
Approximately Correct Label Distribution Learning

Weiwei Li¹ Haitao Wu^{2,1} Yunan Lu^{2,3} Xiuyi Jia²

Abstract

Label distribution learning (LDL) is a powerful learning paradigm that emulates label polysemy by assigning label distributions over the label space. However, existing LDL evaluation metrics struggle to capture meaningful performance differences due to their insensitivity to subtle distributional changes, and existing LDL learning objectives often exhibit biases by disproportionately emphasizing a small subset of samples with extreme predictions. As a result, the LDL metrics lose their discriminability, and the LDL objectives are also at risk of overfitting. In this paper, we propose DeltaLDL, a percentage of predictions that are approximately correct within the context of LDL, as a solution to the above problems. DeltaLDL can serve as a novel evaluation metric, which is parameter-free and reflects more on real performance improvements. DeltaLDL can also serve as a novel learning objective, which is differentiable and encourages most samples to be predicted as approximately correct, thereby mitigating overfitting. Our theoretical analysis and empirical results demonstrate the effectiveness of the proposed solution.

1. Introduction

Label polysemy is a common phenomenon in real-world scenarios, traditionally emulated by multi-label learning (Zhang & Zhou, 2013), which treats labels in a binary fashion. In contrast, label distribution learning (LDL) (Geng, 2016) offers a more refined framework by answering the question: “To what extent does each label y describe the instance x ?” This is achieved through the concept of a

¹College of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics, Nanjing, China ²School of Computer Science and Engineering, Nanjing University of Science and Technology, Nanjing, China ³Department of Computing, The Hong Kong Polytechnic University, Hong Kong, China. Correspondence to: Xiuyi Jia <jiaxy@njust.edu.cn>.

label distribution d , represented as a probability simplex, where each label is associated with a real-valued *description degree* d_x^y . By quantifying label polysemy, LDL enables a broader spectrum of applications, including tasks like object detection (Xu et al., 2023), age estimation (Wen et al., 2023), sentiment analysis (Chen et al., 2020), and diagnostic assistance (Li et al., 2023), etc.

For years, there are some deep-rooted problems in the field of LDL: 1) existing evaluation metrics for LDL are based on some distance/similarity metrics, e.g., Kullback-Leibler divergence (KLD), which have very small values and thus have poor discriminability between superior and inferior models, since the label distribution is subject to two constraints, non-negativity (i.e., $d_x^y \geq 0$) and sum-to-one (i.e., $\sum_{y \in \mathcal{Y}} d_x^y = 1$); 2) existing LDL algorithms mainly focus on minimizing the average of some kind of measurement, usually the KLD as well, which may lead to overfitting due to excessive preference for some samples, resulting even worse performance than that calculated from uniform vectors.

To kill two birds with one stone, the idea of approximately correct prediction can be introduced. On the one hand, it enables the design of percentage-based metrics, quantifying what proportion of the samples are predicted approximately correctly, which is intuitive and semantically meaningful; on the other hand, it enables the design of noise-robust objectives, allowing the model to sacrifice samples that are particularly difficult to learn, ensuring that most samples are predicted approximately correctly, thereby avoiding overfitting to a specific subset of data.

Contributions and organizational structure

In this paper: 1) we conduct a theoretical analysis of the KLD to demonstrate its unsuitability as an evaluation metric/learning objective for LDL (Section 2.1); 2) we propose DeltaLDL, a function that calculates the percentage of label distributions that are approximately correctly predicted, exhibiting desirable properties, making it a promising candidate for both metrics and objectives (Section 2.2); 3) based on DeltaLDL, we propose a novel evaluation metric, the μ metric, which possesses superior discriminative power according to our analysis (Section 3); 4) based on DeltaLDL, we propose a novel learning objective and encapsulate it in a new LDL algorithm, named δ -LDL, which can achieve highly compet-

itive performance according to our experiments (Sections 4 and 5). Related work and limitations are discussed in Sections 6 and 7, respectively.

2. Theoretical analysis

In this section, we conduct a theoretical analysis to argue that the KLD is not suitable as an evaluation metric/learning objective for LDL. Then, we propose our new solution, named DeltaLDL.

Preliminaries Let $\mathbf{x} \in \mathcal{X} = \mathbb{R}^q$ denote the feature of the instance and $\mathbf{d} = (d_{\mathbf{x}}^{y_j})_{j=1}^c \in \Delta^{c-1}$ denote the label distribution, where c is the number of the labels, y_j is the j -th label in the label space $\mathcal{Y} = \{y_j\}_{j=1}^c$, and

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

is the $(k-1)$ -dimensional probability simplex. LDL's goal is to find a mapping $f: \mathcal{X} \mapsto \Delta^{c-1}$. More specifically, given a training set $\{\mathbf{x}_i, \mathbf{d}_i\}_{i=1}^m$ and an LDL model $f(\cdot; \Theta)$, the loss function can be defined as

$$\ell(\Theta) = \frac{1}{m} \sum_{i=1}^m \text{dist}(f(\mathbf{x}_i; \Theta), \mathbf{d}_i), \quad (2)$$

where Θ is the model parameters and $\text{dist}(\cdot, \cdot)$ is a distance metric, usually the KLD. The optimal parameters Θ^* can be obtained by minimizing Equation (2), i.e.,

$$\Theta^* = \arg \min_{\Theta} \ell(\Theta). \quad (3)$$

2.1. Issues of existing metrics: KLD as an example

KLD is defined as

$$\text{KL}(\mathbf{u} \parallel \mathbf{v}) = \sum_{j=1}^c u_j \log \frac{u_j}{v_j}, \quad (4)$$

which exposes shortcomings in two aspects. First, KLD is *not bounded and is only sensitive to small probability values of the prediction*, i.e., $\text{KL}(\mathbf{u} \parallel \mathbf{v}) \in [0, +\infty)$ and $\text{KL}(\mathbf{u} \parallel \mathbf{v}) \rightarrow \infty$ if $v_j \rightarrow 0^+$, $\exists j \in [c]$. Most existing LDL methods are based on maximum entropy models or neural networks with softmax as the final mathematical processing, hence it is difficult to trigger its sensitivity. Second, KLD *does not adequately capture the changing nature of probability simplexes*, i.e., $\text{KL}(\mathbf{u} \parallel \mathbf{v})$ lacks the significance of changes to tell the difference between superior and inferior models when $\mathbf{u}, \mathbf{v} \in \Delta^{c-1}$. To illustrate the second aspect, we introduce the following assumption and proposition.

Assumption 2.1. The ground-truth label distribution \mathbf{d} follows a Dirichlet distribution, i.e., the probability density function of \mathbf{d} is given by

$$p(\mathbf{d}_i; \alpha) = \frac{1}{\mathbf{B}(\alpha)} \prod_{j=1}^c (d_{\mathbf{x}_i}^{y_j})^{\alpha_j - 1}, \quad \mathbf{B}(\alpha) = \frac{\prod_{j=1}^c \Gamma(\alpha_j)}{\Gamma(\sum_{j=1}^c \alpha_j)}, \quad (5)$$

Table 1. Case analysis of two LDL datasets

Case	Dataset	$\mathbb{E}(\text{KL}(\mathbf{u} \parallel \mathbf{v}))$	SA-BFGS	LRR
1)	SBU_3DFE	.0707	.0627	.0571
2)	M ² B	.6335	.6878	.6664

where $\Gamma(\cdot)$ is the gamma function. Assumption 2.1 is currently the most common assumption in the literature of LDL (Lu & Jia, 2022; He et al., 2024). Based on Equation (5), we have the following proposition.

Proposition 2.2. The expectation of the KLD between a Dirichlet-distributed random vector \mathbf{u} and a uniform vector

$$\mathbf{v} = \underbrace{\left(\frac{1}{c}, \frac{1}{c}, \dots, \frac{1}{c}\right)}_c \quad (6)$$

is given by

$$\mathbb{E}(\text{KL}(\mathbf{u} \parallel \mathbf{v})) = \sum_{j=1}^c \left(\frac{\alpha_j}{\sum_{k=1}^c \alpha_k} \left(\psi(\alpha_j + 1) - \psi\left(\sum_{k=1}^c \alpha_k + 1\right) \right) + \frac{\alpha_j}{\sum_{k=1}^c \alpha_k} \log c \right), \quad (7)$$

where $\psi(\cdot)$ is the digamma function. $\mathbb{E}(\text{KL}(\mathbf{u} \parallel \mathbf{v}))$ can be regarded as the theoretical stochastic KLD performance of an LDL model.

Proof. See Appendix A for details. \square

Remark 2.3. To estimate α from given data $(\mathbf{d}_i)_{i=1}^m$, one can employ maximum likelihood estimation (MLE), i.e., minimizing:

$$\begin{aligned} \ell_{\text{MLE}}(\alpha) = & - \sum_{i=1}^m \log p(\mathbf{d}_i; \alpha) = \\ & - m \left(\ln \Gamma\left(\sum_{j=1}^c \alpha_j\right) - \sum_{j=1}^c \ln \Gamma(\alpha_j) \right) \\ & - \sum_{i=1}^m \sum_{j=1}^c (\alpha_j - 1) \ln d_{\mathbf{x}_i}^{y_j}, \end{aligned} \quad (8)$$

the closed-form solution of which does not exist. Thus, one can employ numerical optimization methods, e.g., L-BFGS, to solve the optimization problem.

