desired goal. Below, we introduce the forward correction loss, before which we need to define
 118 inner risk and proper loss. For this purpose we focus on loss functions of the form $L : \Delta^C \times \mathcal{Y} \rightarrow \mathbb{R}$.

119 **Definition 1.** Let $L : \Delta^C \times \mathcal{Y} \rightarrow \mathbb{R}$ be a loss function. The **inner L -risk** at x with probability
 120 measure P is $\mathcal{C}_{L,P,x} : \Delta^C \rightarrow \mathbb{R}$, $\mathcal{C}_{L,P,x}(q) := \mathbb{E}_{Y \sim P(\cdot | x)} L(q, Y)$. The **minimal inner L -risk** at x
 121 with a probability measure P is $\mathcal{C}_{L,P,x}^* := \inf_{q \in \Delta^C} \mathcal{C}_{L,P,x}(q)$.

122 **Definition 2.** $\ell : \Delta^C \times \mathcal{Y} \rightarrow \mathbb{R}$ is a **proper loss** if \forall probability measures P on $\mathcal{X} \times \mathcal{Y}$, $\forall x \in$
 123 \mathcal{X} , $\mathcal{C}_{\ell,P,x}^* = \mathcal{C}_{\ell,P,x}(\eta(x))$, and a proper loss is called **strictly proper** if the minimizer of $\mathcal{C}_{\ell,P,x}$ is
 124 unique for all $x \in \mathcal{X}$.

125 Commonly used proper losses include the *log loss* $\ell^{log}(q, c) = -\log q_c$, the *square loss* $\ell^{sq}(q, c) =$
 126 $\sum_{c'=1}^C (\mathbb{1}_{c=c'} - q_{c'})^2$, and the 0-1 loss $\ell^{01}(q, c) = \mathbb{1}_{c \neq \min\{\arg \max_j q_j\}}$, among which only the log
 127 loss and the square loss are strictly proper [49]. Here $\mathbb{1}$ denotes the indicator function. Note that it is
 128 common to compose proper losses with inverted link functions, leading to familiar losses like the
 129 cross-entropy loss. Such losses are discussed further in Section 4.

130 We are now ready to introduce the forward correction loss.

131 **Definition 3.** Let ℓ be a strictly proper loss and let T be a noise transition matrix. Define the **forward**
 132 **correction loss** of ℓ as $\ell_T : \Delta^C \times \mathcal{Y} \rightarrow \mathbb{R}$, $\ell_T(q, c) := \ell(Tq, c)$.

133 It follows from the definition that, if T is invertible, then the inner ℓ_T -risk under the distribution P_T
 134 has a unique minimizer $\eta(x)$. Next we introduce the definition of L -risk and L -Bayes risk associated
 135 with a loss function L .

136 **Definition 4.** Let $L : \Delta^C \times \mathcal{Y} \rightarrow \mathbb{R}$ and P be a probability measure. Define the **L -risk** of f with
 137 distribution P to be $\mathcal{R}_{L,P} : \mathcal{F} \rightarrow \mathbb{R}$, $\mathcal{R}_{L,P}(f) := \mathbb{E}_P [L(f(X), Y)]$ and the **L -Bayes risk** to be
 138 $\mathcal{R}_{L,P}^* := \inf_{f \in \mathcal{F}} \mathcal{R}_{L,P}(f)$.

139 We call $\mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^*$ the *excess L-risk* of f under distribution P . Given a proper loss ℓ , Theorem
140 2 of Patrini et al. [30] establishes Fisher consistency of the FC loss, meaning the minimizer of ℓ -risk
141 under the clean distribution P is the same as the minimizer of ℓ_T -risk under noisy distribution P_T :
142 $\arg \min_{f \in \mathcal{F}} \mathcal{R}_{L,P}(f) = \arg \min_{f \in \mathcal{F}} \mathcal{R}_{\ell_T, P_T}(f)$. Next, we present a stronger result relating the
143 excess ℓ_T -risk under the noisy distribution P_T to the excess 0-1 risk under the clean distribution P .

144 3 Calibration Analysis for the Forward Correction Loss

145 Our objective in this section is to show that when L is the 0-1 loss and ℓ is a continuous strictly
146 proper surrogate loss, there exists a strictly increasing, invertible function θ with $\theta(0) = 0$ such that
147 $\forall f \in \mathcal{F}$ and \forall distributions P , $\theta(\mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^*) \leq \mathcal{R}_{\ell_T, P_T}(f) - \mathcal{R}_{\ell_T, P_T}^*$. Given such a bound,
148 it follows that consistency *w.r.t* the surrogate risk implies consistency *w.r.t* the target risk. The
149 results in this section are standalone results for the FC loss that may be of independent interest, and
150 will be extended in the next section in relation to LLP.

151 The following theorem guarantees the existence of such function θ , given that T is invertible.

152 **Theorem 5.** *Let ℓ be a continuous strictly proper loss and T be an invertible column-stochastic*
153 *matrix. Let L be the 0-1 loss. Assume $\mathcal{R}_{\ell_T, P_T}^* < \infty$. Then $\exists \theta : [0, 1] \rightarrow [0, \infty]$ that is*
154 *strictly increasing and continuous, satisfying $\theta(0) = 0$, such that $\forall f \in \mathcal{F}$, $\mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^* \leq$*
155 *$\theta^{-1}(\mathcal{R}_{\ell_T, P_T}(f) - \mathcal{R}_{\ell_T, P_T}^*)$.*

156 The function θ depends on ℓ and T . The following proposition provides a convex lower bound on θ
157 for the commonly used log loss $\ell^{\log}(q, c) = -\log q_c$. Let $M \in \mathbb{R}^{C \times C}$ be a matrix and let $\|\cdot\|$ be a
158 norm on \mathbb{R}^C . The *subordinate matrix norm* induced by $\|\cdot\|$ is $\|M\| := \sup_{x \in \mathbb{R}^C : x \neq 0} \frac{\|Mx\|}{\|x\|}$. When
159 $\|\cdot\|$ is the 1-norm on \mathbb{R}^C , the induced norm is denoted $\|M\|_1$, referred to as the matrix 1-norm, and
160 can be computed as $\|M\|_1 = \max_{1 \leq j \leq C} \sum_{i=1}^C |M(i, j)|$ [10].

161 **Proposition 6.** *Let $T \in \mathbb{R}^{C \times C}$ be an invertible, column-stochastic matrix. Define $\underline{\theta}_T : [0, \infty] \rightarrow$*
162 *$[0, \infty]$ by $\underline{\theta}_T(\epsilon) = \frac{1}{2} \frac{\epsilon^2}{\|T^{-1}\|_1^2}$. If L is the 0/1 loss, ℓ is the log loss, then for all $f \in \mathcal{F}$ and distributions*
163 *P , $\mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^* \leq \underline{\theta}_T^{-1}(\mathcal{R}_{\ell_T, P_T}(f) - \mathcal{R}_{\ell_T, P_T}^*) = \sqrt{2} \|T^{-1}\|_1 \sqrt{\mathcal{R}_{\ell_T, P_T}(f) - \mathcal{R}_{\ell_T, P_T}^*}$*

164 The factor $\|T^{-1}\|_1$ may be viewed as a constant that captures the overall amount of label noise.
165 The more noise, the larger the constant. For example, let I and N be the identity and the all $1/C$'s
166 matrices, respectively. Let $\alpha \in [0, 1]$ and $T = (1 - \alpha)I + \alpha N$. Thus, $\alpha = 0$ represents the noise-free
167 case and $\alpha = 1$ the noise-only case. It is easy to verify that $T^{-1} = (1 - \alpha)^{-1}(I - \alpha N)$ and
168 $\|T^{-1}\|_1 = (1 - \alpha)^{-1}(1 + (1 - 2/C)\alpha)$.

