LEARNING LABEL DISTRIBUTION WITH SUBTASKS

Anonymous authors

Paper under double-blind review

ABSTRACT

Label distribution learning (LDL) is a novel learning paradigm that emulates label polysemy by assigning label distributions over the label space. However, recent LDL work seems to exhibit a notable contradiction: 1) some existing LDL methods employ auxiliary tasks to enhance performance, which narrows their focus to specific domains, thereby lacking generalization capability; 2) conversely, LDL methods without auxiliary tasks rely on losses tailored solely to label distributions of the primary task, lacking additional supervised information to guide the learning process. In this paper, we propose S-LDL, a novel and minimalist solution that partitions the label distribution of the primary task into subtask label distributions, i.e., a form of pseudo-supervised information, to reconcile the above contradiction. S-LDL encompasses two key aspects: 1) an algorithm capable of generating subtasks without any extra knowledge, with subtasks deemed valid and reconstructable via our analysis; and 2) a plug-and-play framework seamlessly compatible with existing LDL methods, and even adaptable to derivative tasks of LDL. Experiments demonstrate that S-LDL is effective and efficient. To the best of our knowledge, this represents the first endeavor to address LDL via subtasks. The code will soon be available on GitHub to facilitate reproducible research.

000

001

004

006 007

008 009

010

011

012

013

014

015

016

017

018

019

021

023

1 INTRODUCTION

028 Multi-label learning (MLL) (Zhang and Zhou, 2013) handles label polysemy in a binary manner, 029 whereas label distribution learning (LDL) (Geng, 2016) offers a more nuanced perspective by answering: "How much does each label y describe the instance x?". This is accomplished through 031 the concept of a *label distribution d*, which is a form of probability simplex that assigns a real value (i.e., description degree d_{α}^{y}) to each label of each instance. This form introduces a quantitative manner 033 to address label polysemy and extends LDL's practical applications to a wider range, e.g., counting 034 (or grading) (Geng et al., 2013; Wu et al., 2019), sentiment analysis (Chen et al., 2020; Le et al., 2023), segmentation (Gao et al., 2017; Li et al., 2023b), etc. Concurrently, more and more derivative tasks of LDL (González et al., 2021b; Lu and Jia, 2022; Wang, Jing and Geng, Xin, 2019; Xu and 037 Zhou, 2017; Xu et al., 2019) are emerging to offer assistance in various real-world dilemmas.

However, LDL encounters a spectrum of challenges: 1) label distributions are bound by two constraints, non-negativity (i.e., $d_x^y \ge 0$) and sum-to-one (i.e., $\sum_{y \in \mathcal{Y}} d_x^y = 1$), and are often formed from mixture distributions, posing significant hurdles for fitting, particularly when employing a maximum entropy model (Shen et al., 2017); 2) label distribution matrices are usually obtained via crowdsourcing, which is time-consuming and labor-intensive, so one often copes with scarce and low-quality datasets (Wang et al., 2023). These two key issues stand as formidable barriers to performance improvement in LDL.

With the widespread use of multi-task learning, some LDL work tries to compensate for performance from the perspective of *auxiliary tasks*, which are learned concurrently alongside the primary task, thereby refining its representations and ultimately boosting performance. Unfortunately, though these methods can exploit additional supervised information, they 1) do not address the first key issue mentioned above; and 2) require extra knowledge (e.g., facial characteristics (Chen et al., 2020), pathology criteria (Wu et al., 2019), emotion wheel theory in psychology (Yang et al., 2017a), etc.)
or similar domain-specific data (Zhao et al., 2023b), limiting their generalization capability to those corresponding specific domains. Conversely, LDL methods that do not take advantage of auxiliary tasks, despite their efforts in loss function engineering and network structure design, they 1) do not address the second key issue mentioned above; and 2) focus solely on one aspect of label correlations

(e.g., correlation of local instances (Jia et al., 2019), ranking relation (Jia et al., 2023), suboptimal label (Wang, Jing and Geng, Xin, 2019), etc.), each with its own set of limitations.

The generalizability across various domains appears to conflict with the ability to exploit additional data, so benefiting from both simultaneously seems elusive. However, we can still see the light from some MLL methods, which partition the label space and apply operations on these subspaces (Tsoumakas et al., 2008; 2010). These methods construct *subtasks* without involving extra knowledge and exhibit applicability across various domains. Intuitively, in the context of LDL, reliable supervised information can be generated from these subtasks, which can eventually be aggregated and reconstructed to the information of the primary task via ensemble strategies. Although existing label distribution ensemble practices demonstrate promising performance (González et al., 2021a; Shen et al., 2017), they focus only on the supervised information of the primary task.

In this paper, we introduce S-LDL, a novel and minimalist label distribution learning algorithm that constructs and exploits subtasks, to reconcile the contradiction between the generalizability across various domains and the ability to exploit additional data. Serving as auxiliary tasks, subtasks 1)
provide different views of the primary task distribution, rendering the mixture of distributions more traceable (i.e., the key issue one); 2) furnish additional supervised data to mitigate the scarcity and ambiguity inherent in LDL datasets (i.e., the key issue two); 3) require no extra knowledge from specific domains; and 4) emphasize various label correlations via partitioning of the label space.

The main contributions of this paper are outlined below: 1) we propose S-LDL, which is considered the first endeavor to address LDL via subtasks; 2) our analysis shows the validity and reconstructability of these subtasks; 3) we present a plug-and-play framework seamlessly compatible with existing LDL methods, and adaptable to derivative tasks of LDL; and 4) the code will be available on GitHub soon, facilitating reproducible research endeavors.

077 078

2 RELATED WORK

079

LDL Our work is mainly related to LDL. Initially employed to tackle age estimation (Geng et al., 081 2013), LDL has evolved into a novel machine learning paradigm (Geng, 2016), which is supported 082 by theoretical underpinnings (Wang and Geng, 2019) and features various derivative tasks (González 083 et al., 2021b; Lu and Jia, 2022; Wang, Jing and Geng, Xin, 2019; Xu and Zhou, 2017; Xu et al., 2019). 084 Most methods focus on improving performance via loss function engineering (Jia et al., 2019; 2023; 085 Ren et al., 2019; Wen et al., 2023) or efficient model structures (González et al., 2021a; Jin et al., 086 2024; Shen et al., 2017; Yang et al., 2017b), while some work is dedicated to practical application scenarios (Gao et al., 2017; Li et al., 2023a; Shirani et al., 2019; Wu et al., 2019). However, the 087 880 scarcity of label distribution datasets and the complexity of the label distribution itself make it difficult to further improve performance, at which point one may think of leveraging auxiliary tasks. 089

090

LDL with auxiliary tasks While there are LDL methods that leverage auxiliary tasks to enhance 091 performance, they often rely on knowledge from disparate domains, extending beyond the scope of 092 the LDL task. For example, LDL-ALSG (Chen et al., 2020) designs auxiliary tasks dedicated to facial 093 emotion recognition, necessitating the use of external tools to extract facial points and action units 094 from human faces. Wu et al. (2019) exploit the Hayashi criterion, a rule for counting and grading in 095 acne lesions, which results in their method being only applicable in a small branch of the dermatology 096 field. Yang et al. (2017a) employ a multi-task framework for image emotion classification, designing 097 constraints inspired by Mikel's wheel, a psychological emotion model, which also suffers from 098 similar limitations. As LDL methods of transfer learning, GLDL (Zhao et al., 2023b) utilizes data 099 from one or more source domains, which is not easy to obtain in practical applications. The need for 100 specific extra knowledge significantly narrows the application scenarios of these methods.

101

MLL with partitioning of the label space For reference, there exist MLL methods based on partitioning of the label space, which can construct multi-label subtasks without involving additional knowledge and can be widely used in various domains. The most classic related work is that of HOMER (Tsoumakas et al., 2008) and RAkEL (Tsoumakas et al., 2010), the former forms a hierarchy of label subspaces while the latter randomly selects label subspaces. Many subsequent papers have been inspired by them (Prabhu et al., 2018; Read et al., 2013; Wang et al., 2021). Read et al. (2014) present a general framework of label subspaces and provide some theoretical justification for it. Since

Table 1: Key notation and terminology in this paper									
Symbol	Description	Example							
$\mathcal{Y} = \{y_j\}_{j=1}^L$	Label space (L labels)	$\mathcal{Y} = \{\texttt{HA}, \texttt{SA}, \texttt{SU}, \texttt{AN}, \texttt{DI}, \texttt{FE}\}^1$							
${\cal Y}^\circ$	Subtask label spaces	$\{\cdots, \mathcal{Y}^{(t)} = \{ extsf{HA}, extsf{SA}, extsf{SU}, extsf{FE}\}, \cdots\}$							
$d_{oldsymbol{x}_i}^{y_j}$	Description degree of x_i about y_j	$d^{y_0}_{oldsymbol{x}_i}=0.4,$ i.e., HA describes $oldsymbol{x}_i$ by 0.4							
$\boldsymbol{d}_i = (d_{\boldsymbol{x}_i}^{y_j})_{j=1}^L$	Label distribution of x_i	$\boldsymbol{d}_i = (0.4, 0.05, 0.3, 0.1, 0.1, 0.05)$							
$oldsymbol{D} = (oldsymbol{d}_i)_{i=1}^N = (oldsymbol{d}_{ullstyle j})_{j=1}^L$	Distribution matrix (N samples)	$(\cdots, oldsymbol{d}_i, \cdots)$							
$d_i^{(t)} = (d_{x_i}^{(t)y_j})_{j=1}^{ \mathcal{Y}^{(t)} }$	Subtask label distribution	$\boldsymbol{d}_{i}^{(t)} = (0.5, 0.0625, 0.375, 0.0625)$							
\mathcal{D}°	Subtask distribution matrices	$\{\cdots, oldsymbol{D}^{(t)}, \cdots\}$							
$\boldsymbol{M} = (\boldsymbol{m}_t)_{t=1}^T = (M_{tj})$	Mask matrix (T anticipated tasks)	$(\cdots, \boldsymbol{m}_t = (1, 1, 1, 0, 0, 1), \cdots)$							

108

label distribution contains rich knowledge, we can follow the patterns of these methods to construct label distribution subtasks.

122 123

LDL with ensemble strategy It is imperative to aggregate the output of subtasks. Fortunately, 124 ensemble-based LDL methods have demonstrated promising performance. For instance, LDLFs 125 (Shen et al., 2017) learns different label distributions on the leaf nodes of differentiable decision 126 trees and learns weights that aggregate these label distributions. DF-LDL (González et al., 2021a) 127 aggregates the label distribution of output of multiple base models by simple averaging, while Zhai 128 et al. (2018) focus on aggregating the results of various neural networks via a combining learner. 129 However, 1) the above methods are not suitable for incomplete label spaces (i.e., subtask label 130 spaces); and 2) none of them involve the partitioning of the label space, therefore no extra supervised 131 information of label distributions is constructed.

Drawing from the analysis of the aforementioned related work, we introduce S-LDL, which leverages
 pseudo-supervised information from subtasks to eliminate reliance on additional knowledge from
 disparate domains, and facilitates the creation of a novel knowledge dimension in a generic framework.

136 137

138 139

3 SUBTASK CONSTRUCTION

3.1 PRELIMINARY

Notation Vectors are denoted by lowercase bold letters, e.g., 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., A. The row vector a_i indicates its i-th row and the column vector $a_{\bullet j}$ indicates its j-th column. A_{ij} is the element in i-th row and j-th column of A. The superscript (t) indicates that a symbol corresponds to the t-th subtask. Table 1 outlines the key notation in this paper.

146 **Problem definition** Let $x \in \mathbb{R}^P$ denote the feature of the instance and $d \in \Delta^{L-1}$ denote the label 147 distribution, where $\Delta^{k-1} \triangleq \{v \in \mathbb{R}^k \mid \mathbf{1}v^T = 1, v \ge 0\}$ is the (k-1)-dimensional probability 148 simplex. The goal of LDL is to find a mapping $\zeta : x \mapsto d$. In this paper, we partition the label 149 space \mathcal{Y} corresponding to d to obtain the subtask label space set \mathcal{Y}° , then accordingly generate 150 pseudo-supervised information, i.e., subtask distribution matrix set \mathcal{D}° , to guide the learning of ζ .

151

Technical challenges Our first challenge arises from *the exponential growth in partitions* as the 152 number of labels increases (Tsoumakas et al., 2010). When generating T tasks from a label space 153 with L labels, the number of unique partitions is given by $(2^L - L - 2)!/(T!(2^L - L - 2 - T)!)$. This makes it 154 impractical to calculate metric for each case to select subtasks. We tackle this challenge in a mask 155 matrix learning manner. The second challenge lies in discerning reasonable partitions. Since the 156 label distribution matrix is usually imbalanced in average description degree (Zhao et al., 2023a), 157 some partitions exhibit unreasonable local ignorance. As a result, the corresponding spaces struggle 158 to handle the majority of instances, because 1) theoretically, there is no objective standard for the 159 degree of negative correlation; and 2) empirically, weakly or negatively correlated information is

¹HA, SA, SU, AN, DI, FE, and NE represent the seven common emotions in sentiment analysis datasets, namely happiness, sadness, surprise, anger, disgust, fear, and neutral, respectively.



Figure 1: (a) is sourced from the emotion6 (Yang et al., 2017b) dataset, which has only 7 labels, but
the number of potential partitions is huge. (b) exemplifies a subtask label space {FE, AN, DI}, which
is challenging to describe (a). (c)'s two subtask label spaces encompass all descriptive information,
meaning no new knowledge is generated about (a). This example vividly illustrates the limitations
that may result from local ignorance and lack of diversity in subtask label spaces.

easily overlooked by human annotators in crowdsourced datasets. To mitigate this, we incorporate the description degree as a metric for the reliability of supervised information in guiding the generation of subtask masks. The third challenge is *avoiding analogous pseudo-supervised information*, i.e., to generate label distributions containing new knowledge (González et al., 2021b). This necessitates fostering richness and diversity in both subtask label spaces and distributions. To achieve this objective, we 1) minimize pairwise similarity among subtask masks; and 2) normalize each subtask label distribution to yield brand new insights distinct from the label distribution of the primary task. Fig. 1 portraits an illustrative example of these challenges.

3.2 LEARNING SUBTASK MASKS

Let $M \in \{0, 1\}^{T \times L}$ denote the subtask mask matrix, where T represents the number of anticipated tasks. To ensure that the subtask label spaces contain as reliable information as possible, the learning of the subtask mask matrix can be converted into this problem: $\arg \max_M \|DM^{\top}\|_{F}$.