With Proposition 2.2 and Remark 2.3, one can calculate $\mathbb{E}(\text{KL}(\mathbf{u} \parallel \mathbf{v}))$ for each LDL dataset. Here we conduct a case analysis of two LDL datasets, SBU_3DFE and M²B, and compare the KLD performance of two learned LDL methods, SA-BFGS (Geng, 2016) and LRR (Jia et al., 2023b), both of which exploit average KLD as the learning objective.

The results are shown in Table 1. In case 1), since label distributions in SBU_3DFE are non-sparse, the calculated values are very small, lacking discriminability, i.e., considerable performance improvements are only reflected via these nuanced changes, which is counterintuitive; in case 2), results of SA-BFGS and LRR are worse than the stochastic performance, indicating that minimizing average KLD sometimes does not work as expected, which may be attributed to the negative contribution of some extremely predicted samples.

2.2. The proposed solution: DeltaLDL

On the one hand, Proposition 2.2 fully illustrates the non-significance of KLD as an evaluation metric/learning objective for LDL; on the other hand, it provides a baseline of performance improvements to get a glimpse of *real* performance improvements, which will be discussed in Section 3. Taking inspiration from the analysis of Table 1, we define DeltaLDL as our new solution.

Definition 2.4 (Approximately correct LDL; DeltaLDL). Suppose that an LDL model $f(\cdot; \Theta)$ is evaluated on $\{\mathbf{x}_i, \mathbf{d}_i\}_{i=1}^n$; DeltaLDL of a distance metric $\text{dist}(\cdot, \cdot)$ is defined as the following function:

$$\mathfrak{D}(\text{dist}, \delta; f) \triangleq \frac{1}{n} \sum_{i=1}^n \mathbb{I}(\text{dist}(\mathbf{d}_i, f(\mathbf{x}_i; \Theta)) \leq \delta), \quad (9)$$

where $\mathbb{I}(\cdot)$ is the indicator function to count the number of approximately correct predictions. DeltaLDL of a similarity metric $\text{sim}(\cdot, \cdot)$ is defined as the following function:

$$\mathfrak{D}(\text{sim}, \delta; f) \triangleq \frac{1}{n} \sum_{i=1}^n \mathbb{I}(\text{sim}(\mathbf{d}_i, f(\mathbf{x}_i; \Theta)) \geq \delta). \quad (10)$$

Remark 2.5. The following properties hold for \mathfrak{D} :

- $\mathfrak{D}(\cdot, \cdot; \cdot) \in [0, 1]_{\mathbb{R}}$;
- $\mathfrak{D}(\text{dist}, \sup_{\mathbf{u}, \mathbf{v} \in \Delta^{c-1}} \text{dist}(\mathbf{u}, \mathbf{v}); \cdot) = 1$;
- $\mathfrak{D}(\text{sim}, \inf_{\mathbf{u}, \mathbf{v} \in \Delta^{c-1}} \text{sim}(\mathbf{u}, \mathbf{v}); \cdot) = 1$;
- $\mathfrak{D}(\text{dist}, \delta; \cdot) \leq \mathfrak{D}(\text{dist}, \delta'; \cdot)$, if $\delta \leq \delta'$;
- $\mathfrak{D}(\text{sim}, \delta; \cdot) \geq \mathfrak{D}(\text{sim}, \delta'; \cdot)$, if $\delta \leq \delta'$;
- $\text{dist}(\mathbf{d}_i, f(\mathbf{x}_i)) = \inf_{\mathbf{u}, \mathbf{v} \in \Delta^{c-1}} \text{dist}(\mathbf{u}, \mathbf{v})$, $\forall i \in [n]$, if and only if $\mathfrak{D}(\text{dist}, \inf_{\mathbf{u}, \mathbf{v} \in \Delta^{c-1}} \text{dist}(\mathbf{u}, \mathbf{v}); f) = 1$;
- $\text{sim}(\mathbf{d}_i, f(\mathbf{x}_i)) = \sup_{\mathbf{u}, \mathbf{v} \in \Delta^{c-1}} \text{sim}(\mathbf{u}, \mathbf{v})$, $\forall i \in [n]$, if and only if $\mathfrak{D}(\text{sim}, \sup_{\mathbf{u}, \mathbf{v} \in \Delta^{c-1}} \text{sim}(\mathbf{u}, \mathbf{v}); f) = 1$.

As stated in Definition 2.4, any distance/similarity metric can be transformed into a percentage of predictions that are approximately correct within the context of LDL. Furthermore, Remark 2.5 shows some good properties of DeltaLDL, which are intuitive and easy to proof by readers. These properties indicate that DeltaLDL has excellent potential as an evaluation metric/learning objective, and also highlight the relationship between DeltaLDL and the distance/similarity

metric employed. In this paper, we focus primarily on $\mathfrak{D}(\text{KL}, \cdot; \cdot)$, since the KLD is the most widely used evaluation metric in the field of LDL. Following the definition of DeltaLDL, one may wonder how it can be utilized as an evaluation metric/learning objective.

3. DeltaLDL as an evaluation metric

To transform $\mathfrak{D}(\text{KL}, \cdot; \cdot)$ into an evaluation metric, the following challenges must be addressed: 1) the metric should be parameter-free, as δ varies with different label distributions and is difficult to determine; 2) the metric should effectively reflect real performance improvements, enabling significant differentiation between superior and inferior models. To this end, we propose the following definition.

Definition 3.1 (The discriminative power μ). Suppose that an LDL model is evaluated on the test set, where the ground-truth label distribution matrix is $\mathbf{D} = (\mathbf{d}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{D})_{i=1}^n$, and the predicted label distribution matrix is $\tilde{\mathbf{D}} = (\tilde{\mathbf{d}}_i)_{i=1}^n$, which is given by the LDL model f ; our new evaluation metric based on $\mathfrak{D}(\text{KL}, \delta; f)$, denoted as μ , is defined as

$$\begin{aligned} \mu(\mathbf{D}, \tilde{\mathbf{D}}) &\triangleq \frac{1}{\delta_0} \int_0^{\delta_0} \mathfrak{D}(\text{KL}, \delta; f) d\delta \\ &= \frac{1}{\delta_0} \int_0^{\delta_0} \frac{1}{n} \sum_{i=1}^n \mathbb{I}(\text{KL}(\tilde{\mathbf{d}}_i \parallel \mathbf{d}_i) \leq \delta) d\delta \\ &\in [0, 1]_{\mathbb{R}}, \end{aligned} \quad (11)$$

where

$$\delta_0 = \mathbb{E}(\text{KL}(\mathbf{u} \parallel \mathbf{v})), \quad \mathbf{u} \sim \mathcal{D}, \quad (12)$$

and \mathbf{v} is a uniform vector (the same as Equation (6)).

Since the test set is available in the evaluation phase, Equation (12) can be calculated as:

$$\delta_0 \approx \delta'_0 = \frac{1}{n} \sum_{i=1}^n \text{KL}(\mathbf{d}_i \parallel \mathbf{v}) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^c d_{\mathbf{x}_i}^{y_j} \ln d_{\mathbf{x}_i}^{y_j} c. \quad (13)$$

According to Definition 3.1, μ 's semantic essence lies in the overall performance of the predictions being approximately correct across different values of δ ; mathematically, μ is defined as the *ratio* of the area enclosed by the function $\mathfrak{D}(\text{KL}, \delta; \cdot)$ w.r.t. δ and the coordinate axes, where δ_0^{-1} serves as the scaling factor. By adopting an integral form, μ eliminates dependency on specific δ , addressing the aforementioned challenge 1; the introduction of δ_0 establishes a baseline via the value based on uniform vectors, exclusively capturing performance improvements above this baseline, thus effectively resolving the challenge 2.

Discriminability analysis For better comparison, we define the *real* performance improvement based on the $\text{KLD}/\mu(\cdot, \cdot)$, with a little symbol abuse.

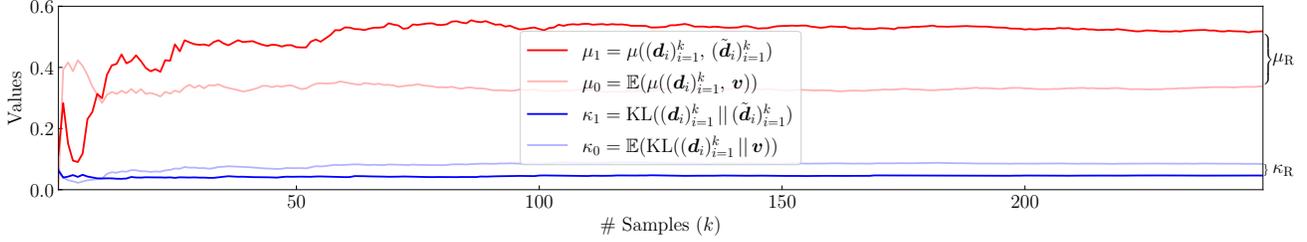


Figure 1. κ_1 , κ_0 , μ_1 , and μ_0 of DF-LDL w.r.t. the number of samples in the test set of SBU_3DFE.