169 4 Learning with Multiple Noise Transition Matrices

170 Our algorithm for LLP, formally stated in subsection 5.4, reduces the problem of LLP to LLN by
171 partitioning bags into groups and modeling each group as an LLN problem. Since each group has its
172 own noise transition matrix, this leads to a new problem that we refer to as learning with multiple
173 noise transition matrices (LMNTM). In this section, we show how to extend the calibration analysis
174 of section 3 to this setting. In addition, we offer a generalization error bound that justifies an empirical
175 risk minimization learning procedure based on a weighted sum of FC losses.

176 4.1 Learning with Multiple Noise Transition Matrices

177 We first define the LMNTM problem formally. For all $n \in \mathbb{N}$, denote $\mathbb{N}_n = \{1, 2, \dots, n\}$. Consider
178 a clean distribution P on $\mathcal{X} \times \mathcal{Y}$ and a finite sequence of noise transition matrices T_1, T_2, \dots, T_N .
179 For each i we denote the noisy prior as the $\alpha_i \in \hat{\Delta}^C$ where, $\forall c \in \mathcal{Y}$, $\alpha_i(c) = P_{T_i}(\tilde{Y} = c)$.
180 We assume the α_i 's are known for theoretical analysis. In practice, α_i is estimable as discussed
181 below. In LMNTM, we observe data points $S = \{X_{i,c,j} : i \in \mathbb{N}_N, c \in \mathcal{Y}, j \in \mathbb{N}_{n_{i,c}}\}$ where
182 $X_{i,c,j} \stackrel{i.i.d.}{\sim} P_{T_i}(\cdot | c)$, and $n_{i,c} \in \mathbb{N}$ is the number of data points drawn from the class conditional

183 distribution $P_{T_i}(\cdot | c)$. Assume all $X_{i,c,j}$'s are mutually independent. We make additional remarks
 184 on this setting in subsection [D.1](#) in the appendix.

185 4.2 An Risk for LMNTM

186 The following result extends Theorem [5](#) to LMNTM. It establishes that the risk $\tilde{R}_{\ell,P,\mathcal{T}}$, which can be
 187 estimated from LMNTM training data, is a valid surrogate risk. This type of result is not needed for
 188 the backward correction approach of Scott and Zhang [\[39\]](#).

189 **Theorem 7.** *Let L be the 0-1 loss and $N \in \mathbb{N}$. Consider a sequence of invertible column-stochastic*
 190 *matrices $\mathcal{T} = \{T_i\}_{i=1}^N$ and a continuous strictly proper loss function ℓ . Let $w = (w_i)_{i=1}^N \in \Delta^N$.*
 191 *Define $\tilde{R}_{\ell,P,\mathcal{T}} : \mathcal{F} \rightarrow \mathbb{R}$ by $\tilde{R}_{\ell,P,\mathcal{T}}(f) := \sum_{i=1}^N w_i \mathcal{R}_{\ell_{T_i}, P_{T_i}}(f)$ and $\tilde{R}_{\ell,P,\mathcal{T}}^* = \inf_{f \in \mathcal{F}} \tilde{R}_{\ell,P,\mathcal{T}}(f)$.*
 192 *Assume $\forall i \in \{1, 2, \dots, N\}, \mathcal{R}_{\ell_{T_i}, P_{T_i}}^* < \infty$. Then \exists a strictly increasing continuous function $\theta :$
 193 $[0, 1] \rightarrow [0, \infty]$ with $\theta(0) = 0$ s.t. for all $P, \forall f \in \mathcal{F}, \theta(\mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^*) \leq \tilde{R}_{\ell,P,\mathcal{T}}(f) - \tilde{R}_{\ell,P,\mathcal{T}}^*$.*

194 The weights w_i allow the user flexibility, for example, to place different weights on noisier or larger
 195 subsets of data. Unlike Scott and Zhang [\[39\]](#), however, because the weights appear in both our excess
 196 risk bound and generalization error bound, it is not straightforward to optimize them a priori.

197 4.3 Generalization Error Bound

198 The aggregate risk $\tilde{R}_{\ell,P,\mathcal{T}}$ is desirable because it can naturally be estimated from the given data. We
 199 propose the empirical risk

$$\hat{\mathcal{R}}_{w,S}(f) = \sum_{i=1}^N w_i \sum_{c=1}^C \frac{\alpha_i(c)}{n_{i,c}} \sum_{j=1}^{n_{i,c}} \ell_{T_i}(f(X_{i,c,j}), c). \quad (1)$$

200 It should be noted that $\hat{\mathcal{R}}_{w,S}(f)$ is an unbiased estimate of $\tilde{R}_{\ell,P,\mathcal{T}}(f)$. Here we establish a general-
 201 ization error bound for this estimate, that is, we show that the empirical risk concentrates around the
 202 risk uniformly. Our results build on Rademacher complexity analysis.

203 To state the bound, we must first introduce the notion of a *proper composite loss* [\[49\]](#). This stems
 204 from the fact that in practice, a function f outputting values in Δ^C is typically obtained by composing
 205 a \mathbb{R}^C -valued function (such as a neural network with C output layer nodes), with another function
 206 $\mathbb{R}^C \rightarrow \Delta^C$ such as the softmax function. Thus, let $\psi : \mathcal{U} \subset \Delta^C \rightarrow \mathcal{V}$ be an invertible function where
 207 \mathcal{V} is a subset of a normed space, referred to as an *invertible link function*. Consider $\mathcal{G} \subset \psi \circ \mathcal{F} :=$
 208 $\{\psi \circ f : f \in \mathcal{F}\}$, and observe that $\forall g \in \mathcal{G}, \psi^{-1} \circ g \in \mathcal{F}$. In practice, ψ is fixed and we seek to learn
 209 $g \in \mathcal{G}$ that leads to an $f \in \mathcal{F}$ with a risk close to the Bayes risk. An example of ψ^{-1} is the softmax
 210 function so that $\psi : \mathcal{U} \rightarrow \mathcal{V}, \psi_i(p) = \log p_i - \frac{1}{C} \sum_{k=1}^C \log p_k, (\psi^{-1})_i(s) = \frac{e^{s_i}}{\sum_{k=1}^C e^{s_k}}$ where \mathcal{U} is
 211 the interior of Δ^C and $\mathcal{V} = \{s \in \mathbb{R}^C : \sum_{i=1}^C s_i = 0\}$. This motivates the following definition.

212 **Definition 8** (Proper Composite Loss). *Given an invertible link function $\psi : \mathcal{U} \subset \Delta^C \rightarrow \mathcal{V}$, we define*
 213 *the **proper composite loss** λ_ℓ of a proper loss $\ell : \Delta^C \times \mathcal{Y} \rightarrow \mathbb{R}$ to be $\lambda_\ell : \mathcal{V} \times \mathcal{Y} \rightarrow \mathbb{R}, \lambda_\ell(v, c) =$
 214 $\ell(\psi^{-1}(v), c)$.*

215 For example, when ℓ is the log loss and ψ^{-1} is the softmax function, λ_ℓ is the cross-entropy (or
 216 multinomial logistic) loss. With this notation, we are now able to state our generalization error bound
 217 for LMNTM. We study two popular choices of function classes, the reproducing kernel Hilbert space
 218 (RKHS) and the multilayer perceptron (MLP).

