

LEARNING GENERALIZED LABEL DISTRIBUTIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Label ambiguity is pervasive in supervised learning, motivating a variety of representations beyond the traditional single-label setting. While label distribution (LD) provides a probabilistic description and has attracted increasing attention, we reveal its inherent limitations, including inconsistency with raw data, distortion of inter-sample order, and limited applicability. To address these issues, we introduce generalized label distribution (GLD), a unified representation that can perfectly recover raw data while preserving inter-sample order consistency, transform into existing forms of label representations without information loss, and capture out-of-distribution samples as well as negative label correlations. We further develop GLD learning algorithms and demonstrate their effectiveness through both theoretical analysis and extensive experiments.

1 INTRODUCTION

In supervised learning, the *single-label* (SL) setting has been the most widely adopted, where a sample is associated with a categorical variable. Growing recognition of label ambiguity has motivated the use of *logical label* (LL), i.e., a binary vector, to describe a sample, giving rise to multi-label learning (MLL) (Tsoumakas et al., 2010; Zhang & Zhou, 2013). In recent years, *label ranking* (LR) (Brinker et al., 2006; Lu & Jia, 2022) and *ternary label* (TL) (Lu & Jia, 2024) have also been introduced as alternative ways to model label ambiguity. Following this direction, *label distribution* (LD) has attracted increasing attention as a more fine-grained representation, using a probability distribution to describe a sample, leading to label distribution learning (LDL) (Geng, 2016).

However, in this paper, we argue that LD suffers from inherent limitations, including inconsistency with the raw data, disruption of inter-sample order, and limited applicability. These limitations not only challenge the claim that “*LDL represents a generalized form of MLL*” (Geng, 2016), but also expose critical shortcomings in current derivative tasks of LDL, e.g., label enhancement (LE) (Xu et al., 2019), joint LDL & LE (Liu et al., 2021; Jia et al., 2024), and classification-oriented LDL (Wang & Geng, 2019a; Wang et al., 2021), whose methodologies exhibit room for a more substantial refinement. Addressing these issues is both urgent and necessary, as they hinder further progress in learning with label ambiguity. In response, we propose a new form of representation, namely generalized label distribution (GLD), which extends LD to provide a more versatile and comprehensive characterization of label ambiguity. Specifically, GLD offers the following advantages:

- GLD serves as a truly unified form of label polysemy representation: it can be mapped upward to recover the raw data perfectly, and downward to derive other existing forms of label representations (Xu et al., 2020) (e.g., LD, LL, TL, LR, and SL) without any information loss.
- GLD can naturally characterize out-of-distribution (OOD) samples (Ren et al., 2019; Fort et al., 2021) and explicitly capture fine-grained negative correlations among labels (Xu et al., 2019).
- GLD does not counterintuitively distort the order relations across samples (Bergeron et al., 2008), making it more consistent with both machine learning models and human perception.

Contributions & organizational structure In this paper: (1) we conduct a thorough theoretical analysis of the inherent limitations of LD (Section 2); (2) we introduce GLD as a unified representation that addresses these limitations (Section 3.1); (3) we further propose several algorithms for effectively learning GLDs, which can also leverage and adapt existing LDL methods to extend their applicability to GLD (Section 3.2); (4) we provide both theoretical analyses and extensive experiments to demonstrate the effectiveness of the proposed GLD learning framework (Sections 3.3 and 4). The code will be available on Github soon, facilitating reproducibility and further research.

2 LIMITATIONS OF THE LABEL DISTRIBUTION

2.1 PRELIMINARIES

Notation Vectors are denoted by lowercase bold letters, e.g., \mathbf{v} , and the corresponding regular letter with subscript i , i.e., v_i , indicates its i -th element. Matrices are denoted by uppercase bold letters, e.g., \mathbf{A} , with \mathbf{a}_i as the i -th *column* and $\mathbf{a}_{\bullet j}$ as the j -th *row*. $[k]$ is the set containing all integers from 1 up to k . $\llbracket \cdot \rrbracket$ is the Iverson bracket, which equals 1 if the condition is true and 0 otherwise. \odot and \oslash denote the element-wise multiplication and division, respectively.

Assumption 2.1 (Projection-based label distribution). Let $\mathbf{q} \in \prod_{j \in [c]} [a_j, b_j]_{\mathbb{R}}$ denote the raw data, where c is the number of labels and $[a_j, b_j]_{\mathbb{R}}$ is the value range of each label j . The corresponding LD $\mathbf{d} \in \Delta^{c-1}$ is obtained via a projection operator $\text{proj}(\cdot)$, i.e., $\mathbf{d} = \text{proj}(\mathbf{q})$, where

$$\Delta^{c-1} \triangleq \{\mathbf{v} \in \mathbb{R}^c \mid \mathbf{1}^\top \mathbf{v} = 1, \mathbf{v} \geq 0\} \quad (1)$$

is the $(c-1)$ -dimensional probability simplex and $\mathbf{1}$ is a c -dimensional vector of all ones. Each d_j of \mathbf{d} , namely *description degree*, indicates the degree to which the j -th label describes the sample.

Assumption 2.2 (Proportion-based label distribution). Building on Assumption 2.1, the raw data \mathbf{q} satisfies $\forall_{j \in [c]} (q_j \geq 0)$ and $\exists_{j \in [c]} (q_j > 0)$, then the LD is obtained by $\mathbf{d} = \mathbf{q} / \sum_{j \in [c]} q_j$.

Definition 2.3 (Logical label). The LL is a binary vector $\mathbf{l} \in \{0, 1\}^c$, where $l_j = 1$ indicates that the j -th label $y_j \in \mathcal{Y}$ is relevant to the sample, and $l_j = 0$ otherwise.

2.2 LACK OF FIDELITY TO THE RAW DATA

In this subsection, we show that the LD representation is not strictly consistent with the raw data, which can lead to substantial information distortion. To illustrate this issue, we consider the conversion from an LD to its corresponding LL as an example (Kou et al., 2024). Specifically, we introduce the subset error rate based on the definition of subset accuracy (Zhang & Zhou, 2013). Given two LLs, \mathbf{l}^* and \mathbf{l} , the subset error rate between them is defined as

$$\text{S.Err.}(\mathbf{l}^*, \mathbf{l}) \triangleq \llbracket \forall_{j \in [c]} ((\mathbf{l}^*)_j \neq (\mathbf{l})_j) \rrbracket. \quad (2)$$

According to Kou et al., there exist two binarization algorithms to convert an LD to its corresponding LL, i.e., top- k -based and τ -thresholding-based, denoted by $\text{bin}_1(\cdot; k)$ and $\text{bin}_2(\cdot; \tau)$, respectively, where $k \in [c]$ and $\tau \in [0, 1]_{\mathbb{R}}$ are hyperparameters. The two algorithms can be formulated as follows:

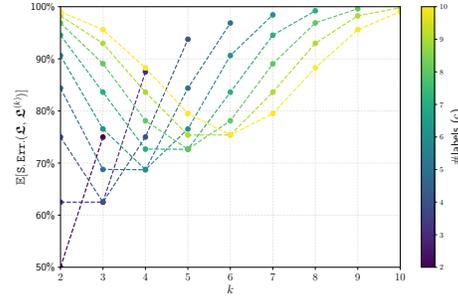
$$\text{bin}_1(\mathbf{d}; k) \triangleq \left\llbracket y_j \in \arg \text{top-}k(\mathbf{d}) \right\rrbracket_{y \in \mathcal{Y}}, \quad (3)$$

$$\text{bin}_2(\mathbf{d}; \tau) \triangleq \text{bin}_1(\mathbf{d}; k^*), \quad (4)$$

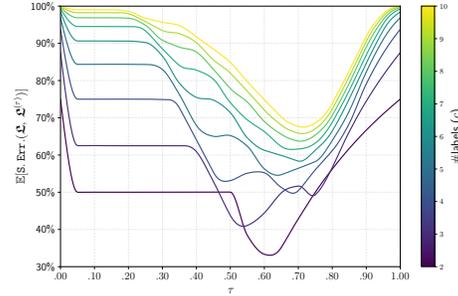
where $k^* = \text{card}(\min\{k \in [c] \mid \sum_{j \in [k]} \rho_j \geq \tau\})$, $\boldsymbol{\rho}$ is the sorted LD in *descending order*. For the top- k -based method, we have the following proposition.

Proposition 2.4. Let $\mathbf{L}^{(k)}$ denote the LL derived from the LD via the conversion rule in Equation (3); let \mathbf{L}^* denote the LL directly derived from the raw data random vector $\boldsymbol{\mathfrak{P}}$ by thresholding at τ' , i.e., $\mathbf{L}_j = \llbracket \boldsymbol{\mathfrak{P}}_j \geq \tau' \rrbracket$ for all $j \in [c]$. The raw data are independently drawn from a uniform distribution over $[a, b]_{\mathbb{R}}$, and the LD is obtained by its projection. Then, we have the following expectation:

$$\mathbb{E}[\text{S.Err.}(\mathbf{L}^*, \mathbf{L}^{(k)})] = 1 - \binom{c}{k} \left(\frac{b - \tau'}{b - a} \right)^k \left(\frac{\tau' - a}{b - a} \right)^{c-k} \geq 50\%. \quad (5)$$



(a) $\mathbb{E}[\text{S.Err.}(\mathbf{L}^*, \mathbf{L}^{(k)})]$ w.r.t. k & c .



(b) $\mathbb{E}[\text{S.Err.}(\mathbf{L}^*, \mathbf{L}^{(\tau)})]$ w.r.t. τ & c .

Figure 1: Visualization of the subset error rate between the LL derived from the raw data and that derived from the LD via two conversion methods: (a) top- k -based and (b) τ -thresholding-based.

The equality holds if and only if $c = 2$ and $k = 1$.

Proof sketch. \mathfrak{L}^* arises as a Bernoulli vector determined by the uniform thresholding, while $\mathfrak{L}^{(k)}$ corresponds to selecting the k largest order statistics, which directly yield the stated formula. \square

Corollary 2.5. *The following equation holds:*

$$\arg \min_k \mathbb{E}[\text{S. Err.}(\mathfrak{L}^*, \mathfrak{L}^{(k)})] = \left\lfloor \frac{c}{2} \right\rfloor. \quad (6)$$

Let $\mathfrak{L}^{(\tau)}$ denote the LL derived from the LD by Equation (4). For the τ -thresholding-based method, obtaining a closed-form expression of the expectation of the error rate is intractable. Therefore, we resort to Monte Carlo estimation under the similar setting in Proposition 2.4, with a sampling size of 10^6 and varying the number of labels $c \in [2, 10]_{\mathbb{Z}}$. The results are visualized in Figure 1b, where we also include the visualization for the top- k -based method for comparison, as shown in Figure 1a.

Remark 2.6. As established in Proposition 2.4 and corroborated by the empirical evidence in Figure 1, it is observed that the top- k -based method consistently incurs an error rate above 50%, even when using the recommended choice of k , i.e., $\lfloor c/2 \rfloor$ specified in Corollary 2.5; the τ -thresholding-based method achieves comparatively lower error rates, yet they remain exceeding 30%, with τ optimally selected within the interval $[0.5, 0.75]_{\mathbb{R}}$. For both methods, the minimum attainable error rate grows unfavorably with the number of labels c .

2.3 DISRUPTION OF INTER-SAMPLE ORDER

In this subsection, we show that LD can counterintuitively distort the order relations across samples, which is undesirable in many scenarios. To quantify this effect, we adopt Kendall’s tau, a widely used metric to measure the ordinal association between two sequences. Its type-A definition is as follows:

$$\text{Ken.}(u, v) \triangleq \frac{2 \sum_{i < j} \text{sign}(u_i - u_j) \text{sign}(v_i - v_j)}{n(n-1)}. \quad (7)$$

Proposition 2.7. *Assume that the raw data matrix $\mathbf{Q} \in [0, 1]_{\mathbb{R}}^{c \times n}$ has entries drawn independently from the uniform distribution. The corresponding LD matrix is obtained by $\mathbf{D} = \text{proj}(\mathbf{Q})$. Treating each row as a sequence across samples, the inter-sample order consistency can be quantified by*

$$\begin{aligned} \mathbb{E}_{\mathbf{Q}}[\text{Ken.}(\mathbf{q}_{\bullet}, \mathbf{d}_{\bullet})] &= \mathbb{E}_{\mathbf{Q}}[\text{sign}((\mathbf{q}_{\bullet} - \mathbf{q}'_{\bullet})^{\top} (\mathbf{q}_{\bullet} \otimes \mathbf{s} - \mathbf{q}'_{\bullet} \otimes \mathbf{s}))] \\ &= 2\mathbb{P}((\mathbf{q}_{\bullet} - \mathbf{q}'_{\bullet})^{\top} (\mathbf{q}_{\bullet} \otimes \mathbf{s} - \mathbf{q}'_{\bullet} \otimes \mathbf{s}) > 0) - 1, \end{aligned} \quad (8)$$

where \mathbf{q}_{\bullet} and \mathbf{q}'_{\bullet} are two independent sequences drawn from the same distribution as the rows of \mathbf{Q} , and $s_i = \sum_{j \in [c]} (\mathbf{q}_i)_j$ is the LD normalization factor for each sample i .