Definition 3.2 (The *real* performance improvement based on the KLD). Let $\kappa_0 = \mathbb{E}(\text{KL}(\mathbf{D} \parallel \mathbf{v}))$, where \mathbf{v} is a uniform vector; suppose that an LDL model is evaluated on the test set and the KLD is κ_1 ; the *real* performance improvement based on the KLD is defined as $\kappa_R \triangleq \kappa_1 - \kappa_0$.

Definition 3.3 (The *real* performance improvement based on the μ metric). Let $\mu_0 = \mathbb{E}(\mu(\mathbf{D}, \mathbf{v}))$, where \mathbf{v} is a uniform vector; suppose that an LDL model is evaluated on the test set and the value of the μ metric is μ_1 ; the *real* performance improvement based on the μ metric is defined as $\mu_R \triangleq \mu_1 - \mu_0$.¹

Similarly, readers can get the *real* performance improvement based on other evaluation metrics. We train a DF-LDL model on the training set of the SBU_3DFE dataset, and evaluate it on the test set by calculating κ_1 , κ_0 , μ_1 , and μ_0 with varying numbers of test samples. The results are shown in Figure 1. As can be seen, μ_R is much larger than κ_R , indicating that the μ metric is more sensitive to performance improvements than the KLD, which is consistent with our theoretical analysis. One limitation worth noting is that, when there are few samples for testing, the μ metric is unstable, which is the nature characteristic of percentage-based metrics and can be eliminated by increasing the number of samples. Since DeltaLDL has been successfully transformed into an evaluation metric, the μ metric, one may wonder how it can be utilized as a learning objective.

4. DeltaLDL as a learning objective

To transform $\mathfrak{D}(\text{KL}, \cdot, \cdot)$ into a learning objective, the following challenges must be addressed: 1) the objective should not contain the indicator function $\mathbb{I}(\cdot)$, as it is difficult to optimize; 2) the objective should be to sacrifice a small number of samples that are difficult to learn and ensure that most samples can be predicted as approximately correct. To solve the challenge 1, we introduce a smoothing function $\sigma(\cdot)$ to replace $\mathbb{I}(\cdot)$. Suppose that an LDL model is training on the training set, where the ground-truth label

¹For readers who find these definitions obscure, we hereby provide a metaphor: κ_0 and μ_0 are like the accuracy of random selection in multi-classification problems.

distribution matrix is $\mathbf{D} = (\mathbf{d}_i^{\text{i.i.d.}})_{i=1}^m$, and the predicted label distribution matrix is $\tilde{\mathbf{D}} = (\tilde{\mathbf{d}}_i)_{i=1}^m$, which is given by the LDL model f ; then $\mathfrak{D}(\text{KL}, \delta; f)$ can be transformed into the following basic loss function:

$$\ell_0(\delta; \mathbf{D}, \tilde{\mathbf{D}}) = \frac{1}{m} \sum_{i=1}^m \sigma \left(\left(\sum_{j=1}^c \tilde{d}_{\mathbf{x}_i}^{y_j} \ln \frac{\tilde{d}_{\mathbf{x}_i}^{y_j}}{d_{\mathbf{x}_i}^{y_j}} \right) - \delta \right), \quad (14)$$

which is differentiable if $\sigma(\cdot)$ is appropriate. The above challenge 2 lies in ensuring that $\mathfrak{D}(\text{KL}, \delta; \cdot)$ should be as large as possible while keeping δ as small as possible. (Note that δ is difficult to determine, as mentioned in Section 3.) This requirement can be understood geometrically as minimizing the area enclosed by the function $\ell_0(\delta; \cdot, \cdot)$ w.r.t. δ and the coordinate axes. We can use adaptive Simpson’s rule (ASR) (McKeeman, 1962), a numerical integration method, for loss function calculation, to solve this challenge.

Remark 4.1. From Simpson’s rule, we have

$$\begin{aligned} \int_a^b \ell_0(\delta; \cdot, \cdot) d\delta &\approx \varphi_0(a, b) \\ &= \frac{b-a}{6} \left(\ell_0(a) + 4\ell_0\left(\frac{a+b}{2}\right) + \ell_0(b) \right). \end{aligned} \quad (15)$$

In adaptive Simpson’s rule, the term $v_{a,b}$ quantifies the error of the basic Simpson’s approximation by considering the subdivision at the midpoint. It is defined as:

$$v_{a,b} = \varphi_0(a, c) + \varphi_0(c, b) - \varphi_0(a, b), \quad c = \frac{a+b}{2}. \quad (16)$$

The adaptive Simpson’s rule refines the integration by recursively subdividing the interval until the error is below a specified tolerance ε .

Remark 4.2. From the adaptive Simpson’s rule, we have

$$\begin{aligned} \int_a^b \ell_0(\delta; \cdot, \cdot) d\delta &\approx \varphi(a, b, \varepsilon) = \\ &\begin{cases} \varphi_0(a, c) + \varphi_0(c, b) + v_{a,b}/15, & |v_{a,b}| \leq 15\varepsilon, \\ \varphi(a, c, \varepsilon/2) + \varphi(c, b, \varepsilon/2), & \text{o/w.} \end{cases} \end{aligned} \quad (17)$$

Algorithm 1 Adaptive Simpson’s rule: ASR

Input: Interval bounds a, b ; error tolerance ε ; current integral estimate s ; maximum recursion depth ξ .

Output: Integral estimate s' .

- 1: $c \leftarrow (a+b)/2$;
- 2: $l \leftarrow \varphi_0(a, c)$; ▷ Equation (15)
- 3: $r \leftarrow \varphi_0(c, b)$;
- 4: **if** $|l + r - s| \leq 15\varepsilon$ **or** $t \leq 0$ **then**
- 5: **return** $l + r + (l+r-s)/15$; ▷ Equation (17)
- 6: **end if**
- 7: **return** $\text{ASR}(a, c, \varepsilon/2, l, \xi - 1) + \text{ASR}(c, b, \varepsilon/2, r, \xi - 1)$;

Algorithm 2 Our proposed algorithm: δ -LDL

Input: Training set $\{(\mathbf{x}_i, \mathbf{d}_i)\}_{i=1}^m$, maximum recursion depth ξ , test sample \mathbf{x}' .

Output: Label distribution \mathbf{d}' .

- 1: $\text{EPS} \leftarrow 10^{-7}$;
- 2: $\Theta \leftarrow$ random initialization;
- 3: $\delta'_0 \leftarrow 1/m \sum_{i=1}^m \sum_{j=1}^c d_{\mathbf{x}_i}^{y_j} \ln d_{\mathbf{x}_i}^{y_j} c$; ▷ Equation (19)
- 4: **for** each epoch **to** anticipated epoch number **do**
- 5: $\mathcal{B} \leftarrow$ randomly split $\{(\mathbf{x}_i, \mathbf{d}_i)\}_{i=1}^m$ into batches;
- 6: **for** each batch **in** \mathcal{B} **do**
- 7: $\tilde{\mathbf{D}} \leftarrow f(\mathbf{x}; \Theta), \forall \mathbf{x} \in$ the batch; ▷ Equation (20)
- 8: $s_0 \leftarrow \varphi_0(0, \delta'_0)$; ▷ Equation (15)
- 9: $\ell \leftarrow \text{ASR}(0, \delta'_0, \text{EPS}, s_0, \xi)$; ▷ Algorithm 1
- 10: $\Theta \leftarrow$ update by Adam, tracing ℓ ; ▷ Equation (21)
- 11: **end for**
- 12: **end for**
- 13: $\mathbf{d}' \leftarrow f(\mathbf{x}'; \Theta)$; ▷ Equation (20)
- 14: **return** \mathbf{d}' ;

The rule is summarized in Algorithm 1. Finally, the loss function can be defined as:

$$\ell(\mathbf{D}, \tilde{\mathbf{D}}) = \varphi(0, \delta'_0, \text{EPS}). \quad (18)$$

Implementation details In Equation (14), $\sigma(\cdot)$ is implemented as the ReLU function; in Algorithm 2, ξ is set to 5; in Equation (18),

$$\delta'_0 = \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^c d_{\mathbf{x}_i}^{y_j} \ln d_{\mathbf{x}_i}^{y_j} c, \quad (19)$$

serving as a performance baseline, and EPS can be a small positive number, e.g., 10^{-7} . The above loss function can be optimized by gradient descent methods, such as Adam (Kingma & Ba, 2015). The LDL model f is implemented by a naive network with the parameter matrix $\Theta \in \mathbb{R}^{q \times c}$:

$$f(\mathbf{x}; \Theta) = \varsigma(\mathbf{x}\Theta), \quad (20)$$

where $\varsigma(\cdot)$ is the softmax function as the final mathematical processing. The overall algorithm is summarized in Algorithm 2, denoted as δ -LDL.

Remark 4.3. Let $\ell_0(\delta; \mathbf{D}, f(\mathbf{X}; \Theta))$ be $\ell_0(\Theta)$; assume that $\sigma(\cdot)$ is the ReLU function; $\tilde{\mathbf{D}} = f(\mathbf{X}; \Theta)$; $\ell_0(\Theta)$ is

differentiable and its gradient w.r.t. Θ is given by

$$\frac{\partial \ell_0(\Theta)}{\partial \Theta} = \mathbb{I}(\text{KL}(\mathbf{D} \parallel \tilde{\mathbf{D}}) > \delta) \odot (\mathbf{X}^\top (\tilde{\mathbf{D}} - \mathbf{D})), \quad (21)$$

where \odot is the element-wise product.