219 **Definition 9.** *Let k be a symmetric positive definite (SPD) kernel, and let \mathcal{H} be the associated*
 220 *reproducing kernel Hilbert space (RKHS). Assume k is bounded by K , meaning $\forall x, \|k(\cdot, x)\|_{\mathcal{H}} \leq K$.*
 221 *Let $\mathcal{G}_{K,R}^k$ denote the ball of radius R in \mathcal{H} . Define $\mathcal{G}_1 = \mathcal{G}_{K,R}^k \times \mathcal{G}_{K,R}^k \times \dots \times \mathcal{G}_{K,R}^k$ (C copies).*

222 We follow Zhang et al. [\[54\]](#) and define real-valued MLPs inductively:

223 **Definition 10.** *Define $\mathcal{N}_1 = \{x \rightarrow \langle x, v \rangle : v \in \mathbb{R}^d, \|v\|_2 \leq \beta\}$, and for $m > 2$, inductively define*
 224 $\mathcal{N}_m = \left\{ x \rightarrow \sum_{j=1}^d v_j \mu(f_j(x)) : v \in \mathbb{R}^d, \|v\|_1 \leq \beta, f_j \in \mathcal{N}_{m-1} \right\}$, *where $\beta \in \mathbb{R}_+$ and μ is a 1-*
 225 *Lipschitz activation function. Define an MLP which outputs a vector in \mathbb{R}^C by $\mathcal{G}_2 = \mathcal{N}_m \times \mathcal{N}_m \times$
 226 $\dots \times \mathcal{N}_m$ (C copies). We additionally assume that the choice of μ satisfies $\forall m \in \mathbb{N}, 0 \in \mu \circ \mathcal{N}_m$.*

Theorem 11. Let T_1, T_2, \dots, T_N be a finite sequence of invertible column-stochastic matrices. Let ℓ be a proper loss such that $\forall i, c$ the function $\lambda_{\ell_{T_i}}(\cdot, c)$ is Lipschitz continuous w.r.t. the 2-norm. Let S be the set of data points as defined in Section 4.1. Assume $\sup_{x \in \mathcal{X}, g \in \mathcal{G}_q} \|g(x)\|_2 \leq A_q$ for some constant $A_q, \forall q \in \{1, 2\}$. Let $\hat{\mathcal{R}}_{w,S}(g) := \sum_{i=1}^N w_i \sum_{c=1}^C \frac{\alpha_i(c)}{n_{i,c}} \sum_{j=1}^{n_{i,c}} \lambda_{\ell_{T_i}}(g(X_{i,c,j}), c)$, $\tilde{\mathcal{R}}(g) := \tilde{R}_{\ell,P,T}(\psi^{-1} \circ g) = \mathbb{E}_S[\hat{\mathcal{R}}_{w,S}(g)]$. Then for each $q \in \{1, 2\}, \forall \delta \in [0, 1]$, with probability at least $1 - \delta$,

$$\sup_{g \in \mathcal{G}_q} \left| \hat{\mathcal{R}}_{w,S}(g) - \tilde{\mathcal{R}}(g) \right| \leq (\max_i (|\lambda_{\ell_{T_i}}| A_q + |\lambda_{\ell_{T_i}}|_0) \sqrt{2 \log \frac{2}{\delta}} + C B_q \max_i |\lambda_{\ell_{T_i}}|) \sqrt{\sum_{i=1}^N \sum_{c=1}^C \frac{w_i^2}{n_{i,c}}}.$$

where B_q is a constant depending on $\mathcal{G}_q, |\lambda_{\ell_{T_i}}|_0 = \max_c |\lambda_{\ell_{T_i}}(0, c)|$, and $|\lambda_{\ell_{T_i}}|$ is the smallest real number such that it is a Lipschitz constant of $\lambda_{\ell_{T_i}}(\cdot, c), \forall i, c$.

It should be noted that Theorem 11 is a special case of of Lemma 26 which extends the notion of Rademacher complexity to the LMNTM setting and applies to arbitrary function classes. Lemma 26 is presented in the appendix.

Let HM_i denote the harmonic mean of $n_{i,1}, \dots, n_{i,C}$, i.e., $HM_i = \frac{C}{\sum_{c=1}^C \frac{1}{n_{i,c}}}$. The term $\sqrt{\sum_{i=1}^N \sum_{c=1}^C \frac{w_i^2}{n_{i,c}}}$ could be written as $\sqrt{C \sum_{i=1}^N \frac{w_i^2}{HM_i}}$ and is optimized by $w_i = \frac{HM_i}{\sum_{m=1}^N HM_m}$, leading to $\sqrt{\sum_{i=1}^N \sum_{c=1}^C \frac{w_i^2}{n_{i,c}}} = \sqrt{\frac{C}{\sum_{i=1}^N HM_i}}$. The term $\sqrt{\frac{C}{\sum_{i=1}^N HM_i}}$ vanishes (needed to establish consistency) when N goes to infinity, or when $\exists i$ s.t. $\forall c, n_{i,c}$ goes to infinity. For the special case where all bags have the same size n and all weights w_i are $1/N$, $\sqrt{\sum_{i=1}^N \sum_{c=1}^C \frac{w_i^2}{n_{i,c}}} = \sqrt{\frac{C}{Nn}}$. Thus, consistency is possible even if bag size remains bounded.

Assuming ℓ is the log loss and ψ^{-1} is the softmax function, we next study the constants $|\lambda_{\ell_{T_i}}|$ and $|\lambda_{\ell_{T_i}}|_0$.

Proposition 12. Let ℓ be the log loss, ψ^{-1} be the softmax function, and T be a column-stochastic matrix. Then $|\lambda_{\ell_T}| \leq \sqrt{2}$.

The constant $|\lambda_{\ell_T}|_0 = \max_c |\lambda_{\ell_T}(0, c)| = \max_c -\log\left(\frac{1}{C} \sum_{j=1}^C t_{c,j}\right)$. The invertibility of T guarantees $\sum_{j=1}^C t_{c,j}$ is positive and hence the finiteness of $|\lambda_{\ell_T}|_0$. However, if we have a “bad” T , $\sum_{j=1}^C t_{c,j}$ could be arbitrarily close to 0 leading to a large $|\lambda_{\ell_T}|_0$.

Following Theorem 11, if the function class \mathcal{G} has a universal approximation property, such as an RKHS associated to a universal kernel, or an MLP with increasing number of nodes, consistency for LMNTM via (regularized) minimization of $\hat{\mathcal{R}}_{w,S}(g)$ can be shown by leveraging standard techniques, provided $N \rightarrow \infty$ (bag size may remain bounded). Then the excess risk bound in Theorem 7 would automatically imply consistency with respect to 0-1 loss.

5 The LLPFC algorithms

In this section, we define a probabilistic model for LLP, show how LLP reduces to LMNTM, and introduce algorithms that we refer to as the LLPFC algorithms.

5.1 Probabilistic Model for LLP

Given a measure P on the space $\mathcal{X} \times \mathcal{Y}$, let $\{P_c : c \in \mathcal{Y}\}$ denote the class-conditional distributions of \mathcal{X} , i.e., \forall events $\mathcal{A} \subset \mathcal{X}, P_c(\mathcal{A}) = P(\mathcal{A} | Y = c)$. Let $\sigma(c) = P(Y = c), \forall c \in \mathcal{Y}$ and call $\sigma = (\sigma(1), \dots, \sigma(C))$ the clean prior. Assume $\forall c \in \mathcal{Y}, \sigma(c) \neq 0$. Given $z = (z(1), \dots, z(C)) \in \Delta^C$, let P_z be the probability measure on $\mathcal{X} \times \mathcal{Y}$ s.t. \forall events $\mathcal{A} \subset \mathcal{X}, \forall i \in \mathcal{Y}, P_z(\mathcal{A} \times \{i\}) = z(i)P_i(\mathcal{A})$. Thus P_z has the same class-conditional distributions as P but a variable prior z .