Obviously, a senseless solution is $m_t = 1$ where $t = 1, \dots, T$, i.e., all pseudo-supervised information is equivalent to the primary task information. Therefore, solving the above problem alone is inappropriate. To address this, we consider pairwise similarity among subtask masks. We also employ exponential tricks to convert maximization into minimization. Finally, M is calculated as

197

203

204

205

206 207

208 209

210

211 212

213 214

215

178

187

188

$$\boldsymbol{M}^{*} = \operatorname*{arg\,min}_{\boldsymbol{M}} \left(\frac{1}{NT} \sum_{t=1}^{T} \sum_{i=1}^{N} \exp\left(-\boldsymbol{d}_{i}\boldsymbol{m}_{t}^{\top}\right) + \frac{2\lambda}{(T(T-1))} \sum_{i,j,i\neq j} \frac{\boldsymbol{m}_{i}\boldsymbol{m}_{j}^{\top}}{\|\boldsymbol{m}_{i}\| \|\boldsymbol{m}_{j}\|} \right), \quad (1)$$

s.t.
$$M_{tj} \in \{0, 1\}; t = 1, \cdots, T; j = 1, \cdots, L,$$

where λ is a trade-off parameter. Eq. (1) is slightly more complicated than conventional integer programming. For convenience, we solve it using the stochastic gradient descent (SGD) method, with its constraint enforced via sigmoid (a conversion threshold is set, where outputs greater than it are set to 1, while those below are set to 0). Refer to Section 4.1 for an analysis of the validity of Eq. (1).

3.3 GENERATING SUBTASK DISTRIBUTIONS

We slice the label distribution matrix according to the subtask label space. To generate diversified subtask label distributions, we perform normalization on each subtask distribution with

$$[\mathcal{N}_{\text{SUM}}(\boldsymbol{v})]_j = \frac{v_j}{\sum_{i=1}^{|\boldsymbol{v}|} v_i}.$$
(2)

²Despite the discrete label space, in the field of LDL, the label distribution is intentionally plotted as a curve, to distinguish it from the logical labels.

216 Algorithm 1 Subtask construction 217 **Input**: Input matrix D, trade-off parameter λ , anticipated number of subtasks T. 218 **Output**: Subtask distribution matrices \mathcal{D}° (with corresponding subtask label spaces \mathcal{Y}°). 219 1: Initialization: $\mathcal{Y}^{\circ} = \{\emptyset\}, \mathcal{D}^{\circ} = \{\emptyset\};$ 220 2: Calculate *M* using SGD; ▷ (Eq. (1)) 3: for t = 1 to T do 4: $\mathcal{Y}^{(t)} \leftarrow \{y_j\}$ if $M_{tj} = 1$; 221 222 if $|\mathcal{Y}^{(t)}| = L$ or $|\mathcal{Y}^{(t)}| \leq 1$ then 5: 223 6: continue; /* Ignore invalid subtask masks. */ 224 7: end if 225 /* The clip(x, a, b) function limits x to be within [a, b]. */ 8: 226 $D^{(t)} \leftarrow \operatorname{clip}(d_{\bullet i}, \varepsilon, 1)$ if $y_i \in \mathcal{Y}^{(t)}$; /* ε is a very small positive number. */ 9: 227 10: for i = 1 to N do Normalization: $d_i^{(t)} \leftarrow \mathcal{N}_{\text{SUM}}(d_i^{(t)});$ ⊳ (Eq. (2)) 11: 228 end for 12: 229 $\mathcal{Y}^{\circ} = \mathcal{Y}^{(t)} \cup \mathcal{Y}^{\circ}, \mathcal{D}^{\circ} = \boldsymbol{D}^{(t)} \cup \mathcal{D}^{\circ};$ 13: 230 14: **end for** 231 232 Algorithm 2 S-LDL (shallow regime) 233 **Input**: Feature matrix X, label distribution matrix D, testing instance x'. 234 **Output**: Predicted label distribution d' for instance x'. 235 1: Initialize parameter of each estimator; 236 2: $\mathcal{D}^{\circ} \leftarrow SC(\boldsymbol{D});$ 237 3: for t = 1 to $|\mathcal{D}^{\circ}|$ do 238 Fit an estimator $f^{(t)}$ on dataset $\{X, D^{(t)}\}$; 4: 239 $\boldsymbol{d}^{(t)\prime} \leftarrow f^{(t)}(\boldsymbol{x}');$ 5: 6: end for 7: Concatenate X and all of the $D^{(t)}$ s to get Z, where $t = 1, \dots, |\mathcal{D}^{\circ}|$; 241 8: Fit an estimator f on dataset $\{Z, D\}$; 242 9: Concatenate x' and all of the $d^{(t)'}$ s to get z', where $t = 1, \dots, |\mathcal{D}^{\circ}|$; 243 10: $\boldsymbol{d}' \leftarrow f(\boldsymbol{z}');$

244 245 246

247

248

249 250

251

253

254

255

256

257

258

259

The rationale for utilizing N_{SUM} as the normalization function can be found in Section 4.2. The overall subtask construction process, denoted by SC, is illustrated in Alg. 1. Then, one can naturally come up with an adaptive LDL pipeline based on the shallow regime, as depicted in Alg. 2.

4 ANALYSIS ABOUT SUBTASK CONSTRUCTION

In this section, we analyze the subtask construction algorithm SC by studying the following questions:

- Q_1 : Are the subtask spaces provided by Eq. (1) valid for performance improvement? Can one configure λ and T in Eq. (1) without any prior knowledge?
- Q₂: Are the subtask label distributions provided by Eq. (2) reconstructable? Can one replace Eq. (2) with other normalization functions?
- Q_3 : What is the overall time complexity of SC? Is it practical for large-scale datasets?

The validity, reconstructability, and complexity analysis are conducted for Q_1 , Q_2 , and Q_3 , respectively.

262

264

4.1 VALIDITY ANALYSIS

Eq. (1) manages the intricate task of selecting subtask spaces via λ and T. On the one hand, we strive to explain that it is useful for performance improvement to suppress local ignorance and increase diversity of each subtask label space simultaneously. On the other hand, we seek to determine the appropriate λ and T without any prior knowledge. To this end, we design the following two metrics.

Definition 1 (Information rate). We call it informative if $M_{tj} = 1$ where $t = 1, \dots, T$ and $j = 1, \dots, L$. Let I be the summation of information; we define the information rate as I/(TL).



Figure 2: Visualized results of the validity analysis. Results of (a) are the average of experiments on all datasets (introduced in the appendix), while results of (b) and (c) are on the emotion6 dataset. The blue lines in (b) and (c) represent performance without auxiliary tasks.

Definition 2 (Mask valid rate). Let $\delta(\cdot, \cdot)$ be the Kronecker delta function. For all $t = 1, \dots, T$, the following are considered counting of invalid masks: 1) $\delta(\mathbf{m}_t, \mathbf{1})$, or 2) $\delta(|\mathbf{m}_t|!, 1)$, or, 3) excluding masks in cases 1) and 2), for any remaining mask index i, $\delta(\mathbf{m}_i, \mathbf{m}_j)$, where $j = 1, \dots, i$.³ Let S be the summation of all invalid subtasks; we define the mask valid rate as 1 - S/T.

290 While it is unknown which metric is more important, we intuitively claim that higher values for both 291 metrics are likely to lead to better performance. First, we calculate the average of these two metrics 292 for all datasets with varying λ , results of which are shown in Fig. 2(a). The grey area depicts the 293 average of the two metrics, implying that the appropriate value of λ may be greater than 0.1.

294 It is important to highlight that Alg. 2 relies on naive concatenation operations and is not tied to 295 representation learning. Consequently, any performance improvement over the base estimator is 296 solely attributed to the effects of the subtask label distributions. Therefore, we employ Alg. 2 as the 297 "scaffolding" for our analysis. With λ varying and T fixed at 10, we conduct ten-fold experiments 298 repeated 10 times on the emotion 6 dataset using Alg. 2. Here, $f^{(t)}$ s and f are implemented by a 299 representative LDL method, LDSVR (Geng and Hou, 2015). We record the average Spearman's coefficient (the higher the better). The results, which are shown in Fig. 2(b), support our claim. The 300 panels from left to right display examples of subtask label spaces when the λ is 0.01, 0.05, 0.2, 1, 301 and 10, respectively. When λ is suitable, label spaces are diverse and do not have excessive local 302 ignorance; as λ decreases, label spaces tend to be homogeneous, and invalid masks account for the 303 majority; as λ increases, the local ignorance of each label space becomes significant. 304

Besides, we also study the parameter sensitivity of T with λ fixed at 0.2. Results are shown in Fig. 2(c), illustrating that having a plethora of auxiliary tasks are detrimental to performance, which may be due to overfitting.

The validity analysis demonstrates that simultaneously avoiding local ignorance and homogeneity can lead to more efficient subtask label spaces, thereby improving performance. Without any prior knowledge, λ and T are recommended to be set to 0.2 and 10, respectively. Since the validity of Eq. (1) is ensured, one may wonder about the rationality and necessity of Eq. (2).

312

281

282

283

284 285

287

288

289

4.2 RECONSTRUCTABILITY ANALYSIS314

We strive to choose a normalization function so that subtask label distributions retain more information, even efficacious enough to reconstruct the label distribution of the primary task. Theorem 1 illustrates that Eq. (2) is the only possibility.

Theorem 1. Let each subtask label space form a connected graph with its each label as a node. Then merge these graphs according to their respective labels to form \mathcal{G} . If and only if \mathcal{N}_{SUM} is used for normalization, the primary label distribution can be reconstructed from these subtask label distributions, when the following conditions are satisfied: 1) \mathcal{G} is connected; 2) \mathcal{G} covers all labels in the label space, and 3) corresponding description degrees of all cut vertices of \mathcal{G} are not zero.

³These three cases correspond to 1) masks that are exactly the same as the primary task; 2) masks that fail to form label distributions; and 3) duplicate masks among the remaining masks, respectively.

324 *Proof.* We solely discuss the extreme case where two subtask label spaces overlap with just one label. 325 Further specialized cases can be deduced by the reader via induction. With a little bit of symbol abuse, 326 let the general normalization function be defined as $\mathcal{N}(v) \triangleq \frac{p(v)}{q(v)}$. Assume that there is a label 327 distribution $d = (d_1, \dots, d_L)$ and its corresponding label space is $\mathcal{Y} = \{y_1, \dots, y_L\}$. The two 328 decompositions of \mathcal{Y} are $\mathcal{Y}_{a} = \{y_{1}, \dots, y_{k}\}$ and $\mathcal{Y}_{b} = \{y_{k}, \dots, y_{L}\}$, respectively. It is obvious that $\mathcal{Y}_a \cup \mathcal{Y}_b = \mathcal{Y}$ and $\mathcal{Y}_a \cap \mathcal{Y}_b = \{y_k\}$. Let the subspace label distribution corresponding to these two decompositions be $a = (a_1, \dots, a_k)$ and $b = (b_k, \dots, b_L)$. According to our assumptions, 330 $d_k \neq 0$. Then, for any integer $j \in [1, k]$, we have 331

$$\frac{a_j}{a_k} = \frac{[\mathcal{N}(d)]_j}{[\mathcal{N}(d)]_k} = \frac{[p(d)]_j}{[q(d)]_j} \frac{[q(d)]_k}{[p(d)]_k}.$$
(3)

 $a_k \quad [\mathcal{N}(d)]_k \quad [q(d)]_j \ [p(d)]_k$. Typically, for most normalization functions, $q(\cdot)$ is a normalizing constant, i.e., $[q(d)]_j = [q(d)]_k$. 334 335 Thus Eq. (3) can be rewritten into $a_j[p(d)]_k = a_k[p(d)]_j$. Plug it into $\sum_{j=1}^k a_j = 1$, and do the 336 same for **b** as well, and get 337

$$\frac{a_k \sum_{j=1}^k [p(d)]_j}{[p(d)]_k} = 1, \quad \frac{b_k \sum_{j=k}^L [p(d)]_j}{[p(d)]_k} = 1.$$
(4)

Add these two equations together, we have

332

333

338 339 340

345 346

347

354

355

361

362

371 372

373

$$[p(d)]_k + \sum_{j=1}^{L} [p(d)]_j = \frac{[p(d)]_k}{a_k} + \frac{[p(d)]_k}{b_k}.$$
(5)

Eq. (5) implies that $\sum_{i=1}^{L} [p(d)]_j$ must be given, and $[p(d)]_k$ is related to d_k , and only d_k . To make it possible, the only thing we can exploit is the sum-to-one constraint of d, i.e., $\sum_{j=1}^{L} d_j = 1$. Therefore $[p(v)]_j = v_j$. Since $\sum_i^{|v|} [\mathcal{N}(v)]_i = 1$, we have $q(v) = \sum_i^{|v|} v_i$, i.e., the finally deduced normalization function is Eq. (2). In this case, for any integer $j \in [1, L]$, we have

$$d_{j} = \begin{cases} \frac{a_{j}b_{k}}{a_{k} + b_{k} - a_{k}b_{k}}, & j = 1, \cdots, k\\ \frac{a_{k}b_{j}}{a_{k} + b_{k} - a_{k}b_{k}}, & j = k + 1, \cdots, L \end{cases},$$
(6)

which illustrates that the original label distribution d can be reconstructed by subtask label distributions a and b. This is possible thanks to the use of \mathcal{N}_{SUM} .

356 Theorem 1 also states that it is not appropriate 357 to replace Eq. (2) with the min-max or softmax 358 function because doing so destroys the recon-359 struction information. 360

4.3 COMPLEXITY ANALYSIS

363 The overall time cost of SC is primarily influ-364 enced by the calculation of M and the normal-365 ization process. The time complexity of com-366 puting and updating M are $\mathcal{O}(L(TN + T^2))$ 367 and $\mathcal{O}(LT)$, respectively. The time complexity of the normalization process is $\mathcal{O}(LTN)$. The 368



Figure 3: The overview of S-LDL (deep regime). White, red and gray highlight our proposed, existing methods, and loss functions, respectively.

overall time complexity of each iteration of SC is $\mathcal{O}(L(TN+T^2))$, which is linear with respect to 369 the number of instances and labels. Therefore, it is clear that SC can be applied to large-scale datasets. 370

5 S-LDL of the deep regime

374 The aforementioned analysis has exposed the problems of the shallow regime: 1) shallow methods as 375 base estimators have low potential in themselves; 2) there is a training gap between the primary task and subtasks, i.e., no representation learning is involved. Therefore, it is necessary to introduce our 376 proposed S-LDL of the deep regime, the overview of which is illustrated in Fig. 3. We illustrate our 377 framework by introducing the learnable parts one by one.