Time complexity analysis Let $\mathbf{X} \in \mathbb{R}^{m \times q}$ denote the feature matrix and $\mathbf{D} \in \mathbb{R}^{m \times c}$ denote the ground-truth label distribution matrix; $\tilde{\mathbf{D}} = f(\mathbf{X}; \Theta)$; the overall time cost of δ -LDL (Algorithm 2) is primarily influenced by the following calculations: computing δ'_0 in Equation (19) has a complexity of $\mathcal{O}(mc)$; the forward pass step in Equation (20) requires $\mathcal{O}(mqc)$; each calculation of the KLD, $\kappa = \mathbf{D} \odot (\log \mathbf{D} - \log \tilde{\mathbf{D}})$, has a complexity of $\mathcal{O}(mc)$; calculating the gradient, $\partial \kappa / \partial \Theta = \mathbf{X}^\top (\tilde{\mathbf{D}} - \mathbf{D})$, involves a complexity of $\mathcal{O}(mqc)$. Since the parameter ξ is fixed, ℓ_0 is computed a limited number of times and this number remains relatively small. The overall time complexity of each iteration of δ -LDL is $\mathcal{O}(kmqc)$, where k is the number of calculations of Equation (14), attributed to Algorithm 1. It is linear w.r.t. the number of samples m /labels c , clarifying that our proposed δ -LDL can be applied to large-scale datasets. Then, one may wonder how δ -LDL performs in practice.

5. Experiments

In this section, extensive experiments are conducted to illustrate the superiority of the μ metric and δ -LDL. Details of all implementations are openly accessible at GitHub.²

5.1. Experimental setup

Metrics & datasets We mainly evaluate LDL models’ performance via the $\mu \uparrow$ metric proposed in Section 3. We also represent the same metrics suggested in (Jia et al., 2023b), which are Cheby. \downarrow (Chebyshev distance), Clark \downarrow (Clark distance), Can. \downarrow (Canberra distance), KLD \downarrow (Kullback-Leibler divergence), Cosine \uparrow (cosine similarity), Int. \uparrow (intersection similarity) and Spear. \uparrow (Spearman’s coefficient), respectively. Here \downarrow (\uparrow) indicates “the lower (higher) the better”. We adopt several widely used label distribution datasets, including: M²B (Nguyen et al., 2012), fbp5500 (Liang et al., 2018), RAF_ML (Li & Deng, 2019), SBU_3DFE (Yin et al., 2006), Natural_Scene (Geng et al., 2021), Music (Lee et al., 2021) and Painting (Machajdik & Hanbury, 2010).

Comparison methods & methodology We compare δ -LDL with several LDL methods. On the one hand, methods based on the optimization of the average KLD are mainly compared, including: SA-BFGS (Geng, 2016), LDLLC (Jia et al., 2018), LDLSF (Ren et al., 2019a), LCLR (Ren et al., 2019b), SCL (Jia et al., 2019), DPA (Jia et al., 2023a) and

²<https://github.com/SpriteMisaka/PyLDL>

Table 2. Experimental results on M^2B , fbp5500, and SBU_3DFE, formatted as “mean \pm std”

Algorithms	Cheby. \downarrow	Clark \downarrow	Can. \downarrow	KLD \downarrow	Cosine \uparrow	Int. \uparrow	Spear. \uparrow	μ \uparrow
M^2B								
PT-Bayes	• .5394 \pm .043	• 1.7487 \pm .064	• 3.6172 \pm .169	• 1.7422 \pm .393	• .5272 \pm .053	• .4498 \pm .042	• .5758 \pm .052	• 25.76% \pm .043
SCL	• .4788 \pm .013	• 1.6358 \pm .010	• 3.3495 \pm .032	• .7773 \pm .027	• .6842 \pm .014	• .4741 \pm .010	• .5752 \pm .110	• 26.67% \pm .012
LDLSF	• .3865 \pm .018	• 1.5441 \pm .031	• 3.0624 \pm .083	• 2.0495 \pm .497	• .7363 \pm .022	• .5991 \pm .019	• .6405 \pm .026	• 32.95% \pm .030
LDLF	• .4171 \pm .024	• 1.5266 \pm .048	• 3.1013 \pm .114	• .5688 \pm .051	• .7584 \pm .020	• .5604 \pm .028	• .6984 \pm .031	• 44.91% \pm .054
LCLR	• .3867 \pm .022	• 1.2516 \pm .036	• 2.4208 \pm .089	• .7303 \pm .098	• .7152 \pm .026	• .6040 \pm .022	• .6862 \pm .026	• 52.13% \pm .031
DPA	• .3847 \pm .022	• 1.1734 \pm .030	• 2.2681 \pm .071	• .7013 \pm .089	• .7187 \pm .026	• .6075 \pm .022	• .6936 \pm .024	• 52.95% \pm .031
LDLLC	• .3845 \pm .022	• 1.1733 \pm .030	• 2.2683 \pm .070	• .7003 \pm .089	• .7189 \pm .026	• .6077 \pm .022	• .6939 \pm .025	• 52.99% \pm .031
SA-BFGS	• .3826 \pm .021	• 1.2228 \pm .028	• 2.3520 \pm .066	• .6878 \pm .091	• .7231 \pm .025	• .6080 \pm .021	• .6886 \pm .026	• 53.06% \pm .029
LRR	• .3793 \pm .021	• 1.1675 \pm .028	• 2.2525 \pm .066	• .6664 \pm .088	• .7271 \pm .025	• .6127 \pm .021	• .6986 \pm .024	• 53.90% \pm .029
AA- k NN	• .3721 \pm .020	◦ 1.0566 \pm .031	◦ 1.9578 \pm .067	• .6179 \pm .074	• .7426 \pm .023	• .6191 \pm .020	• .7022 \pm .026	• 55.21% \pm .028
DF-LDL	• .3699 \pm .019	1.1011 \pm .025	2.0883 \pm .058	• .5743 \pm .072	• .7457 \pm .023	• <u>.6219</u> \pm .019	.7049 \pm .022	• <u>55.36%</u> \pm .027
δ -LDL	.3599 \pm .017	<u>1.0966</u> \pm .023	<u>2.0629</u> \pm .051	.4882 \pm .053	.7768 \pm .020	.6333 \pm .017	.6943 \pm .023	58.25% \pm .025
fbp5500								
PT-Bayes	• .6569 \pm .011	• 1.6829 \pm .010	• 3.1196 \pm .032	• 15.2111 \pm .362	• .4779 \pm .015	• .3413 \pm .011	• .6987 \pm .010	• .08% \pm .001
SCL	• .2822 \pm .005	• 1.4191 \pm .007	• 2.6654 \pm .022	• .4370 \pm .010	• .8442 \pm .006	• .6339 \pm .004	• .7314 \pm .053	• 38.39% \pm .011
LDLSF	• .1450 \pm .003	◦ 1.0405 \pm .016	◦ 1.6992 \pm .031	• .3391 \pm .030	• .9515 \pm .002	• .8381 \pm .003	• .8876 \pm .005	• 69.42% \pm .013
AA- k NN	• .1515 \pm .004	◦ <u>1.0482</u> \pm .015	◦ <u>1.7383</u> \pm .031	• .1819 \pm .013	• .9418 \pm .003	• .8315 \pm .004	• .8860 \pm .006	• 78.19% \pm .010
LCLR	• .1386 \pm .003	• 1.2978 \pm .009	• 2.2253 \pm .022	• .1114 \pm .005	• .9523 \pm .003	• .8475 \pm .004	• .9009 \pm .005	• 84.36% \pm .007
SA-BFGS	• .1369 \pm .003	• 1.2945 \pm .009	• 2.2146 \pm .022	• .1084 \pm .005	• .9536 \pm .003	• .8494 \pm .004	• .9025 \pm .004	• 84.77% \pm .006
LDLF	• .1348 \pm .004	1.2786 \pm .011	2.1763 \pm .028	• .1058 \pm .006	• .9549 \pm .003	• .8503 \pm .005	• .9017 \pm .005	• 85.18% \pm .008
DF-LDL	• .1346 \pm .003	• 1.2898 \pm .010	• 2.2008 \pm .023	• .1055 \pm .005	• .9548 \pm .002	• .8518 \pm .003	• .9042 \pm .004	• 85.18% \pm .006
DPA	• .1347 \pm .003	• 1.2844 \pm .010	• 2.1895 \pm .024	• .1053 \pm .005	• .9548 \pm .003	• .8512 \pm .003	• .9038 \pm .004	• 85.19% \pm .006
LDLLC	• .1348 \pm .003	• 1.2844 \pm .010	• 2.1894 \pm .024	• .1052 \pm .005	• .9549 \pm .003	• .8511 \pm .003	• .9036 \pm .004	• 85.20% \pm .006
LRR	• <u>.1343</u> \pm .003	• 1.2837 \pm .010	• 2.1871 \pm .024	• <u>.1046</u> \pm .004	• <u>.9552</u> \pm .002	• .8516 \pm .003	• .9040 \pm .004	• <u>85.29%</u> \pm .006
δ -LDL	.1300 \pm .003	1.2783 \pm .010	2.1685 \pm .023	.0997 \pm .005	.9576 \pm .002	.8560 \pm .003	.9061 \pm .004	85.99% \pm .006
SBU_3DFE								
PT-Bayes	• .1289 \pm .004	• .4019 \pm .009	• .8606 \pm .020	• .0772 \pm .004	• .9247 \pm .003	• .8452 \pm .004	• .1444 \pm .035	• 33.40% \pm .016
AA- k NN	• .1272 \pm .004	• .4001 \pm .009	• .8281 \pm .020	• .0801 \pm .004	• .9217 \pm .004	• .8489 \pm .004	• .2053 \pm .030	• 36.18% \pm .016
LDLSF	• .1091 \pm .003	• .4020 \pm .011	• .8224 \pm .021	• .0874 \pm .020	• .9402 \pm .003	• .8591 \pm .004	• .3358 \pm .025	• 37.61% \pm .020
SA-BFGS	• .1151 \pm .004	• .3706 \pm .009	• .7827 \pm .020	• .0627 \pm .003	• .9384 \pm .003	• .8599 \pm .004	• .3001 \pm .029	• 39.68% \pm .018
LCLR	• .1061 \pm .004	• .3627 \pm .009	• .7531 \pm .020	• .0579 \pm .003	• .9438 \pm .003	• .8669 \pm .004	• .3573 \pm .025	• 43.42% \pm .020
LDLLC	• .1078 \pm .003	• .3581 \pm .009	• .7476 \pm .018	• .0573 \pm .003	• .9439 \pm .003	• .8671 \pm .003	• .3441 \pm .025	• 43.49% \pm .016
DPA	• .1077 \pm .003	• .3580 \pm .009	• .7474 \pm .019	• .0572 \pm .003	• .9440 \pm .003	• .8672 \pm .003	• .3446 \pm .027	• 43.56% \pm .017
LRR	• .1078 \pm .003	• .3576 \pm .009	• .7470 \pm .018	• .0571 \pm .003	• .9440 \pm .003	• .8672 \pm .003	• .3421 \pm .027	• 43.56% \pm .017
SCL	• .1111 \pm .004	• .3512 \pm .008	• .7523 \pm .018	• .0573 \pm .003	• .9436 \pm .003	• .8658 \pm .003	• .2952 \pm .027	• 43.61% \pm .017
LDLF	• .1023 \pm .003	.3310 \pm .008	• .6940 \pm .017	• .0512 \pm .003	• .9492 \pm .003	• .8760 \pm .003	• .4045 \pm .026	• 48.82% \pm .018
DF-LDL	• <u>.1001</u> \pm .003	<u>.3276</u> \pm .007	<u>.6817</u> \pm .016	.0493 \pm .003	<u>.9511</u> \pm .003	<u>.8780</u> \pm .003	<u>.4273</u> \pm .022	• <u>49.93%</u> \pm .018
δ -LDL	.0986 \pm .004	.3272 \pm .008	.6794 \pm .018	<u>.0494</u> \pm .003	.9513 \pm .003	.8793 \pm .003	.4345 \pm .029	50.73% \pm .018

