



LSTC: When label-specific features meet third-order label correlations

Xing-Yi Zhang^{a,b}, Fan Min^{a,b,e,*}, Guojie Song^c, Hong Yu^d

^a School of Computer Science, Southwest Petroleum University, Chengdu 610500, China

^b Lab of Machine Learning, Southwest Petroleum University, Chengdu 610500, China

^c School of Sciences, Southwest Petroleum University, Chengdu 610500, China

^d College of Computer Science and Technology, Chongqing University of Posts and Telecommunications, Chongqing 400065, China

^e Institute for Artificial Intelligence, Southwest Petroleum University, Chengdu 610500, China

ARTICLE INFO

Keywords:

Label correlations
Label-specific features
Multi-label learning
Neural network

ABSTRACT

Multi-label learning (MLL) copes with applications such as text classification and image recognition where each instance is associated with multiple labels. Currently, some label-specific feature conversion approaches such as LIFT serve each individual label well, however do not take advantage of label correlation. Some label correlation exploitation approaches such as GLOCAL have independent model prediction and label correlation stages, thus limiting their learning ability. This paper introduces LSTC, a new algorithm that generates label-specific features and exploits third-order label correlation. On the one hand, for each label, we choose the same number of positive and negative representative instances inspired by density peak. These instances are used to generate a new data matrix with label-specific features. On the other hand, we train a paired output prediction network for each label based on these matrices of its own and that of two auxiliary labels. In this way, third-order label correlations are implicitly exploited. Specifically, the two auxiliary labels are the most similar and least similar, respectively, to avoid getting stuck in local optima. Experiments are conducted on seventeen benchmark datasets in comparison with ten popular algorithms. Results on eight measures demonstrate the superiority of LSTC on data from various domains except the text-domain. The source code is available at github.com/fansmale/lstc.

1. Introduction

Multi-label learning (MLL) [21] has become a popular machine learning task due to its wide range of applications. In text classification [17,21,26], MLL can simultaneously model the topical aspects of the feature data in the text, such as papers, authors, and publication venues. It provides services such as recommendation and information retrieval. In image classification [18], MLL can easily use more labels to embody a large amount of information in image compared with multi-classification. It also utilizes the connection between latent information and supervision information to achieve more accurate predictions. In video annotation [30], the use of MLL improves the reliability of video browsing, searching, and navigation. Specifically, by label correlations, the complex concepts in video can be better inferred based on the correlations with the other concepts.

* Corresponding author at: School of Computer Science, Southwest Petroleum University, Chengdu 610500, China.
E-mail address: minfan@swpu.edu.cn (F. Min).