Proposition 2.7 follows directly since the $\text{sign}(\cdot)$ function always yields ± 1 in the absence of ties. For the special case $c = 2$, the probability term equals $3/4$ and the expected Kendall’s tau is 0.5. For $c > 2$, however, the expectation is analytically intractable due to the high-dimensional integrals involved. Therefore, we employ Monte Carlo estimation, fixing the number of samples at 10^6 and varying the number of labels $c \in [2, 10]_{\mathbb{Z}}$. The results are visualized in Figure 2.

Remark 2.8. As demonstrated in Proposition 2.7 and Figure 2, the expected Kendall’s tau increases with the number of labels c , yet it remains significantly below 1, indicating that LD can substantially distort the inter-sample order. Note that this apparent increase with c is primarily due to the uniform distribution spreading the values more evenly across labels, which reduces the likelihood of inconsistencies on average; values across different c may not be directly comparable.

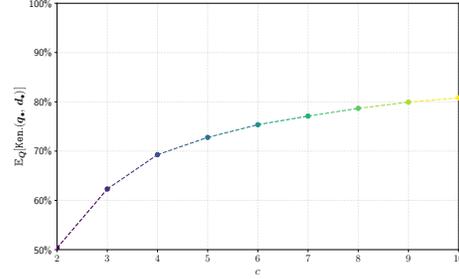


Figure 2: Visualization of Monte Carlo estimation of inter-sample order consistency, i.e., $\mathbb{E}_{\mathbf{Q}}[\text{Ken.}(\mathbf{q}_{\bullet}, \mathbf{d}_{\bullet})]$ w.r.t. c .

2.4 LIMITED PRACTICAL APPLICABILITY

In this subsection, we simply discuss the flaws of LD in practical applications using some examples.

Remark 2.9. Under Assumption 2.1, the practical limitations of LD can be summarized as follows:

- Inability to reconstruct the raw data: LD can *not* recover the original information, such as absolute magnitude and label annotation (Xu et al., 2020). For example, different colors (e.g., Magenta $\mathbf{q} = \langle 255, 0, 255 \rangle^\top$ and Patriarch $\mathbf{q}' = \langle 128, 0, 128 \rangle^\top$ in the RGB color space) may share the same LD ($\mathbf{d} = \langle 0.5, 0, 0.5 \rangle^\top$) in an LD color space.
- Lack of mechanism to handle OOD data: LD assumes that every label is relevant to the sample to some extent, which may *not* hold in real-world scenarios (Wu et al., 2025). For example, in an emotion recognition task, an image (e.g., a neutral face) may not convey any of the predefined emotions (e.g., the primary emotions like happiness, sadness, anger, fear, surprise, and disgust), but LD would still assign non-zero description degrees to all emotions.

Given the aforementioned limitations of LD, a more expressive and flexible representation is needed. This motivates the concept of generalized label distribution, which we present in the next section.

3 GENERALIZED LABEL DISTRIBUTION

3.1 FORMULATION OF GENERALIZED LABEL DISTRIBUTION

Remark 3.1. Assume a GLD mapping $\eta(\cdot)$. To address the limitations inherent in traditional LD, a GLD should satisfy the following properties:

- Regarding Remark 2.6, GLD should allow conversion to other forms of label representations, e.g., LD, LL, TL, LR, and SL, without any information loss. For example, in the case of LL, there should exist a mapping function $\text{GLD2LL}(\cdot)$ such that, for any LL \mathbf{l}^* derived from the raw data \mathbf{q} , the equation $\text{S.Err.}(\mathbf{l}^*, \text{GLD2LL}(\eta(\mathbf{q}))) = 0$ holds.
- Regarding Remark 2.8, GLD should preserve the inter-sample order, i.e., for the j -th row of any raw data matrix \mathbf{Q} , the equation $\text{Ken.}(\mathbf{q}_{\bullet j}, \eta(\mathbf{Q})_{\bullet j}) = 1$ holds.
- Regarding Remark 2.9, GLD should be capable of handling OOD data while providing a clear characterization of both positive and negative correlations; $\eta(\cdot)$ should be bijective to ensure that the raw data can be accurately reconstructed, i.e., $\mathbf{q} = \eta^{-1}(\eta(\mathbf{q}))$ for any raw data \mathbf{q} .

With these properties in mind, we present the formal definition of GLD as follows:

Definition 3.2 (Generalized label distribution). Given the raw data $\mathbf{q} \in \prod_{j \in [c]} [a_j, b_j]_{\mathbb{R}}$, the GLD is defined as

$$[-1, 1]_{\mathbb{R}}^c \ni \mathbf{g} = \eta(\mathbf{q}; \mathbf{a}, \mathbf{b}) \triangleq 2(\mathbf{q} - \mathbf{a}) \odot (\mathbf{b} - \mathbf{a}) - 1. \quad (9)$$

Conversely, the raw data can be recovered from the GLD by

$$\prod_{j \in [c]} [a_j, b_j]_{\mathbb{R}} \ni \mathbf{q} = \eta^{-1}(\mathbf{g}; \mathbf{a}, \mathbf{b}) \triangleq \frac{(\mathbf{g} + 1) \odot (\mathbf{b} - \mathbf{a})}{2} + \mathbf{a}. \quad (10)$$

Based on Definition 3.2, we have the following theorem regarding the expressiveness of GLD.

Theorem 3.3. Let $\mathfrak{P} \sim \mathcal{Q}$ be the random vector of the raw data, $\mathfrak{D} = \text{proj}(\mathfrak{P})$, and $\mathfrak{G} = \eta(\mathfrak{P})$. Under Assumption 2.1, we have

$$I(\mathfrak{P}; \mathfrak{D}) \leq I(\mathfrak{P}; \mathfrak{G}), \quad (11)$$

where $I(\cdot; \cdot)$ is the mutual information. The equality holds if $\mathcal{Q} \subseteq \psi(\Delta^{c-1})$, where $\psi(\cdot)$ is an invertible transformation and \mathcal{Q} is the distribution of the raw data.

Proof. See Appendix A.1 for details. \square

In practical, the mapping $\psi(\Delta^{c-1})$ can be a scaled simplex $\kappa\Delta^{c-1}$ ($\kappa > 0$), a.k.a. the *Aitchison simplex* (Aitchison, 1994), i.e., the raw data are compositional data constrained to a fixed sum κ . It should be noted that this is a rather stringent condition, rarely satisfied in general scenarios. Theorem 3.3 indicates that GLD provides a more expressive and faithful representation compared to conventional LD, supporting tasks that rely on nuanced label correlations.

Underlying philosophy The essence of GLD lies in the concept of *net probability* from a statistical perspective. It reflects the intrinsic association between samples and labels, analogous to *odds*, but within a more interpretable range that can express both positive and negative correlations. Whereas previous work in LD focused on the description degrees, i.e., answering “*To what extent does a label describe a sample?*”, GLD acts as a set of correlation coefficients, i.e., answering “*To what extent is a sample associated with a label?*”. In this sense, we name the values in GLD as *relative degrees*. Table 1 illustrates the conversion capabilities among various label ambiguity representations, showing that GLD can be mapped to all other forms seamlessly and losslessly. Additional insights of Table 1 are provided in Appendix B.

Table 1: Conversion capabilities among label ambiguity representations. \checkmark indicates that the conversion can be performed losslessly.

From	To					
	GLD	LD	TL	LL	LR	SL
GLD		\checkmark	\checkmark	\checkmark	\checkmark	\checkmark
LD					\checkmark	\checkmark
TL						
LL						
LR						\checkmark
SL						

Given the numerous benefits of GLD, one may wonder how to design a learning framework capable of effectively learning GLDs. This is the focus of the next subsection.

3.2 LEARNING GENERALIZED LABEL DISTRIBUTIONS

Problem formulation Let $\mathcal{X} = \mathbb{R}^m$ denote the input space, and $\mathcal{G} \subseteq [-1, 1]_{\mathbb{R}}^c$ the output space. Given the training set $\mathcal{S} = \{(\mathbf{x}_i, \mathbf{g}_i)\}_{i \in [n]}$, the goal is to find a mapping $f : \mathcal{X} \mapsto \mathcal{G}$, where $\mathbf{x}_i \in \mathcal{X}$ and $\mathbf{g}_i \in \mathcal{G}$ for all $i \in [n]$. The function f is not only required to accurately predict the GLD \mathbf{g} for any unseen instance \mathbf{x} , but also to perform well on other paradigms, e.g., LDL, MLL, etc.

Inspired by Geng, we discuss three strategies to learn GLDs, i.e., problem transformation, algorithm adaptation, and specialized algorithm. Note that these following methods are designed to facilitate fair comparisons with their respective reference models.

3.2.1 PROBLEM TRANSFORMATION & ALGORITHM ADAPTATION

Problem transformation The problem of learning GLDs can be reformulated as a series of constrained regression tasks, which bears similarity to classical methods in LDL and multi-target regression (Liu et al., 2009; Geng & Hou, 2015). Like them, we employ the support vector regression (SVR) algorithm as the underlying model. Note that we do *not* adopt the same problem transformation strategy used in (Geng, 2016), as it would compromise the intrinsic structure of the target data. To address the bounded nature of GLDs, we incorporate a forward-regressor-inverse transformation framework. Specifically, before training, the target values are first mapped to the real line via $\text{arctanh}(\text{clip}(\cdot; -1 + \varepsilon, 1 - \varepsilon))$, where $\text{clip}(\cdot; a, b)$ restricts the input to the specified range $[a, b]_{\mathbb{R}}$, and ε is a small positive constant to avoid numerical issues. After training, the predictions are mapped back to the original space using $\tanh(\cdot)$. The resulting method, referred to as GLD-SVR, is straightforward and intuitive.

Algorithm adaptation We can also naturally adapt some existing algorithms to deal with GLDs, among which we consider the k -nearest neighbor algorithm as an example, since there are already some extended variants for LDL and MLL (Zhang & Zhou, 2007; Geng, 2016). For a given test sample \mathbf{x} , we first determine the set $\mathcal{N}(\mathbf{x})$ of its k nearest neighbors from the training set. The GLD prediction is then obtained by averaging the GLDs of these neighbors, i.e., $1/k \sum_{\mathbf{x}_i \in \mathcal{N}(\mathbf{x})} \mathbf{g}_i$. The resulting method, referred to as GLD- k NN, is also simple yet effective.

3.2.2 SPECIALIZED ALGORITHM

Different from the above two strategies, we can also design specialized algorithms to directly match the GLD learning problem. Here, we consider a simple network-based model with a single linear layer followed by a $\tanh(\cdot)$ activation function, i.e., $f(\mathbf{x}) = \tanh(\mathbf{W}\mathbf{x})$, where $\mathbf{W} \in \mathbb{R}^{c \times m}$ is the weight matrix to be learned.

Remark 3.4. Consider the GLD model where each observation is assumed to be corrupted by multivariate Gaussian noise, i.e., $\mathbf{g}_i = f(\mathbf{x}_i) + \varepsilon_i$ with $\varepsilon_i \sim \text{Gaussian}(\mathbf{0}, \Sigma)$ for all $i \in [n]$, where $\Sigma \in \mathbb{R}^{c \times c}$ is the covariance matrix.