We first define a model for a single bag. Given $z \in \Delta^C$, we say that bag b is governed by $z \in \Delta^C$ if b is a collection of feature vectors $\{X_j : j \in \mathbb{N}_{|b|}\}$ annotated by label proportion

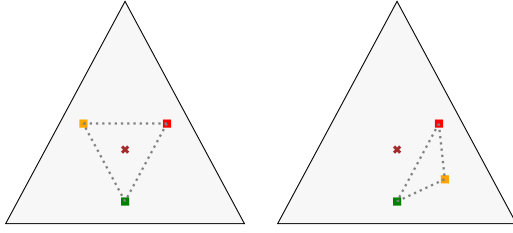


Figure 1: The gray triangle represents the probability simplex Δ^3 . The squares represent γ_1, γ_2 , and γ_3 . The cross is σ . The ternary graph on the left visualizes an example where Assumption 13 holds. The one on the right visualizes an example where Assumption 13 fails.

261 $\hat{z} = (\hat{z}(1), \hat{z}(2), \dots, \hat{z}(C))$, where $|b|$ denotes the cardinality of the bag, each X_j is paired with
 262 an unobserved label Y_j s.t. $(X_j, Y_j) \stackrel{iid}{\sim} P_z$, and $\hat{z}(c) = \frac{1}{|b|} \sum_{j=1}^{|b|} \mathbb{1}_{Y_j=c}$. Note $\mathbb{E}_{P_z}[\hat{z}] = z$ and
 263 $P_z(Y_j = c) = z(c)$. We think of z as the true label proportion and \hat{z} as the empirical label proportion.
 264 Using this model for individual bags, we now formally state a model for LLP. Given bags $\{b_k\}$,
 265 let each b_k be governed by γ_k . Each b_k is a collection of feature vectors $\{X_j^k : j \in \mathbb{N}_{|b_k|}\}$ where
 266 $(X_j^k, Y_j^k) \stackrel{i.i.d.}{\sim} P_{\gamma_k}$ and Y_j^k is unknown. Further assume the X_j^k 's are independent for all k and j .
 267 In practice, γ_k is unknown and we observe $\hat{\gamma}_k$ with $\hat{\gamma}_k(c) = \frac{1}{|b_k|} \sum_{j=1}^{|b_k|} \mathbb{1}_{Y_j^k=c}$ instead. The goal is
 268 learn an f that minimizes the risk $\mathcal{R}_{L,P} = \mathbb{E}_{(X,Y) \sim P}[L(f(X), Y)]$ where L is the 0-1 loss, given
 269 access to the training data $\{(b_k, \hat{\gamma}_k)\}$.

270 5.2 The Case of C Bags: Reduction to LLN

271 To explain our reduction of LLP to LLN, we first consider the case of exactly C bags b_1, b_2, \dots, b_C ,
 272 governed by respective (unobserved) $\gamma_1, \dots, \gamma_C \in \Delta^C$, and annotated with label proportions
 273 $\hat{\gamma}_1, \dots, \hat{\gamma}_C$. Define $\Gamma \in \mathbb{R}^{C \times C}$ by $\Gamma(i, j) = \gamma_i(j)$, and let Γ^{tr} denote the transpose of Γ . Re-
 274 call that σ is the class prior associated to P . To model LLP with C bags as an LLN problem, we
 275 make the following assumption on Γ and σ :

276 **Assumption 13.** \exists unique $\alpha \in \Delta^C$ s.t. $\Gamma^{tr} \alpha = \sigma$.

277 We write $\alpha = (\alpha(1), \dots, \alpha(C))$. Assumption 13 is equivalent to: $\{\gamma_1, \dots, \gamma_C\}$ is a linearly
 278 independent set and σ is in the interior of the convex hull of $\{\gamma_1, \dots, \gamma_C\}$. Ternary plots in Figure
 279 5.2 visualize examples where assumption 13 holds and fails when $C = 3$. Intuitively, assumption
 280 13 is more likely to hold when $\{\gamma_i : i \in \mathbb{N}_C\}$ are more “spread out” in Δ^C , in which case it is more
 281 likely for σ to reside in the convex hull of $\{\gamma_i : i \in \mathbb{N}_C\}$.

282 To reduce LLP with C bags to LLN, we simply propose to assign the “noisy label” $\tilde{Y} = i$ to all
 283 elements of bag b_i and to construct a noise transition matrix T with $T(i, j) = \frac{\gamma_i(j)\alpha(i)}{\sigma(j)}$. Assumption
 284 13 ensures T is indeed a column-stochastic matrix. Thus, the probability measure \bar{P}_T on $\mathcal{X} \times \mathcal{Y} \times \mathcal{Y}$
 285 satisfies $\alpha(i) = \bar{P}_T(\tilde{Y} = i)$ and $P_{\gamma_i}(\cdot) = \bar{P}_T(\cdot | \tilde{Y} = i)$, which further implies $\gamma_i(c) = \bar{P}_T(Y = c |$
 286 $\tilde{Y} = i)$. We confirm these facts in Section E in the appendix. Such construction transforms LLP with
 287 C bags into LLN with an estimable noise transition matrix T . Each element of a bag can then be
 288 viewed as a triplet (X, Y, \tilde{Y}) , with Y unobserved, such that (X, Y) is drawn from $P_{\gamma_{\tilde{Y}}}$.

289 After assigning the noisy labels, we have a dataset $\bigcup_{c=1}^C \{(X_j^c, c) : j \in \mathbb{N}_{|b_c|}\}$ along with the
 290 noise transition matrix T . This allows us to leverage the forward correction loss ℓ_T to mini-
 291 mize the objective $\mathcal{R}_{\ell_T, P_T}(f) = \mathbb{E}_{P_T}[\ell_T(f(X), \tilde{Y})]$ which can be estimated by the empirical
 292 risk $\sum_{c=1}^C \frac{\alpha(c)}{|b_c|} \sum_{j=1}^{|b_c|} \ell_T(f(X_j^c), c)$.

293 5.3 The General Case: Reduction to LMNTM

294 More generally, consider LLP with NC bags, $N \in \mathbb{N}$. We propose to randomly partition the bags
 295 into N groups, each with C bags indexed from 1 to C . Let $k_{i,c}$ denote the index of the c -th bag in the

Algorithm 1 LLPFC-ideal

1: **Input:** $\{(b_k, \gamma_k)\}_{k=1}^{NC}$ and $w \in \Delta^N$ where $b_k = \{X_j^k : j \in \mathbb{N}_{|b_k|}\}$.
 2: Randomly partition the bags into N groups $\{G_i\}_{i=1}^N$ s.t. $G_i = \{(b_{k_{i,c}}, \gamma_{k_{i,c}}) : c \in \mathcal{Y}\}$ and $\{k_{i,c} : i \in \mathbb{N}_N, c \in \mathcal{Y}\} = \mathbb{N}_{NC}$.
 3: **for** $i = 1 : N$ **do**
 4: $\Gamma_i \leftarrow [\gamma_{k_{i,1}}, \gamma_{k_{i,2}}, \dots, \gamma_{k_{i,C}}]^{tr}$
 5: $\alpha_i \leftarrow \Gamma_i^{-tr} \sigma$
 6: **for** $c_1 = 1 : C, c_2 = 1 : C$ **do**
 7: $T_i(c_1, c_2) \leftarrow \frac{\gamma_{k_{i,c_1}}(c_2) \alpha_i(c_1)}{\sigma(c_2)}$
 8: **end for**
 9: **end for**
 10: Train f with the empirical objective [\(1\)](#)
 11: **Return:** the trained model f .