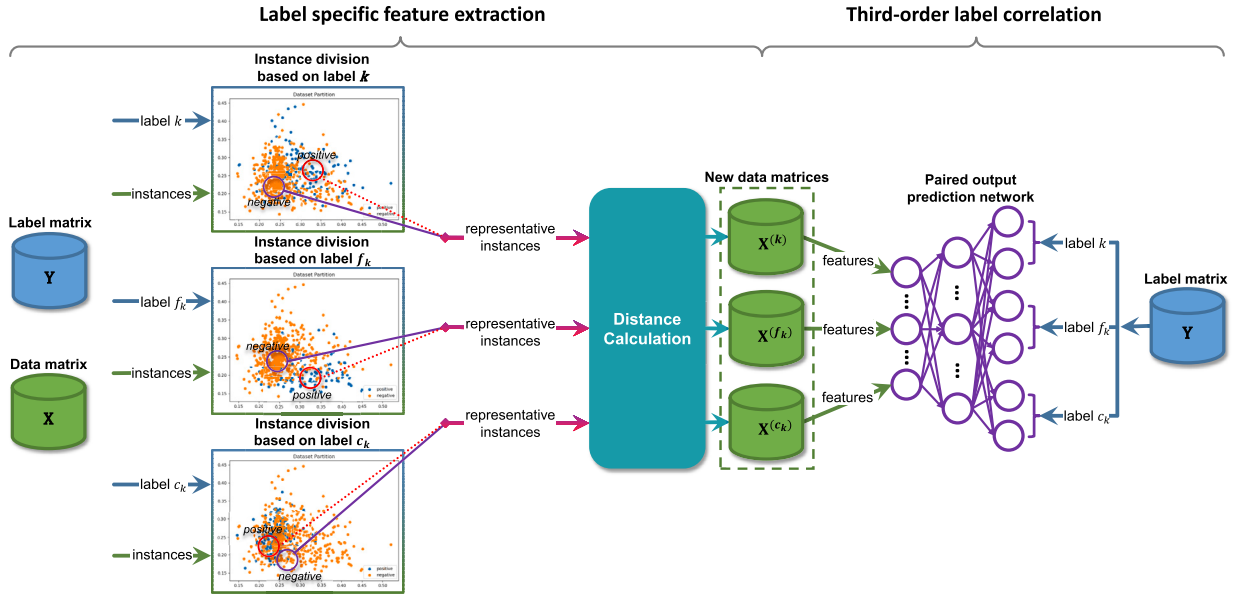


Fig. 1. The framework of LSTC. The left part is feature conversion through representative instance selection and distance calculation; the right part is the neural network training using data from the current label and two auxiliary labels.

There are various methods that deal with different characteristics of MLL tasks. When an instance contains a label (e.g., there is a cat in a given picture), we call them a positive instance and a positive label, respectively. Otherwise, we call them a negative instance and a negative label, respectively. For label sparsity [4,16,31], i.e., few (often < 5%) positive labels, LIFT [45] extracts the same number of features for positive and negative instances. Deep pre-trained transformer models [4] are also good at this issue. For data sparsity, i.e., most data features are missing, CNN models [17] are often effective. For missing labels [31], DM2L [19] imposes a high-rank structure on predictions from instances with different labels to provide better discriminability. For noisy labels, [38] adopts HISC technique to handle partial multi-label without recovering ground-truth labels. For label correlation exploitation [11], ACKEL [41] transforms MLL into a set of multi-class learning problems. Meanwhile, BP-MLL [46] and MASP [22] use neural network to embody this correlation.

LIFT [45] is an original method with two main advantages. One is that it opens up a new idea of multi-label classification with embedded features, i.e., *label-specific features*. This is done by clustering the positive and negative instances of each label and then forming new features based on the cluster centers. Another is to handle label sparsity (or equivalently, imbalanced data). This is done by building an equal number of features from positive/negative instances of the imbalanced label [24]. However, LIFT does not consider label correlation, which leads to the overlook of important information.

In this paper, we propose LSTC, a new algorithm that generates label-specific features and exploits third-order label correlations. Fig. 1 presents the main framework of LSTC, which follows the feature conversion and model training scheme of LIFT. In the first stage, each label independently guides the partitioning of the data matrix X . We select a set of representative instances in each partition through a clustering algorithm. Then, the distance between each instance in X and those representative instances is calculated to construct a new data matrix $X^{(k)}$. In the second stage, the features of these new data matrices are fed into the network in parallel for training. Meanwhile, the new data matrices (denoted by $X^{(f_k)}$ and $X^{(c_k)}$) of two auxiliary labels participate in the network training of the current label. The contribution of this algorithm is three-fold.

- We adopt the density peak under Gaussian kernel to select representative instances, which are employed further to construct new features. Compared with the k -Means clustering technique employed by LIFT [45], our representative instances are not influenced by outliers. Consequently, our new features are more stable. Regarding algorithm implementation, we design techniques to obtain a trade-off between time and space complexity.
- We use data from two auxiliary labels to help train the predictive neural network for each label. In this way, the third-order label correlations are implicitly exploited. We require auxiliary labels to be the most and least similar according to Manhattan distance. The idea behind is to avoid getting stuck in local optima. Also, the network has paired outputs for more stable results. Compared with other popular label correlation approaches that convert MLL to multi-class learning [3,39,41] or consider association rules [8], our algorithm does not change the relationship between data and labels.
- For prediction, we use the softmax function to obtain numerical outputs from paired neural network outputs. Compared with simple normalization, this function is more reasonable especially when the output values are close. With numerical output, our algorithm can be compared with other ranking-based measures that are more appropriate than accuracy or F_1 -measure for MLL tasks.