Let $\varphi(\mathbf{x}, \mathbf{g}; \Sigma)$ denote the probability density function, the corresponding negative log-likelihood can be expressed as

$$-\log \prod_{i \in [n]} \varphi(\mathbf{x}_i, \mathbf{g}_i; \Sigma) = \underbrace{\frac{nc}{2} \log(2\pi) + \frac{n}{2} \log \det(\Sigma)}_{\text{constant}} + \frac{1}{2} \sum_{i \in [n]} \underbrace{(\mathbf{g}_i - f(\mathbf{x}_i))^\top \Sigma^{-1} (\mathbf{g}_i - f(\mathbf{x}_i))}_{\text{squared Mahalanobis distance}}. \quad (12)$$

Equation (12) reveals that, apart from an additive constant, the negative log-likelihood naturally reduces to the sum of squared Mahalanobis distances between the predicted outputs and the observed GLDs. Consequently, adopting the Mahalanobis distance as a loss function is theoretically well justified. Formally, the loss function can be defined as follows:

$$\ell(\mathbf{W}; \{(\mathbf{x}_i, \mathbf{g}_i)\}_{i \in [n]}) = \frac{1}{n} \sum_{i \in [n]} (\mathbf{g}_i - \tanh(\mathbf{W} \mathbf{x}_i))^\top \Sigma^{-1} (\mathbf{g}_i - \tanh(\mathbf{W} \mathbf{x}_i)). \quad (13)$$

Optimization We can learn the optimal model parameters \mathbf{W}^* by minimizing the empirical loss, i.e., $\mathbf{W}^* = \arg \min_{\mathbf{W}} (\ell)$, which can be solved using gradient-based optimization methods.

Remark 3.5. ℓ is differentiable and its gradient w.r.t. \mathbf{W} is given by

$$\frac{\partial \ell(\mathbf{W})}{\partial \mathbf{W}} = -\frac{2}{n} \sum_{i \in [n]} ((\Sigma^{-1} (\mathbf{g}_i - \tanh(\mathbf{W} \mathbf{x}_i))) \odot \text{sech}^2(\mathbf{W} \mathbf{x}_i)) \mathbf{x}_i^\top, \quad (14)$$

derivation of which is provided in Appendix A.2, where $\text{sech}(\cdot)$ is the hyperbolic secant function.

In practice, Σ needs to be updated iteratively. However, in the early stages of training, the residuals are not yet reliable, hence retaining \mathbf{I} during the first half of the training process is advisable. Let t' denote the current iteration and t the maximum number of iterations. In the latter half of the training, we estimate Σ' using the Ledoit-Wolf method (Ledoit & Wolf, 2004) when $t' \equiv 0 \pmod{r}$, where r is a predefined frequency. The inverse covariance matrix is then computed as $\alpha(\Sigma')^{-1} + (1 - \alpha)\mathbf{I}$, where $\alpha = t'/t$ is an annealing factor. This helps avoid abrupt gradient fluctuations.

The resulting method is referred to as GLD-BFGS when employing the L-BFGS-B algorithm (Zhu et al., 1997) for optimization, which is a quasi-Newton method that approximates the Hessian matrix using gradient evaluations. Importantly, this specialized algorithm strategy is not only limited to GLD-BFGS itself, but also can be integrated into and adapted from existing LDL methods, extending their capability to support GLD learning. Further details are provided in Section 4.1.

Time complexity analysis Let \mathbf{X} denote the feature matrix, \mathbf{G} the GLD matrix. The overall time cost of GLD-BFGS is primarily influenced by the following calculations: the forward pass step $\tanh(\mathbf{W} \mathbf{X})$ requires $\mathcal{O}(nmc)$; the squared Mahalanobis distance computation in Equation (13), i.e., $(\mathbf{G} - \tanh(\mathbf{W} \mathbf{X})) \odot \Sigma^{-1} (\mathbf{G} - \tanh(\mathbf{W} \mathbf{X}))$, has a complexity of $\mathcal{O}(nc^2)$; the gradient calculation in Equation (14), vectorized as $\Sigma^{-1} ((\mathbf{G} - \tanh(\mathbf{W} \mathbf{X})) \odot (\mathbf{1} - \tanh^2(\mathbf{W} \mathbf{X}))) \mathbf{X}^\top$, involves a complexity of $\mathcal{O}(npc + nc^2)$. The calculation of Σ^{-1} has a complexity of $\mathcal{O}(nc^2)$, but this is performed infrequently. Therefore, the overall time complexity of GLD-BFGS is $\mathcal{O}(tnc^2 + tnpc)$.

3.3 THEORETICAL ANALYSIS

Theorem 3.6. Let \mathfrak{D}^* denote the underlying LD derived from the raw data random vector under Assumption 2.2; let \mathfrak{D}' and \mathfrak{D} be the predictions of GLD-kNN and LD-kNN, respectively. Suppose both \mathfrak{D} and \mathfrak{D}' are unbiased estimators of \mathfrak{D}^* . Then, the variance of \mathfrak{D}' with respect to \mathfrak{D}^* is no larger than that of \mathfrak{D} , i.e., $\mathbb{V}[\mathfrak{D}'_j - \mathfrak{D}^*_j] \leq \mathbb{V}[\mathfrak{D}_j - \mathfrak{D}^*_j]$ for all $j \in [c]$. The equality holds under the same condition as that required for the equality case in Theorem 3.3.

Theorem 3.7. Let \mathfrak{L}^* denote the underlying LL derived from the raw data random vector \mathfrak{P} by thresholding at τ' ; let \mathfrak{L}' and \mathfrak{L} be the predictions of GLD-kNN and LL-kNN, respectively. Suppose

both \mathfrak{L} and \mathfrak{L}' are unbiased estimators of \mathfrak{L}^* . Then, for the j -th label, $\mathbb{V}[\mathfrak{L}'_j - \mathfrak{L}^*_j] \leq \mathbb{V}[\mathfrak{L}_j - \mathfrak{L}^*_j]$ holds when $\sqrt{\mathbb{V}[\mathfrak{P}_j]} \leq 2|\tau'_j| \sqrt{\mathbb{P}(\mathfrak{P}_j \geq \tau'_j)(1 - \mathbb{P}(\mathfrak{P}_j \geq \tau'_j))}$ holds.

Theorem 3.8. Let $\hat{\mathfrak{R}}_n$ be the empirical Rademacher complexity w.r.t. \mathcal{S} with n samples; \mathcal{H} be a family of functions; and \mathcal{H}_j be the j -th component of the output of \mathcal{H} for all $j \in [c]$. For the squared Mahalanobis distance loss function $M^2(\cdot, \cdot)$ with covariance matrix Σ , we have

$$\hat{\mathfrak{R}}_n(M^2 \circ \tanh \circ \mathcal{H} \circ \mathcal{S}) \leq 4\sqrt{2c}\|\Sigma^{-1}\| \sum_{j \in [c]} \hat{\mathfrak{R}}_n(\mathcal{H}_j \circ \mathcal{S}). \quad (15)$$

Theorem 3.9. Define a family of functions \mathcal{H} and the j -th component of the output $\mathcal{H}_j = \{\mathbf{x} \mapsto \mathbf{w}_j \cdot \mathbf{x}\}$, where $\|\mathbf{w}_j\|_2 \leq \xi$ for all $j \in [c]$. Let \mathcal{F} be the family of functions for methods corresponding to Section 3.2.2, for any $\delta > 0$, with probability at least $1 - \delta$, for all $f \in \mathcal{F}$, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [M^2(\mathbf{g}_x, f(\mathbf{x}))] &\leq \frac{1}{n} \sum_{i \in [n]} M^2(\mathbf{g}_i, f(\mathbf{x}_i)) + 12c\|\Sigma^{-1}\| \sqrt{\frac{\log(2/\delta)}{2n}} + \\ &\frac{8\xi c\sqrt{2c}\|\Sigma^{-1}\|}{\sqrt{n}} \max_{i \in [n]} \|\mathbf{x}_i\|, \end{aligned} \quad (16)$$

where \mathcal{D} is the underlying distribution of the input space, and \mathbf{g}_x is the ground-truth GLD of \mathbf{x} .

Proof. Proofs of these above theorems are provided in Appendices A.3 to A.6, respectively. \square

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Datasets We construct an artificial dataset `artf` similar to (Geng, 2016), with only a limited number of training samples to simulate a challenging learning scenario; details are given in Appendix C. The real-world datasets have been examined in previous research (Spyromitros-Xioufis et al., 2016; González et al., 2021), including: a facial emotion recognition dataset of Japanese female faces, denoted as `jaf`; a geochemical composition dataset collected from the topsoil of the Swiss Jura region, denoted as `jura`; a water quality dataset `wq`; two datasets for price prediction, `atp` and `scm`; and a dataset on heating/cooling loads of buildings (i.e., energy efficiency), denoted as `enb`.

Evaluation metrics We do *not* directly evaluate the GLD prediction task itself, as its practical application is unexplored. Instead, we assess performance through metrics in these following versatile tasks, where \downarrow (\uparrow) indicates the smaller (larger) the better; see Appendix C for more details.

- OOD classification: For GLD, if all relative degree values are negative, the sample is regarded as an OOD instance. While LD cannot provide explicit OOD information, the situation where all labels are related to a sample is uncommon. Therefore, by leveraging the uniform distribution, we design a new metric `OOD Err.` \downarrow to make a relatively fair comparison.
- Intra/inter-sample ranking: The Spearman rank correlation coefficient `Spear.` \uparrow is commonly used for evaluating label ranking performance (Jia et al., 2023). In addition, we compute Kendall’s tau for each *row* of the predicted/ground-truth matrix to assess the consistency of inter-sample ordering, and denote this metric as `Ken.` \uparrow .
- LL prediction: We convert the model outputs into LLs and evaluate them using the Hamming distance `Ham.` \downarrow and subset accuracy `S. Acc.` \uparrow (Zhang & Zhou, 2013). For LD, the binarization algorithm `bin2` is adopted.
- LD prediction: We evaluate using the Clark distance `Clark` \downarrow (Geng, 2016) and a “regularized” K-L divergence (KLD) denoted as `μ_{KLD}` \uparrow (Li et al., 2025).

Baselines GLD-SVR, GLD- k NN, and GLD-BFGS in Section 3.2 correspond to baselines proposed in (Geng, 2016), which we denote as LD-SVR, LD- k NN, and LD-BFGS for clarity. In addition, we design the following GLD methods based on the specialized algorithm strategy:

- GLD-DF: This method is based on an ensemble method, denoted as LD-DF, which assigns different data batches to different label pairs (González et al., 2021). In our adaptation, we replace its base estimators with GLD-BFGS models, and substitute the built-in LD- k NN with GLD- k NN.

Table 2: Experimental results on datasets artf, jura, and wq formatted as (mean \pm std).