LRR (Jia et al., 2023b); on the other hand, methods based on specified structures are also compared, including: PT-Bayes (Geng, 2016), AA- k NN (Geng, 2016), LDLF (Shen et al., 2017) and DF-LDL (González et al., 2021). 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.

5.2. Results and discussion

Label distribution prediction Table 2 shows the representative label distribution prediction results. • (◦) indicates “ δ -LDL is statistically superior (inferior) to the comparing methods” (pairwise t -test at 0.05 significance level); there is no significant if neither • nor ◦ is shown; the best and second-best results are highlighted in **bold** and underline, respectively. In terms of performance on M^2B , a considerable number of state-of-the-art methods are even inferior to

AA- k NN, exposing the limitations of minimizing average KLD; in contrast, δ -LDL achieves the best performance on five metrics, significantly outperforming the other methods. Note that the KLD evaluation of δ -LDL on M^2B is 0.4882, which is lower than the stochastic performance, i.e., 0.6335 in Table 1, which fully demonstrates the effectiveness of our proposed δ -LDL. The results on fbp5500 and SBU_3DFE are similar, which implies that δ -LDL is effective on different datasets. Remaining results can be found in Appendix B.

Regarding the μ metric, it is always significant on t -test as our expectation. Note that the ranking derived from the μ metric may be different from that derived from the KLD. An obvious example lies in the comparison between LDLF and DF-LDL on M^2B , where the KLD evaluation of LDLF is 0.5688, which is lower than that of DF-LDL, i.e., 0.5743, while the μ metric shows that DF-LDL is superior to LDLF. This is because the μ metric focuses on the percentage of

predictions that are approximately correct, which is less susceptible to extreme predictions. Compared with other metrics, the μ metric has explicit range, is more discriminable between superior and inferior models, and is more intuitive, reliable, and interpretable.

Ablation study & analysis on parameters Here we are interested in the importance of each processing step in δ -LDL, thus an ablation study on M^2B is conducted: 1) we replace the loss function with solely $\ell_0(\delta'_0; \cdot, \cdot)$ to examine the importance of ℓ , and this model is denoted as δ -LDL $_{\ell_0}$; 2) we replace the ASR with solely Simpson’s rule by setting $\xi = 0$, and this model is denoted as δ -LDL w/o ASR. The results are shown in Figure 2(a), which confirms that each processing step in δ -LDL is indispensable. We conduct an analysis on parameters on M^2B as well, where the KLD and training time w.r.t. ξ are also shown in Figure 2(a). The training time initially increases sharply with the growth of ξ , but gradually levels off. This is because, as ξ increases, more samples are able to meet the conditions for halting the computation of the ASR loss in advance. For the balance between training time and performance, ξ is fixed at 5 in our proposed δ -LDL.

Though δ -LDL is δ -parameter-free, we can still analyze the impact of δ on $\mathcal{D}(\cdot, \delta; \cdot)$ with predictions of an experiment. We analyze various $\mathcal{D}(\cdot, \delta; \cdot)$ s w.r.t. δ with predictions of one experiment on M^2B , results of which are shown in Figure 2(b) and Appendix B. The gray lines indicate the integration interval. In terms of our proposed δ -LDL, $\mathcal{D}(\text{KL}, \delta; \cdot)$ is larger than that of other methods when δ is small, which makes the area under the curve larger. It is worth noting that $\mathcal{D}(\text{KL}, \delta; \cdot)$ is a monotonically increasing function of δ , which is consistent with our theoretical analysis.

Robustness testing We now evaluate the robustness of the μ metric, exploiting the artificial dataset introduced in (Geng, 2016). First, we generate the artificial dataset with 200×200 samples. Then we continuously apply noise along a certain axis of the features, which is visualized in Figure 3(a). The KLD and μ metric w.r.t. the noise level are shown in Figure 3(b). The KLD lacks the significance of changes when the noise level is low. Then, it sharply increases and becomes highly volatile once the noise level surpasses a certain threshold. This is because when the noise increases to a certain extent, the label distribution will become sparse, triggering the sensitivity of the KLD, which is consistent with our theoretical analysis. The early tiny changes/the later substantial fluctuations in the KLD are detrimental to the discriminable/reliable evaluation of LDL models. In contrast, the μ metric remains stable and gradually decreases with the noise level. This is because the μ metric is based on the percentage of predictions that are approximately correct, which is less sensitive to noise.

Next, we evaluate the robustness of δ -LDL and other meth-

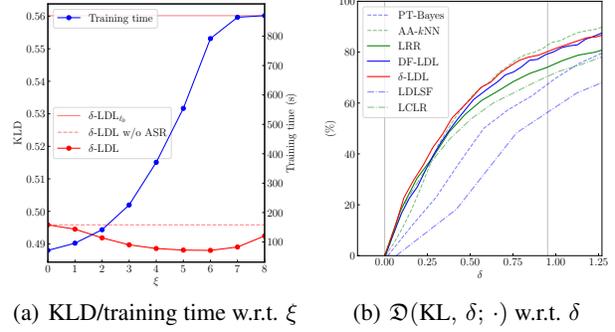


Figure 2. (Best viewed in color) Visualized results of the analysis on parameters.

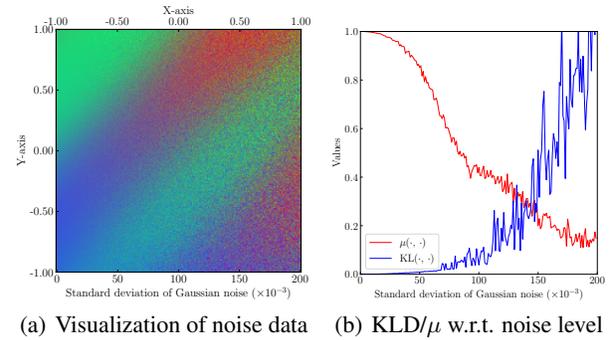


Figure 3. (Best viewed in color) Visualized results of the robustness testing of the μ metric.

ods, exploiting the same artificial dataset mentioned above. From bottom to top, from left to right, the four types of processing in Figure 4(a) 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 4(b) to 4(g), and the ground truth, i.e., the original artificial data without noise, is shown in Figure 4(h) for comparison. The consistency of each prediction result reflects the robustness of each method. Although it is hard to visualize the so-called “decision boundary” in the field of LDL, Figure 4 shows that our method predicts accurately and the prediction results are not easily affected by noise.