Table 2: Modifications of different task adaptations

Туре	Subtask construction	ℓ _{PRI}	ℓ_{SUB}
Vanilla LDL	$(\boldsymbol{D}^{(1)},\cdots) \leftarrow \mathtt{SC}(\boldsymbol{D})$	$\ell(oldsymbol{D}, ilde{oldsymbol{D}})\in\mathcal{L}_{ ext{LDL}}$	$\ell_{ ext{SUB}}(oldsymbol{D}^{(1)},\cdots; ilde{oldsymbol{D}}^{(1)},\cdots)$
LDL4C	$(\boldsymbol{D}^{(1)},\cdots) \leftarrow \mathtt{SC}(\boldsymbol{D})$	$\ell(oldsymbol{D}, ilde{oldsymbol{D}}, ilde{oldsymbol{D}})\in\mathcal{L}_{ ext{LDL4C}}$	$\ell_{ ext{SUB}}(oldsymbol{D}^{(1)},\cdots; ilde{oldsymbol{D}}^{(1)},\cdots)$
IncomLDL	$(\boldsymbol{D}_{\Omega}^{(1)},\cdots) \leftarrow \mathtt{SC}(\mathcal{R}_{\Omega}(\boldsymbol{D}))$	$\ell(\mathcal{R}_{\Omega}(\boldsymbol{D}), \mathcal{R}_{\Omega}(\tilde{\boldsymbol{D}})) \in \mathcal{L}_{\mathrm{IncomLDL}}$	$\ell_{\text{SUB}}(\boldsymbol{D}_{\Omega}^{(1)},\cdots;\mathcal{R}_{\Omega}(\tilde{\boldsymbol{D}}^{(1)}),\cdots)$
LE	$(\boldsymbol{L}^{(1)}, \cdots) \leftarrow \mathtt{SC}(\boldsymbol{L})$	$\ell(oldsymbol{L}, ilde{oldsymbol{D}})\in\mathcal{L}_{ ext{LE}}$	$\ell_{ ext{SUB}}(m{L}^{(ilde{1})},\cdots; ilde{m{D}}^{(1)},\cdots)$

• $\varphi(\cdot)$ is guided by subtasks to learn a powerful representation, i.e., $\mathbf{R} = \varphi(\mathbf{X})$.

• $\psi(\cdot)$ is responsible for predicting subtask label distributions, i.e., $(\tilde{D}^{(1)}, \cdots) = \psi(R)$. To ensure the precise prediction of subtask label distributions for reconstruction, we employ the mean absolute error function for subtask learning. The loss is weighted by the summation of the description degrees corresponding to the primary tasks, allowing more reliable label spaces to receive more attention. The subtask learning loss has the following form:

$$\ell_{\text{SUB}}\left(\mathcal{D}^{\circ};\,\tilde{\mathcal{D}}^{\circ}\right) = \frac{1}{N\left|\mathcal{Y}^{\circ}\right|} \sum_{\mathcal{Y}^{(t)}\in\mathcal{Y}^{\circ}} \sum_{i=1}^{N} \left(\sum_{y_{k}\in\mathcal{Y}^{(t)}} d_{\boldsymbol{x}_{i}}^{y_{k}}\right) \sum_{j=1}^{\left|\mathcal{Y}^{(t)}\right|} \left|d_{\boldsymbol{x}_{i}}^{(t)y_{j}} - \tilde{d}_{\boldsymbol{x}_{i}}^{(t)y_{j}}\right|.$$
(7)

• $\omega(\cdot)$ can be any existing method that can be expressed as a network structure theoretically. Since the concatenation of the representation and subtask label distributions, we have $\mathbf{Z} = (\mathbf{R}, \psi(\mathbf{R}))$ and $D = \omega(Z)$. In the case of the primary task being vanilla LDL, the primary task loss ℓ_{PRI} can be

$$\ell_{\mathrm{KL}}\left(\boldsymbol{D},\,\tilde{\boldsymbol{D}}\right) = \frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{L} d_{\boldsymbol{x}_{i}}^{y_{j}} \ln \frac{d_{\boldsymbol{x}_{i}}^{y_{j}}}{\tilde{d}_{\boldsymbol{x}_{i}}^{y_{j}}}, \quad \ell_{\mathrm{KL}} \in \mathcal{L}_{\mathrm{LDL}}.$$
(8)

Note that ℓ_{PRI} changes as the primary task changes. Finally, we can learn the model parameters Θ by

$$\Theta^* = \arg\min(\ell_{\text{PRI}} + \alpha \ell_{\text{SUB}}),\tag{9}$$

where α is a trade-off parameter. Compared with the shallow regime, S-LDL of the deep regime 405 has the following advantages: 1) There is no two-stage training gap, which makes the representation contain insights from both the primary task and the subtasks; 2) the framework not only serves LDL, but can also be directly applied to derivative tasks of LDL, e.g., LDL for classification (LDL4C) 408 (Wang, Jing and Geng, Xin, 2019), incomplete LDL (IncomLDL) (Xu and Zhou, 2017), label 409 enhancement (LE) (Xu et al., 2019). The modifications involved are shown in Table 2, where \mathcal{L}_X 410 indicates the set of losses for adaptable methods in the task of type "X". Special mathematical procedures of LDL4C and IncomLDL are defined as 412

$$\begin{bmatrix} \bar{\boldsymbol{D}}^{(t)} \end{bmatrix}_{ij} \triangleq \begin{cases} 1, & \text{if } y_j = \arg\max_{\bar{y} \in \mathcal{Y}^{(t)}} d_{\boldsymbol{x}_i}^{\bar{y}} \\ 0, & \text{otherwise} \end{cases}, \quad \begin{bmatrix} \mathcal{R}_{\Omega} \left(\boldsymbol{D} \right) \end{bmatrix}_{ij} \triangleq \begin{cases} \begin{bmatrix} \boldsymbol{D} \end{bmatrix}_{ij}, & \text{if } (i, y_j) \in \Omega \\ 0, & \text{otherwise} \end{cases}, \quad (10)$$

respectively, where $[\cdot]_{ii}$ represents the element in *i*-th row of the matrix corresponding to label y_i , 415 and Ω represents observed elements sampled uniformly at random from D in IncomLDL. Such 416 modifications are rational since: 1) targets of LDL4C and IncomLDL, i.e., D and $\mathcal{R}_{\Omega}(D)$, are 417 essentially different forms of degradation of the label distribution matrix; and 2) the target of LE is a 418 logical label matrix L, the same as the target of MLL, which is actually a special case of LDL. 419

420

421 422

423

424

EXPERIMENTS 6

In this section, we evaluate S-LDL of the deep regime. Due to page limitations, datasets, comparison methods, and their parameter settings are introduced in the appendix.

425 **Metrics** For LDL, we use the same metrics suggested by Jia et al. (2023). Due to page limitations, 426 we only present results on Cheby. \downarrow (Chebyshev distance), Clark \downarrow (Clark distance), Cosine \uparrow 427 (cosine similarity), and Spear. \uparrow (Spearman's coefficient) in the main paper, where \downarrow (\uparrow) indicates 428 "the lower (higher) the better". Note that these metrics are *not* as intuitive as accuracy or error rate, i.e., 429 small changes can mean large performance differences. For LDL4C, objective of which is different 430 from LDL, we use $0/1 \log \downarrow$ (zero one loss) and Err. prob. \downarrow (error probability) as metrics (Wang, 431 Jing and Geng, Xin, 2019).

386

387

388

389

390

391

392 393 394

396

397

398

403 404

406

407

411

Table 3: Experime	ental results of LDL on	JAFFE and Yeast	diau formatted as	$(\texttt{mean} \pm \texttt{std})$	rank))

1			—	((
Algorithms	JAFFE (Lyon	is et al., 1998)	Algorithms	Yeast_diau (Geng, 2016)		
Algonullis	Clark \downarrow	Cosine ↑	Aigonullins	Cheby.↓	Spear. ↑	
LDSVR (Geng and Hou, 2015)	.3280 ±.027 (6)	.9549 ±.010 (7)	CPNN (Geng et al., 2013)	$.0385 \pm .001$ (9)	$.2962 \pm .034$ (10)	
AA-kNN (Geng, 2016)	.3483 ±.032 (8)	.9497 ±.010 (9)	AA-kNN	$.0385 \pm .001$ (9)	.3674 ±.029 (9)	
LDLFs (Shen et al., 2017)	$.3637 \pm .032 (10)$	$.9494 \pm .009 (10)$	LDLFs	$.0371 \pm .001$ (8)	$.4088 \pm .021$ (8)	
DF-BFGS (González et al., 2021a)	.3062 ±.025 (3)	.9633 ±.007 (2)	DF-BFGS	$.0368 \pm .001$ (5)	.4161 ±.027 (5)	
KLD (Geng, 2016) •	.3608 ±.031 (9)	.9538 ±.008 (8)	LRR •	$.0370 \pm .001$ (7)	$.4154 \pm .023$ (6)	
S-KLD	.3007 ±.032 (2)	.9625 ±.009 (3)	S-LRR	.0366 ±.001 (1)	.4198 ±.023 (2)	
SCL (Jia et al., 2019) •	.3358 ±.024 (7)	$.9592 \pm .006$ (6)	QFD^2 (Wen et al., 2023) •	$.0369 \pm .001$ (6)	.4118 ±.025 (7)	
S-SCL	$.3184 \pm .025$ (4)	$.9604 \pm .008 (5)$	S-QFD ²	.0366 ±.001 (1)	$.4203 \pm .021 (1)$	
LRR (Jia et al., 2023) •	$.3230 \pm .027$ (5)	.9616 ±.008 (4)	CJS (Wen et al., 2023)	$.0367 \pm .001$ (4)	$.4164 \pm .025$ (4)	
S-LRR	.2934 ±.028 (1)	.9635 ±.008 (1)	S-CJS	.0366 ±.001 (1)	$.4198 \pm .024$ (2)	
					× /	

Table 4: Experimental results of LDL4C on sBU_3DFE and Flickr formatted as (mean ± std(rank))

Algorithms	sBU_3DFE (sBU_3DFE (Geng, 2016)		Flickr (Yang et al., 2017b)		
Algoriums	0/1loss↓ Err.prob.↓		Aigoritimis	$0/1 \log \downarrow$	Err.prob.↓	
LDL4C (Wang, Jing and Geng, Xin, 2019)	.5578 ±.028 (6)	.7671 ±.007 (5)	LDL4C	.8971 ±.008 (6)	$.8884 \pm .004$ (6)	
S-LDL4C	.5526 ±.025 (5)	.7686 ±.006 (6)	S-LDL4C	.8705 ±.138 (5)	.8702 ±.100 (5)	
HR (Wang and Geng, 2021a) •	.5167 ±.027 (3)	.7596 ±.006 (2)	HR •	.4513 ±.015 (4)	$.5823 \pm .007$ (4)	
S-HR	$.5069 \pm .025$ (2)	.7598 ±.006 (3)	S-HR	.4219 ±.015 (1)	.5639 ±.007 (1)	
LDLM (Wang and Geng, 2021b) •	.5258 ±.034 (4)	.7619 ±.009 (4)	LDLM •	$.4384 \pm .014$ (3)	.5740 ±.007 (3)	
S-LDLM	.4809 ±.024 (1)	.7524 ±.005 (1)	S-LDLM	.4321 ±.016 (2)	.5667 ±.007 (2)	

Results and discussion We apply S-LDL to 454 existing methods to demonstrate performance 455 improvements. For each dataset we conduct 456 ten-fold experiments repeated 10 times, and the 457 average performance is recorded. Tables 3 to 4 show representative results and the remainder 458 are in the appendix, where \bullet (\circ) indicates that 459 more than half of the metrics support that "S-X" 460 is statistically superior (inferior) to the corre-461 sponding methods "X" (pairwise t-test at 0.05 462 significance level); there is no significant if nei-463 ther \bullet nor \circ is shown. LRR focuses on the label 464 ranking relationship, which is also emphasized 465 by each subtask. We believe this is why S-LDL 466 and LRR fit so well. Note that our method has 467 the least improvement in SCL, which may be



Figure 4: Visualized results of the ablation study and the parameter sensitivity analysis, which is on the Natural_Scene (Geng, 2016) dataset.

468attributed to its reliance on shallow regime methods in the prediction phase. It is also worth noting469that the improvement in vanilla KLD is considerable, which just illustrates the limitations of loss470function engineering that considers label correlation one-sidedly. QFD² and CJS are tailored for471ordinary LDL, and may have better results than LRR on this regard. Powered by S-LDL, these472methods can all achieve better level. For LDL4C, S-LDL significantly improves both HR and LDLM.473However, it can be observed that S-LDL4C is unstable on the Flickr dataset, which is not surprising474since LDL4C itself fails on it. We believe this is caused by the combined effect of the sparsity of the475dataset and the information entropy operation involved in LDL4C.

475

Parameter sensitivity We check the sensitivity of the trade-off parameter α on the LDL task with the Natural_Scene dataset by varying the parameter in {0.01, 0.05, 0.1, 0.5, 1, 5}. Results are shown in Fig. 4. Spearman's coefficient of S-LDL first increases and then decreases as α varies, demonstrating a desirable bell-shaped curve. This justifies our motivation of jointly learning the primary task and subtasks, as a good trade-off between them can enhance the performance.

481

Ablation study Here we are interested in the importance of each part of S-LDL, thus an ablation study is performed with S-KLD: 1) we replace \mathcal{N}_{SUM} in SC with the min-max function to examine the importance of the subtask distribution reconstruction, and this model is denoted as S-KLD (min-max); 2) we remove the identity mapping in Fig. 3 to examine the importance of the prediction via subtask representation, and this model is denoted as S-KLD w/o id.; 3) we train without the term of ℓ_{SUB} (i.e.,

451 452

432 433



486 setting $\alpha = 0$) to examine the importance of subtask learning, and this model is denoted as S-KLD 487 w/o ℓ_{SUB} . Results are also shown in Fig. 4, which confirms that each part of S-LDL contributes as 488 long as there is a good trade-off.

489 490

491 492

493

495

497

503

521

526

527

528

7 LIMITATIONS AND CONCLUSION

Limitations First, S-LDL of the shallow regime is proposed out of intuition, and in Section 5, we have discussed its limitations, which are addressed via the designing of S-LDL of the deep regime. 494 Second, when the label space is large, especially when labels are continuous and result in unimodal label distributions (e.g., age estimation), our proposed cannot be rationally applied. Fortunately, one 496 possible workaround is to use a binning tricks for preprocessing, and then construct subtasks.