296 i -th group. Thus, $b_{k_{i,c}}$ is the c -th bag in the i -th group and it is governed by $\gamma_{k_{i,c}}$. For $i \in \mathbb{N}_N$, define
 297 the matrix $\Gamma_i \in \mathbb{R}^{C \times C}$ by $\Gamma_i(c_1, c_2) = \gamma_{k_{i,c_1}}(c_2), \forall c_1, c_2 \in \mathcal{Y}$. We make the following assumption
 298 on the Γ_i 's and σ :

299 **Assumption 14.** For each $i \in \mathbb{N}_N, \exists$ unique $\alpha_i \in \hat{\Delta}^C$ s.t. $\Gamma_i^{tr} \alpha_i = \sigma$.

300 Thus, every group i can be modeled as above as an LLN problem with noise transition matrix T_i
 301 where $T_i(c_1, c_2) = \frac{\gamma_{k_{i,c_1}}(c_2) \alpha_i(c_1)}{\sigma(c_2)}$. Data points in the bag assigned with noisy label c in the i -th
 302 group can be viewed as drawn *i.i.d.* from the class conditional distribution $P_{T_i}(\cdot | c)$. This problem
 303 now maps directly to LMNTM as described in Section [4.1](#) and satisfies the associated performance
 304 guarantees. In the next subsection, we spell out the associated algorithm.

305 5.4 Algorithm

306 As above, assume we have NC bags where $N \in \mathbb{N}$. Let each bag b_k be governed by $\gamma_k \in \Delta^C$ and
 307 be annotated by label proportion $\hat{\gamma}_k$. Denote the size of b_k by $|b_k|$. We first present the LLPFC-ideal
 308 algorithm in an ideal setting where σ , the γ_k 's and the α_i 's are known precisely and Assumption [14](#)
 309 holds. We then present the real-world adaptations LLPFC-uniform and LLPFC-approx in practical
 310 settings where these assumptions might not hold.

311 The LLPFC-ideal algorithm is presented in Algorithm [1](#). We follow the idea in section [5.3](#) to partition
 312 the bags into N groups of C bags, and model each group as an LLN problem. In Algorithm [1](#), we
 313 assume γ_k and σ are known and Assumption [14](#) holds so that $\alpha_i \in \hat{\Delta}^C$ for every i . The theoretical
 314 analysis in Section [4](#) is immediately applicable to the LLPFC-ideal algorithm. We partition the
 315 bags by uniformly randomly partitioning the set of indices \mathbb{N}_{NC} into disjoint subsets $\{k_{i,c} : c \in \mathcal{Y}\}$,
 316 $i \in \mathbb{N}_N$, where $k_{i,c}$ denotes the index of the c -th bag in the i -th group. Hence, the c -th bag in the i -th
 317 group is denoted as $b_{k_{i,c}}$ and it is governed by $\gamma_{k_{i,c}}$. Following our convention, in Algorithms [1](#), [2](#),
 318 and [3](#) we denote the element at c_1 -th row and c_2 -th column in a matrix T by $T(c_1, c_2)$ and the c_1 -th
 319 element in vector v by $v(c_1)$. We denote the inverse transpose of Γ_i by Γ_i^{-tr} .

320 In practice, when γ_k is unknown, we replace γ_k with $\hat{\gamma}_k$ as a plug-in method. Hence, we work with
 321 $\hat{\sigma} = \frac{\sum_{k=1}^{NC} |b_k| \hat{\gamma}_k}{\sum_{k=1}^{NC} |b_k|}$ and $\hat{\Gamma}_i = [\hat{\gamma}_{k_{i,1}}, \hat{\gamma}_{k_{i,2}}, \dots, \hat{\gamma}_{k_{i,C}}]^{tr}$ instead of σ and Γ_i in Algorithm [1](#), respectively.
 322 Here $\hat{\sigma}$ is the label proportion of all training data points and we use it as an estimate of the clean prior σ .
 323 Likewise, $\alpha_i = \Gamma_i^{-tr} \sigma$ in Algorithm [1](#) should be replaced with $\hat{\alpha}_i = \hat{\Gamma}_i^{-tr} \hat{\sigma}$ and we would like to use
 324 $\hat{\Gamma}_i, \hat{\sigma}$, and $\hat{\alpha}_i$ to calculate \hat{T}_i as an estimate of T_i . For this to make sense, we need $\hat{\alpha}_i = \hat{\Gamma}_i^{-tr} \hat{\sigma} \in \hat{\Delta}^C$,
 325 which is equivalent to $\hat{\sigma}$ being in the interior of the convex hull of $\{\hat{\gamma}_{k_{i,c}} : c \in \mathcal{Y}\}$ for all i . However,
 326 this may not be the case in practice. Thus, we consider two heuristics to estimate \hat{T}_i as real-world
 327 adaptations of the LLPFC-ideal algorithm. The first, called LLPFC-uniform, is presented in Algorithm
 328 [2](#) in the Appendix which sets $\hat{\alpha}_i$ by counting the occurrences of the noisy labels. This is motivated by
 329 our model wherein α_i is the noisy class prior for the i -th group. The second, called LLPFC-approx, is
 330 presented in Algorithm [3](#) in the Appendix and sets $\hat{\alpha}_i$ to be the solution of $\arg \min_{\alpha \in \Delta^C} \|\hat{\sigma} - \hat{\Gamma}_i \alpha\|_2^2$.
 331 Both real-world adaptations perform reasonably well in experiments.

332 6 Experiments

333 We compare against three previous works that have studied LLP applying deep learning to image
334 data: Dulac-Arnold et al. [9] study the KL loss described in the introduction, and a novel loss
335 based on optimal transport. They find that KL performs just as well as the novel loss. Liu et al.
336 [22] employ the KL loss within a generative adversarial framework (LLPGAN). Tsai and Lin [46]
337 propose augmenting the KL loss with a regularizer from semi-supervised learning and show improved
338 performance (LLPVAT). We compare both LLPFC-uniform and LLPFC-approx against the KL loss,
339 LLPGAN, and LLPVAT to clearly establish which empirical objective is better. Recent papers on
340 multiclass LLP for which code is not available were not included [23, 16].

341 We generate bags with fixed, equal sizes in $\{32, 64, 128, 256, 512, 1024, 2048\}$. To generate each
342 bag, we first sample a label proportion γ from the uniform distribution on Δ^C . Then we sample data
343 points from a benchmark dataset without replacement using a multinomial distribution with parameter
344 γ . It should be noted that Tsai and Lin [46], Dulac-Arnold et al. [9], and Liu et al. [22] generate bags
345 by shuffling all data points and making every B data points a bag where B is a fixed bag size. Their
346 method is equivalent to sampling data points without replacement using a multinomial distribution
347 with a fixed parameter $\gamma = \{\frac{1}{C}, \frac{1}{C}, \dots, \frac{1}{C}\}$. As noted by Scott and Zhang [39], this leads to bags
348 with very similar label proportions which makes the learning task much more challenging. All models
349 are trained on a single Nvidia Tesla v100 GPU with 16GB RAM. We repeat each experiment 5 times
350 and report the mean test accuracy and standard deviation.

351 For the comparison against KL and LLPVAT, we perform experiments on three benchmark image
352 datasets: the “letter” split of EMNIST [6], SVHN [29], and CIFAR10 [17]. To show that our approach
353 is robust to the choice of architecture, we experiment with three different networks: Wide ResNet-16-
354 4 [52], ResNet18 [11], and VGG16 [42]. We train these networks with the parameters suggested in
355 the original papers. The test accuracies are reported in Tables 1, 2, and 3 in the appendix.

356 Since convergence in the GAN framework is sensitive to the choice of architecture and hyperparam-
357 eters, we compare LLPFC against LLPGAN using the architecture proposed in the original paper
358 along with the hyperparameters suggested in their code¹. It should be noted that for LLPFC we only
359 use the discriminator for classification and did not use the generator to augment data. Since Liu
360 et al. [22] only provide hyperparameters for colored images, we perform experiments on SVHN and
361 CIFAR10 only. The test accuracies are reported in Table 4 in the appendix.