Table 1
Notations.

Notation	Meaning
$\mathbf{X} = [x_{ij}]_{n \times d}$	Data matrix
$\mathbf{x}_i = [x_{i1}, \dots, x_{id}]$	The i -th instance
$\mathbf{Y} = [y_{ik}]_{n \times q}$	Label matrix
$\mathbf{y}_i = [y_{i1}, \dots, y_{iq}]$	Labels of \mathbf{x}_i
$\mathbf{X}^{(k)} = [x_{ij}^{(k)}]_{n \times 2m_k}$	Data matrix after k -th label feature conversion
$\mathbf{P}^{(k)}$	The set of positive instances wrt. the k -th label
$\mathbf{N}^{(k)}$	The set of negative instances wrt. the k -th label
Θ_k	Network model for the k -th label

Experiments were undertaken on seventeen datasets on label-based measures including Micro-AUC and macro-AUC, sample-based measures including Hamming Loss, ranking-based measures including Peak F_1 -score, NDCG, One-Error, Coverage and Ranking Loss. Results show that: 1) All components of LSTC are appropriate; 2) LSTC outperforms ten state-of-the-art algorithms on twelve datasets from various domains; and 3) LSTC is inappropriate for five text-domain data.

The rest of this paper is organized as follows. Section 2 introduces some preliminary knowledge and related work. Section 3 presents our new algorithm with complexity analysis. Section 4 presents experimental results. Finally, Section 5 draws our conclusions and outlines further research directions.

2. Related work

This section first introduces the data model used by LSTC. Then we will discuss some techniques for dealing with MLL, including feature conversion and label correlation exploitation. Some of them will be compared with LSTC in the experimental section.

2.1. Data model

Table 1 lists some important notations used throughout the paper. The data matrix is denoted by $\mathbf{X} = [x_{ij}]_{n \times d} \in \mathbb{R}^{n \times d}$, where n and d are the number of instances and features, respectively. The i -th row, denoted by \mathbf{x}_i , represents an instance, e.g., a piece of music or an image, which is described by a set of features with distinct semantics. The label matrix is denoted by $\mathbf{Y} = [y_{ik}]_{n \times q} \in \{-1, +1\}^{n \times q}$, where q is the number of labels. The i -th row, denoted by \mathbf{y}_i , is the label array of \mathbf{x}_i . $y_{ik} = +1$ indicates \mathbf{x}_i has the k -th label, while $y_{ik} = -1$ indicates it does not have this label. Obviously, an instance \mathbf{x}_i has q labels, and these q labels form \mathbf{y}_i .

2.2. Feature conversion and label correlation approaches to MLL

Feature conversion for MLL is more important than that of traditional classification (e.g., single-label learning). The first category is to learn label-specific features. The original method is LIFT [45], which is based on the idea that different labels are determined by different information. This is the first study that introduces label-specific feature conversion. However, LIFT does not make full use of the information between labels, nor does it remove redundant information for extracted features. LLSF [13] considers second-order label correlation by analyzing the features of two labels. The shared features of two strongly correlated labels extend the information that a classifier can learn. But LLSF weakens the possible contribution of irrelevant labels to learned information. FRS-SS-LIFT [43] adopts a fuzzy rough set and sample selection to reduce unnecessary interference in feature space. It removes the redundant information of extracted features, but does not fully consider the correlation between different labels. LSML [10] considers missing labels through learning a $q \times q$ supplementary matrix for high-order label correlations. It addresses inaccurate label correlation exploitation due to missing labels. However, it only considers positive label correlations, which may cause limitations for some datasets [5]. The second category is closely related to data clustering. CBMLC [27] decomposes a dataset into several disjoint data clusters, and extracts features for each cluster. It improves the learning speed when the number of labels is large, but reduces the data in each classifier, so the performance might degrade [25]. OMLC [28] uses a decay mechanism during training and changes features during evolution. It combines multi-label and clustering algorithms under the framework of an incremental learning approach. IMCC [34] applies data augmentation techniques to obtain virtual centers of clusters so that local label correlation can be learned. But the limited number of virtual centers obtained by clustering may limit the overall performance.