Algorithms	Clark \downarrow	$\mu_{\text{KLD}} \uparrow$	Ham. \downarrow	S. Acc. \uparrow	Spear. \uparrow	Ken. \uparrow	OOD Err. \downarrow	Avg. Rank
artf								
LD-SVR	.0307 \pm .005	<u>98.58%</u> \pm .005	• .2338 \pm .046	• 40.85% \pm .097	• .9640 \pm .028	• .5083 \pm .065	• 12.70% \pm .076	6.14
GLD-SVR	.0307 \pm .006	98.57% \pm .005	<u>.0143</u> \pm .014	<u>95.70%</u> \pm .041	.9763 \pm .022	.9502 \pm .014	<u>03.35%</u> \pm .041	<u>2.86</u>
LD-kNN	.0475 \pm .008	96.40% \pm .011	• .2340 \pm .047	• 40.65% \pm .102	.9333 \pm .036	• .4942 \pm .064	• 13.95% \pm .071	9.43
GLD-kNN	.0476 \pm .008	96.41% \pm .011	.0547 \pm .031	84.10% \pm .088	.9328 \pm .037	.9046 \pm .024	08.25% \pm .056	6.86
LD-BFGS	• .0783 \pm .010	• 90.90% \pm .022	• .2532 \pm .048	• 36.85% \pm .094	• .9470 \pm .036	• .5063 \pm .065	• 09.30% \pm .063	10.3
GLD-BFGS	.0444 \pm .009	96.55% \pm .010	.0670 \pm .030	80.40% \pm .088	.9565 \pm .030	<u>.9589</u> \pm .014	04.25% \pm .042	5.29
LD-DF	• .0492 \pm .007	• 96.38% \pm .010	• .2497 \pm .047	• 37.10% \pm .096	• .9200 \pm .046	• .5103 \pm .066	• 16.35% \pm .082	10.1
GLD-DF	.0237 \pm .004	99.05% \pm .003	.0118 \pm .013	96.45% \pm .038	<u>.9758</u> \pm .021	.9648 \pm .011	02.45% \pm .031	1.14
LD-LRR	• .0783 \pm .010	• 90.90% \pm .022	• .2532 \pm .048	• 36.85% \pm .094	• .9470 \pm .036	• .5062 \pm .065	• 09.30% \pm .063	10.4
GLD-LRR	.0444 \pm .009	96.55% \pm .010	.0670 \pm .030	80.40% \pm .088	.9565 \pm .030	.9589 \pm .014	04.25% \pm .042	5.14
LD-Delta	<u>.0285</u> \pm .012	98.41% \pm .014	• .2413 \pm .048	• 38.95% \pm .098	• .9683 \pm .027	• .5103 \pm .066	• 10.75% \pm .069	6.43
GLD-Delta	.0299 \pm .006	98.56% \pm .006	.0303 \pm .020	90.90% \pm .059	.9590 \pm .027	.9559 \pm .014	06.60% \pm .058	3.86
jura								
LD-SVR	.2709 \pm .027	• 91.90% \pm .018	• .4687 \pm .033	• 02.23% \pm .023	.9217 \pm .030	• .3160 \pm .053	• 61.78% \pm .090	9.29
GLD-SVR	.2665 \pm .025	92.55% \pm .016	.1137 \pm .027	68.27% \pm .069	.9235 \pm .030	.5419 \pm .063	17.16% \pm .059	5.71
LD-kNN	.2833 \pm .027	• 91.29% \pm .019	• .4752 \pm .036	• 03.26% \pm .025	.9177 \pm .031	• .2806 \pm .050	• 61.44% \pm .091	10.4
GLD-kNN	.2768 \pm .029	90.39% \pm .024	.1184 \pm .029	68.04% \pm .073	.9188 \pm .031	.4981 \pm .060	20.61% \pm .059	8.00
LD-BFGS	• .3009 \pm .031	92.60% \pm .015	• .4803 \pm .033	• 03.12% \pm .026	<u>.9348</u> \pm .029	• .3120 \pm .050	• 57.12% \pm .097	8.71
GLD-BFGS	.2530 \pm .024	92.86% \pm .016	.1155 \pm .027	68.38% \pm .070	.9301 \pm .030	<u>.5772</u> \pm .056	<u>15.43%</u> \pm .052	<u>4.29</u>
LD-DF	• .2562 \pm .025	93.56% \pm .014	• .4795 \pm .036	• 03.26% \pm .029	.9308 \pm .028	• .3149 \pm .050	• 59.05% \pm .095	7.43
GLD-DF	.2394 \pm .023	<u>93.67%</u> \pm .014	.1151 \pm .027	<u>69.42%</u> \pm .068	.9291 \pm .030	.6028 \pm .051	15.91% \pm .052	2.86
LD-LRR	.2595 \pm .025	• 93.59% \pm .014	• .4792 \pm .034	• 02.93% \pm .027	.9336 \pm .031	• .3125 \pm .050	• 59.21% \pm .097	7.57
GLD-LRR	.2733 \pm .029	90.21% \pm .029	.1130 \pm .026	68.91% \pm .070	.9189 \pm .070	.5553 \pm .099	15.38% \pm .051	5.86
LD-Delta	.2494 \pm .024	• 94.33% \pm .014	• .4720 \pm .037	• 03.71% \pm .031	• .9403 \pm .030	• .3264 \pm .052	• 58.96% \pm .093	5.00
GLD-Delta	<u>.2483</u> \pm .022	93.42% \pm .017	.1102 \pm .027	70.55% \pm .068	.9319 \pm .030	<u>.5950</u> \pm .054	18.08% \pm .057	2.86
wq								
LD-SVR	3.4504 \pm .030	• 00.48% \pm .006	• .1844 \pm .009	• 05.51% \pm .021	• .3400 \pm .021	• .0809 \pm .024	• 78.89% \pm .036	10.4
GLD-SVR	3.6075 \pm .010	00.77% \pm .006	.1649 \pm .010	<u>11.50%</u> \pm .032	.3643 \pm .021	.1110 \pm .022	70.23% \pm .046	8.00
LD-kNN	• <u>2.7350</u> \pm .038	20.28% \pm .025	• .2317 \pm .009	• 02.46% \pm .015	.3628 \pm .023	• .2388 \pm .022	• 75.84% \pm .036	7.57
GLD-kNN	2.7234 \pm .038	20.55% \pm .024	.1644 \pm .010	11.02% \pm .031	.3631 \pm .026	.2306 \pm .023	69.30% \pm .043	5.71
LD-BFGS	• 3.1244 \pm .032	• 27.57% \pm .014	• .2646 \pm .011	• 01.20% \pm .010	• .3901 \pm .022	• .2516 \pm .019	• 77.31% \pm .039	7.57
GLD-BFGS	3.1396 \pm .037	<u>25.83%</u> \pm .034	<u>.1640</u> \pm .010	<u>10.70%</u> \pm .030	<u>.3706</u> \pm .032	<u>.2142</u> \pm .030	<u>68.62%</u> \pm .047	<u>6.57</u>
LD-DF	• 3.1184 \pm .032	<u>29.36%</u> \pm .014	• .2556 \pm .009	• 01.58% \pm .011	• .4092 \pm .022	• .2730 \pm .019	• 77.18% \pm .040	5.43
GLD-DF	3.1278 \pm .031	29.07% \pm .015	.1623 \pm .010	10.99% \pm .029	.3985 \pm .023	.2508 \pm .018	<u>68.02%</u> \pm .048	<u>4.29</u>
LD-LRR	3.1235 \pm .032	• 27.41% \pm .014	• .2679 \pm .010	• 01.22% \pm .011	• .3893 \pm .022	• .2469 \pm .020	• 77.68% \pm .038	8.14
GLD-LRR	3.1272 \pm .033	26.81% \pm .014	.1650 \pm .010	10.61% \pm .028	.3679 \pm .023	.2183 \pm .018	68.76% \pm .045	6.71
LD-Delta	3.1277 \pm .031	• 30.24% \pm .017	• .2396 \pm .011	• 02.23% \pm .012	• .4162 \pm .022	• .2778 \pm .020	• 76.26% \pm .041	5.29
GLD-Delta	3.1221 \pm .032	29.10% \pm .016	.1607 \pm .010	11.87% \pm .031	.4025 \pm .024	.2540 \pm .019	66.33% \pm .044	2.29

- GLD-LRR: The reference method LD-LRR designs a loss function that incorporates label ranking relationships (Jia et al., 2023). We convert the predicted GLD into LD to make it compatible with this loss, and replace the original KLD with the squared Mahalanobis distance.
- GLD-Delta: The corresponding method, denoted as LD-Delta, optimizes a “regularized” KLD to ensure that the majority of samples are approximately predicted, thereby mitigating strong interference from outlier data (Li et al., 2025). To adapt it for GLD, both the loss and the worst-case expectation are computed using the squared Mahalanobis distance.

Methodology To ensure a fair comparison, for each dataset and for each method we conduct ten-fold experiments repeated 10 times, and the average performance is recorded. Representative results are in Table 2. The best and second-best results are highlighted in **bold** and underline, respectively. • (◦) indicates “GLD-X is statistically superior (inferior) to the comparing methods LD-X” (pairwise t -test at 0.05 significance level); if neither • nor ◦ is present, there is no significant difference.

4.2 DISCUSSION

The artificial dataset The first observation is that, across almost all evaluation metrics, the proposed GLD-based methods demonstrate clear advantages over their LD counterparts. In particular, the improvements are most pronounced in LL prediction, inter-sample ranking, and OOD classification.¹ These three aspects correspond directly to the limitations of LD discussed in Remarks 2.6,

¹It should be noted that the observed performance gap primarily stems from the representational limitations of LD itself. The conversion strategies (e.g., bin₂) have already been applied in the most effective way possible.

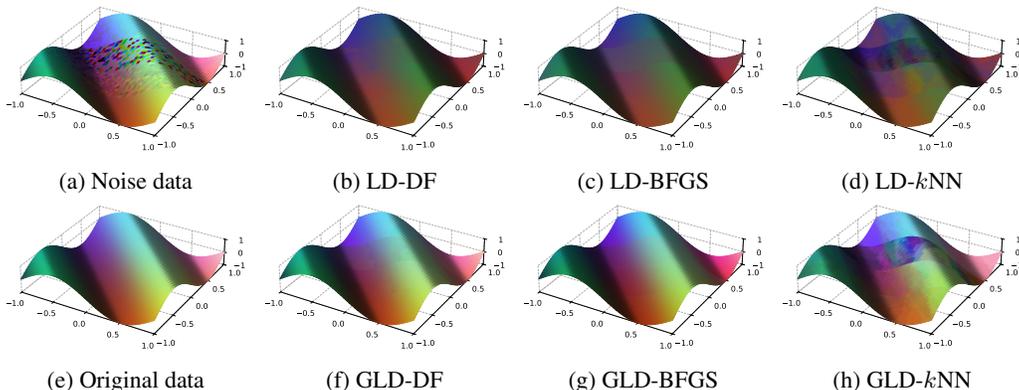


Figure 3: (Best viewed in color) Visualized results of the robustness testing.

2.8 and 2.9, respectively. GLD effectively overcomes these deficiencies, yielding substantial performance gains in these tasks, while at the same time maintaining high accuracy in LD prediction.

Secondly, regarding the adapted algorithms, although their relative performance on the LDL task can be roughly ordered as LD-Delta \succ LD-DF \succ LD-LRR, the ranking changes after our adaptation, becoming GLD-DF \succ GLD-Delta \succ GLD-LRR. A possible explanation is that we do not introduce any specific optimizations tailored for GLD prediction. As a result, the more general ensemble framework demonstrates stronger adaptability and delivers superior performance.

The real-world datasets In most cases, GLD methods still outperform their corresponding LD baseline methods. Regarding the methods in Section 3.2, although they achieve impressive results on the artificial dataset, their performance on real-world datasets is notably less satisfactory. For example, on the *artf* dataset, GLD-SVR achieves an average rank of 2.86 among 12 methods, whereas on *jura* and *wq* it only reaches 5.71 and 8.00, respectively. On the *wq* dataset, the performance of inter-sample ranking is poor, due to the discrete nature of the target data values, which leads to misjudgments in tie cases. The remaining experimental results are shown in Appendix C.

Robustness analysis We use the squared Mahalanobis distance to guide the learning of GLD, where the covariance matrix not only captures pairwise label correlations but also, to some extent, serves as a noise-robust regularization term. This is an aspect that needs to be verified. Therefore, we conduct robustness tests similar to (Li et al., 2025), exploiting the same artificial dataset mentioned above. The difference is that we use GLDs as the original data for GLD methods.² From bottom to top, from left to right, the four types of processing in Figure 3a are: (1) no treatment; (2) applying Gaussian noise (He et al., 2024); (3) randomly setting description degrees to zero (Xu & Zhou, 2017); (4) randomly emphasizing description degrees (Kou et al., 2023). These treatments correspond to different possible types of noise in label distributions. The prediction results of representative methods are shown in Figures 3b to 3d and 3f to 3h, and the ground truth, i.e., the original artificial data without noise, is shown in Figure 3e for comparison. The consistency of each prediction result reflects the robustness of each method. Figure 3 shows that methods like GLD-BFGS and GLD-DF predict accurately and the prediction results are not easily affected by noise.

5 LIMITATIONS & CONCLUSION

Limitations Despite these promising results, several limitations remain. In particular, the scenario when the semantics of the label space are continuous rather than discrete has not been explored. More applications of GLD are yet to be explored.

Conclusion In this paper, we systematically analyze the limitations of LD, and propose GLD as a more *unified*, *versatile*, and *faithful* representation of label ambiguity. GLD can recover raw data while preserving inter-sample order consistency, derive other label forms without information loss, and naturally capture out-of-distribution samples and negative label correlations.