362 LLPFC-uniform and LLPFC-approx substantially outperform the competitors in a clear majority
363 of settings. In a few settings, the competitors perform better, but LLPFC is not far behind. On the
364 other hand, in many settings, the competitors perform very poorly. Quantitatively, comparing mean
365 accuracies across all models and bag sizes, $\min(\text{llpfc-uniform}, \text{llpfc-approx}) - \max(\text{kl}, \text{llpgan}, \text{llpvat})$
366 $> 22\%$ on CIFAR10, $> 27\%$ on SVHN, and $> 28\%$ on EMNIST. In addition, all three competitors
367 perform gradient descent with minibatches of bags and the GPU could potentially run out of memory
368 when the bag size is large. Our implementation, which also uses stochastic optimization (with random
369 regrouping of bags at every epoch), does not suffer from this phenomenon. Full experimental details
370 are in an appendix.

371 7 Conclusions and Future Work

372 We propose a theoretically supported approach to LLP by reducing it to learning with label noise
373 and using the forward correction (FC) loss. An excess risk bound and generalization error analysis
374 are established. Our approach outperforms leading existing methods in deep learning scenarios
375 across multiple datasets and architectures. A limitation of our approach is that the theory makes
376 an assumption that may not be valid in practice. Future research directions include optimizing the
377 grouping of bags and adapting LLPFC to other objectives beyond accuracy.

¹<https://github.com/liujiabin008/LLP-GAN>

378 **References**

- 379 [1] Ehsan M. Ardehaly and Aron Culotta. “Co-Training for Demographic Classification Using
380 Deep Learning from Label Proportions”. In: *2017 IEEE International Conference on Data
381 Mining Workshops (ICDMW)*. 2017, pp. 1017–1024.
- 382 [2] Gerda Bortsova, Florian Dubost, Silas Ørting, Ioannis Katramados, Laurens Hogeweg, Laura
383 Thomsen, Mathilde Wille, and Marleen de Bruijne. “Deep Learning from Label Proportions
384 for Emphysema Quantification”. In: *Medical Image Computing and Computer Assisted Inter-
385 vention – MICCAI 2018*. Ed. by Alejandro F. Frangi, Julia A. Schnabel, Christos Davatzikos,
386 Carlos Alberola-López, and Gabor Fichtinger. Cham: Springer International Publishing, 2018,
387 pp. 768–776. ISBN: 978-3-030-00934-2.
- 388 [3] S. Chen, B. Liu, M. Qian, and C. Zhang. “Kernel K-means Based Framework for Aggregate
389 Outputs Classification”. In: *2009 IEEE International Conference on Data Mining Workshops*.
390 2009, pp. 356–361.
- 391 [4] Tao Chen, Felix X. Yu, Jiawei Chen, Yin Cui, Yan-Ying Chen, and Shih-Fu Chang. “Object-
392 Based Visual Sentiment Concept Analysis and Application”. In: *Proceedings of the 22nd ACM
393 International Conference on Multimedia*. MM ’14. Orlando, Florida, USA: Association for
394 Computing Machinery, 2014, pp. 367–376. ISBN: 9781450330633. DOI: [10.1145/2647868.
395 2654935](https://doi.org/10.1145/2647868.2654935). URL: <https://doi.org/10.1145/2647868.2654935>.
- 396 [5] Zhensong Chen, Zhiquan Qi, Bo Wang, Limeng Cui, Fan Meng, and Yong Shi. “Learning
397 with label proportions based on nonparallel support vector machines”. In: *Knowledge-Based
398 Systems* 119 (2017), pp. 126–141.
- 399 [6] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and André van Schaik. “EMNIST: an exten-
400 sion of MNIST to handwritten letters”. In: *ArXiv abs/1702.05373* (2017).
- 401 [7] Lucio Mwinmaarong Dery, Benjamin Nachman, Francesco Rubbo, and Ariel Schwartzman.
402 “Weakly Supervised Classification For High Energy Physics”. In: *Journal of Physics: Con-
403 ference Series* 1085 (Sept. 2018), p. 042006. DOI: [10.1088/1742-6596/1085/4/042006](https://doi.org/10.1088/1742-6596/1085/4/042006)
404 URL: <https://doi.org/10.1088/1742-6596/1085/4/042006>.
- 405 [8] Yongke Ding, Yuanxiang Li, and W. Yu. “Learning from label proportions for SAR image
406 classification”. In: *EURASIP Journal on Advances in Signal Processing* 2017 (2017), pp. 1–12.
- 407 [9] Gabriel Dulac-Arnold, Neil Zeghidour, Marco Cuturi, Lucas Beyer, and Jean-Philippe Vert.
408 “Deep multi-class learning from label proportions”. In: *ArXiv abs/1905.12909* (2019).
- 409 [10] Jean Gallier and Jocelyn Quaintance. *Linear Algebra and Optimization with Applications to
410 Machine Learning*. World Scientific, 2020.
- 411 [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. “Deep Residual Learning for
412 Image Recognition”. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition
413 (CVPR)*. 2016, pp. 770–778. DOI: [10.1109/CVPR.2016.90](https://doi.org/10.1109/CVPR.2016.90).
- 414 [12] J. Hernández-González, I. Inza, Lorena Crisol-Ortíz, M. A. Guedebe, M. J. Iñarra, and J. A.
415 Lozano. “Fitting the data from embryo implantation prediction: Learning from label propor-
416 tions”. In: *Statistical Methods in Medical Research* 27.4 (2018), pp. 1056–1066.
- 417 [13] J. Hernández-González, I. Inza, and J. A. Lozano. “Learning Bayesian network classifiers
418 from label proportions”. In: *Pattern Recognition* 46.12 (2013), pp. 3425–3440.
- 419 [14] Olav Kallenberg. *Foundations of Modern Probability*. Springer, 2002.
- 420 [15] Diederik P. Kingma and Jimmy Ba. “Adam: A Method for Stochastic Optimization”. In: *3rd
421 International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May
422 7-9, 2015, Conference Track Proceedings*. Ed. by Yoshua Bengio and Yann LeCun. 2015. URL:
423 <http://arxiv.org/abs/1412.6980>.
- 424 [16] Ryoma Kobayashi, Yusuke Mukuta, and Tatsuya Harada. *Risk Consistent Multi-Class Learning
425 from Label Proportions*. 2022. DOI: [10.48550/ARXIV.2203.12836](https://doi.org/10.48550/ARXIV.2203.12836) URL: [https://arxiv
426 org/abs/2203.12836](https://arxiv.org/abs/2203.12836).
- 427 [17] Alex Krizhevsky. “Learning Multiple Layers of Features from Tiny Images”. In: *University of
428 Toronto* (May 2012).
- 429 [18] Hendrik Kück and Nando de Freitas. “Learning about Individuals from Group Statistics”. In:
430 *Proceedings of the Twenty-First Conference on Uncertainty in Artificial Intelligence*. 2005,
431 pp. 332–339.