Label correlation exploitation is essential to MLL. The first category is to analyze the label matrix independently. GLOCAL [48] decomposes the label matrix into two low-rank subspaces and then builds a linear prediction model for the first subspaces. It is the first MLL algorithm to exploit global and local correlations of the label matrix, but it is unsuitable for handling asymmetric correlations. k -Means based Apriori multi-label classification [8] uses k -Means to divide pre-processing space of Apriori [1] in advance, thereby generating association rules for each cluster to independently predict each label. This algorithm shows that association rules are also good strategies for extracting information from the label matrix. LDML [32] transforms logical label into label distribution with an approach of random variable distribution with granular computation. It can use mutual information and label enhancement to eliminate redundant and irrelevant information. But the discretization of features may also lead to the loss of key information [37]. LFCMLL [12] learns a sparse weight parameter vector for each label through a linear regression model and directly constrains the label correlation in output. This algorithm combines the reconstruction coefficient matrix with the label-specific features. It is further

demonstrated that the information matrix facilitates label correlation. The second category is to consider a set of labels in each learning subproblem. LP [3] considers a label subset each time, therefore converting MLL into a set of multi-class learning problems. But LP lacks optimization for label subsets, resulting in data starvation and high runtime. RAKEL [39] considers the runtime of LP, thus it decomposes the initial label set into length- k subsets randomly. Although the performance has been significantly improved, the class-imbalance [24] in label subset still exists. As an extension of RAKEL, ACKEL [41] considers factors such as separability and class balance in the process of label subset construction. Although ACKEL generally outperforms RAKEL, its label selection is time-consuming.

3. The proposed approach

In this section, we present our approach, which consists of feature conversion, network training, and label prediction.

3.1. Feature conversion

Similar to LIFT [45], we generate new features for each label. This is done by selecting a few representative instances. The distances between all instances and all representative instances will form a new data matrix $\mathbf{X}^{(k)}$. In this way, the number of new features is equal to the number of representative instances. Therefore, we are faced with two sub-problems: 1) how to determine the number of representative instances, and 2) how to calculate the representativeness of instances?

3.1.1. Determination of the number of representative instances

When we are concerned with the k -th label, the data matrix \mathbf{X} can be split in two as

$$\mathbf{P}^{(k)} = \{\mathbf{x}_i \mid y_{ik} = +1\} \quad (1)$$

and

$$\mathbf{N}^{(k)} = \{\mathbf{x}_i \mid y_{ik} = -1\}. \quad (2)$$

Due to label sparsity, in many datasets $|\mathbf{P}^{(k)}| \ll |\mathbf{N}^{(k)}|$. We will select the same number of representative instances from $\mathbf{P}^{(k)}$ and $\mathbf{N}^{(k)}$ to embody the equal importance of positive and negative instances. We set this number to

$$m_k = \min \left\{ R, \left\lceil r |\mathbf{P}^{(k)}| \right\rceil, \left\lceil r |\mathbf{N}^{(k)}| \right\rceil \right\}, \quad (3)$$

where R is the user-specified upper bound, r is the control ratio with domain $[0, 1]$. The idea behind is that the number of representative positive and negative instances should not exceed their respective totals, nor should they have excessive features. Consequently, the total number of representative instances is $2m_k \leq 2R$.

3.1.2. Instance representativeness calculation

We compute instance representativeness inspired by the density peak clustering algorithm [33]. There are two factors [40] that characterize the instance representativeness: 1) its density, and 2) its independency (i.e., its distance to master).

When we consider all instances in \mathbf{X} , we obtain the *global density* under the Gaussian kernel as

$$\rho_i = \sum_{j \neq i} e^{-\left(\frac{d_{ij}}{d_c}\right)^2}, \quad (4)$$

where d_{ij} is the Euclidean distance between \mathbf{x}_i and \mathbf{x}_j , and d_c is a custom threshold describing the degree of distance decay.