498 **Conclusion** We propose S-LDL, a subtask learning framework nested into LDL. S-LDL is generic: 499 it generates pseudo-supervised information via subtask construction without any extra knowledge; S-LDL is minimalist: it can be attached to existing methods and handle derivative tasks; S-LDL 500 is efficient: it captures a wide variety of label correlations. The analysis shows the validity and 501 reconstructability of subtasks, and experiments show the superiority of our framework. 502

- 504 REFERENCES
- 505 Shikai Chen, Jianfeng Wang, Yuedong Chen, Zhongchao Shi, Xin Geng, and Yong Rui. Label distri-506 bution learning on auxiliary label space graphs for facial expression recognition. In Proceedings 507 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 13984–13993, 508 2020. 509
- Bin-Bin Gao, Chao Xing, Chen-Wei Xie, Jianxin Wu, and Xin Geng. Deep label distribution learning 510 with label ambiguity. IEEE Transactions on Image Processing, 26(6):2825–2838, 2017. 511
- 512 Xin Geng. Label distribution learning. IEEE Transactions on Knowledge and Data Engineering, 28 513 (7):1734-1748, 2016.
- 514 Xin Geng and Peng Hou. Pre-release prediction of crowd opinion on movies by label distribution 515 learning. In Proceedings of the 24th International Joint Conference on Artificial Intelligence, pages 516 3511-3517, 2015. 517
- 518 Xin Geng, Chao Yin, and Zhi-Hua Zhou. Facial age estimation by learning from label distributions. 519 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(10):2401–2412, 2013.
- 520 Manuel González, Germán González-Almagro, Isaac Triguero, José-Ramón Cano, and Salvador García. Decomposition-fusion for label distribution learning. Information Fusion, 66:64-75, 522 2021a. 523
- Manuel González, Julián Luengo, José-Ramón Cano, and Salvador García. Synthetic sample 524 generation for label distribution learning. Information Sciences, 544:197–213, 2021b. 525
 - Xiuyi Jia, Zechao Li, Xiang Zheng, Weiwei Li, and Sheng-Jun Huang. Label distribution learning with label correlations on local samples. *IEEE Transactions on Knowledge and Data Engineering*, 33(4):1619-1631, 2019.
- 529 Xiuyi Jia, Xiaoxia Shen, Weiwei Li, Yunan Lu, and Jihua Zhu. Label distribution learning by 530 maintaining label ranking relation. IEEE Transactions on Knowledge and Data Engineering, 35 531 (02):1695-1707, 2023.532
- 533 Yufei Jin, Richard Gao, Yi He, and Xingquan Zhu. Gldl: Graph label distribution learning. In 534 Proceedings of the AAAI Conference on Artificial Intelligence, pages 12965–12974, 2024.
- Diederik P. Kingma and Jimmy Lei Ba. Adam: A method for stochastic optimization. In Proceedings 536 of the 3rd International Conference on Learning Representations, 2015. 537
- Nhat Le, Khanh Nguyen, Quang Tran, Erman Tjiputra, Bac Le, and Anh Nguyen. Uncertainty-aware 538 label distribution learning for facial expression recognition. In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pages 6088–6097, 2023.

540 541 542	Xiangyu Li, Xinjie Liang, Gongning Luo, Wei Wang, Kuanquan Wang, and Shuo Li. Ambiguity- aware breast tumor cellularity estimation via self-ensemble label distribution learning. <i>Medical</i> <i>Image Analysis</i> , 90:102944, 2023a.
543 544	Xiangyu Li, Gongning Luo, Wei Wang, Kuanquan Wang, and Shuo Li. Curriculum label distribution learning for imbalanced medical image segmentation. <i>Medical Image Analysis</i> , 89:102911, 2023b.
545 546 547	Lingyu Liang, Luojun Lin, Lianwen Jin, Duorui Xie, and Mengru Li. Scut-fbp5500: A diverse benchmark dataset for multi-paradigm facial beauty prediction. In <i>Proceedings of the 24th</i>
549 550	Yunan Lu and Xiuyi Jia. Predicting label distribution from multi-label ranking. In <i>Proceedings of the</i>
551 552	36th Annual Conference on Neural Information Processing Systems, pages 36931–36943, 2022. Michael Lyons, Shigeru Akamatsu, Miyuki Kamachi, and Jiro Gyoba. Coding facial expressions
553 554 555	with gabor wavelets. In <i>Proceedings of the 3rd IEEE International Conference on Automatic Face and Gesture Recognition</i> , pages 200–205, 1998.
556 557 558	Yashoteja Prabhu, Anil Kag, Shrutendra Harsola, Rahul Agrawal, and Manik Varma. Parabel: Partitioned label trees for extreme classification with application to dynamic search advertising. In <i>Proceedings of the 2018 World Wide Web Conference</i> , pages 993–1002, 2018.
559 560	Jesse Read, Concha Bielza, and Pedro Larrañaga. Multi-dimensional classification with super-classes. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 26(7):1720–1733, 2013.
561 562 563	Jesse Read, Antti Puurula, and Albert Bifet. Multi-label classification with meta-labels. In <i>Proceed</i> - ings of the 2014 IEEE International Conference on Data Mining, pages 941–946, 2014.
564 565 566	Tingting Ren, Xiuyi Jia, Weiwei Li, Lei Chen, and Zechao Li. Label distribution learning with label-specific features. In <i>Proceedings of the 28th International Joint Conference on Artificial Intelligence</i> , pages 3318–3324, 2019.
567 568 569	Wei Shen, Kai Zhao, Yilu Guo, and Alan Yuille. Label distribution learning forests. In <i>Proceedings</i> of the 31st Annual Conference on Neural Information Processing Systems, pages 834–843, 2017.
570 571 572 573	Amirreza Shirani, Franck Dernoncourt, Paul Asente, Nedim Lipka, Seokhwan Kim, Jose Echevarria, and Thamar Solorio. Learning emphasis selection for written text in visual media from crowd-sourced label distributions. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1167–1172, 2019.
574 575 576 577	Grigorios Tsoumakas, Ioannis Katakis, and Ioannis Vlahavas. Effective and efficient multilabel classification in domains with large number of labels. In <i>ECML/PKDD 2008 Workshop on Mining Multidimensional Data</i> , volume 21, pages 53–59, 2008.
578 579	Grigorios Tsoumakas, Ioannis Katakis, and Ioannis Vlahavas. Random k-labelsets for multilabel classification. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 23(7):1079–1089, 2010.
580 581	Jing Wang and Xin Geng. Theoretical analysis of label distribution learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 33, pages 5256–5263, 2019.
582 583 584	Jing Wang and Xin Geng. Learn the highest label and rest label description degrees. In <i>Proceedings</i> of the 30th International Joint Conference on Artificial Intelligence, pages 3097–3103, 2021a.
585 586	Jing Wang and Xin Geng. Label distribution learning machine. In <i>Proceedings of the 38th Interna-</i> <i>tional Conference on Machine Learning</i> , pages 10749–10759. PMLR, 2021b.
587 588 589	Ke Wang, Ning Xu, Miaogen Ling, and Xin Geng. Fast label enhancement for label distribution learning. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 35(02):1502–1514, 2023.
590 591	Ran Wang, Sam Kwong, Xu Wang, and Yuheng Jia. Active k-labelsets ensemble for multi-label classification. <i>Pattern Recognition</i> , 109:107583, 2021.
592 593	Wang, Jing and Geng, Xin. Classification with label distribution learning. In <i>Proceedings of the 28th International Joint Conference on Artificial Intelligence</i> , pages 3712–3718, 2019.

- 594 Changsong Wen, Xin Zhang, Xingxu Yao, and Jufeng Yang. Ordinal label distribution learning. In 595 Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 23481–23491, 596 2023. 597 Xiaoping Wu, Ni Wen, Jie Liang, Yukun Lai, Dongyu She, Mingming Cheng, and Jufeng Yang. Joint 598 acne image grading and counting via label distribution learning. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 10642–10651, 2019. 600 601 Miao Xu and Zhi-Hua Zhou. Incomplete label distribution learning. In Proceedings of the 26th 602 International Joint Conference on Artificial Intelligence, pages 3175–3181, 2017. 603 Ning Xu, Yunpeng Liu, and Xin Geng. Label enhancement for label distribution learning. IEEE 604 *Transactions on Knowledge and Data Engineering*, 33(4):1632–1643, 2019. 605 606 Ning Xu, Jun Shu, Renyi Zheng, Xin Geng, Deyu Meng, and Min-Ling Zhang. Variational label 607 enhancement. IEEE Transactions on Pattern Analysis and Machine Intelligence, 45(5):6537–6551, 2023. 608 609 Jufeng Yang, Dongyu She, and Ming Sun. Joint image emotion classification and distribution learning 610 via deep convolutional neural network. In Proceedings of the 26th International Joint Conference 611 on Artificial Intelligence, pages 3266-3272, 2017a. 612 Jufeng Yang, Ming Sun, and Xiaoxiao Sun. Learning visual sentiment distributions via augmented 613 conditional probability neural network. In Proceedings of the AAAI Conference on Artificial 614 Intelligence, pages 224–230, 2017b. 615 616 Yansheng Zhai, Jianhua Dai, and Hong Shi. Label distribution learning based on ensemble neural 617 networks. In Proceedings of the 25th International Conference on Neural Information Processing, 618 pages 593–602, 2018. 619 Min-Ling Zhang and Zhi-Hua Zhou. A review on multi-label learning algorithms. IEEE Transactions 620 on Knowledge and Data Engineering, 26(8):1819–1837, 2013. 621 622 Xingyu Zhao, Yuexuan An, Ning Xu, Jing Wang, and Xin Geng. Imbalanced label distribution learning. In Proceedings of the AAAI Conference on Artificial Intelligence, pages 11336–11344, 623 2023a. 624 625 Xingyu Zhao, Lei Qi, Yuexuan An, and Xin Geng. Generalizable label distribution learning. In 626 Proceedings of the 31st ACM International Conference on Multimedia, pages 8932–8941, 2023b. 627 Qinghai Zheng, Jihua Zhu, and Haoyu Tang. Label information bottleneck for label enhancement. In 628 Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 629 7497–7506, 2023. 630 631 632 **APPENDIX: DETAILS OF EXPERIMENTS** Α 633 634 Here we try our best to provide as much information as possible for reproducible research. 635 636 A.1 METRICS 637 638 For LDL, IncomLDL, and LE, we use the same metrics suggested by Geng (2016), which are 639 Cheby. \downarrow (Chebyshev distance), Clark \downarrow (Clark distance), Can. \downarrow (Canberra distance), KLD \downarrow 640 (Kullback-Leibler divergence), Cosine \uparrow (cosine similarity), and Int. \uparrow (intersection similarity), 641 respectively. Here \downarrow (\uparrow) indicates "the lower (higher) the better". For LDL and IncomLDL, we
- For LDL, IncomLDL, and LE, we use the same metrics suggested by Geng (2016), which are Cheby. \downarrow (Chebyshev distance), Clark \downarrow (Clark distance), Can. \downarrow (Canberra distance), KLD \downarrow (Kullback-Leibler divergence), Cosine \uparrow (cosine similarity), and Int. \uparrow (intersection similarity), respectively. Here \downarrow (\uparrow) indicates "the lower (higher) the better". For LDL and IncomLDL, we additionally use two ranking metrics: Spear. \uparrow (Spearman's coefficient) and Ken. \uparrow (Kendall's coefficient) (Jia et al., 2023). Note that these metrics are *not* as intuitive as accuracy or error rate, i.e., *small changes can mean large performance differences*. For LDL4C, objective of which is different from LDL, we use $0/1 \text{ loss } \downarrow$ (zero one loss) and Err. prob. \downarrow (error probability) as metrics (Wang, Jing and Geng, Xin, 2019). Let the real distribution be denoted by $u = \{u_j\}_{j=1}^L$, and the predicted distribution be denoted by $v = \{v_j\}_{j=1}^L$, then the above metrics can be summarized in Table 5, where $\rho(\cdot)$ and $\delta(\cdot, \cdot)$ are the ranking function and the Kronecker delta function, respectively.

Table 5: Summary of the metrics 649 650 Name Formula Name Formula $\operatorname{Sim}_{1}(\boldsymbol{u},\,\boldsymbol{v}) = \frac{\sum_{j=1}^{L} u_{j} v_{j}}{\sqrt{\sum_{j=1}^{L} u_{j}^{2}} \sqrt{\sum_{j=1}^{L} v_{j}^{2}}}$ 651 $\operatorname{Dis}_1(\boldsymbol{u}, \boldsymbol{v}) = \max_j |u_j - v_j|$ Cheby. \downarrow Cosine ↑ 652 653 $Dis_2(u, v) = \sqrt{\sum_{j=1}^{L} \frac{(u_j - v_j)^2}{(u_j + v_j)^2}}$ $\operatorname{Sim}_2(\boldsymbol{u}, \boldsymbol{v}) = \sum_{j=1}^L \min(u_j, v_j)$ $Clark \downarrow$ Int.↑ 654 $\begin{aligned} & \operatorname{Rnk}_{1}(\boldsymbol{u},\,\boldsymbol{v}) = 1 - \frac{6\sum_{j=1}^{L} (\rho(u_{j}) - \rho(v_{j}))^{2}}{L(L^{2} - 1)} \\ & \operatorname{Rnk}_{2}(\boldsymbol{u},\,\boldsymbol{v}) = \frac{2\sum_{j < k} \operatorname{sgn}(u_{j} - u_{k})\operatorname{sgn}(v_{j} - v_{k})}{L(L - 1)} \end{aligned}$ 655 $\operatorname{Dis}_{3}(\boldsymbol{u}, \boldsymbol{v}) = \sum_{j=1}^{L} \frac{|u_{j} - v_{j}|}{u_{j} + v_{j}}$ Spear. \uparrow $Can.\downarrow$ 656 657 $\operatorname{Dis}_4(\boldsymbol{u},\,\boldsymbol{v}) = \sum_{j=1}^L u_j \ln \frac{u_j}{v_j}$ KLD \downarrow Ken.↑ 658 $C_1(\boldsymbol{u}, \boldsymbol{v}) = \delta(\arg \max(\boldsymbol{u}),$ 0/1 Err. 659 $\mathbf{C}_2(\boldsymbol{u}, \boldsymbol{v}) = 1 - u_{\arg\max(\boldsymbol{v})}$ loss↓ $\arg \max(v)$ prob. \downarrow 660

A.2 DATASETS

We adopt several widely used label distribution datasets, including: JAFFE (Lyons et al., 1998);⁴ fbp5500 (Liang et al., 2018);⁵ sBU_3DFE, Movie, Natural_Scene, Yeast_heat, Yeast_diau, Yeast_cold, and Yeast_dtt provided by Geng (2016);⁶ emotion6, Twitter, and Flickr provided by Yang et al. (2017b).⁷ The information of these datasets are summarized in Table 6.