²Differences in color depth of the results also show limitations of the LD representation (Remark 2.9).

REFERENCES

- 486
487
488 John Aitchison. Principles of compositional data analysis. *Lecture Notes-Monograph Series*, 24:
489 73–81, 1994.
- 490 Peter L. Bartlett and Shahar Mendelson. Rademacher and Gaussian complexities: Risk bounds and
491 structural results. *The Journal of Machine Learning Research*, 3:463–482, 2003.
- 492
493 Normand J. Beaudry and Renato Renner. An intuitive proof of the data processing inequality. *Quantum
494 Information and Computation*, 12(5-6):432–441, 2012.
- 495 Charles Bergeron, Jed Zaretski, Curt Breneman, and Kristin P Bennett. Multiple instance ranking.
496 In *Proceedings of the International Conference on Machine Learning*, pp. 48–55, 2008.
- 497
498 Klaus Brinker, Johannes Fürnkranz, and Eyke Hüllermeier. A unified model for multilabel classi-
499 fication and ranking. In *Proceedings of the European Conference on Artificial Intelligence*, pp.
500 489–493, 2006.
- 501 Elijah Cole, Oisín Mac Aodha, Titouan Lorieul, Pietro Perona, Dan Morris, and Nebojsa Jojic.
502 Multi-label learning from single positive labels. In *Proceedings of the IEEE/CVF Conference on
503 Computer Vision and Pattern Recognition*, pp. 933–942, 2021.
- 504
505 Stanislav Fort, Jie Ren, and Balaji Lakshminarayanan. Exploring the limits of out-of-distribution
506 detection. In *Proceedings of the International Conference on Neural Information Processing
507 Systems*, pp. 7068–7081, 2021.
- 508 Wei Gao and Zhi-Hua Zhou. Dropout Rademacher complexity of deep neural networks. *Science
509 China Information Sciences*, 59(7):072104, 2016.
- 510
511 Xin Geng. Label distribution learning. *IEEE Transactions on Knowledge and Data Engineering*, 28
512 (7):1734–1748, 2016.
- 513
514 Xin Geng and Peng Hou. Pre-release prediction of crowd opinion on movies by label distribution
515 learning. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pp.
516 3511–3517, 2015.
- 517 Manuel González, Germán González-Almagro, Isaac Triguero, José-Ramón Cano, and Salvador
518 García. Decomposition-fusion for label distribution learning. *Information Fusion*, 66:64–75,
519 2021.
- 520 Liang He, Yunan Lu, Weiwei Li, and Xiuyi Jia. Generative calibration of inaccurate annotation for
521 label distribution learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp.
522 12394–12401, 2024.
- 523
524 Xiuyi Jia, Xiaoxia Shen, Weiwei Li, Yunan Lu, and Jihua Zhu. Label distribution learning by
525 maintaining label ranking relation. *IEEE Transactions on Knowledge and Data Engineering*, 35
526 (02):1695–1707, 2023.
- 527 Yuheng Jia, Jiawei Tang, and Jiahao Jiang. Label distribution learning from logical label. In *Pro-
528 ceedings of the International Joint Conference on Artificial Intelligence*, pp. 4228–4236, 2024.
- 529
530 Zhiqiang Kou, Jing Wang, Yuheng Jia, Biao Liu, and Xin Geng. Instance-dependent inaccurate
531 label distribution learning. *IEEE Transactions on Neural Networks and Learning Systems*, 36(1):
532 1425–1437, 2023.
- 533
534 Zhiqiang Kou, Jing Wang, Jiawei Tang, Yuheng Jia, Boyu Shi, and Xin Geng. Exploiting multi-label
535 correlation in label distribution learning. In *Proceedings of the International Joint Conference on
Artificial Intelligence*, pp. 4326–4334, 2024.
- 536
537 Olivier Ledoit and Michael Wolf. A well-conditioned estimator for large-dimensional covariance
538 matrices. *Journal of Multivariate Analysis*, 88(2):365–411, 2004.
- 539
Weiwei Li, Haitao Wu, Yunan Lu, and Xiuyi Jia. Approximately correct label distribution learning.
In *Proceedings of the International Conference on Machine Learning*, 2025. In press.

- 540 Guangan Liu, Zhouchen Lin, and Yong Yu. Multi-output regression on the output manifold. *Pattern*
541 *Recognition*, 42(11):2737–2743, 2009.
- 542 Xinyuan Liu, Jihua Zhu, Zhongyu Li, Zhiqiang Tian, Xiuyi Jia, and Lei Chen. Unified framework
543 for learning with label distribution. *Information Fusion*, 75:116–130, 2021.
- 544 Yunan Lu and Xiuyi Jia. Predicting label distribution from multi-label ranking. In *Proceedings of*
545 *the International Conference on Neural Information Processing Systems*, pp. 36931–36943, 2022.
- 546 Yunan Lu and Xiuyi Jia. Predicting label distribution from ternary labels. In *Proceedings of the*
547 *International Conference on Neural Information Processing Systems*, pp. 70431–70452, 2024.
- 548 Andreas Maurer. A vector-contraction inequality for Rademacher complexities. In *Proceedings of*
549 *the International Conference on Algorithmic Learning Theory*, pp. 3–17, 2016.
- 550 Mehryar Mohri, Afshin Rostamizadeh, and Ameeet Talwalkar. *Foundations of Machine Learning*.
551 The MIT Press, 2012.
- 552 Jie Ren, Peter J Liu, Emily Fertig, Jasper Snoek, Ryan Poplin, Mark A DePristo, Joshua V Dillon,
553 and Balaji Lakshminarayanan. Likelihood ratios for out-of-distribution detection. In *Proceedings*
554 *of the International Conference on Neural Information Processing Systems*, pp. 14707–14718,
555 2019.
- 556 Eleftherios Spyromitros-Xioufis, Grigorios Tsoumakas, William Groves, and Ioannis Vlahavas.
557 Multi-target regression via input space expansion: Treating targets as inputs. *Machine Learn-*
558 *ing*, 104(1):55–98, 2016.
- 559 Grigorios Tsoumakas, Ioannis Katakis, and Ioannis Vlahavas. Random k -labelsets for multilabel
560 classification. *IEEE Transactions on Knowledge and Data Engineering*, 23(7):1079–1089, 2010.
- 561 Jing Wang and Xin Geng. Classification with label distribution learning. In *Proceedings of the*
562 *International Joint Conference on Artificial Intelligence*, pp. 3712–3718, 2019a.
- 563 Jing Wang and Xin Geng. Theoretical analysis of label distribution learning. In *Proceedings of the*
564 *AAAI Conference on Artificial Intelligence*, pp. 5256–5263, 2019b.
- 565 Jing Wang, Xin Geng, and Hui Xue. Re-weighting large margin label distribution learning for
566 classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):5445–
567 5459, 2021.
- 568 Haitao Wu, Weiwei Li, and Xiuyi Jia. Divide and conquer: Learning label distribution with subtasks.
569 In *Proceedings of the International Conference on Machine Learning*, 2025. In press.
- 570 Ming-Kun Xie and Sheng-Jun Huang. Partial multi-label learning. In *Proceedings of the AAAI*
571 *Conference on Artificial Intelligence*, pp. 4302–4309, 2018.
- 572 Donna Xu, Yaxin Shi, Ivor W Tsang, Yew-Soon Ong, Chen Gong, and Xiaobo Shen. Survey on
573 multi-output learning. *IEEE Transactions on Neural Networks and Learning Systems*, 31(7):
574 2409–2429, 2020.
- 575 Miao Xu and Zhi-Hua Zhou. Incomplete label distribution learning. In *Proceedings of the 26th*
576 *International Joint Conference on Artificial Intelligence*, pp. 3175–3181, 2017.
- 577 Ning Xu, Yunpeng Liu, and Xin Geng. Label enhancement for label distribution learning. *IEEE*
578 *Transactions on Knowledge and Data Engineering*, 33(4):1632–1643, 2019.
- 579 Min-Ling Zhang and Zhi-Hua Zhou. ML- k NN: A lazy learning approach to multi-label learning.
580 *Pattern recognition*, 40(7):2038–2048, 2007.
- 581 Min-Ling Zhang and Zhi-Hua Zhou. A review on multi-label learning algorithms. *IEEE Transac-*
582 *tions on Knowledge and Data Engineering*, 26(8):1819–1837, 2013.
- 583 Ciyou Zhu, Richard H Byrd, Pei Huang Lu, and Jorge Nocedal. Algorithm 778: L-BFGS-B: Fortran
584 subroutines for large-scale bound-constrained optimization. *ACM Transactions on mathematical*
585 *software*, 23(4):550–560, 1997.

594 A PROOFS AND DERIVATIONS

595 Here, we provide the detailed proofs and derivations omitted in the main paper.

596 A.1 PROOF OF THEOREM 3.3

597 *Proof.* We first introduce the data processing inequality.

598 **Lemma A.1** (Beaudry & Renner (2012)). *Let \mathfrak{A} , \mathfrak{B} , and \mathfrak{C} be random variables. If $\mathfrak{A} \rightarrow \mathfrak{B} \rightarrow \mathfrak{C}$ is*
 599 *a Markov chain, then we have*

$$600 I(\mathfrak{A}; \mathfrak{C}) \leq I(\mathfrak{A}; \mathfrak{B}). \quad (17)$$

601 *The equality holds if $\mathfrak{A} \rightarrow \mathfrak{C} \rightarrow \mathfrak{B}$ is also a Markov chain.*

602 According to Assumption 2.1 and Definition 3.2, we have

$$603 \mathfrak{G} \xrightarrow{\eta^{-1}} \mathfrak{P} \xrightarrow{\text{proj}} \mathfrak{D}. \quad (18)$$

604 According to Lemma A.1, we have

$$605 I(\mathfrak{G}; \mathfrak{D}) \leq I(\mathfrak{G}; \mathfrak{P}). \quad (19)$$

606 $\eta^{-1}(\cdot)$ is bijective. The mutual information is invariant under bijective transformations, i.e.,

$$607 I(\mathfrak{G}; \mathfrak{D}) = I(\mathfrak{P}; \mathfrak{D}). \quad (20)$$

608 The mutual information is symmetric, i.e.,

$$609 I(\mathfrak{G}; \mathfrak{P}) = I(\mathfrak{P}; \mathfrak{G}). \quad (21)$$

610 Combining Equations (19) to (21) gives the desired inequality.

611 If there exists an equality case, $\mathfrak{G} \rightarrow \mathfrak{D} \rightarrow \mathfrak{P}$ is also a Markov chain. Note that $\mathfrak{G} \rightarrow \mathfrak{D}$ has
 612 already been guaranteed by Equation (18). If $\mathfrak{D} \rightarrow \mathfrak{P}$ holds, then $\exists_{\psi}(\forall_{\mathbf{q} \in \mathcal{Q}}(\exists_{\mathbf{d} \in \Delta^{c-1}}(\mathbf{q} = \psi(\mathbf{d}))))$
 613 is a tautology, where $\psi(\cdot)$ is invertible since \mathbf{d} can be determined by \mathbf{q} . The above condition is
 614 equivalent to $\forall_{\mathbf{q} \in \mathcal{Q}}(\mathbf{q} \in \psi(\Delta^{c-1}))$, i.e., $\mathcal{Q} \subseteq \psi(\Delta^{c-1})$. \square

615 A.2 DERIVATION OF EQUATION (14)

616 Let $\mathbf{h}_i = \mathbf{W}\mathbf{x}_i$, $\tilde{\mathbf{g}}_i = \tanh(\mathbf{h}_i)$, and $\mathbf{e}_i = \mathbf{g}_i - \tilde{\mathbf{g}}_i$ for all $i \in [n]$. The quadratic term in the loss is
 617 $\mathbf{e}_i^\top \Sigma^{-1} \mathbf{e}_i$, where Σ^{-1} is symmetric. Its gradient with respect to \mathbf{e}_i is

$$618 \frac{\partial \mathbf{e}_i^\top \Sigma^{-1} \mathbf{e}_i}{\partial \mathbf{e}_i} = (\Sigma^{-1} + (\Sigma^{-1})^\top) \mathbf{e}_i = 2\Sigma^{-1} \mathbf{e}_i. \quad (22)$$