However, when dealing with data conversion for the k -th label, we would like to calculate the density of \mathbf{x}_i considering $\mathbf{P}^{(k)}$ if $y_{ik} = +1$, and $\mathbf{N}^{(k)}$ otherwise. The *local density* is therefore given by

$$\rho_i^{(k)} = \begin{cases} \sum_{\mathbf{x}_j \in \mathbf{P}^{(k)} \setminus \{\mathbf{x}_i\}} e^{-\left(\frac{d_{ij}}{d_c}\right)^2}, & \text{if } y_{ik} = +1; \\ \sum_{\mathbf{x}_j \in \mathbf{N}^{(k)} \setminus \{\mathbf{x}_i\}} e^{-\left(\frac{d_{ij}}{d_c}\right)^2} = \rho_i - \sum_{\mathbf{x}_j \in \mathbf{P}^{(k)}} e^{-\left(\frac{d_{ij}}{d_c}\right)^2}, & \text{otherwise,} \end{cases} \quad (5)$$

where the expression for the second condition is rewritten for reducing time complexity, as will be discussed later.

The independency of \mathbf{x}_i is given by

$$\delta_i^{(k)} = \min_{y_{jk} = y_{ik}, \rho_j^{(k)} > \rho_i^{(k)}} d_{ij}, \quad (6)$$

where $y_{jk} = y_{ik}$ indicates that the instance is only compared these with the same class value in terms of the k -th label. Let

$$z = \operatorname{argmin}_{y_{jk} = y_{ik}, \rho_j^{(k)} > \rho_i^{(k)}} d_{ij}, \quad (7)$$

then \mathbf{x}_z is called the *local master* of \mathbf{x}_i [40].

The representativeness of \mathbf{x}_i wrt. label k is therefore computed by

$$\gamma_i^{(k)} = \rho_i^{(k)} \delta_i^{(k)}. \quad (8)$$

3.1.3. New data matrix construction

According to representativeness, we obtain the most representative m_k positive instances $\mathbf{P}_c^{(k)} = \{\mathbf{p}_1^k, \dots, \mathbf{p}_{m_k}^k\} \subset \mathbf{P}^{(k)}$ and the most representative m_k negative instances $\mathbf{N}_c^{(k)} = \{\mathbf{n}_1^k, \dots, \mathbf{n}_{m_k}^k\} \subset \mathbf{N}^{(k)}$. The new data matrix is given by

$$\mathbf{X}^{(k)} = \begin{pmatrix} \mathbf{x}_{ij}^{(k)} \end{pmatrix}_{n \times 2m_k}, \quad (9)$$

where

$$\mathbf{x}_{ij}^{(k)} = \begin{cases} d(\mathbf{x}_i, \mathbf{p}_j^k), & \text{if } j \leq m_k; \\ d(\mathbf{x}_i, \mathbf{n}_{j-m_k}^k), & \text{otherwise,} \end{cases} \quad (10)$$

where $d(\mathbf{x}_i, \mathbf{p}_j^k)$ stand for Euclidean distance between instance \mathbf{x}_i and representative instances \mathbf{p}_j^k . It is a feasible approach to alleviate the single-label imbalance by replacing features with distances between general instances and positive/negative instances [24,45].

3.1.4. Algorithm description and analysis

Algorithm 1 LSTC feature conversion for all labels.

Input: Data matrix \mathbf{X} , label matrix \mathbf{Y} , target label k .