- 432 [19] K. Lai, F. X. Yu, M. Chen, and S. Chang. “Video Event Detection by Inferring Temporal
433 Instance Labels”. In: *2014 IEEE Conference on Computer Vision and Pattern Recognition*.
434 2014, pp. 2251–2258.
- 435 [20] Kuan-Ting Lai, Felix Yu, Ming-Syan Chen, and S. Chang. “Video Event Detection by Inferring
436 Temporal Instance Labels”. In: June 2014. DOI: [10.1109/CVPR.2014.288](https://doi.org/10.1109/CVPR.2014.288)
- 437 [21] Fan Li and Graham Taylor. “Alter-CNN: An Approach to Learning from Label Proportions
438 with Application to Ice-Water Classification”. In: *Neural Information Processing Systems
439 Workshops (NIPSW) on Learning and privacy with incomplete data and weak supervision*.
440 2015.
- 441 [22] Jiabin Liu, Bo Wang, Zhiquan Qi, Yingjie Tian, and Yong Shi. “Learning from Label Proportions
442 with Generative Adversarial Networks”. In: *Advances in Neural Information Processing
443 Systems*. Ed. by H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and
444 R. Garnett. Vol. 32. Curran Associates, Inc., 2019. URL: [https://proceedings.neurips
445 cc/paper/2019/file/4fc848051e4459b8a6afeb210c3664ec-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/4fc848051e4459b8a6afeb210c3664ec-Paper.pdf)
- 446 [23] Jiabin Liu, Bo Wang, Xin Shen, Zhiquan Qi, and Yingjie Tian. “Two-stage Training for
447 Learning from Label Proportions”. In: *Proceedings of the Thirtieth International Joint Con-
448 ference on Artificial Intelligence, IJCAI-21*. Ed. by Zhi-Hua Zhou. Main Track. International
449 Joint Conferences on Artificial Intelligence Organization, Aug. 2021, pp. 2737–2743. DOI:
450 [10.24963/ijcai.2021/377](https://doi.org/10.24963/ijcai.2021/377). URL: <https://doi.org/10.24963/ijcai.2021/377>
- 451 [24] Nan Lu, Shida Lei, Gang Niu, Issei Sato, and Masashi Sugiyama. “Binary Classification
452 from Multiple Unlabeled Datasets via Surrogate Set Classification”. In: *Proceedings of the
453 38th International Conference on Machine Learning*. Ed. by Marina Meila and Tong Zhang.
454 Vol. 139. Proceedings of Machine Learning Research. PMLR, July 2021, pp. 7134–7144. URL:
455 <https://proceedings.mlr.press/v139/lu21c.html>
- 456 [25] Andreas Maurer. “A vector-contraction inequality for Rademacher complexities”. In: *Internation-
457 al Conference on Algorithmic Learning Theory*. Springer. 2016, pp. 3–17.
- 458 [26] Colin McDiarmid. “On the method of bounded differences”. In: *Surveys in combinatorics
459 141.1* (1989), pp. 148–188.
- 460 [27] Ron Meir and Tong Zhang. “Generalization Error Bounds for Bayesian Mixture Algorithms”.
461 In: *J. Mach. Learn. Res.* 4.null (Dec. 2003), pp. 839–860. ISSN: 1532-4435.
- 462 [28] Nagarajan Natarajan, Inderjit S Dhillon, Pradeep K Ravikumar, and Ambuj Tewari. “Learning
463 with Noisy Labels”. In: *Advances in Neural Information Processing Systems*. Ed. by C. J. C.
464 Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger. Vol. 26. Curran
465 Associates, Inc., 2013. URL: [https://proceedings.neurips.cc/paper/2013/file/
466 3871bd64012152bfb53fdf04b401193f-Paper.pdf](https://proceedings.neurips.cc/paper/2013/file/3871bd64012152bfb53fdf04b401193f-Paper.pdf)
- 467 [29] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng.
468 “Reading Digits in Natural Images with Unsupervised Feature Learning”. In: *NIPS Workshop
469 on Deep Learning and Unsupervised Feature Learning 2011*. 2011. URL: [http://ufldl
470 stanford.edu/housenumbers/nips2011_housenumbers.pdf](http://ufldl.stanford.edu/housenumbers/nips2011_housenumbers.pdf)
- 471 [30] Giorgio Patrini, A. Rozza, A. Menon, R. Nock, and Lizhen Qu. “Making Deep Neural Networks
472 Robust to Label Noise: A Loss Correction Approach”. In: *2017 IEEE Conference on Computer
473 Vision and Pattern Recognition (CVPR)* (2017), pp. 2233–2241.
- 474 [31] Mark S Pinsker. *Information and information stability of random variables and processes*.
475 Holden-Day, 1964.
- 476 [32] R. Poyiadzi, R. Santos-Rodriguez, and N. Twomey. “LABEL PROPAGATION FOR LEARN-
477 ING WITH LABEL PROPORTIONS”. In: *2018 IEEE 28th International Workshop on Ma-
478 chine Learning for Signal Processing (MLSP)*. 2018, pp. 1–6.
- 479 [33] Zhiquan Qi, Bo Wang, Fan Meng, and Lingfeng Niu. “Learning With Label Proportions via
480 NPSVM”. In: *IEEE Transactions on Cybernetics* 47 (2017), pp. 3293–3305.
- 481 [34] Novi Quadrianto, Alex J. Smola, Tiberio S. Caetano, and Quoc V. Le. “Estimating Labels
482 from Label Proportions”. In: *Proceedings of the 25th International Conference on Machine
483 Learning, ICML ’08*. Helsinki, Finland: Association for Computing Machinery, 2008, pp. 776–
484 783. ISBN: 9781605582054. DOI: [10.1145/1390156.1390254](https://doi.org/10.1145/1390156.1390254). URL: [https://doi.org/
485 10.1145/1390156.1390254](https://doi.org/10.1145/1390156.1390254)
- 486 [35] Brendan van Rooyen and Robert C. Williamson. “A Theory of Learning with Corrupted
487 Labels”. In: *Journal of Machine Learning Research* 18.228 (2018), pp. 1–50.