667 668 669

661 662

663 664

665

666

648

A.3 COMPARISON METHODS

670 On the one hand, we apply our pro-671 posed S-LDL to existing methods to 672 demonstrate performance improve-673 ments in the LDL task (denoted 674 by the "S-" prefix). These meth-675 ods are BFGS-LLD (KLD) (Geng, 676 2016), SCL (Jia et al., 2019), LRR (Jia et al., 2023), QFD² (Wen et al., 677 2023), and CJS (Wen et al., 2023) 678 (the losses of these methods consti-679 tute the set \mathcal{L}_{LDL}). On the other 680 hand, we compare S-LDL with 681 methods that have specialized struc-682 ture, which our proposed cannot di-683 rectly adapt to. These methods are

Table 6: Summary of datasets

Dataset	# Instances N	# Features P	# Labels L
JAFFE	213	243	6
sBU_3DFE	2500	243	6
Movie	7755	1869	5
Nature_Scene	2000	294	9
fbp5500	5500	512	5
Yeast_heat	2465	24	6
Yeast_diau	2465	24	7
Yeast_cold	2465	24	4
Yeast_dtt	2465	24	4
emotion6	1980	168	7
Twitter	10045	168	8
Flickr	11150	168	8

⁶⁸⁴ CPNN (Geng et al., 2013), LDSVR (Geng and Hou, 2015), AA-*k*NN (Geng, 2016), LDLFs (Shen et al., 2017), and DF-LDL (denoted by DF-BFGS since we use BFGS-LLDs as base estimators)
⁶⁸⁶ (González et al., 2021a). Moreover, we apply our proposed to derivative tasks of LDL (i.e., LDL4C, IncomLDL, and LE) and the comparison methods involved are LDL4C (Wang, Jing and Geng, Xin, 2019), HR (Wang and Geng, 2021a), LDLM (Wang and Geng, 2021b), IncomLDL (Xu and Zhou, 2017), LP (Xu et al., 2019), GLLE (Xu et al., 2019), LEVI (Xu et al., 2023), and LIBLE (Zheng et al., 2023).

691 692

700

A.4 PARAMETER SETTINGS AND EXPERIMENTAL ENVIRONMENT

The parameter settings of the proposed S-LDL and comparison algorithms are summarized in Table 7. Note that DF-LDL is parameter-free, and we use BFGS-LLDs as its base estimators, parameter settings of which are the same as BFGS-LLD as the comparison algorithm. We use Adam (Kingma and Ba, 2015) for the optimization of S-LDL. For all methods of the deep regime, the learning rate is chosen among $\{1, 2, 5\} \times 10^{\{-4, -3, -2\}}$, and the selection of the number of epochs is nested into

⁶https://palm.seu.edu.cn/xgeng/LDL/download.htm

⁷https://cv.nankai.edu.cn/projects

^{699 &}lt;sup>4</sup>https://zenodo.org/records/3451524

⁵https://github.com/HCIILAB/SCUT-FBP5500-Database-Release

Algorithms	Parameter	Value (Range)
AA-kNN	k: # Neighbors	5
	# Estimators (trees)	5
LDLFs	Depth	6
	Latent units (leaves)	64
PECSUD	ε : Convergence criterion	10^{-6}
BF03-LLD	Max iteration	600
SCI	m: # Clusters	5
SCL	$\lambda_1, \lambda_2, \lambda_3$: Trade-off	$10^{-3}, 10^{-3}, 0.1$
IDD	λ : Trade-off (ranking loss)	$10^{\{-5, -4, -3, -2, -1\}}$
LKK	β : Trade-off (regularization)	$10^{\{-3, -2, -1, 0, 1, 2\}}$
	C_1, C_2 : Balance coefficients	$10^{-2}, 10^{-6}$
LDL4C	ρ : Margin	10^{-2}
Пр	$\lambda_1, \lambda_2, \lambda_3$: Trade-off	$10^{-2}, 10^{-6}$
IIK	ρ : Margin	10^{-2}
	$\lambda_1, \lambda_2, \lambda_3$: Trade-off	$10^{-6}, 10^{\{-3, -2, -1\}}, 10^{\{-3, -2\}}$
LDLM	ρ : Margin	10^{-2}
	ε : Convergence criterion	10^{-6}
IncomLDL	γ : Factor of Lipschitz constant	2
	λ : Trade-off	1
LP	α : Balance coefficient	0.5
CLLE	λ_1, λ_2 : Trade-off	$10^{-2}, 10^{-4}$
OLLE	σ : Width parameter for similarity calculation	10
LEVI	λ : Trade-off	1
LIBLE	α, β : Trade-off	$10^{\{-3, -2, -1, 0, 1, 2\}}$
S-LDL	α, λ, T	0.1, 0.2, 10

a ten-fold cross validation. All the results are obtained on a Linux workstation with Intel Core i9 (3.70GHz), NVIDIA GeForce RTX 3090 (24GB), and 32GB memory.

A.5 FULL EXPERIMENTAL RESULTS

737Here we provide complete results of all conducted experiments. Tables 8 to 19 are results on the LDL
task with different datasets. For IncomLDL, we follow *the incomplete settings* (Xu and Zhou, 2017)
and vary the observed rate ω % from 20% to 40%. Tables 20 to 21 are results on the IncomLDL task.
Tables 22 to 23 are on the LDL4C task. For LE, we follow *the settings of the recovery experiment*
(Xu et al., 2019). Tables 24 to 25 show results on the LE task.

Table 8: Experimental results of LDL on the JAFFE dataset formatted as $(mean \pm std)$

		1						· · ·	/
744	Algorithms	Cheby. \downarrow	$\texttt{Clark}\downarrow$	$\texttt{Can.}\downarrow$	$\texttt{KLD}\downarrow$	$\texttt{Cosine} \uparrow$	Int. \uparrow	Spear. \uparrow	Ken. ↑
745	LDSVR	$.0959 \pm .013$	$.3280 \pm .027$	$.6778 \pm _{.058}$	$.0476 \pm _{.011}$	$.9549 \pm .010$	$.8838 \pm .012$	$.5175 \pm _{.102}$	$.4508 \pm _{.086}$
746	AA-kNN	$.0978 \pm _{.012}$	$.3483 _{ \pm .032 }$	$.7164 _{\pm .066}$	$.0527_{\pm.011}$	$.9497 \pm _{.010}$	$.8766 \pm _{.012}$	$.4111 \pm _{.083}$	$.3514 _{\pm .070}$
747	LDLFs	$.0940 \pm _{.010}$	$.3637 _{\pm .032}$	$.7355 {\scriptstyle \pm .066}$	$.0550 \pm _{.009}$	$.9494 \scriptstyle \pm .009$	$.8766 \pm _{.011}$	$.4364 {\scriptstyle \pm .108}$	$.3749 _{\pm .093}$
748	DF-BFGS	$.0827 \pm _{.009}$	$.3062 \pm _{.025}$	$.6239 _{ \pm .052 }$	$.0388 \pm _{.007}$	$.9633 _{\pm .007}$	$.8944 \pm _{.010}$	$.5244 \pm _{.087}$	$.4493 \pm _{.077}$
749	KLD •	$.0925 \pm _{.010}$	$.3608 \pm _{.031}$	$.7363 \pm _{.064}$	$.0508 \pm _{.009}$	$.9538 \pm _{.008}$	$.8777_{\pm .011}$	$.4572 \pm _{.097}$	$.3873 \pm _{.084}$
750	S-KLD	$.0818 \pm .011$	$.3007 \pm _{.032}$	$.6132 \pm .067$	$.0395 \pm .010$	$.9625 \pm .009$	$.8960 \pm .012$	$\textbf{.5461} \scriptstyle \pm .105$	$.4769 \pm \scriptstyle .096$
750	SCL •	$.0873 \pm _{.008}$	$.3358 \pm _{.024}$	$.6874 \pm _{.051}$	$.0439 \pm _{.006}$	$.9592 \pm .006$	$.8851 \pm _{.009}$	$.4744 \pm .092$	$.4020 \pm _{.080}$
/51	S-SCL	$.0854 {\pm.010}$	$.3184 \pm _{.025}$	$.6526 \pm _{.053}$	$.0420 \pm _{.008}$	$.9604 { \pm .008 }$	$.8896 {\scriptstyle \pm .010}$	$.5110 \pm _{.095}$	$.4388 \pm _{.084}$
752	LRR •	$.0853 \pm _{.010}$	$.3230 _{\pm .027}$	$.6560 \pm _{.055}$	$.0412 \pm _{.008}$	$.9616 {\scriptstyle \pm .008}$	$.8906 \pm _{.010}$	$.5117 \pm _{.094}$	$.4420 \pm _{.084}$
753	S-LRR	$\textbf{.0804} \pm .009$	$\textbf{.2934} \scriptstyle \pm .028$	$\textbf{.5989} \scriptstyle \pm .059$	$\textbf{.0383} \scriptstyle \pm .009$	$\textbf{.9635} \scriptstyle \pm .008$	$\textbf{.8981} \scriptstyle \pm .011$	$.5448 \pm _{.092}$.4819 $_{\pm.084}$
754									

Table 9: Experimental results of LDL on the sBU_3DFE dataset formatted as $(mean \pm std)$

Algorithms	Cheby. \downarrow	$\texttt{Clark}\downarrow$	Can.↓	$\texttt{KLD}\downarrow$	$\texttt{Cosine} \uparrow$	Int. \uparrow	Spear. \uparrow	Ken. ↑
LDSVR	$.1250 \pm .005$	$.3710 _{\pm .010}$	$.8009 _{\pm .021}$	$.0720 \pm _{.004}$	$.9298 \pm _{.004}$	$.8559 _{\pm .004}$	$.3524 \pm _{.031}$	$.3011 _{\pm .026}$
AA-kNN	$.1272 \pm .004$	$.4001 \pm _{.009}$	$.8281 _{\pm .020}$	$.0801 _{\pm .004}$	$.9217 {\scriptstyle \pm .004}$	$.8488 \pm _{.004}$	$.2053 \pm _{.030}$	$.1767 \pm _{.026}$
LDLFs	$.1016 \pm .003$	$.3262 \pm .008$	$.6841 \pm \scriptstyle .017$	$.0504 \pm _{.003}$	$.9499 \pm \scriptstyle .003$	$.8776 \pm _{.003}$	$.4212 \pm .023$	$.3620 {\scriptstyle \pm . 019}$
DF-BFGS	$.1146 \pm .004$	$.3616 \pm .008$	$.7627 \pm .019$	$.0618 \pm .003$	$.9388 \pm .003$	$.8626 \pm \scriptstyle .004$	$.3026 \pm _{.031}$	$.2621 \pm _{.026}$
KLD •	$.1147 \pm .004$	$.3697 \pm .008$	$.7804 \pm .019$	$.0624 \pm .003$	$.9387 \pm .003$	$.8604 \pm \scriptstyle .003$	$.3021 \pm _{.026}$	$.2643 \pm _{.022}$
S-KLD	$.1014 _{\pm .004}$	$.3203 \pm _{.009}$	$.6736 _{\pm .018}$	$.0514 \pm _{.003}$	$.9487 \pm _{.003}$	$.8789 \scriptstyle \pm .004$	$.4334 {\pm.025}$	$.3729 _{\pm . 022}$
SCL •	$.1145 \pm _{.004}$	$.3648 \pm _{.008}$	$.7748 _{\pm .018}$	$.0605 \pm _{.003}$	$.9404 \pm _{.003}$	$.8614 _{\pm .003}$	$.3091 {\scriptstyle \pm .026}$	$.2701 \pm _{.021}$
S-SCL	$.1041 \pm _{.004}$	$.3301 \pm _{.009}$	$.6936 _{\pm .019}$	$.0535 \pm _{.003}$	$.9468 \pm _{.003}$	$.8754 _{\pm .004}$	$.3956 _{\pm .030}$	$.3381 {\scriptstyle \pm .027}$
LRR •	$.1067 _{\pm .003}$	$.3476 _{\pm .008}$	$.7320 _{\pm .017}$	$.0543 \pm _{.003}$	$.9465 \pm _{.003}$	$.8695 {\scriptstyle \pm .003}$	$.3626 \pm _{.026}$	$.3123 \pm _{.022}$
S-LRR	$\textbf{.0996} \pm .004$	$\textbf{.3157} \scriptstyle \pm .008$	$\textbf{.6610} \scriptstyle \pm .017$	$\textbf{.0499} \scriptstyle \pm .003$	$\textbf{.9502} \scriptstyle \pm .003$	$\textbf{.8812} \scriptstyle \pm .003$	$\textbf{.4455} \scriptstyle \pm .026$	$\textbf{.3837} \scriptstyle \pm .023$

Table 10: Experimental results of LDL on the Yeast_heat dataset formatted as (mean \pm std)