6. Related work

Our work is mainly related to LDL. LDL was initially proposed by Geng et al. to tackle the problem of age estimation, and subsequently developed into a novel paradigm of machine learning (Geng, 2016), which is now supported by

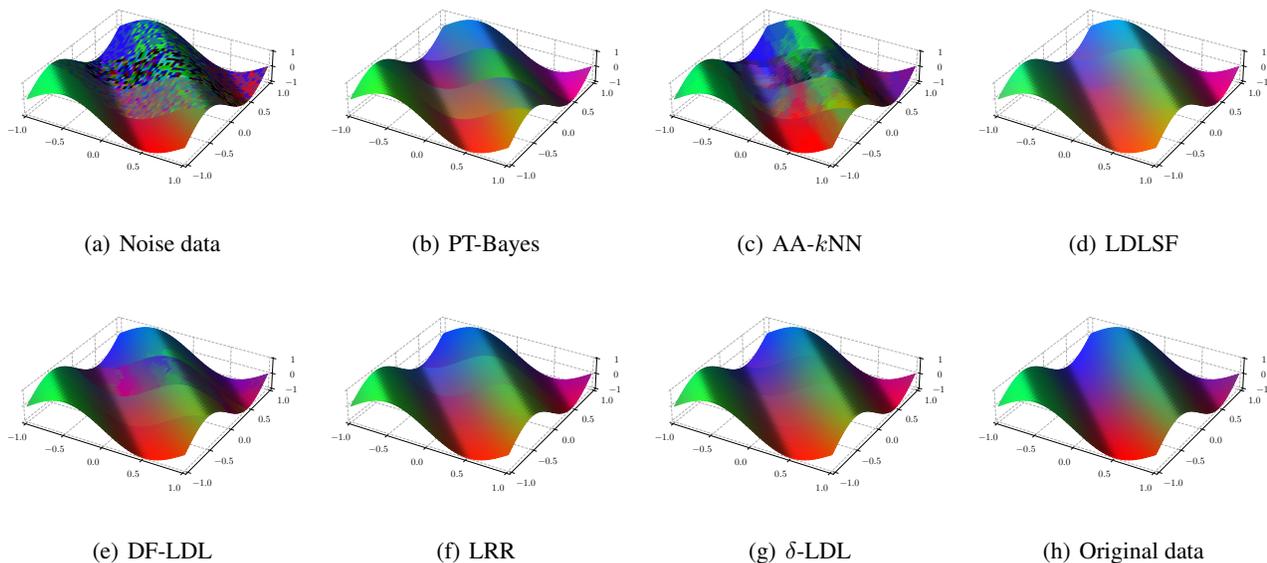


Figure 4. (Best viewed in color) Visualized results of the robustness testing of δ -LDL and other methods.

theoretical underpinnings (Wang & Geng, 2019). Over the years, a large number of specialized algorithms for LDL have emerged (Geng et al., 2013; Geng & Hou, 2015; Shen et al., 2017; González et al., 2021).

Regrading evaluation metrics, the average KLD is the most widely used one and it is first considered by Geng, who advocates hierarchical clustering to select representative metrics, and demonstrates that Clark and Canberra metrics suffer from oversensitivity to small values. Subsequent studies by Xu & Zhou revealed KLD’s unreliability with sparse predictions, prompting calls for alternative measures. Metrics that used in early work therefore directly borrows from existing mathematical foundations, e.g., $\text{Tops}\phi$ (Ren & Geng, 2017), $S\phi$ rensen (Ren et al., 2019a) and Spearman’s coefficient (Jia et al., 2023b), etc. There is little work on designing new metrics until recently, e.g., the DPA proposed by Jia et al., which, however, from the empirical results, still lacks discriminability between superior and inferior models. The CAD, QFD² and CJS proposed by Wen et al. do not address this issue either, and only work with ordinal label distributions. Our proposed μ metric, which is based on the percentage of predictions that are approximately correct, is the first to tackle this issue.

Regrading learning objectives, the average KLD and the MSE are the most widely adopted. For example, AA-BP (Geng, 2016), a simple 3-layer network, minimizes MSE but underperforms due to overfitting. The average KLD is exploited in many LDL work based on loss function engineering via different optimization methods, including BFGS

(Jia et al., 2018; 2023a), ADMM (Ren et al., 2019a;b), and gradient-based methods (Jia et al., 2019; 2023b). These methods 1) implicitly address overfitting through ad-hoc regularization, indicating the field of LDL lacks principled solutions; 2) sometimes are worse than average performance deduced by uniform vectors, which could not even be verified before our work. The essential difference between our proposed δ -LDL and these methods is that it no longer focuses on all samples, but sacrifices a part of them by using an approximately correct interval.

Furthermore, the parameter δ is potentially related to the parameter ρ in the margin theory, which is exploited in some LDL work (Wang & Geng, 2019; 2021a;b). The difference, however, is that: 1) their work focuses on the problem of objective mismatch, using margin theory to improve the performance of classification tasks, destroying the expected shape of the predicted label distributions, while our work still focuses on LDL tasks; 2) the margin ρ is still a hyperparameter that needs to be set in their methods, while our method is δ -parameter-free, which is more convenient for practical scenarios.

7. Limitations and conclusion

Limitations A notable limitation of the proposed δ -LDL algorithm is that it has not yet achieved optimal performance across all metrics, which is close to the theoretical best-case scenario. This is likely due to our current work not considering fitting different shapes of label distributions, which is an area that warrants further exploration.

Conclusion In this paper, we propose DeltaLDL, a percentage of predictions that are approximately correct in LDL. DeltaLDL can serve as a novel evaluation metric, which is parameter-free and reflects real performance improvements; DeltaLDL can also serve as a novel learning objective, which is differentiable and ensures that most samples can be predicted as approximately correct. Our theoretical analysis and empirical results demonstrate the effectiveness of the proposed solution.

Impact statement

This paper presents work whose goal is to advance the field of machine learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

Acknowledgements

This research is partially supported by Shanghai Central Guidance Science and Technology Development Fund (YDZX20233100004008), the Natural Science Foundation of Jiangsu Province (BK20242045), and the National Natural Science Foundation of China (62176123, 62476130).

References

- Chen, S., Wang, J., Chen, Y., Shi, Z., Geng, X., and Rui, Y. Label distribution learning on auxiliary label space graphs for facial expression recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13984–13993, 2020.
- Geng, X. Label distribution learning. *IEEE Transactions on Knowledge and Data Engineering*, 28(7):1734–1748, 2016.
- Geng, X. and Hou, P. Pre-release prediction of crowd opinion on movies by label distribution learning. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pp. 3511–3517, 2015.
- Geng, X., Yin, C., and Zhou, Z.-H. Facial age estimation by learning from label distributions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(10):2401–2412, 2013.
- Geng, X., Zheng, R., Lv, J., and Zhang, Y. Multilabel ranking with inconsistent rankers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):5211–5224, 2021.
- González, M., González-Almagro, G., Triguero, I., Cano, J.-R., and García, S. Decomposition-fusion for label distribution learning. *Information Fusion*, 66:64–75, 2021.
- He, L., Lu, Y., Li, W., and Jia, X. Generative calibration of inaccurate annotation for label distribution learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 12394–12401, 2024.
- Jia, X., Li, W., Liu, J., and Zhang, Y. Label distribution learning by exploiting label correlations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 3310–3317, 2018.
- Jia, X., Li, Z., Zheng, X., Li, W., and Huang, S.-J. Label distribution learning with label correlations on local samples. *IEEE Transactions on Knowledge and Data Engineering*, 33(4):1619–1631, 2019.
- Jia, X., Qin, T., Lu, Y., and Li, W. Adaptive weighted ranking-oriented label distribution learning. *IEEE Transactions on Neural Networks and Learning Systems*, 35(8):11302–11316, 2023a.
- Jia, X., Shen, X., Li, W., Lu, Y., and Zhu, J. Label distribution learning by maintaining label ranking relation. *IEEE Transactions on Knowledge and Data Engineering*, 35(02):1695–1707, 2023b.
- Kingma, D. P. and Ba, J. L. Adam: A method for stochastic optimization. In *Proceedings of the International Conference on Learning Representations*, 2015.
- Kou, Z., Wang, J., Jia, Y., Liu, B., and Geng, X. Instance-dependent inaccurate label distribution learning. *IEEE Transactions on Neural Networks and Learning Systems*, 36(1):1425–1437, 2023.
- Lee, H., Höger, F., Schoenwiesner, M., Park, M., and Jacoby, N. Cross-cultural mood perception in pop songs and its alignment with mood detection algorithms. In *Proceedings of the International Society for Music Information Retrieval Conference*, 2021.
- Li, S. and Deng, W. Blended emotion in-the-wild: Multi-label facial expression recognition using crowdsourced annotations and deep locality feature learning. *International Journal of Computer Vision*, 127(6-7):884–906, 2019.
- Li, X., Luo, G., Wang, W., Wang, K., and Li, S. Curriculum label distribution learning for imbalanced medical image segmentation. *Medical Image Analysis*, 89:102911, 2023.
- Liang, L., Lin, L., Jin, L., Xie, D., and Li, M. SCUT-FBP5500: A diverse benchmark dataset for multi-paradigm facial beauty prediction. In *Proceedings of the International Conference on Pattern Recognition*, pp. 1598–1603, 2018.
- Lu, Y. and Jia, X. Predicting label distribution from multi-label ranking. In *Proceedings of the Annual Conference*