619 Since $\tanh(\cdot)$ is applied element-wise, by the chain rule, we have

$$620 \frac{\partial \mathbf{e}_i}{\partial \tilde{\mathbf{g}}_i} \frac{\partial \tilde{\mathbf{g}}_i}{\partial \mathbf{h}_i} = -\text{diag}(\text{sech}^2(\mathbf{h}_i)). \quad (23)$$

621 Furthermore, the derivative of the loss with respect to the weight matrix \mathbf{W} can be expressed as

$$622 \frac{\partial \ell}{\partial \mathbf{W}} = \frac{1}{n} \sum_{i \in [n]} \left(\frac{\partial \ell}{\partial \mathbf{h}_i} \right) \mathbf{x}_i^\top. \quad (24)$$

623 Combining Equations (22) to (24), we arrive at the final expression:

$$624 \frac{\partial \ell}{\partial \mathbf{W}} = -\frac{2}{n} \sum_{i \in [n]} ((\Sigma^{-1}(\mathbf{g}_i - \tanh(\mathbf{W}\mathbf{x}_i))) \odot \text{sech}^2(\mathbf{W}\mathbf{x}_i)) \mathbf{x}_i^\top. \quad (25)$$

A.3 PROOF OF THEOREM 3.6

Proof. Under Assumption 2.2, the LD prediction of LD- k NN for an unknown sample $\tilde{\mathbf{x}}$ is given by:

$$\tilde{\mathbf{d}} = \frac{1}{k} \sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} \mathbf{d}_i = \frac{1}{k} \sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} \frac{\mathbf{q}_i}{\sum_{j \in [c]} (\mathbf{q}_i)_j}. \quad (26)$$

The LD prediction of GLD- k NN for $\tilde{\mathbf{x}}$ is given by:

$$\tilde{\mathbf{d}}' = \frac{\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} \mathbf{q}_i}{\sum_{j \in [c]} (\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} \mathbf{q}_i)_j}. \quad (27)$$

Define $s_i = \sum_{j \in [c]} (\mathbf{q}_i)_j > 0$ for all $i \in [n]$. Assume that each observed raw vector \mathbf{q}_i can be decomposed into a scaled ground-truth distribution \mathbf{d}^* perturbed by noise:

$$\mathbf{q}_i = s_i \mathbf{d}^* + \boldsymbol{\varepsilon}_i, \quad (28)$$

where $\boldsymbol{\varepsilon}_i$ denotes the noise term. Thus,

$$\mathbf{d}_i = \frac{\mathbf{q}_i}{s_i} = \mathbf{d}^* + \frac{\boldsymbol{\varepsilon}_i}{s_i}. \quad (29)$$

For LD- k NN, we have

$$\tilde{\mathbf{d}} - \mathbf{d}^* = \frac{1}{k} \sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} \left(\mathbf{d}^* + \frac{\boldsymbol{\varepsilon}_i}{s_i} \right) - \mathbf{d}^* = \frac{1}{k} \sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} \frac{\boldsymbol{\varepsilon}_i}{s_i}. \quad (30)$$

For GLD- k NN, we have

$$\tilde{\mathbf{d}}' - \mathbf{d}^* = \frac{\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} (s_i \mathbf{d}^* + \boldsymbol{\varepsilon}_i)}{\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} s_i} - \mathbf{d}^* = \frac{\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} \boldsymbol{\varepsilon}_i}{\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} s_i}. \quad (31)$$

Assume that the noise vectors $\boldsymbol{\varepsilon}_i$ are independent with zero mean and isotropic covariance, i.e., $\mathbb{E}[\boldsymbol{\varepsilon}_i] = \mathbf{0}$ and $\mathbb{V}[\boldsymbol{\varepsilon}_i] = \sigma^2 \mathbf{I}$. The variance of LD- k NN is

$$\mathbb{V}[d_j - d_j^*] = \frac{\sigma^2}{k^2} \sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} \frac{1}{s_i^2}. \quad (32)$$

The variance of GLD- k NN is

$$\mathbb{V}[d_j' - d_j^*] = \frac{k\sigma^2}{\left(\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} s_i \right)^2}. \quad (33)$$

By applying the QM-HM inequality to the sequence $1/s_i$, we have

$$\frac{k^2}{\left(\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} s_i \right)^2} \leq \frac{\left(\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} 1/s_i \right)^2}{k}. \quad (34)$$

By applying the AM-QM inequality to the sequence $1/s_i$, we have

$$\frac{\left(\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} 1/s_i\right)^2}{k^2} \leq \frac{\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} 1/s_i^2}{k}. \quad (35)$$

Combining Equations (34) and (35) yields

$$\frac{k}{\left(\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} s_i\right)^2} \leq \frac{1}{k^2} \sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} \frac{1}{s_i^2}. \quad (36)$$

Substituting Equations (32) and (33) into this relation directly leads to

$$\mathbb{V}[d'_j - d_j^*] \leq \mathbb{V}[d_j - d_j^*]. \quad (37)$$

The equality holds when $\forall \mathbf{x}_i, \mathbf{x}_j \in \mathcal{N}(\tilde{\mathbf{x}}) (s_i = s_j)$ is a tautology.

□

A.4 PROOF OF THEOREM 3.7

Proof. The LL prediction of LL- k NN for an unknown sample $\tilde{\mathbf{x}}$ is given by:

$$\begin{aligned} \tilde{\mathbf{l}} &= \left[\frac{1}{k} \left(\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} l_i \right) \geq \frac{1}{2} \right] \\ &= \left[\frac{1}{k} \left(\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} \mathbb{[q}_i \geq \tau'] \right) \geq \frac{1}{2} \right] \\ &= \left[\left(\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} 2\tau' \odot \mathbb{[q}_i \geq \tau'] \right) \geq k\tau' \right]. \end{aligned} \quad (38)$$

The LL prediction of GLD- k NN for $\tilde{\mathbf{x}}$ is given by:

$$\begin{aligned} \tilde{\mathbf{l}}' &= \left[\frac{1}{k} \left(\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} \mathbf{g}_i \right) \geq \mathbf{0} \right] \\ &= \left[\frac{1}{k} \left(\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} \mathbf{q}_i \right) \geq \tau' \right] \\ &= \left[\left(\sum_{\mathbf{x}_i \in \mathcal{N}(\tilde{\mathbf{x}})} \mathbf{q}_i \right) \geq k\tau' \right]. \end{aligned} \quad (39)$$

The only distinction between the two methods lies in how the target vectors of the k nearest neighbors are aggregated prior to thresholding. If $\mathbb{V}[l'_j - l_j^*] \leq \mathbb{V}[l_j - l_j^*]$, then $\mathbb{V}[q_j] \leq \mathbb{V}[2\tau'_j \mathbb{[q}_j \geq \tau'_j]]$, i.e.,

$$\begin{aligned} \mathbb{V}[q_j] &\leq 4(\tau'_j)^2 \mathbb{P}(q_j \geq \tau'_j)(1 - \mathbb{P}(q_j \geq \tau'_j)) \\ \sqrt{\mathbb{V}[q_j]} &\leq 2|\tau'_j| \sqrt{\mathbb{P}(q_j \geq \tau'_j)(1 - \mathbb{P}(q_j \geq \tau'_j))}. \end{aligned} \quad (40)$$

□

A.5 PROOF OF THEOREM 3.8

Proof. Define $\phi(\mathbf{u}, \mathbf{v})$ as $M^2(\mathbf{u}, \tanh(\mathbf{v}))$. Let $\mathbf{y}_1 = \tanh(\mathbf{v}_1)$ and $\mathbf{y}_2 = \tanh(\mathbf{v}_2)$. Then, we have

$$\begin{aligned} |\phi(\mathbf{u}, \mathbf{v}_1) - \phi(\mathbf{u}, \mathbf{v}_2)| &= |\mathbf{M}(\mathbf{u}, \mathbf{y}_1) - \mathbf{M}(\mathbf{u}, \mathbf{y}_2)| \\ &= |(\mathbf{u} - \mathbf{y}_1)^\top \boldsymbol{\Sigma}^{-1}(\mathbf{u} - \mathbf{y}_1) - (\mathbf{u} - \mathbf{y}_2)^\top \boldsymbol{\Sigma}^{-1}(\mathbf{u} - \mathbf{y}_2)| \\ &= |((\mathbf{u} - \mathbf{y}_1) + (\mathbf{u} - \mathbf{y}_2))^\top \boldsymbol{\Sigma}^{-1}((\mathbf{u} - \mathbf{y}_1) - (\mathbf{u} - \mathbf{y}_2))| \\ &= |(2\mathbf{u} - \mathbf{y}_1 - \mathbf{y}_2)^\top \boldsymbol{\Sigma}^{-1}(\mathbf{y}_2 - \mathbf{y}_1)|. \end{aligned} \quad (41)$$

According to Cauchy-Schwarz inequality,

$$|(2\mathbf{u} - \mathbf{y}_1 - \mathbf{y}_2)^\top \boldsymbol{\Sigma}^{-1}(\mathbf{y}_2 - \mathbf{y}_1)| \leq \|2\mathbf{u} - \mathbf{y}_1 - \mathbf{y}_2\| \cdot \|\boldsymbol{\Sigma}^{-1}\| \cdot \|\mathbf{y}_2 - \mathbf{y}_1\|. \quad (42)$$

And we have

$$\|2\mathbf{u} - \mathbf{y}_1 - \mathbf{y}_2\| \leq 2\|\mathbf{u}\| + \|\mathbf{y}_1\| + \|\mathbf{y}_2\| = 4\sqrt{c}. \quad (43)$$

Combine Equations (41) to (43), which shows that $\phi(\mathbf{u}, \cdot)$ is L -Lipschitz, i.e., for $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^c$,

$$|\phi(\mathbf{u}, \mathbf{v}_1) - \phi(\mathbf{u}, \mathbf{v}_2)| \leq L \|\tanh(\mathbf{v}_1) - \tanh(\mathbf{v}_2)\|, \quad (44)$$

where $L = 4\sqrt{c}\|\boldsymbol{\Sigma}^{-1}\|$. The subsequent proof follows a process similar to (Wang & Geng, 2019b). Recall the definition of Rademacher complexity, i.e.,

$$\hat{\mathfrak{R}}_n(M^2 \circ \tanh \circ \mathcal{H} \circ \mathcal{S}) = \mathbb{E} \left[\sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i \in [n]} \phi(\mathbf{g}_i, h(\mathbf{x}_i)) \nu_i \right], \quad (45)$$

where ν_i are n i.i.d. rademacher random variables with $\mathbb{P}(\nu_i = 1) = \mathbb{P}(\nu_i = -1) = 1/2$ for all $i \in [n]$. And according to (Maurer, 2016), with $\phi(\mathbf{u}, \cdot)$ being $4\sqrt{c}\|\boldsymbol{\Sigma}^{-1}\|$ -Lipschitz, we have

$$\mathbb{E} \left[\sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i \in [n]} \phi(\mathbf{g}_i, h(\mathbf{x}_i)) \nu_i \right] \leq 4\sqrt{2c}\|\boldsymbol{\Sigma}^{-1}\| \mathbb{E} \left[\sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i \in [n]} \sum_{j \in [c]} \nu_{i,j} h_j(\mathbf{x}_i) \right], \quad (46)$$

where $h_j(\cdot)$ is the j -th component of $h(\cdot)$ and $\nu_{i,j}$ are $n \times c$ i.i.d. rademacher random variables. Suppose all \mathcal{H}_j s be classes of functions for network-based architecture, then $\mathcal{H} = \oplus_{j \in [c]} \mathcal{H}_j = \{\mathbf{x} \mapsto h(\mathbf{x})\}$, and we have

$$\begin{aligned} \mathbb{E} \left[\sup_{h \in \oplus_{j \in [c]} \mathcal{H}_j} \frac{1}{n} \sum_{i \in [n]} \sum_{j \in [c]} \nu_{i,j} h_j(\mathbf{x}_i) \right] &\leq \sum_{j \in [c]} \mathbb{E} \left[\sup_{h_j \in \mathcal{H}_j} \frac{1}{n} \sum_{i \in [n]} \nu_{i,j} h_j(\mathbf{x}_i) \right] \\ &\leq \sum_{j \in [c]} \hat{\mathfrak{R}}_n(\mathcal{H}_j \circ \mathcal{S}). \end{aligned} \quad (47)$$

Combine Equations (45) to (47), which finishes proof of Theorem 3.8. \square

A.6 PROOF OF THEOREM 3.9

Proof. We first introduce the following two lemmas.