									-
Algorithms	Cheby. \downarrow	$\texttt{Clark}\downarrow$	$\texttt{Can.}\downarrow$	$\texttt{KLD}\downarrow$	$\texttt{Cosine} \uparrow$	Int. \uparrow	Spear. \uparrow	Ken. ↑	
CPNN	$.0419 _{\pm .001}$	$.1818_{\pm .005}$	$.3633 _{\pm .009}$	$.0125 \pm .001$	$.9881 \pm _{.001}$	$.9404 \pm _{.001}$	$.1507 \pm _{.034}$	$.1221 \pm _{.028}$	
AA-kNN	$.0441 \pm .001$	$.1913 \pm _{.005}$	$.3840 \pm _{.010}$	$.0140 \pm .001$	$.9867 \pm \scriptstyle .001$	$.9370 \pm _{.002}$	$.1678 \pm _{.031}$	$.1384 \pm _{.026}$	
LDLFs	$.0420 \pm .001$	$.1818 \pm .005$	$.3627 \pm .009$	$.0125 \pm .001$	$.9881 \pm .001$	$.9405 \pm .001$	$.1731 \pm _{.032}$	$.1409 \pm _{.026}$	
DF-BFGS	$.0420 \pm _{.001}$	$.1816 \pm _{.005}$	$.3624 _{\pm .009}$	$.0125 \pm .001$	$.9881 \pm _{.001}$	$.9405 \pm _{.001}$.1964 $_{\pm.034}$.1624 $_{\pm .028}$	
LRR •	$.0423 \pm _{.001}$	$.1828 \pm .005$	$.3644 _{\pm .009}$	$.0126 \pm _{.001}$	$.9880 \pm _{.001}$	$.9402 \pm .001$	$.1655 \pm _{.033}$	$.1351 \pm _{.028}$	
S-LRR	$\textbf{.0417} \scriptstyle \pm .001$	$.1806 \pm _{.005}$	$.3609 _{\pm .009}$.0124 $_{\pm.001}$	$\textbf{.9882} \scriptstyle \pm .001$	$.9408 \pm \scriptstyle .001$	$.1882 \pm _{.034}$	$.1548 \pm _{.028}$	
$QFD^2 \bullet$	$.0423 \pm .001$	$.1827 \pm _{.005}$	$.3644 \pm .009$	$.0126 \pm .001$	$.9880 \pm .001$	$.9402 \pm .001$	$.1677 \pm _{.032}$	$.1351 \pm _{.027}$	
S -QFD 2	$\textbf{.0417} \scriptstyle \pm .001$	$.1808 \pm .005$	$.3611 {\scriptstyle \pm .009}$	$\textbf{.0124} \pm .001$	$\textbf{.9882} \pm .001$	$.9408 \pm .001$	$.1880 \pm _{.032}$	$.1544 \pm .027$	
CJS •	$.0423 \pm .001$	$.1827 \pm .005$	$.3643 \pm .009$	$.0126 \pm .001$	$.9880 \pm .001$	$.9402 \pm .001$	$.1632 \pm _{.032}$	$.1329 \pm _{.027}$	
S-CJS	$\textbf{.0417} \scriptstyle \pm .001$	$\textbf{.1804} \scriptstyle \pm .005$	$\textbf{.3603} \scriptstyle \pm .009$.0124 $_{\pm.001}$	$\textbf{.9882} \scriptstyle \pm .001$	$\textbf{.9409} \scriptstyle \pm .001$	$.1940 \pm _{.030}$	$.1589 \pm _{.025}$	

Table 11: Experimental results of LDL on the Yeast_diau dataset formatted as (mean \pm std)

Algorithms	Cheby. \downarrow	$\texttt{Clark}\downarrow$	Can.↓	$\texttt{KLD}\downarrow$	$\texttt{Cosine} \uparrow$	Int. \uparrow	Spear. \uparrow	Ken. ↑
CPNN	$.0385 \pm .001$	$.2069 _{\pm .006}$	$.4439 _{\pm .012}$	$.0138_{\pm.001}$	$.9872 \pm .001$	$.9383 _{\pm .002}$	$.2962 \pm _{.034}$	$.2427 \pm _{.027}$
AA-kNN	$.0385 \pm .001$	$.2085 \pm .006$	$.4487 _ { \pm .014 }$	$.0145 \pm _{.001}$	$.9867 _{\pm .001}$	$.9377_{\pm .002}$	$.3674 _{\pm .029}$	$.2976 _{\pm .024}$
LDLFs	$.0371 _{\pm .001}$	$.2014 \pm .006$	$.4324 {\scriptstyle \pm .012}$	$.0132 \pm .001$	$.9879 _{\pm .001}$	$.9401 \pm _{.002}$	$.4088 \pm _{.021}$	$.3254 {\scriptstyle \pm .018}$
DF-BFGS	$.0368 \pm _{.001}$	$.1999 _{\pm .006}$	$.4294 \pm _{.013}$	$.0131 \pm _{.001}$	$.9879 _{\pm .001}$	$.9405 \pm _{.002}$	$.4161 \pm _{.027}$	$\textbf{.3404} \scriptstyle \pm .022$
LRR •	$.0370 \pm _{.001}$	$.2007 \pm .006$	$.4307 \pm _{.012}$	$.0131_{\pm.001}$	$.9879 _{\pm .001}$	$.9403 \pm _{.002}$	$.4154 \pm _{.023}$	$.3343 _{\pm .020}$
S-LRR	$\textbf{.0366} \pm .001$	$\textbf{.1983} \scriptstyle \pm .006$	$\textbf{.4257} \scriptstyle \pm .012$	$\textbf{.0129} \scriptstyle \pm .001$	$\textbf{.9881} \scriptstyle \pm .001$	$\textbf{.9410} \scriptstyle \pm .002$	$.4198 \pm _{.023}$	$.3389 {\scriptstyle \pm .019}$
$QFD^2 \bullet$	$.0369 _{\pm .001}$	$.2000 \pm .006$	$.4296 _{\pm .012}$	$.0131_{\pm.001}$	$.9879 _{\pm .001}$	$.9404 _{\pm .002}$	$.4118 \pm _{.025}$	$.3326 _{\pm .021}$
S -QFD 2	$\textbf{.0366} \scriptstyle \pm .001$	$.1985 \pm .006$	$.4261 _{\pm .012}$	$\textbf{.0129} \scriptstyle \pm .001$	$\textbf{.9881} \scriptstyle \pm .001$	$.9409 _{\pm .002}$	$\textbf{.4203} \scriptstyle \pm .021$	$.3387 _{\pm .018}$
CJS	$.0367_{\pm.001}$	$.1989 _{\pm .006}$	$.4272 \pm .012$	$.0130_{\pm.001}$	$.9880 _{\pm .001}$	$.9408 \pm _{.002}$	$.4164 \pm _{.025}$	$.3366 _{\pm .021}$
S-CJS	$\textbf{.0366} \scriptstyle \pm .001$	$.1984 \pm _{.006}$	$.4260 \pm _{.012}$	$.0130 \pm _{.001}$	$\textbf{.9881} \scriptstyle \pm .001$	$.9409 \pm _{.002}$	$.4198 \pm _{.024}$	$.3392 {\pm.019}$

Table 12: Experimental results of LDL on the Yeast_cold dataset formatted as $(mean \pm std)$

Algorithms	(h + h	(11-	0	VID	Casina A	T+ A	G +	V ^
Algorithms	Cneby.↓	Clark 4	Can.↓	KLD ↓	Cosine	Int. T	Spear. T	Ken.
CPNN	$\textbf{.0510} \scriptstyle \pm .002$	$.1392 \pm .005$	$.2396 _{\pm .008}$.0121 $_{\pm.001}$.9886 $_{\pm.001}$	$\textbf{.9410} \scriptstyle \pm .002$	$\textbf{.2651} \scriptstyle \pm .036$.2263 $_{\pm.032}$
AA-kNN	$.0542 \pm _{.002}$	$.1476 \pm _{.005}$	$.2549 _{\pm .008}$	$.0135_{\pm.001}$	$.9872 \pm .001$	$.9371 _{\pm .002}$	$.2189 \pm _{.035}$	$.1866 \pm _{.031}$
LDLFs	$.0511_{\pm.002}$	$.1396 _{\pm .005}$	$.2404 \pm _{.009}$	$.0122 \pm .001$	$.9885 {\scriptstyle \pm .001}$	$.9408 \pm _{.002}$	$.2482 \pm _{.038}$	$.2112 \pm _{.033}$
DF-BFGS	$.0514 \pm _{.002}$	$.1404 \pm _{.005}$	$.2424 \pm .008$	$.0123 \pm .001$	$.9885 \pm .001$	$.9403 \pm _{.002}$	$.2581 \pm _{.036}$	$.2190 \pm _{.030}$
LRR	$.0511_{\pm.002}$	$.1395 \pm _{.005}$	$.2402 \pm .009$	$.0122 \pm .001$.9886 $_{\pm.001}$	$.9408 \pm _{.002}$	$.2490 _{ \pm .035 }$	$.2111_{\pm.030}$
S-LRR	$\textbf{.0510} \scriptstyle \pm .002$	$\textbf{.1391} \scriptstyle \pm .005$	$\textbf{.2395} \scriptstyle \pm .009$	$\textbf{.0121} \scriptstyle \pm .001$	$\textbf{.9886} \scriptstyle \pm .001$	$\textbf{.9410} \scriptstyle \pm .002$	$.2618 \pm _{.037}$	$.2238 \pm _{.032}$
QFD^2	$.0513 _{\pm .002}$	$.1401 \pm _{.005}$	$.2413 \pm _{.009}$	$.0123 \pm .001$	$.9885 {\scriptstyle \pm .001}$	$.9405 _{\pm .002}$	$.2534 _{\pm .037}$	$.2158 \pm _{.032}$
S-QFD ²	$\textbf{.0510} \scriptstyle \pm .002$	$\textbf{.1391} \scriptstyle \pm .005$	$.2396 {\scriptstyle \pm .008}$	$\textbf{.0121} \scriptstyle \pm .001$	$\textbf{.9886} \scriptstyle \pm .001$	$\textbf{.9410} \scriptstyle \pm .002$	$.2571 _{\pm .039}$	$.2197 \pm _{.033}$
CJS	$.0513 \pm _{.002}$	$.1401 \pm .005$	$.2412 \pm .008$	$.0123 \pm .001$	$.9884 \pm .001$	$.9406 \pm \scriptstyle .002$	$.2535 \pm .038$	$.2152 \pm .032$
S-CJS	$\textbf{.0510} \scriptstyle \pm .002$	$.1392 \pm .005$	$.2396 _{\pm .009}$	$\textbf{.0121} \scriptstyle \pm .001$	$\textbf{.9886} \scriptstyle \pm .001$	$\textbf{.9410} \scriptstyle \pm .002$	$.2621 _{\pm .037}$	$.2241 \pm _{.031}$

Table 13: Experimental results of LDL on the Yeast_dtt dataset formatted as $(mean \pm std)$

	Algorithms	Cheby. \downarrow	$\texttt{Clark}\downarrow$	$\texttt{Can.}\downarrow$	$\texttt{KLD}\downarrow$	$\texttt{Cosine} \uparrow$	Int. \uparrow	Spear. \uparrow	Ken. ↑
_	CPNN	$.0361 \pm .001$	$.0984 \pm _{.004}$	$.1690 \pm _{.006}$	$.0063 \pm _{.001}$	$.9941 \pm .000$	$.9583 {\scriptstyle \pm .001}$	$.1735 \pm _{.035}$	$.1494 \pm _{.030}$
	AA-kNN	$.0386 \pm _{.001}$	$.1047 \pm _{.004}$	$.1797 \pm _{.006}$	$.0071 \pm _{.001}$	$.9933 _{\pm .000}$	$.9556 {\scriptstyle \pm .001}$	$.1591 \pm _{.033}$	$.1399 {}_{\pm . 030}$
	LDLFs	$.0360 \pm .001$	$.0981 \pm .004$	$.1689 \pm .006$	$.0063 \pm .001$	$\textbf{.9941} \scriptstyle \pm .000$	$.9583 {\scriptstyle \pm .001}$	$.1986 \pm _{.038}$	$.1727 \scriptstyle \pm .034$
	DF-BFGS	$.0365 \pm .001$	$.0995 \pm _{.004}$	$.1712 \pm .006$	$.0064 \pm _{.001}$	$.9939 \pm _{.000}$	$.9578 \pm .001$	$.1804 \pm \scriptstyle .033$	$.1592 \pm _{.030}$
-	LRR	$.0360 \pm _{.001}$	$.0982 \pm _{.004}$	$.1690 \pm _{.006}$	$.0063 _{ \pm .001 }$.9941 $_{\pm.000}$	$.9583 {\scriptstyle \pm .001}$	$.2016 _{\pm .037}$	$.1738_{\pm .032}$
	S-LRR	$\textbf{.0359} \scriptstyle \pm .001$.0977 $_{\pm.004}$	$\textbf{.1680} \scriptstyle \pm .006$.0062 $_{\pm.001}$	$\textbf{.9941} \scriptstyle \pm .000$.9585 $_{\pm.001}$	$.2068 \pm _{.035}$	$.1811 \pm _{.031}$
	QFD^2	$.0362 \pm .001$	$.0986 \pm _{.004}$	$.1696 \pm .006$	$.0063 \pm _{.001}$	$.9940 \pm .000$	$.9582 \pm .001$	$.1917 \scriptstyle \pm .035$	$.1665 _{\pm .031}$
	S-QFD ²	$\textbf{.0359} \scriptstyle \pm .001$.0977 $_{\pm.004}$	$.1681 \pm .006$	$\textbf{.0062} \pm .001$	$\textbf{.9941} \pm .000$	$\textbf{.9585} \scriptstyle \pm .001$	$\textbf{.2086} \scriptstyle \pm .036$	$\textbf{.1822} \scriptstyle \pm .032$
	CJS	$.0361 _{\pm .001}$	$.0984 \pm _{.004}$	$.1692 \pm _{.006}$	$.0063 \pm _{.001}$	$\textbf{.9941} \scriptstyle \pm .000$	$.9582 \pm .001$	$.1975 \pm _{.040}$	$.1722 \pm _{.035}$
	S-CJS	$\textbf{.0359} \scriptstyle \pm .001$	$.0978 \pm _{.004}$	$.1682 \pm .006$.0062 $_{\pm.001}$	$\textbf{.9941} \scriptstyle \pm .000$.9585 $_{\pm.001}$	$.2080 \pm _{.035}$	$.1804 \pm _{.031}$
-									

Table 14: Experimental results of LDL on the emotion6 dataset formatted as (mean \pm std)