- on *Neural Information Processing Systems*, pp. 36931–36943, 2022.
- Machajdik, J. and Hanbury, A. Affective image classification using features inspired by psychology and art theory. In *Proceedings of the ACM International Conference on Multimedia*, pp. 83–92, 2010.
- McKeeman, W. M. Algorithm 145: Adaptive numerical integration by simpson’s rule. *Communications of the ACM*, 5(12):604, 1962.
- Nguyen, T. V., Liu, S., Ni, B., Tan, J., Rui, Y., and Yan, S. Sense beauty via face, dressing, and/or voice. In *Proceedings of the ACM International Conference on Multimedia*, pp. 239–248, 2012.
- Ren, T., Jia, X., Li, W., Chen, L., and Li, Z. Label distribution learning with label-specific features. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pp. 3318–3324, 2019a.
- Ren, T., Jia, X., Li, W., and Zhao, S. Label distribution learning with label correlations via low-rank approximation. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pp. 3325–3331, 2019b.
- Ren, Y. and Geng, X. Sense beauty by label distribution learning. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pp. 2648–2654, 2017.
- Shen, W., Zhao, K., Guo, Y., and Yuille, A. Label distribution learning forests. In *Proceedings of the Annual Conference on Neural Information Processing Systems*, pp. 834–843, 2017.
- Wang, J. and Geng, X. Classification with label distribution learning. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pp. 3712–3718, 2019.
- Wang, J. and Geng, X. Label distribution learning machine. In *Proceedings of the International Conference on Machine Learning*, pp. 10749–10759, 2021a.
- Wang, J. and Geng, X. Learn the highest label and rest label description degrees. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pp. 3097–3103, 2021b.
- Wen, C., Zhang, X., Yao, X., and Yang, J. Ordinal label distribution learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 23481–23491, 2023.
- Xu, H., Liu, X., Zhao, Q., Ma, Y., Yan, C., and Dai, F. Gaussian label distribution learning for spherical image object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1033–1042, 2023.
- Xu, M. and Zhou, Z.-H. Incomplete label distribution learning. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pp. 3175–3181, 2017.
- Yin, L., Wei, X., Sun, Y., Wang, J., and Rosato, M. J. A 3D facial expression database for facial behavior research. In *Proceedings of the International Conference on Automatic Face and Gesture Recognition*, pp. 211–216, 2006.
- Zhang, M.-L. and Zhou, Z.-H. A review on multi-label learning algorithms. *IEEE Transactions on Knowledge and Data Engineering*, 26(8):1819–1837, 2013.

A. Proof of Proposition 2.2

Proof. Let $\mathbf{u} = (u_1, u_2, \dots, u_c)$ be a random vector distributed according to a Dirichlet distribution, i.e., $\mathbf{u} \sim \text{Dir}(\boldsymbol{\alpha})$, with parameters $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_c)$. The KLD between \mathbf{u} and \mathbf{v} is given by $\text{KL}(\mathbf{u} \parallel \mathbf{v}) = \sum_{j=1}^c u_j \log(cu_j)$. Then the expectation of the KLD can be decomposed as

$$\mathbb{E}(\text{KL}(\mathbf{u} \parallel \mathbf{v})) = \mathbb{E}\left(\sum_{j=1}^c u_j \log(cu_j)\right) = \sum_{j=1}^c \mathbb{E}(u_j \log(cu_j)) = \sum_{j=1}^c (\mathbb{E}(u_j \log u_j) + \mathbb{E}(u_j) \log c). \quad (22)$$

From the properties of the Dirichlet distribution, the expected value of u_j is

$$\mathbb{E}(u_j) = \frac{\alpha_j}{\sum_{k=1}^c \alpha_k}. \quad (23)$$

Since $\mathbf{u} \sim \text{Dir}(\boldsymbol{\alpha})$, we have $u_j \sim \text{Beta}(\alpha_j, \sum_{k=1}^c \alpha_k - \alpha_j)$, $\forall j \in [c]$. From the properties of the Beta distribution, we have

$$\mathbb{E}(u_j \log u_j) = \frac{\alpha_j}{\sum_{k=1}^c \alpha_k} \left(\psi(\alpha_j + 1) - \psi\left(\sum_{k=1}^c \alpha_k + 1\right) \right). \quad (24)$$

Combining Equations (22) to (24), the proof is complete. \square

B. Remaining experimental results

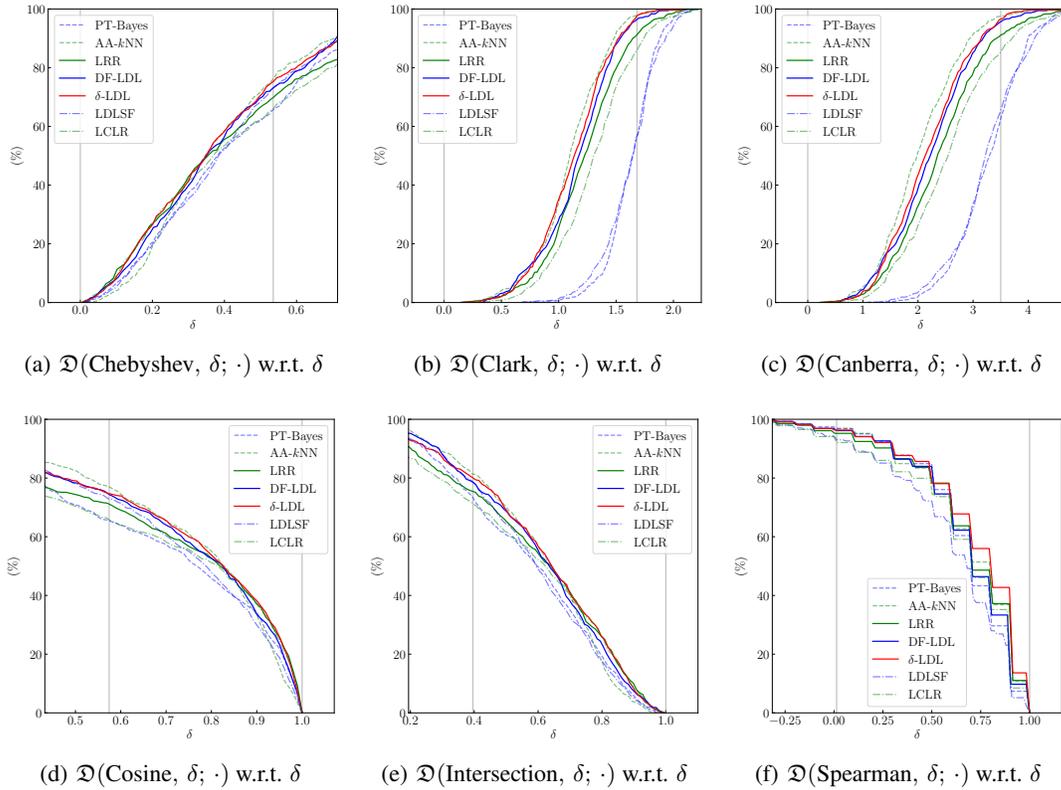


Figure 5. (Best viewed in color) Visualized results of the analysis on parameters.

Approximately Correct Label Distribution Learning

Table 3. Experimental results on RAF_ML, Natural_Scene, Music, and Painting, formatted as “mean \pm std”