Output: The new data matrices for all labels.

```

1: Compute  $\rho_i$  for  $1 \leq i \leq n$  according to Eq. (4);
2: for  $1 \leq k \leq q$  do
3:   Form  $\mathbf{P}^{(k)}$  and  $\mathbf{N}^{(k)}$  based on  $k$ -th label according to Eqs. (1) and (2);
4:   Compute the density  $\rho_i^{(k)}$  for each  $\mathbf{x}_i$  in  $\mathbf{P}^{(k)}$  and  $\mathbf{N}^{(k)}$  according to Eq. (5);
5:   Compute the independence  $\delta_i^{(k)}$  according to Eq. (6);
6:   Compute the representativeness  $\gamma_i^{(k)}$  according to Eq. (8);
7:   Construct the most representative  $m_k$  instances  $\{\mathbf{p}_1^k, \dots, \mathbf{p}_{m_k}^k\}$  in  $\mathbf{P}^{(k)}$  and  $\{\mathbf{n}_1^k, \dots, \mathbf{n}_{m_k}^k\}$  in  $\mathbf{N}^{(k)}$  according to  $\gamma_i^{(k)}$ ;
8:   Construct  $\mathbf{X}^{(k)}$  through all instances in  $\mathbf{X}$  and  $2m_k$  representative instances extracted for the  $k$ -th label according to Eqs. (9) and (10);
9: end for
10: Return  $(\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(q)})$ ;

```

Algorithm 1 illustrates the complete process of LSTC feature conversion for all labels. Step 1 goes through preprocessing to obtain the global density. Step 3 splits the data matrix \mathbf{X} into positive and negative matrices $\mathbf{P}^{(k)}$ and $\mathbf{N}^{(k)}$. Steps 4 to 6 obtain representativeness by sequentially calculating the local density and independence for both matrices $\mathbf{P}^{(k)}$ and $\mathbf{N}^{(k)}$. Steps 7 to 8 calculates the distance between the representative instances and each instance to construct $\mathbf{X}^{(k)}$. Finally, we calculate $\mathbf{X}^{(k)}$ corresponding to each label, obtain $(\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(q)})$ and return.

Now we consider the time and space complexity of Algorithm 1. First, in order to save some time overhead, we will pre-calculate global density ρ_i according to (4). The time complexity of calculating each of instance is $O(dn)$. So the time complexity of global preprocessing is $O(dn^2)$. At the same time, during this process, we calculate the distance between all instances. However, these distance values actually will be reused in subsequent calculations. Therefore, in order to avoid multiple calculation, we will consume $O(n^2)$ space complexity to store them.

For each label k , the time complexity of feature conversion can be described in five parts. For simplicity, in the following description we will replace $|\mathbf{N}^{(k)}|$ and $|\mathbf{P}^{(k)}|$ with s_n and s_p .

- 1) According to Eqs. (1) and (2), the time complexity of dividing the data matrix \mathbf{X} based on positive/negative labels is $O(n)$.
- 2) After dividing the positive/negative matrices, the representative selection strategy is performed on the two matrices respectively.
 - a) For the negative matrix $\mathbf{N}^{(k)}$, we can assist the local density $\rho_i^{(k)}$ computation with global density ρ_i according to Eq. (5). Meanwhile, we also can know that the time complexity of computing its all local densities $\rho_i^{(k)}$ is $O(s_n s_p)$.
 - b) For the positive matrix $\mathbf{P}^{(k)}$, since the number of instances is sufficiently small compared to n . Therefore, all local densities can be calculated by traversing the matrix $\mathbf{P}^{(k)}$ twice. Then the time complexity is $O(s_p^2)$.
- 3) According to Eq. (6), determining the master of any instance \mathbf{x}_i also requires traversing the entire matrix with time complexity $O(dn)$. Therefore, the total time complexity of determining each $\delta_i^{(k)}$ of $\mathbf{P}^{(k)}$ and $\mathbf{N}^{(k)}$ is $O(s_p^2 + s_n^2)$.
- 4) The time complexity of obtaining a representativeness array is $O(s_p + s_n) = O(n)$ through Eq. (8). We can then use a *simple selection sort* with complexity $O(R(s_p + s_n)) = O(Rn)$ to filter out the most representative R instances.
- 5) Then according to Eqs. (9) and (10), the time complexity of constructing the new matrix $\mathbf{X}^{(k)}$ is $O(2Rn)$.

Let $K = \max_{1 \leq k \leq q} s_p$, in case of sparse labels, $K \ll n$. Similarly, because of label sparsity, we can also approximate $s_p^2 + s_n^2$ as n^2 . The total time complexity is shown in Table 2.

Table 2
Time complexity of Algorithm 1.