	-						,	
Algorithms	Cheby. \downarrow	$\texttt{Clark}\downarrow$	$\texttt{Can.}\downarrow$	$\texttt{KLD}\downarrow$	$\texttt{Cosine} \uparrow$	Int. \uparrow	Spear. \uparrow	Ken.↑
LDSVR	$.3152 \pm .010$	$1.8217 \pm _{.020}$	$4.1452 \pm .064$	$1.0744 \pm _{.081}$	$.6906 \pm _{.015}$	$.5773 \pm _{.012}$	$.3915 \pm _{.030}$	$.3235 \pm _{.025}$
AA-kNN	$.3288 {\scriptstyle \pm .011}$	$1.7116_{\pm .026}$	$3.8757_{\pm .076}$	$.9512 {\scriptstyle \pm .115}$	$.6632 \pm .013$	$.5564 { }_{ \pm .010 }$	$.2920 \pm _{.027}$	$.2401 \pm _{.022}$
LDLFs	$.3120 {\scriptstyle \pm .010}$	$1.6625_{\pm .026}$	$3.7330_{\pm .075}$	$.5871 _{\pm .024}$	$.7143 \pm .011$	$.5802 \pm _{.010}$	$.3631 _{\pm .029}$	$.3025 \pm _{.024}$
DF-BFGS	$.3026 _{\pm .010}$	$1.6765_{\pm .025}$	$3.7675_{\pm .071}$	$.5805 \pm _{.026}$	$.7206 \pm .013$	$.5909 _{\pm .010}$	$.3940 _{\pm .027}$	$.3256 \pm _{.022}$
KLD •	$.3037 _{\pm .010}$	$1.6774_{\pm .025}$	$3.7729_{\pm .074}$	$.5863 { }_{ \pm .027 }$	$.7191 \pm .013$	$.5897 {\scriptstyle \pm .011}$	$.3959 \pm _{.028}$	$.3259 {\scriptstyle \pm .023}$
S-KLD	$.3024 \pm .010$	$1.6548 \pm _{.026}$	$3.6984 \pm .075$	$.5631 {\scriptstyle \pm .024}$	$.7282 \pm .012$	$.5926 \pm .010$	$.4063 \pm _{.027}$	$.3361 \pm _{.023}$
SCL •	$.3020 \pm .010$	$1.6750 \pm _{.025}$	$3.7642 \pm _{.073}$	$.5803 \pm _{.027}$	$.7219 \pm .013$	$.5917 \pm _{.011}$	$.4003 \pm _{.028}$	$.3299 \pm _{.023}$
S-SCL	$\textbf{.3018} \scriptstyle \pm .010$	$1.6554_{\pm .026}$	$3.6993 _{\pm .076}$	$.5631 _{\pm .025}$	$.7281 _{\pm .012}$	$\textbf{.5936} \scriptstyle \pm .010$	$\textbf{.4089} \scriptstyle \pm .027$	$\textbf{.3383} \scriptstyle \pm .023$
LRR •	$.3030 _{\pm .010}$	$1.6736_{\pm .025}$	$3.7601 _{\pm .073}$	$.5804 { }_{ \pm .026 }$	$.7212 \pm .013$	$.5899 _{\pm .010}$	$.3941 _{ \pm .027 }$	$.3243 \pm _{.023}$
S-LRR	$.3028 \pm _{.009}$	$\textbf{1.6524} \scriptstyle \pm .026$	$\textbf{3.6923} \scriptstyle \pm .074$.5607 $_{\pm.023}$.7299 $_{\pm.011}$	$.5923 \pm _{.010}$	$.4078 \pm _{.027}$	$.3373 _{\pm .022}$

Table 15: Experimental results of LDL on the Twitter dataset formatted as $(mean \pm std)$

	Algorithms	Cheby. \downarrow	$\texttt{Clark}\downarrow$	$\texttt{Can.}\downarrow$	$\texttt{KLD}\downarrow$	$\texttt{Cosine} \uparrow$	$\texttt{Int.} \uparrow$	Spear. \uparrow	Ken.↑
-	LDSVR	$.4236 \pm .008$	$2.6722 \pm .002$	$7.3015 \pm .009$	$5.0018 \pm .115$	$.7627 \pm .008$	$.5761 \pm _{.008}$	$.5237 \pm .008$	$.4246 \pm \scriptstyle .007$
	AA-kNN	$.3172 \pm .004$	$\textbf{2.0142} \scriptstyle \pm .012$	$\textbf{4.5597} \scriptstyle \pm .043$	$3.1429 \pm _{.148}$	$.7926 \pm _{.006}$	$.6024 \pm _{.005}$	$.5014 \pm _{.009}$	$.4432 \pm .008$
	LDLFs	$.4035 _{\pm .014}$	$2.5461 _{\pm .010}$	$6.8269 \pm _{.040}$	$1.6884_{\pm.115}$	$.6756 \pm _{.018}$	$.5318 \pm .013$	$.4164 \pm _{.013}$	$.3349 \pm \scriptstyle .010$
	DF-BFGS	$.2982 \pm _{.004}$	$2.4025 \pm .005$	$6.2416_{\pm .020}$	$.6304 \pm _{.012}$	$.8250 \pm _{.006}$	$.6220 _{\pm .004}$	$.5467 \pm _{.008}$	$.4454 {\scriptstyle \pm .007}$
-	KLD 0	$.2966 \pm _{.004}$	$2.4059 \pm .005$	$6.2558 \pm _{.020}$	$.6307 \pm _{.013}$	$.8243 \pm _{.006}$	$.6249 \pm _{.005}$	$.5470 \pm _{.008}$	$.4456 \pm \scriptstyle .007$
	S-KLD	$.2995 \pm _{.005}$	$2.4112 \pm .005$	$6.2883 \pm .020$	$.6491 \pm .013$	$.8203 \pm .006$	$.6205 \pm _{.005}$	$.5384 \pm .009$	$.4385 \pm _{.008}$
	SCL •	$.2977_{\pm.004}$	$2.4028 \pm .005$	$6.2435_{\pm .021}$	$.6262 \pm _{.013}$	$.8256 \pm _{.006}$	$.6233 _ { \pm .005 }$	$.5488 \pm _{.008}$	$.4471 \pm .007$
	S-SCL	$.2940 \pm _{.005}$	$2.4059 \pm .006$	$6.2589 _{\pm .023}$	$.6203 \pm _{.013}$	$.8268 \pm _{.006}$	$.6281 \pm _{.006}$	$.5518 \pm _{.008}$	$.4497 \pm \scriptstyle .007$
	LRR •	$.2984 \pm _{.004}$	$2.4046 \pm _{.005}$	$6.2525_{\pm .019}$	$.6351 _{\pm .012}$	$.8232 \pm .006$	$.6220 _{\pm .004}$	$.5443 \pm _{.008}$.4636 $_{\pm.007}$
	S-LRR	$\textbf{.2937} \pm .004$	$2.4056 \pm .005$	$6.2589 \pm .019$.6189 $_{\pm.013}$	$\textbf{.8271} \scriptstyle \pm .006$	$\textbf{.6283} \pm .005$	$\textbf{.5519} \scriptstyle \pm .008$	$.4498 \pm .007$

Table 16: Experimental results of LDL on the Flickr dataset formatted as (mean \pm std)

								· · · · ·
Algorithms	Cheby. \downarrow	$\texttt{Clark}\downarrow$	$\texttt{Can.}\downarrow$	$\texttt{KLD}\downarrow$	$\texttt{Cosine} \uparrow$	Int. \uparrow	Spear. \uparrow	Ken. ↑
LDSVR	$.5174 _{\pm .006}$	$2.6364_{\pm .002}$	$7.2094_{\pm .011}$	$5.0366_{\pm .086}$	$.6636 _{\pm .008}$	$.4683 _{\pm .006}$	$.4622 \pm _{.009}$	$.3811_{\pm .008}$
AA-kNN	$.3286 _{\pm .005}$	$\textbf{2.0685} \scriptstyle \pm .009$	$\textbf{4.9363}_{\pm.033}$	$2.2172 \pm .107$	$.7200 \pm .006$	$.5582 \pm _{.005}$	$.4265 \pm _{.009}$	$.3465 \scriptstyle \pm .007$
LDLFs	$.4051 \pm .011$	$2.4012 \pm .012$	$6.3262 \pm .050$	$1.4274 \pm _{.077}$	$.6073 \pm _{.015}$	$.4822 \pm .011$	$.3478 \pm \scriptstyle .014$	$.2847 \pm _{.012}$
DF-BFGS	$.3007 \pm _{.005}$	$2.1995 \pm .007$	$5.4900 \pm .025$	$.6309 _{\pm .011}$	$.7801 \pm _{.005}$	$.5979 \pm \scriptstyle .004$	$.5102 \pm .009$	$.4226 \pm \scriptstyle .008$
KLD 0	$.3015 _{\pm .005}$	$2.2008 \pm .007$	$5.4969 _{\pm .025}$	$.6348 \pm _{.012}$	$.7787 \pm _{.005}$	$.5973 _{\pm .004}$	$.5113 _{ \pm .009 }$	$.4234 \pm _{.008}$
S-KLD	$.3052 \pm _{.005}$	$2.2044 \pm .007$	$5.5222_{\pm .026}$	$.6485 \pm _{.012}$	$.7720 {\scriptstyle \pm .005}$	$.5926 _{\pm .004}$	$.5030 \pm _{.009}$	$.4166 \pm _{.008}$
SCL •	$.3280 \pm .013$	$2.2986 _{\pm .024}$	$5.9247_{\pm .099}$	$.8301 {\scriptstyle \pm .055}$	$.7268 \pm .018$	$.5713 _{\pm .015}$	$.4566 \pm _{.024}$	$.3748 \pm _{.022}$
S-SCL	.2929 $_{\pm.005}$	$2.2045 \pm .007$	$5.5289_{\pm .028}$	$.6113 \pm _{.012}$	$.7862 \pm .005$	$\textbf{.6070} \scriptstyle \pm .005$	$\textbf{.5265} \scriptstyle \pm .008$.4373 $_{\pm.007}$
LRR •	$.3057 \pm _{.005}$	$2.1969 \pm .007$	$5.4763 \pm _{.025}$	$.6431 \pm .012$	$.7752 \pm .006$	$.5929 \pm \scriptstyle .004$	$.5047 \pm _{.009}$	$.4229 \pm .008$
S-LRR	$.2938 \pm _{.005}$	$2.2013_{\pm .006}$	$5.5138_{\pm .023}$.6105 $_{\pm.012}$.7864 $_{\pm.005}$	$.6058 \pm _{.005}$	$.5261 \pm _{.009}$	$.4369 _{\pm .008}$

Table 17: Experimental results of LDL on the Natural_Scene dataset formatted as (mean \pm std)

	Algorithms	Cheby. \downarrow	$\texttt{Clark}\downarrow$	$\texttt{Can.}\downarrow$	$\texttt{KLD} \downarrow$	$\texttt{Cosine} \uparrow$	Int. \uparrow	Spear. \uparrow	Ken.↑
_	LDSVR	$.4899 \pm _{.016}$	$2.0831 \pm .025$	$5.7724_{\pm .092}$	$2.0862 \pm .085$	$.5740 \pm _{.017}$	$.4430 \pm .015$	$.4997 \pm _{.015}$	$.3695 \pm _{.012}$
	AA-kNN	$.3113 \pm .014$	$\textbf{1.9066} \scriptstyle \pm .034$	$\textbf{4.5413}_{\pm.110}$	$1.0874_{\pm .082}$	$.7113 \pm .015$	$.5636 _{\pm .013}$	$.4921 \pm _{.021}$	$.3518 _{\pm .016}$
	LDLFs	$.2808 \pm _{.034}$	$2.4329 \pm _{.024}$	$6.6027_{\pm.108}$.6464 $_{\pm.118}$	$.7679 _{\pm .046}$	$.5839 _{\pm .043}$	$.5406 \pm _{.058}$	$.4072 \pm _{.045}$
	DF-BFGS	$.3074 \pm \scriptstyle .013$	$2.4126 \pm .017$	$6.5896 \pm _{.072}$	$.7603 \pm _{.033}$	$.7381 \pm .013$	$.5568 \pm _{.011}$	$.5110 \pm .017$	$.3837 \pm _{.013}$
-	KLD •	$.3201 \pm .013$	$2.4242 \pm .017$	$6.6560 \pm .070$	$.8285 \pm _{.044}$	$.7172 \pm .015$	$.5485 \pm .011$	$.4958 \pm .016$	$.3715 \pm _{.012}$
	S-KLD	$.2743 \pm _{.013}$	$2.3866 _{\pm .020}$	$6.4733 _{\pm .077}$	$.6608 \pm .039$.7751 $_{\pm.014}$	$.6133 _{\pm .012}$	$.5592 {\scriptstyle \pm .017}$	$.4221 \pm .014$
	SCL •	$.3379 _{\pm .014}$	$2.4800 \pm .018$	$6.8659_{\pm .076}$	$.8867 {\scriptstyle \pm .035}$	$.7014 _{\pm .014}$	$.4801 \pm _{.014}$	$.4109 \pm _{.018}$	$.3025 \pm .013$
	S-SCL	$\textbf{.2733} \scriptstyle \pm .013$	$2.3734_{\pm .018}$	$6.4376 _{\pm .072}$	$.6703 \pm _{.043}$	$.7744 {\scriptstyle \pm .015}$	$.6156 \pm _{.013}$	$.5573 _{\pm .018}$	$.4207 \pm .014$
	LRR •	$.3138 \pm .013$	$2.4469 \pm .018$	$6.7118 \pm .074$	$.7703 \pm _{.032}$	$.7363 \pm _{.013}$	$.5456 \pm .011$	$.5056 \pm _{.016}$	$.3782 \pm .012$
	S-LRR	$.2740 \pm _{.018}$	$2.3461 _{\pm .023}$	$6.3467 _{\pm .087}$	$.6867 _{\pm .070}$	$.7715 \pm _{.021}$.6199 $_{\pm.017}$	$\textbf{.5595} \scriptstyle \pm .023$	$\textbf{.4228} \scriptstyle \pm .018$

Table 18: Experimental results of LDL on the Movie dataset formatted as $(mean \pm std)$