Lemma A.2 (Bartlett & Mendelson (2003); Mohri et al. (2012)). *Let \mathcal{H} be a family of functions. For a loss function ℓ bounded by μ , then for any $\delta > 0$, with probability at least $1 - \delta$, for all $h \in \mathcal{H}$ such that*

$$\vartheta_{\mathcal{D}}(h) \leq \hat{\vartheta}_{\mathcal{S}}(h) + 2\hat{\mathfrak{R}}(\ell \circ \mathcal{H} \circ \mathcal{S}) + 3\mu\sqrt{\frac{\log(2/\delta)}{2n}}. \quad (48)$$

Lemma A.3 (Bartlett & Mendelson (2003); Gao & Zhou (2016)). *Define class of functions $\mathcal{H}_j = \{\mathbf{x} \mapsto \mathbf{w}_j \cdot \mathbf{x}\}$, where $\|\mathbf{w}_j\|_2 \leq \xi$ for all $j \in [c]$. We have*

$$\hat{\mathfrak{R}}_n(\mathcal{H}_j \circ \mathcal{S}) \leq \frac{\xi \max_{i \in [n]} \|\mathbf{x}_i\|}{\sqrt{n}}. \quad (49)$$

Substituting Theorem 3.8 and Lemma A.3 into Lemma A.2, and Theorem 3.9 follows directly. \square

B MORE INSIGHTS ABOUT TABLE 1

A more detailed version of Table 1 is presented in Table 3.

Table 3: Conversion capabilities among label ambiguity representations.

From	To					
	GLD	LD	TL	LL	LR	SL
GLD	† ₁	Eq. (50)	Eq. (51)	Eq. (52)	Eq. (53)	Eq. (54)
LD		† ₂	† ₃	† ₄	Eq. (55)	Eq. (56)
TL		† ₅		† ₆	† ₇	† ₈
LL		† ₉		† ₁₀	† ₁₁	† ₁₂
LR		† ₁₃				Eq. (57)
SL				† ₁₄		† ₁₅

GLD can be transformed into other label ambiguity representations through Equations (50) to (54). The conversion from GLD to LD relies on the capability of GLD to recover the raw data:

$$\text{GLD2LD}(\mathbf{g}) = \text{proj} \left(\frac{(\mathbf{g} + \mathbf{1}) \odot (\mathbf{b} - \mathbf{a})}{2} + \mathbf{a} \right). \quad (50)$$

Definition B.1 (Ternary label). The TL is a ternary vector $\mathbf{t} \in \{-1, 0, 1\}^c$, where $t_j = 1$ indicates that the j -th label $y_j \in \mathcal{Y}$ is positively associated with the sample, $t_j = -1$ indicates a negative association, and $t_j = 0$ means the label is irrelevant.

Based on Definition B.1 and Definition 2.3, GLD can be naturally converted into TL and LL by

$$\text{GLD2TL}(\mathbf{g}; \boldsymbol{\tau}^{(+)}, \boldsymbol{\tau}^{(-)})_j = \begin{cases} \llbracket g_j > \tau_j^{(+)} \rrbracket, & g_j \geq \tau_j^{(-)}, \\ -1, & g_j < \tau_j^{(-)}, \end{cases} \quad (51)$$

$$\text{GLD2LL}(\mathbf{g}; \boldsymbol{\tau})_j = \llbracket g_j \geq \tau_j \rrbracket, \quad (52)$$

respectively. To ensure consistency with the raw data, the threshold vectors are generally set as $\boldsymbol{\tau}^{(+)} = 1/3\mathbf{1}$, $\boldsymbol{\tau}^{(-)} = -1/3\mathbf{1}$, and $\boldsymbol{\tau} = \mathbf{0}$. Next, we define LR in a simplified form.

Definition B.2 (Label ranking). The LR is defined as a label vector $\mathbf{r} \in \mathcal{Y}^c$, where each element corresponds to a label $y \in \mathcal{Y}$ and the order reflects their relative preference or relevance to the sample, and $r_j \neq r_k$ for all $j \neq k$.

Both GLD and LD can then be converted into LR and SL via argsort and arg max , since the necessary information is preserved in their respective data structures:

$$\text{GLD2LR}(\mathbf{g}) = \underset{y \in \mathcal{Y}}{\text{argsort}}(\mathbf{g}), \quad (53)$$

$$\text{GLD2SL}(\mathbf{g}) = \underset{y \in \mathcal{Y}}{\text{arg max}}(\mathbf{g}), \quad (54)$$

$$\text{LD2LR}(\mathbf{d}) = \underset{y \in \mathcal{Y}}{\text{argsort}}(\mathbf{d}), \quad (55)$$

$$\text{LD2SL}(\mathbf{d}) = \underset{y \in \mathcal{Y}}{\text{arg max}}(\mathbf{d}). \quad (56)$$

Moreover, LR can be further converted into SL by selecting the top-ranked label:

$$\text{LR2SL}(\mathbf{r}) = r_1. \quad (57)$$

The remaining white area in Table 3 corresponds to lossy/impossible conversions. For example, the conversion from LD to LL, i.e., \dagger_4 , has been proven to be inconsistent with the raw data, as detailed in Section 2. Similarly, \dagger_3 also represents a lossy conversion. In the case of \dagger_6 , irrelevant labels cannot be reassigned in a reasonable manner. Likewise, conversions such as \dagger_7 , \dagger_8 , \dagger_{11} and \dagger_{12} are not feasible, since labels with the same discrete value do not possess an inherent order.

Next, we provide further insights into Table 3. The gray area reveals potential learning paradigms. The diagonal entries correspond to standard learning paradigms. For example, \dagger_1 corresponds precisely to the learning paradigm proposed in this paper, \dagger_2 refers to LDL (Geng, 2016) or incomplete/inaccurate LDL (Xu & Zhou, 2017; Kou et al., 2023), \dagger_{10} refers to MLL (Zhang & Zhou, 2013) or partial multi-label learning (Xie & Huang, 2018), and \dagger_{15} is for single-label classification and its variants. The lower triangular section comprises recovery/enhancement paradigms. For example, \dagger_9 represents label enhancement (Xu et al., 2019), \dagger_{14} can be regarded as single positive label learning (Cole et al., 2021), and \dagger_5 and \dagger_{13} correspond to (Lu & Jia, 2024) and (Lu & Jia, 2022), respectively.

Table 3 not only demonstrates the advantage of GLD in being fully compatible with existing label ambiguity representations, but also illustrates how our work links prior related studies.

C MORE DETAILS OF EXPERIMENTS

The artificial dataset In this dataset, \mathbf{x} is of three-dimensional and there are three labels, i.e., $m = 3$ and $c = 3$. The corresponding GLD \mathbf{g} of \mathbf{x} is generated in the following ways:

$$t_k = ax_k + bx_k^2 + cx_k^3 + d \quad (k = 1, 2, 3), \quad (58)$$

$$\varpi_1 = \tanh(\mathbf{w}_1^\top \mathbf{t}), \quad (59)$$

$$\varpi_2 = (1 - \lambda_1) \tanh(\mathbf{w}_2^\top \mathbf{t}) + \lambda_1 \varpi_1, \quad (60)$$

$$\varpi_3 = (1 - \lambda_2) \tanh(\mathbf{w}_3^\top \mathbf{t}) + \lambda_2 \varpi_2, \quad (61)$$

$$g_j = \frac{\varpi_j}{\varpi_1 + \varpi_2 + \varpi_3} \quad (j = 1, 2, 3). \quad (62)$$

Note that Equations (60) and (61) deliberately make the relative degree of one label depend on those of other labels. The parameters in Equations (58) to (61) are set as $a = 2$, $b = 4$, $c = 1$, $d = 1$,

$\mathbf{w}_1 = \langle 0.4, 0.2, -1.0 \rangle^\top$, $\mathbf{w}_2 = \langle 0.2, 0.1, 0.4 \rangle^\top$, $\mathbf{w}_3 = \langle -0.1, 0.4, 0.2 \rangle^\top$, and $\lambda_1 = \lambda_2 = 0.01$. To generate the dataset, each component of \mathbf{x} is uniformly sampled within the range $[-1, 1]_{\mathbb{R}}$, and then the GLD \mathbf{g} corresponding to each \mathbf{x} is calculated via Equations (58) to (62). In total, there are 200 data generated in this way.

The real-world datasets According to Definition 3.2 and Equation (50), we preprocess the raw data in $\prod_{j \in [c]} [a_j, b_j]_{\mathbb{R}}$ into GLDs in $[-1, 1]_{\mathbb{R}}^c$ and LDs in Δ^{c-1} , thereby constructing real-world datasets for experiments, where \mathbf{a} and \mathbf{b} serve only as meta-data. Additional experimental results are provided in Table 4. It is important to note that for the `jaf` dataset, the raw data is shifted from the original $[1, 5]_{\mathbb{R}}$ interval to $[0, 4]_{\mathbb{R}}$ to ensure greater reasonability. Consequently, the reported metric values are not directly comparable to those presented in previous literature.

Evaluation metrics The Clark distance is one of the earliest metrics used for evaluating the task of LDL (Geng, 2016). It is sensitive to values approaching zero and, when the ground-truth LD is sparse, can reflect the sparsity of the predicted LD to some extent. Its definition is as follows:

$$\text{Clark}(\mathbf{u}, \mathbf{v}) \triangleq \sqrt{\sum_{j \in [c]} \frac{(u_j - v_j)^2}{(u_j + v_j)^2}}. \quad (63)$$

With a slight abuse of notation, we use $\text{Clark}(\mathbf{D}, \tilde{\mathbf{D}}) \triangleq 1/n \sum_{i \in [n]} \text{Clark}(\mathbf{d}_i, \tilde{\mathbf{d}}_i)$ to evaluate the performance on the test set, where \mathbf{D} and $\tilde{\mathbf{D}}$ are the matrices of the ground-truth and predicted LDs for all n test samples, respectively. The μ metric is proposed in (Li et al., 2025). It can be combined with existing metrics to quantify the extent to which samples are predicted as approximately correct. Its definition, when combined with the KLD, is as follows:

$$\mu_{\text{KLD}}(\mathbf{D}, \tilde{\mathbf{D}}) \triangleq \frac{1}{\delta_0} \int_0^{\delta_0} \frac{1}{n} \sum_{i \in [n]} \mathbb{I}[\text{KLD}(\mathbf{d}_i \parallel \tilde{\mathbf{d}}_i) \leq \delta] d\delta, \quad (64)$$

where

$$\text{KLD}(\mathbf{u} \parallel \mathbf{v}) \triangleq \sum_{j \in [c]} u_j \log \frac{u_j}{v_j}, \quad (65)$$

and

$$\delta_0 = \frac{1}{n} \sum_{i \in [n]} \sum_{j \in [c]} (\mathbf{d}_i)_j \log(c(\mathbf{d}_i)_j) \quad (66)$$

is the worst-case KLD, calculated between the ground-truth LDs and the uniform vector.

For the MLL task, Hamming distance and subset accuracy are commonly used evaluation metrics (Zhang & Zhou, 2013). The Hamming distance is defined as

$$\text{Ham.}(\mathbf{u}, \mathbf{v}) \triangleq \frac{1}{c} \sum_{j \in [c]} \mathbb{I}[u_j \neq v_j]. \quad (67)$$

We further evaluate on the test set using $\text{Ham.}(\mathbf{L}, \tilde{\mathbf{L}}) \triangleq 1/n \sum_{i \in [n]} \text{Ham.}(\mathbf{l}_i, \tilde{\mathbf{l}}_i)$, where \mathbf{L} and $\tilde{\mathbf{L}}$ are the matrices of the ground-truth and predicted LLs for all n test samples, respectively. The definition of subset error rate is given in Equation (2), and thus the subset accuracy on the test set is calculated as $\text{S. Acc.}(\mathbf{L}, \tilde{\mathbf{L}}) \triangleq 1/n \sum_{i \in [n]} 1 - \text{S. Err.}(\mathbf{l}_i, \tilde{\mathbf{l}}_i)$.