- 488 [36] Stefan Rueding. “SVM Classifier Estimation from Group Probabilities”. In: *Proceedings of the*
489 *27th International Conference on International Conference on Machine Learning*. ICML’10.
490 Haifa, Israel: Omnipress, 2010, pp. 911–918. ISBN: 9781605589077.
- 491 [37] Rishi Saket. “Learnability of Linear Thresholds from Label Proportions”. In: *Advances in*
492 *Neural Information Processing Systems*. Ed. by A. Beygelzimer, Y. Dauphin, P. Liang, and
493 J. Wortman Vaughan. 2021. URL: <https://openreview.net/forum?id=5BnaKeEwuYk>
- 494 [38] Rishi Saket, Aravindan Raghuv eer, and Balaraman Ravindran. “On Combining Bags to Better
495 Learn from Label Proportions”. In: *Proceedings of The 25th International Conference on*
496 *Artificial Intelligence and Statistics*. Ed. by Gustau Camps-Valls, Francisco J. R. Ruiz, and
497 Isabel Valera. Vol. 151. Proceedings of Machine Learning Research. PMLR, Mar. 2022,
498 pp. 5913–5927. URL: <https://proceedings.mlr.press/v151/saket22a.html>.
- 499 [39] Clayton Scott and Jianxin Zhang. “Learning from Label Proportions: A Mutual Contamination
500 Framework”. In: *Advances in Neural Information Processing Systems*. Ed. by H. Larochelle,
501 M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin. Vol. 33. Curran Associates, Inc., 2020,
502 pp. 22256–22267. URL: <https://proceedings.neurips.cc/paper/2020/file/fcde14913c766cf307c75059e0e89af5-Paper.pdf>
- 503 [40] Yong Shi, Limeng Cui, Zhensong Chen, and Zhiquan Qi. “Learning from label proportions
504 with pinball loss”. In: *International Journal of Machine Learning and Cybernetics* 10 (2017),
505 pp. 187–205.
- 506 [41] Yong Shi, Jiabin Liu, Zhiquan Qi, and Bo Wang. “Learning from label proportions on high-
507 dimensional data”. In: *Neural Networks* 103 (2018), pp. 9–18. ISSN: 0893-6080. DOI: <https://doi.org/10.1016/j.neunet.2018.03.004>, URL: <https://www.sciencedirect.com/science/article/pii/S0893608018300893>
- 508 [42] Karen Simonyan and Andrew Zisserman. “Very Deep Convolutional Networks for Large-Scale
509 Image Recognition”. In: *International Conference on Learning Representations*. 2015.
- 510 [43] Ingo Steinwart. “How to Compare Different Loss Functions and Their Risks”. In: *Constructive*
511 *Approximation* 26 (2007), pp. 225–287.
- 512 [44] Marco Stolpe and Katharina Morik. “Learning from Label Proportions by Optimizing Cluster
513 Model Selection”. In: *Machine Learning and Knowledge Discovery in Databases*. Ed. by
514 Dimitrios Gunopulos, Thomas Hofmann, Donato Malerba, and Michalis Vazirgiannis. Berlin,
515 Heidelberg: Springer Berlin Heidelberg, 2011, pp. 349–364. ISBN: 978-3-642-23808-6.
- 516 [45] Tao Sun, Dan Sheldon, and Brendan O’Connor. “A Probabilistic Approach for Learning with
517 Label Proportions Applied to the US Presidential Election”. In: *2017 IEEE International*
518 *Conference on Data Mining (ICDM)*. 2017, pp. 445–454. DOI: [10.1109/ICDM.2017.54](https://doi.org/10.1109/ICDM.2017.54)
- 519 [46] Kuen-Han Tsai and Hsuan-Tien Lin. “Learning from Label Proportions with Consistency Reg-
520 ularization”. In: *Proceedings of The 12th Asian Conference on Machine Learning, ACML 2020,*
521 *18-20 November 2020, Bangkok, Thailand*. Ed. by Sinno Jialin Pan and Masashi Sugiyama.
522 Vol. 129. Proceedings of Machine Learning Research. PMLR, 2020, pp. 513–528. URL:
523 <http://proceedings.mlr.press/v129/tsai20a.html>
- 524 [47] B. Wang, Z. Chen, and Z. Qi. “Linear Twin SVM for Learning from Label Proportions”. In:
525 *2015 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent*
526 *Technology (WI-IAT)*. Vol. 3. 2015, pp. 56–59.
- 527 [48] Nik Weaver. *Lipschitz Algebras*. WORLD SCIENTIFIC, 1999. DOI: [10.1142/4100](https://doi.org/10.1142/4100) eprint:
528 <https://www.worldscientific.com/doi/pdf/10.1142/4100>, URL: <https://www.worldscientific.com/doi/abs/10.1142/4100>
- 529 [49] Robert C. Williamson, Elodie Vernet, and Mark D. Reid. “Composite Multiclass Losses”. In:
530 *Journal of Machine Learning Research* 17.222 (2016), pp. 1–52. URL: <http://jmlr.org/papers/v17/14-294.html>
- 531 [50] Felix Yu, Dong Liu, Sanjiv Kumar, Jebara Tony, and Shih-Fu Chang. “ ∞ SVM for Learning
532 with Label Proportions”. In: *Proceedings of the 30th International Conference on Machine*
533 *Learning*. Ed. by Sanjoy Dasgupta and David McAllester. Vol. 28. Proceedings of Machine
534 Learning Research 3. Atlanta, Georgia, USA: PMLR, June 2013, pp. 504–512. URL: <http://proceedings.mlr.press/v28/yu13a.html>
- 535 [51] Felix X. Yu, Krzysztof Choromanski, Sanjiv Kumar, Tony Jebara, and Shih-Fu Chang. *On*
536 *Learning from Label Proportions*. Tech. rep. arXiv:1402.5902. 2015.

- 543 [52] Sergey Zagoruyko and Nikos Komodakis. “Wide Residual Networks”. In: *Proceedings of the*
544 *British Machine Vision Conference (BMVC)*. Ed. by Edwin R. Hancock Richard C. Wilson and
545 William A. P. Smith. BMVA Press, Sept. 2016, pp. 87.1–87.12. ISBN: 1-901725-59-6. DOI:
546 [10.5244/C.30.87](https://dx.doi.org/10.5244/C.30.87). URL: <https://dx.doi.org/10.5244/C.30.87>
- 547 [53] Mingyuan Zhang, Jane Lee, and Shivani Agarwal. “Learning from Noisy Labels with No
548 Change to the Training Process”. In: *Proceedings of the 38th International Conference on*
549 *Machine Learning*. Ed. by Marina Meila and Tong Zhang. Vol. 139. Proceedings of Machine
550 Learning Research. PMLR, July 2021, pp. 12468–12478. URL: [https://proceedings.mlr](https://proceedings.mlr.press/v139/zhang21k.html)
551 [press/v139/zhang21k.html](https://proceedings.mlr.press/v139/zhang21k.html)
- 552 [54] Yuchen Zhang, Jason Lee, Martin Wainwright, and Michael I. Jordan. “On the learnability of
553 fully-connected neural networks”. In: *Proceedings of the 20th International Conference on*
554 *Artificial Intelligence and Statistics*. Ed. by Aarti Singh and Jerry Zhu. Vol. 54. Proceedings of
555 Machine Learning Research. Fort Lauderdale, FL, USA: PMLR, Apr. 2017, pp. 83–91. URL:
556 <http://proceedings.mlr.press/v54/zhang17a.html>.

557 Checklist

- 558 1. For all authors...
- 559 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
560 contributions and scope? [Yes]
- 561 (b) Did you describe the limitations of your work? [Yes] Limitations are described in
562 Section [7](#).
- 563 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 564 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
565 them? [Yes]
- 566 2. If you are including theoretical results...
- 567 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 568 (b) Did you include complete proofs of all theoretical results? [Yes] All proofs are included
569 in the appendix.
- 570 3. If you ran experiments...
- 571 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
572 perimental results (either in the supplemental material or as a URL)? [Yes] Code is
573 included in the supplemental material. Datasets are public and we provide code to
574 download them. We include a README file with instructions on how to reproduce
575 experimental results.
- 576 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
577 were chosen)? [Yes] A concise description of experiments is in Section [6](#) in the main
578 paper with full details in section [B](#) in the appendix.
- 579 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
580 ments multiple times)? [Yes] In table [1](#), [2](#), [3](#) and [4](#), we report the mean and standard
581 deviation for 5 trials with different random seeds.
- 582 (d) Did you include the total amount of compute and the type of resources used (e.g.,
583 type of GPUs, internal cluster, or cloud provider)? [No] We provide the hardware
584 information in Section [6](#) but not the amount of compute. The computational time varies
585 for different experimental settings. It would be too exhaustive to present them given
586 that we have nearly 300 different settings.
- 587 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 588 (a) If your work uses existing assets, did you cite the creators? [Yes] Yes. Creators of code,
589 data, and models are all cited.
- 590 (b) Did you mention the license of the assets? [N/A] We did not directly run experiments
591 using others’ code. However, we implement their algorithms with their code as a
592 reference.
- 593 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
594 Our experiment code is included in the supplemental materials.

- 595 (d) Did you discuss whether and how consent was obtained from people whose data you're
596 using/curating? [N/A]
- 597 (e) Did you discuss whether the data you are using/curating contains personally identifiable
598 information or offensive content? [N/A]
- 599 5. If you used crowdsourcing or conducted research with human subjects...
- 600 (a) Did you include the full text of instructions given to participants and screenshots, if
601 applicable? [N/A]
- 602 (b) Did you describe any potential participant risks, with links to Institutional Review
603 Board (IRB) approvals, if applicable? [N/A]
- 604 (c) Did you include the estimated hourly wage paid to participants and the total amount
605 spent on participant compensation? [N/A]