Algorithms	Cheby. \downarrow	Clark \downarrow	Can. \downarrow	KLD \downarrow	Cosine \uparrow	Int. \uparrow	Spear. \uparrow	μ \uparrow
RAF_ML								
SCL	.3566 \pm .006	1.5754 \pm .012	3.4061 \pm .031	.6360 \pm .011	.7080 \pm .006	.5337 \pm .004	.3372 \pm .025	20.95% \pm .006
PT-Bayes	.3482 \pm .022	1.7253 \pm .033	3.7506 \pm .106	.8545 \pm .123	.7346 \pm .024	.5983 \pm .023	.5575 \pm .028	32.15% \pm .039
LDLSF	.1885 \pm .004	1.5429 \pm .013	3.0619 \pm .034	.8256 \pm .069	.9102 \pm .005	.7672 \pm .005	.7713 \pm .010	48.58% \pm .017
AA-kNN	.2643 \pm .005	1.4882 \pm .014	3.0326 \pm .036	.4073 \pm .014	.8221 \pm .006	.6737 \pm .006	.6041 \pm .013	49.37% \pm .012
LDLF	.1974 \pm .008	1.4077 \pm .017	2.7417 \pm .045	.2444 \pm .014	.9071 \pm .006	.7569 \pm .009	.7494 \pm .012	68.14% \pm .017
DF-LDL	.1546 \pm .005	1.4196 \pm .015	2.7375 \pm .036	.1957 \pm .011	.9273 \pm .005	.8087 \pm .005	.7817 \pm .010	75.68% \pm .010
SA.BFGS	.1538 \pm .005	1.4199 \pm .015	2.7385 \pm .036	.1957 \pm .011	.9277 \pm .005	.8096 \pm .005	.7840 \pm .010	75.77% \pm .010
LCLR	.1538 \pm .005	1.4199 \pm .015	2.7386 \pm .036	.1957 \pm .011	.9277 \pm .005	.8096 \pm .005	.7839 \pm .010	75.77% \pm .010
LDLLC	.1538 \pm .005	1.4178 \pm .015	2.7335 \pm .036	.1956 \pm .011	.9276 \pm .005	.8096 \pm .005	.7842 \pm .010	75.79% \pm .010
LRR	.1538 \pm .005	1.4185 \pm .015	2.7356 \pm .037	.1957 \pm .011	.9276 \pm .005	.8096 \pm .005	.7841 \pm .010	75.79% \pm .010
DPA	.1538 \pm .005	1.4185 \pm .016	2.7355 \pm .037	.1957 \pm .011	.9276 \pm .005	.8096 \pm .005	.7841 \pm .010	75.80% \pm .010
δ -LDL	.1527 \pm .005	1.4159 \pm .016	2.7267 \pm .038	.1929 \pm .011	.9284 \pm .005	.8108 \pm .006	.7847 \pm .010	76.06% \pm .010
Natural_Scene								
PT-Bayes	.7469 \pm .019	2.9025 \pm .008	8.5810 \pm .033	19.3376 \pm .971	.3446 \pm .022	.2416 \pm .019	.3467 \pm .024	10.19% \pm .017
SCL	.3486 \pm .017	2.4728 \pm .020	6.8978 \pm .083	.9959 \pm .028	.6647 \pm .008	.4180 \pm .010	.3368 \pm .041	25.11% \pm .007
LDLSF	.3185 \pm .012	2.4948 \pm .017	6.8440 \pm .071	2.6764 \pm .348	.7193 \pm .012	.5493 \pm .011	.4718 \pm .016	35.02% \pm .020
LCLR	.3590 \pm .016	2.3688 \pm .018	6.5106 \pm .072	1.1351 \pm .093	.6681 \pm .018	.5258 \pm .014	.4697 \pm .017	35.33% \pm .024
SA.BFGS	.3323 \pm .015	2.3921 \pm .019	6.5725 \pm .073	.9239 \pm .070	.7016 \pm .016	.5480 \pm .014	.4898 \pm .016	39.00% \pm .025
DF-LDL	.3067 \pm .013	2.3698 \pm .019	6.4633 \pm .074	.7828 \pm .041	.7323 \pm .013	.5670 \pm .012	.5092 \pm .016	42.55% \pm .022
AA-kNN	.2888 \pm .014	1.7886 \pm .038	4.1708 \pm .110	1.0107 \pm .081	.7391 \pm .015	.5969 \pm .014	.5319 \pm .020	42.71% \pm .024
LDLLC	.2994 \pm .012	2.4214 \pm .018	6.6168 \pm .076	.7327 \pm .035	.7504 \pm .013	.5718 \pm .011	.5251 \pm .015	44.30% \pm .018
DPA	.2994 \pm .012	2.4207 \pm .019	6.6152 \pm .079	.7327 \pm .035	.7500 \pm .012	.5722 \pm .011	.5252 \pm .015	44.31% \pm .018
LRR	.2991 \pm .012	2.4210 \pm .018	6.6161 \pm .076	.7322 \pm .035	.7499 \pm .013	.5722 \pm .011	.5254 \pm .015	44.33% \pm .018
LDLF	.2842 \pm .012	2.4188 \pm .018	6.5960 \pm .074	.6787 \pm .037	.7706 \pm .014	.5899 \pm .012	.5473 \pm .017	47.47% \pm .021
δ -LDL	.2760 \pm .012	2.4123 \pm .017	6.5501 \pm .075	.6586 \pm .033	.7807 \pm .012	.5974 \pm .011	.5553 \pm .017	48.96% \pm .018
Music								
PT-Bayes	.1733 \pm .026	1.1169 \pm .101	2.8681 \pm .292	3.290 \pm .082	.7818 \pm .040	.6930 \pm .034	.2079 \pm .105	04.85% \pm .026
LDLSF	.0946 \pm .006	.9498 \pm .072	2.4003 \pm .194	.6168 \pm .306	.8718 \pm .018	.7678 \pm .019	.3123 \pm .072	23.89% \pm .045
SA.BFGS	.1020 \pm .008	.8802 \pm .060	2.2770 \pm .172	.1912 \pm .030	.8665 \pm .018	.7673 \pm .018	.3030 \pm .072	25.57% \pm .047
DF-LDL	.1022 \pm .008	.8489 \pm .057	2.1694 \pm .160	.1719 \pm .026	.8726 \pm .018	.7763 \pm .018	.2688 \pm .080	26.18% \pm .045
LCLR	.1000 \pm .008	.8644 \pm .064	2.2295 \pm .182	.1821 \pm .034	.8738 \pm .019	.7733 \pm .020	.3673 \pm .073	26.23% \pm .049
DPA	.0971 \pm .007	.8532 \pm .059	2.2104 \pm .166	.1758 \pm .027	.8746 \pm .017	.7741 \pm .018	.3224 \pm .072	27.04% \pm .047
LRR	.0969 \pm .007	.8519 \pm .058	2.2061 \pm .166	.1751 \pm .027	.8750 \pm .017	.7745 \pm .018	.3214 \pm .071	27.11% \pm .047
LDLLC	.0968 \pm .007	.8518 \pm .057	2.2063 \pm .160	.1751 \pm .027	.8751 \pm .017	.7745 \pm .017	.3212 \pm .070	27.12% \pm .046
SCL	.0776 \pm .004	.7676 \pm .050	1.9339 \pm .131	.1157 \pm .013	.9150 \pm .008	.8065 \pm .012	.4544 \pm .055	32.12% \pm .028
LDLF	.0775 \pm .005	.7532 \pm .054	1.9050 \pm .140	.1134 \pm .014	.9161 \pm .009	.8096 \pm .014	.4538 \pm .078	34.01% \pm .045
AA-kNN	.0784 \pm .006	.7499 \pm .061	1.9309 \pm .172	.1258 \pm .021	.9067 \pm .015	.8051 \pm .018	.3989 \pm .081	36.75% \pm .045
δ -LDL	.0733 \pm .005	.7185 \pm .055	1.8221 \pm .150	.1061 \pm .016	.9218 \pm .011	.8186 \pm .015	.4942 \pm .058	39.35% \pm .036
Painting								
PT-Bayes	.7070 \pm .056	2.5747 \pm .056	7.0292 \pm .221	7.4050 \pm 1.412	.3337 \pm .055	.2275 \pm .044	.1484 \pm .075	00.07% \pm .003
LDLSF	.2853 \pm .023	1.8290 \pm .067	4.3005 \pm .204	4.1107 \pm 1.205	.6565 \pm .033	.5507 \pm .025	.2047 \pm .067	11.89% \pm .040
LCLR	.2979 \pm .028	1.8632 \pm .070	4.5308 \pm .229	.9065 \pm .215	.6579 \pm .036	.5538 \pm .028	.2454 \pm .066	16.23% \pm .045
SA.BFGS	.2965 \pm .028	1.8559 \pm .073	4.5050 \pm .237	.8997 \pm .233	.6596 \pm .036	.5547 \pm .028	.2363 \pm .069	16.37% \pm .044
DF-LDL	.2714 \pm .021	1.8051 \pm .064	4.3184 \pm .202	.6950 \pm .096	.6772 \pm .032	.5746 \pm .023	.2310 \pm .073	17.13% \pm .044
SCL	.2579 \pm .019	1.6915 \pm .067	3.9339 \pm .193	.5574 \pm .042	.7223 \pm .017	.5940 \pm .019	.1987 \pm .078	19.36% \pm .024
AA-kNN	.2572 \pm .021	1.7553 \pm .060	4.1427 \pm .182	.7100 \pm .161	.6987 \pm .027	.5888 \pm .021	.2241 \pm .072	19.76% \pm .041
DPA	.2585 \pm .020	1.7561 \pm .069	4.1452 \pm .215	.5999 \pm .067	.7116 \pm .028	.5951 \pm .022	.2780 \pm .069	20.48% \pm .041
LRR	.2587 \pm .019	1.7556 \pm .068	4.1422 \pm .212	.5988 \pm .064	.7118 \pm .027	.5954 \pm .022	.2782 \pm .068	20.48% \pm .040
LDLLC	.2586 \pm .020	1.7560 \pm .069	4.1438 \pm .216	.5994 \pm .066	.7120 \pm .027	.5953 \pm .022	.2774 \pm .068	20.48% \pm .041
LDLF	.2485 \pm .019	1.7059 \pm .066	3.9700 \pm .200	.5343 \pm .046	.7336 \pm .023	.6075 \pm .021	.3054 \pm .088	22.54% \pm .036
δ -LDL	.2484 \pm .020	1.7019 \pm .067	3.9520 \pm .201	.5288 \pm .048	.7377 \pm .023	.6102 \pm .022	.3146 \pm .092	23.05% \pm .037