To evaluate order consistency, we consider both intra-sample and inter-sample ranking metrics. For intra-sample ranking, we adopt Spearman’s rank correlation coefficient (Jia et al., 2023), defined as

972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000
 1001
 1002
 1003
 1004
 1005
 1006
 1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025

$$\text{Spear.}(\mathbf{u}, \mathbf{v}) \triangleq 1 - \frac{6 \sum_{j \in [c]} (\rho(u_j) - \rho(v_j))^2}{c(c-1)}, \quad (68)$$

where $\rho(\cdot)$ denotes the ranking function. The overall performance on the test set is then measured by $\text{Spear.}(\mathbf{U}, \mathbf{V}) \triangleq 1/n \sum_{i \in [n]} \text{Spear.}(\mathbf{u}_i, \mathbf{v}_i)$.

For inter-sample ranking, we employ an averaged Kendall’s tau statistic (type-A) across labels, given by $\text{Ken.}'(\mathbf{U}, \mathbf{V}) \triangleq 1/c \sum_{j \in [c]} \text{Ken.}(\mathbf{u}_{\bullet j}, \mathbf{v}_{\bullet j})$ where $\text{Ken.}(\cdot, \cdot)$ is defined in Equation (7).

For the task of OOD classification, we propose a newly designed evaluation metric, defined as

$$\text{OOD Err.}(\mathbf{g}, \mathbf{v}) \triangleq \begin{cases} \llbracket \forall_{j \in [c]} (v_j \leq 0) \rrbracket, & \text{if } \mathbf{v} \text{ is GLD,} \\ \llbracket \sum_{j \in [c]} v_j \log(cv_j) \leq \delta_0/2 \rrbracket, & \text{if } \mathbf{v} \text{ is LD,} \\ \llbracket \arg \max_j(\mathbf{g}) = \arg \max_j(\mathbf{v}) \rrbracket, & \text{o/w,} \end{cases} \quad \forall_{j \in [c]} (g_j \leq 0), \quad (69)$$

where δ_0 is defined in Equation (66), and \mathbf{g} is the ground-truth GLD. We further evaluate on the test set using $\text{OOD Err.}(\mathbf{G}, \mathbf{V}) \triangleq 1/n \sum_{i \in [n]} \text{OOD Err.}(\mathbf{g}_i, \mathbf{v}_i)$.

Table 4: Experimental results on datasets scm1d, jaf, enb, and atp1d formatted as (mean \pm std).

Algorithms	Clark \downarrow	μ KLD \uparrow	Ham. \downarrow	S. Acc. \uparrow	Spear. \uparrow	Ken. \uparrow	OOD Err. \downarrow	Avg. Rank
scm1d								
LD-SVR	.0955 \pm .002	90.50% \pm .004	.4544 \pm .004	00.09% \pm .001	.9247 \pm .003	.3421 \pm .007	49.31% \pm .016	7.57
GLD-SVR	.0845 \pm .002	92.06% \pm .004	.0731 \pm .003	45.09% \pm .015	.9369 \pm .003	.7150 \pm .006	37.23% \pm .014	2.43
LD-kNN	.0654 \pm .002	94.05% \pm .003	.4484 \pm .004	00.12% \pm .001	.9540 \pm .002	.3617 \pm .007	38.23% \pm .014	4.29
GLD-kNN	.0653 \pm .002	94.07% \pm .003	.0488 \pm .002	56.45% \pm .016	.9541 \pm .002	.7843 \pm .006	32.13% \pm .014	1.00
LD-BFGS	.1059 \pm .002	88.05% \pm .005	.4604 \pm .004	00.09% \pm .001	.9104 \pm .004	.3165 \pm .007	48.83% \pm .015	10.1
GLD-BFGS	.0991 \pm .002	89.34% \pm .004	.1088 \pm .003	36.49% \pm .012	.9182 \pm .003	.6627 \pm .006	43.51% \pm .016	6.86
LD-DF	.0959 \pm .002	89.87% \pm .004	.4586 \pm .004	00.11% \pm .001	.9260 \pm .003	.3290 \pm .007	45.62% \pm .014	7.43
GLD-DF	.0923 \pm .002	90.56% \pm .004	.1002 \pm .003	38.52% \pm .012	.9265 \pm .003	.6880 \pm .006	41.34% \pm .015	4.00
LD-LRR	.1148 \pm .002	86.32% \pm .005	.4615 \pm .005	00.07% \pm .001	.8985 \pm .004	.3036 \pm .007	50.94% \pm .017	12.0
GLD-LRR	.1083 \pm .002	87.72% \pm .005	.1339 \pm .004	29.88% \pm .014	.9058 \pm .004	.6327 \pm .007	49.79% \pm .019	8.86
LD-Delta	.0965 \pm .002	89.71% \pm .004	.4602 \pm .004	00.09% \pm .001	.9209 \pm .003	.3347 \pm .007	46.72% \pm .016	8.57
GLD-Delta	.0960 \pm .002	89.91% \pm .004	.0815 \pm .003	43.25% \pm .015	.9222 \pm .003	.6737 \pm .006	41.09% \pm .015	4.86
jaf								
LD-SVR	.6414 \pm .061	50.14% \pm .060	.3246 \pm .040	09.99% \pm .062	.4821 \pm .109	.3538 \pm .073	51.41% \pm .103	9.71
GLD-SVR	.6189 \pm .056	51.40% \pm .069	.2177 \pm .037	23.74% \pm .084	.4936 \pm .107	.3915 \pm .077	44.37% \pm .101	5.43
LD-kNN	.6185 \pm .053	50.98% \pm .070	.3142 \pm .043	08.54% \pm .058	.5064 \pm .093	.3919 \pm .078	50.94% \pm .101	7.86
GLD-kNN	.6145 \pm .058	52.37% \pm .066	.2146 \pm .038	25.23% \pm .083	.5280 \pm .091	.3883 \pm .075	49.38% \pm .108	4.71
LD-BFGS	.6297 \pm .047	51.43% \pm .089	.2856 \pm .047	18.31% \pm .094	.5084 \pm .100	.3984 \pm .079	46.19% \pm .110	5.86
GLD-BFGS	.6781 \pm .054	47.83% \pm .091	.2134 \pm .039	26.74% \pm .092	.5116 \pm .096	.4094 \pm .069	45.51% \pm .116	5.00
LD-DF	.5740 \pm .049	57.70% \pm .075	.2913 \pm .039	15.88% \pm .074	.5365 \pm .084	.4287 \pm .075	40.79% \pm .116	3.43
GLD-DF	.5951 \pm .054	56.25% \pm .076	.1967 \pm .037	32.26% \pm .104	.5427 \pm .095	.4337 \pm .064	44.56% \pm .101	1.71
LD-LRR	.6993 \pm .051	42.23% \pm .089	.3137 \pm .048	16.06% \pm .083	.4137 \pm .103	.3195 \pm .085	51.86% \pm .109	10.7
GLD-LRR	.7596 \pm .212	39.83% \pm .135	.2488 \pm .085	25.31% \pm .118	.3969 \pm .149	.3255 \pm .125	51.28% \pm .138	9.57
LD-Delta	.6152 \pm .052	52.75% \pm .090	.3114 \pm .048	14.05% \pm .084	.4704 \pm .110	.3734 \pm .101	44.26% \pm .114	6.43
GLD-Delta	.6864 \pm .084	45.68% \pm .092	.2243 \pm .039	28.71% \pm .086	.4482 \pm .102	.3569 \pm .078	47.42% \pm .089	7.57
atp1d								
LD-SVR	.1051 \pm .011	77.83% \pm .035	.4793 \pm .016	01.28% \pm .019	.8200 \pm .039	.1885 \pm .054	52.17% \pm .069	6.57
GLD-SVR	.1021 \pm .012	78.12% \pm .035	.0321 \pm .019	89.97% \pm .051	.8102 \pm .038	.5191 \pm .050	09.00% \pm .042	1.43
LD-kNN	.1072 \pm .013	76.44% \pm .035	.4818 \pm .015	00.86% \pm .015	.8262 \pm .039	.1755 \pm .050	51.84% \pm .069	7.71
GLD-kNN	.1072 \pm .013	76.52% \pm .035	.0382 \pm .021	89.01% \pm .051	.8232 \pm .038	.4782 \pm .058	11.26% \pm .048	3.43
LD-BFGS	.1462 \pm .016	64.06% \pm .051	.4815 \pm .018	01.91% \pm .024	.7008 \pm .061	.1542 \pm .050	49.80% \pm .076	10.1
GLD-BFGS	.1392 \pm .013	65.81% \pm .048	.0467 \pm .026	86.34% \pm .057	.7280 \pm .051	.4079 \pm .053	11.28% \pm .051	7.14
LD-DF	.1224 \pm .013	71.26% \pm .042	.4801 \pm .016	01.82% \pm .022	.7549 \pm .053	.1718 \pm .049	51.68% \pm .081	9.43
GLD-DF	.1119 \pm .012	75.01% \pm .040	.0421 \pm .023	87.65% \pm .055	.7926 \pm .041	.4932 \pm .053	10.93% \pm .048	4.57
LD-LRR	.1163 \pm .012	74.04% \pm .038	.4784 \pm .016	01.99% \pm .022	.7586 \pm .050	.1763 \pm .046	47.57% \pm .072	7.71
GLD-LRR	.2967 \pm .055	17.03% \pm .123	.1200 \pm .046	81.70% \pm .080	.4057 \pm .163	.1043 \pm .079	15.78% \pm .087	9.43
LD-Delta	.1053 \pm .011	77.64% \pm .036	.4771 \pm .017	01.81% \pm .025	.7926 \pm .046	.1903 \pm .051	51.48% \pm .074	6.14
GLD-Delta	.1144 \pm .018	74.45% \pm .064	.0353 \pm .024	89.01% \pm .050	.7812 \pm .052	.4958 \pm .055	09.36% \pm .046	4.29
enb								
LD-SVR	.0209 \pm .001	84.06% \pm .032	.5306 \pm .013	00.26% \pm .005	.7157 \pm .076	.0309 \pm .010	58.31% \pm .049	8.14
GLD-SVR	.0242 \pm .003	80.31% \pm .030	.1235 \pm .028	80.24% \pm .047	.6891 \pm .078	.4556 \pm .054	14.19% \pm .039	8.00
LD-kNN	.0148 \pm .001	90.22% \pm .022	.5172 \pm .013	01.60% \pm .013	.7760 \pm .069	.0490 \pm .011	57.68% \pm .053	4.71
GLD-kNN	.0160 \pm .002	89.10% \pm .019	.0564 \pm .021	89.57% \pm .038	.8226 \pm .058	.7060 \pm .039	12.61% \pm .038	1.86
LD-BFGS	.0213 \pm .002	82.48% \pm .030	.5293 \pm .014	00.39% \pm .007	.7111 \pm .083	.0315 \pm .010	61.37% \pm .047	8.71
GLD-BFGS	.0238 \pm .002	79.13% \pm .035	.0724 \pm .024	88.26% \pm .038	.7309 \pm .072	.5535 \pm .049	11.30% \pm .039	5.14
LD-DF	.0203 \pm .002	83.31% \pm .032	.5172 \pm .013	01.60% \pm .013	.7624 \pm .077	.0359 \pm .011	58.66% \pm .053	6.00
GLD-DF	.0217 \pm .002	81.80% \pm .030	.0691 \pm .020	86.59% \pm .039	.7010 \pm .078	.5669 \pm .058	11.73% \pm .038	5.86
LD-LRR	.0213 \pm .002	82.49% \pm .030	.5293 \pm .014	00.39% \pm .007	.7111 \pm .083	.0315 \pm .010	61.37% \pm .047	8.57
GLD-LRR	.0272 \pm .024	77.68% \pm .089	.0729 \pm .027	88.18% \pm .047	.7187 \pm .084	.5553 \pm .049	11.34% \pm .039	6.14
LD-Delta	.0211 \pm .002	82.30% \pm .032	.5297 \pm .014	00.35% \pm .007	.7146 \pm .087	.0285 \pm .011	62.89% \pm .054	9.43
GLD-Delta	.0232 \pm .002	80.12% \pm .039	.0744 \pm .025	88.18% \pm .039	.7222 \pm .079	.5688 \pm .048	11.58% \pm .039	5.43