Algorithms	Cheby. \downarrow	$\texttt{Clark}\downarrow$	$\texttt{Can.}\downarrow$	$\texttt{KLD}\downarrow$	$\texttt{Cosine} \uparrow$	$\texttt{Int.} \uparrow$	Spear. \uparrow	Ken. ↑
CPNN	$.1337_{\pm.003}$	$.5639 _{\pm .010}$	$1.0746_{\pm .020}$	$.1191 _{ \pm .005 }$	$.9194_{\pm .003}$	$.8164 _{\pm .004}$	$.6610 _{ \pm .013 }$.7080 $_{\pm.002}$
AA-kNN	$.1223 \pm _{.002}$	$.5451 _{\pm .009}$	$1.0445_{\pm.018}$	$.1129 _{\pm .004}$	$.9254 _{\pm .003}$	$.8250 \pm .003$	$.6557_{\pm.011}$	$.5710 _{\pm .010}$
LDLFs	$.1172 \pm .003$	$.5233 _{ \pm .013 }$	$1.0134_{\pm .026}$	$.1086 _{\pm .006}$	$.9305 \pm _{.003}$	$.8324 \pm _{.004}$	$.6929 _{\pm .013}$	$.6051 _ { \pm .012 }$
DF-BFGS	$.1210 \pm .002$	$.5282 \pm .009$	$1.0158 \pm .019$	$.1084 \pm _{.005}$	$.9289 \pm .003$	$.8301 \pm _{.003}$	$.6848 \pm _{.012}$	$.5963 \pm _{.012}$
LRR •	$.1135 \pm .002$	$.5101 \pm .009$	$.9770 \pm _{.018}$	$.0957 \pm _{.004}$	$.9369 \pm _{.002}$	$.8385 \pm .003$	$.7119 \pm .011$	$.6203 \pm .011$
S-LRR	$.1125 \pm .002$	$.5086 _{\pm .009}$	$.9717_{\pm .018}$.0945 $_{\pm.004}$	$\textbf{.9376} \scriptstyle \pm .002$	$.8398 \pm _{.003}$.7126 $\scriptstyle \pm .011$	$.6227 \pm .011$
$QFD^2 \bullet$	$.1159 \pm _{.002}$	$.5200 \pm .009$	$.9920 _{\pm .018}$	$.0975 _{\pm .004}$	$.9355_{\pm .002}$	$.8357_{\pm.003}$	$.7075 \pm _{.011}$	$.6158_{\pm .011}$
S -QFD 2	$\textbf{.1123} \scriptstyle \pm .002$	$.5073 _{ \pm .009 }$	$.9700 \pm _{.018}$.0945 $_{\pm.004}$.9376 $_{\pm.002}$	$\textbf{.8401} \scriptstyle \pm .003$	$.7125 \pm .011$	$.6224 _{\pm .011}$
CJS •	$.1153 \pm _{.002}$	$.5127 _{\pm .009}$	$.9845 \pm _{.019}$	$.0984 \scriptstyle \pm .004$	$.9352 \pm .002$	$.8368 \pm _{.003}$	$.7103 \pm _{.012}$	$.6178 \pm .011$
S-CJS	$\textbf{.1123} \scriptstyle \pm .002$	$\textbf{.5072} \scriptstyle \pm .009$	$\textbf{.9699} \scriptstyle \pm .018$	$\textbf{.0945} \scriptstyle \pm .004$	$\textbf{.9376} \scriptstyle \pm .002$	$\textbf{.8401} \scriptstyle \pm .003$	$.7125 \pm .011$	$.6223 \pm .011$

Table 19: Experimental results of LDL on the fbp5500 dataset formatted as $(mean \pm std)$

Algorithms	Cheby. \downarrow	$\texttt{Clark}\downarrow$	$\texttt{Can.}\downarrow$	$\texttt{KLD}\downarrow$	$\texttt{Cosine} \uparrow$	Int. \uparrow	Spear. \uparrow	Ken.↑
CPNN	$.1864 \pm _{.005}$	$1.3367 \pm .009$	$2.3604 \pm .020$	$.1664 \scriptstyle \pm .005$	$.9281 \pm .004$	$.7958 \pm .005$	$.8688 \pm _{.005}$	$.7831 \pm .007$
AA-kNN	$.1515_{\pm.004}$	$\textbf{1.0443} \scriptstyle \pm .015$	$\textbf{1.7295} \scriptstyle \pm .031$	$.1846 _{\pm .016}$	$.9419 _{\pm .004}$	$.8317 {\scriptstyle \pm .005}$	$.8865 {\scriptstyle \pm .006}$	$.8123 \pm _{.008}$
LDLFs	$.1307 \pm _{.003}$	$1.2787_{\pm.010}$	$2.1703 \pm _{.024}$	$.1002 \pm .005$	$.9575 \pm _{.003}$	$.8552 \pm .004$	$.9060 \pm _{.005}$	$.8352 \pm \scriptstyle .007$
DF-BFGS	$.1341 _{ \pm .003 }$	$1.2889_{\pm.010}$	$2.1982 \pm .023$	$.1050 _{\pm .005}$	$.9551 _{\pm .003}$	$.8523 \pm _{.004}$	$.9047 \pm _{.005}$	$.8337 { }_{ \pm .007 }$
LRR	$.1312 _{\pm .003}$	$1.2767_{\pm.010}$	$2.1655_{\pm .024}$	$.1004 \pm _{.004}$	$.9575 \pm _{.002}$	$.8547 _{\pm .003}$	$.9059 _{\pm .004}$	$.8350 {\scriptstyle \pm .006}$
S-LRR	$\textbf{.1302} \scriptstyle \pm .003$	$1.2796 \pm .010$	$2.1717 \pm .024$	$\textbf{.0997} \pm .005$	$\textbf{.9576} \scriptstyle \pm .002$	$.8558 \pm .003$	$.9063 \pm _{.004}$	$.8425 \pm .006$
$QFD^2 \bullet$	$.1380 \pm .003$	$1.2803 \pm .010$	$2.1858 \pm _{.024}$	$.1084 \pm _{.005}$	$.9535 \pm .003$	$.8476 \pm _{.004}$	$.9021 \pm .004$	$.8297 \pm \scriptstyle .006$
S -QFD 2	$.1321 \pm .004$	$1.2811 \pm .010$	$2.1779 \pm _{.024}$	$.1027 \pm .006$	$.9561 {\scriptstyle \pm .003}$	$.8537 \pm \scriptstyle .004$	$.9044 \pm .005$	$.8330 \pm \scriptstyle .007$
CJS •	$.1343 \pm _{.003}$	$1.3057_{\pm.010}$	$2.2374 \pm _{.024}$	$.1084 \pm _{.005}$	$.9544 {\scriptstyle \pm .003}$	$.8527 \pm _{.004}$	$.9044 { }_{ \pm .004 }$	$.8334 \pm _{.006}$
S-CJS	$\textbf{.1302} \scriptstyle \pm .003$	$1.2802 \pm .010$	$2.1731 \pm _{.024}$.0997 $_{\pm.005}$	$.9575 \pm _{.002}$	$\textbf{.8559} \scriptstyle \pm .003$	$\textbf{.9066} \scriptstyle \pm .004$	$\textbf{.8429} \scriptstyle \pm .007$

Table 20: Experimental results of IncomLDL on the JAFFE dataset formatted as $(mean \pm std)$

Algorithms				$\omega =$	= 20%			
Augoriumis	Cheby. \downarrow	$\texttt{Clark}\downarrow$	$\mathtt{Can.}\downarrow$	$\texttt{KLD}\downarrow$	$\texttt{Cosine} \uparrow$	$\texttt{Int.} \uparrow$	Spear. \uparrow	Ken. ↑
IncomLDL \bullet	$.0898 \pm _{.010}$	$.3304 _{\pm .024}$	$.6742 \pm _{.049}$.0425 $_{\pm.007}$.9598 $_{\pm.007}$	$.8861 _{\pm .009}$	$.4742 \pm _{.094}$	$.4017_{\pm.082}$
\mathcal{S} -IncomLDL	$\textbf{.0863} \scriptstyle \pm .013$	$\textbf{.3179} \scriptstyle \pm .036$	$\textbf{.6525} \scriptstyle \pm .077$	$.0433 \pm _{.012}$	$.9590 \pm _{.011}$	$\textbf{.8893} \pm .014$	$\textbf{.5034} \pm .114$	$\textbf{.4401} \scriptstyle \pm .100$
				(<i>u</i>) =	40%			
Δlgorithms					1070			
Algorithms	Cheby. \downarrow	$\texttt{Clark}\downarrow$	$\texttt{Can.}\downarrow$	$\texttt{KLD}\downarrow$	Cosine ↑	Int. \uparrow	Spear. \uparrow	Ken. ↑
Algorithms	Cheby.↓ .0946±.010	$\begin{array}{c} \texttt{Clark} \downarrow \\ .3454 \scriptstyle \pm .026 \end{array}$	Can.↓ .7073±.053	$KLD \downarrow$.0465 ±.007	Cosine ↑ .9558 ±.007	Int.↑ .8801±.010	Spear.↑ .4231±.086	Ken. ↑ .3534±.074

Table 21: Experimental results of IncomLDL on the SBU_3DFE dataset formatted as (mean \pm std)

Algorithms		$\omega = 20\%$							
Augonumis	Cheby. \downarrow	$\texttt{Clark}\downarrow$	$\texttt{Can.}\downarrow$	$\texttt{KLD}\downarrow$	$\texttt{Cosine} \uparrow$	$\texttt{Int.} \uparrow$	Spear. \uparrow	Ken. ↑	
IncomLDL •	$.1088 \pm .003$	$.3586 _{\pm .008}$	$.7586 {\scriptstyle \pm .017}$	$.0574 _{\pm .003}$	$.9439 _{\pm .003}$	$.8655_{\pm .003}$	$.3171 _{\pm .027}$	$.2746 _{\pm .023}$	
S-IncomLDL	.1014 $_{\pm.004}$	$\textbf{.3208} \scriptstyle \pm .009$.6727 $_{\pm.019}$	$\textbf{.0516} \scriptstyle \pm .003$	$\textbf{.9485} \scriptstyle \pm .003$	$\textbf{.8790} \scriptstyle \pm .004$	$\textbf{.4255} \scriptstyle \pm .026$	$\textbf{.3679} \scriptstyle \pm .023$	
Algorithms				$\omega =$	= 40%				
Algorithms	Cheby.↓	$Clark\downarrow$	Can.↓	$\omega =$ KLD \downarrow	= 40% Cosine ↑	Int. ↑	Spear. ↑	Ken. ↑	
Algorithms	Cheby.↓ .1104±.003	Clark \downarrow .3621 \pm .008	Can.↓ .7673±.017	$\omega =$ KLD \downarrow .0586 \pm .003	= 40% Cosine ↑ .9426±.003	Int.↑ .8638±.003	Spear. ↑ .3003 ±.024	Ken. ↑ .2595±.020	

Table 22: Experimental results of LDL4C on JAFFE and Twitter formatted as (mean \pm std)

Algorithms	JA	FFE	Algorithms	Twitter		
Augoriumis	$0/1 \log \downarrow$	Err.prob. \downarrow	Augonumis	$0/1 \texttt{loss} \downarrow$	Err.prob. \downarrow	
LDL4C •	.4973 $_{\pm.108}$.7665 ±.020	LDL4C	.9081 ±.009	.8846 $_{\pm.005}$	
S-LDL4C	.4453 $_{\pm.102}$.7600 ±.019	S-LDL4C	.8714 $_{\pm.207}$.8729 $_{\pm.156}$	
LDL-HR	.4786 $_{\pm.097}$.7676 ±.020	LDL-HR •	.3656 $_{\pm.017}$.4928 $_{\pm.011}$	
S-HR	.4653 $_{\pm.105}$.7655 ±.019	S-HR	.2753 $_{\pm.013}$.4250 $_{\pm.008}$	
LDLM	.4787 $_{\pm.109}$.7687 $_{\pm.021}$	LDLM •	.2814 $\pm .013$.4291 $_{\pm.008}$	
S-LDLM	.4737 $_{\pm.097}$.7689 ±.019	S-LDLM	.2753 $\scriptstyle \pm .014$.4250 $\pm .008$	

Table 23: Experimental results of LDL4C on sBU_3DFE and Flickr formatted as (mean \pm std)

Algorithms	sBU	_3DFE	Algorithms	Flickr		
Aigonumis	$0/1 \texttt{loss} \downarrow$	Err.prob. \downarrow	Augoritimis	$0/1 \texttt{loss} \downarrow$	Err.prob. \downarrow	
LDL4C	.5578 ±.028	.7671 ±.007	LDL4C	.8971 ±.008	.8884 $\pm .004$	
S-LDL4C	.5526 $_{\pm.025}$.7686 ±.006	S-LDL4C	.8705 ±.138	.8702 $_{\pm.100}$	
LDL-HR •	.5167 ±.027	.7596 ±.006	LDL-HR •	.4513 $_{\pm.015}$.5823 $\pm .007$	
S-HR	.5069 $_{\pm.025}$.7598 ±.006	S-HR	.4219 $_{\pm.015}$.5639 $\pm .007$	
LDLM •	.5258 $_{\pm.034}$.7619 ±.009	LDLM •	.4384 $_{\pm.014}$.5740 $_{\pm.007}$	
S-LDLM	.4809 $_{\pm.024}$.7524 $_{\pm.005}$	S-LDLM	.4321 $_{\pm.016}$.5667 $_{\pm.007}$	

Table 24: Experimental results of LE on the JAFFE dataset formatted as $(mean \pm std)$

Algorithms	Cheby. \downarrow	$\texttt{Clark}\downarrow$	$\mathtt{Can.}\downarrow$	$\texttt{KLD}\downarrow$	$\texttt{Cosine} \uparrow$	$\texttt{Int.} \uparrow$
LP	.0812 $_{\pm.001}$.3446 $\pm .002$.7125 ±.005	.0424 $\pm .001$.9618 $_{\pm.001}$.8808 $\pm .001$
GLLE	.0821 $_{\pm.002}$.3196 $_{\pm.013}$.6518 $_{\pm.028}$.0386 $_{\pm.003}$.9638 $_{\pm.002}$.8901 $_{\pm.004}$
LEVI	.0787 $_{\pm.003}$.3316 $_{\pm.013}$.6864 $_{\pm .028}$.0391 $_{\pm.003}$.9649 $_{\pm.002}$	$.8860 {~\pm .004}$
LIBLE •	.0813 ±.006	.3106 ±.020	.6358 $_{\pm.044}$.0370 ±.005	.9652 $_{\pm.005}$.8929 ±.008
S-LIBLE	.0770 $_{\pm.003}$.2942 $_{\pm.007}$.5997 $_{\pm.016}$.0332 $_{\pm.002}$.9685 $_{\pm.002}$.8987 $_{\pm.003}$

Table 25: Experimental results of LE on the Yeast_heat dataset formatted as (mean \pm std)

-						
Algorithms	Cheby. \downarrow	$\texttt{Clark}\downarrow$	$\mathtt{Can.}\downarrow$	$\texttt{KLD}\downarrow$	$\texttt{Cosine} \uparrow$	Int. \uparrow
LP	.0421 ±.000	.2148 $\pm .000$.4711 ±.001	.0153 ±.000	.9860 ±.000	.9235 ±.000
GLLE	.0481 $\pm .001$.2114 $\pm .005$.4282 $_{\pm.011}$.0168 $_{\pm.001}$.9842 $\pm .001$.9298 $_{\pm.002}$
LEVI	.0494 $_{\pm.007}$.2125 $_{\pm.027}$.4307 ±.056	.0169 $_{\pm.004}$.9838 $_{\pm.004}$.9289 $_{\pm.009}$
LIBLE •	.0453 ±.000	.1973 ±.001	.3982 ±.003	.0148 $_{\pm.000}$.9859 ±.000	.9346 $_{\pm.000}$
S-LIBLE	.0445 $_{\pm.000}$.1901 $_{\pm.002}$.3790 $_{\pm.005}$.0137 $\scriptstyle \pm .000$.9869 $_{\pm.000}$.9376 $_{\pm.001}$