Part	Complexity	Description
Line 1	$O(dn^2)$	Compute the global density
Line 2	$O(q)$	Enumerate all labels
Line 3	$O(n)$	Divide positive/negative based on k -th label
Line 4	$O(nK)$	Density calculation
Line 5	$O(n^2)$	Independence calculation
Line 6	$O(n)$	Representativeness calculation
Line 7	$O(Rn)$	Representative instances selection
Line 8	$O(2Rn)$	Construct $\mathbf{X}^{(k)}$
Total	$O(dn^2 + q(n + nK + n^2 + n + Rn + 2Rn) = O((d + q)n^2)$	

In general, the setting of R is small compared to dn and ds_p . Therefore, the time complexity can be simplified to $O((d + q)n^2)$ on feature conversion.

Regarding the analysis of space complexity, we can find that the largest storage overhead in the whole process is to use the space complexity of $O(n^2)$ to store the distance between all instances. Although for some large datasets, such space overhead will take up too much memory. However, in this algorithm, this space-for-time scheme can greatly save the program runtime.

In general, the process of converting the original feature matrix \mathbf{X} into $\mathbf{X}^{(k)}$ is similar to the LIFT algorithm, but the difference between this feature conversion approach from LIFT are two-folds.

- 1) LSTC adopts a more efficient and convenient density peak clustering algorithm to compute representative instances. First, the density peak algorithm is not only appropriate for spherical data. Second, the clustering center is a concrete point, not a virtual center points. Third, it is insensitive to outliers. Last, its complexity is smaller than traditional clustering schemes like k -Means.
- 2) Because LSTC considers the correlation of multiple labels, the final number of features input to the network is the sum of converted features of three labels, thus LSTC adopts a different control Eq. (3) for the number of clusters than LIFT. The details will be explained in Section 3.2.

3.2. Model training

We use two auxiliary matrices and the respective labels to help train the model of each label. Specifically, for label k , let the most similar and least similar labels be c_k and f_k , respectively. One is the most similar, given by

$$c_k = \operatorname{argmin}_{j \neq k} \sum_{i=1}^n |y_{ik} - y_{ij}|, \quad (11)$$

while the other is the least similar, given by

$$f_k = \operatorname{argmax}_{j \neq k} \sum_{i=1}^n |y_{ik} - y_{ij}|. \quad (12)$$

It is a natural idea that the most similar labels serve as auxiliary information. We could also use the more similar labels, such as 2NN, however the enhancement is trivial. The use of the least similar is to avoid overfitting and getting stuck in local optima. We will show that the effectiveness of technique over the 2NN technique.

The input of the Θ_k is given by $\mathbf{X}^{(k)}$, $\mathbf{X}^{(c_k)}$, and $\mathbf{X}^{(f_k)}$. Therefore, Θ_k has $2(m_k + m_{c_k} + m_{f_k})$ input nodes.

Fig. 2 illustrates an example of our prediction network. There are three parts of the input, one for each label. Then current instance is $\mathbf{x}_i^{(k)} = [x_{i,1}^{(k)}, \dots, x_{i,2m_k}^{(k)}, x_{i,1}^{(c_k)}, \dots, x_{i,2m_{c_k}}^{(c_k)}, x_{i,1}^{(f_k)}, \dots, x_{i,2m_{f_k}}^{(f_k)}]$, which determines the features of the three related labels. Meanwhile, in order to fully embody the correlation between the k -th label and its two auxiliary labels during the training process, the output of Θ_k will be composed of three pairs of ports with domain $[0, 1]$, corresponding to label k , label c_k , and label f_k respectively.

Therefore the network Θ_k has six output nodes. The network is formally denoted by a mapping

$$F_k : \mathbb{R}^{2(m_k + m_{c_k} + m_{f_k})} \rightarrow [0, 1]^6. \quad (13)$$

Let $\hat{g}_k = (\hat{g}_k^-, \hat{g}_k^+)$ denote the two prediction outputs for the negative and positive classes, respectively. During training, the relationship between $\hat{g}_k^-(\mathbf{x}_i^{(k)})$ and $\hat{g}_k^+(\mathbf{x}_i^{(k)})$ can define $\hat{y}_{ik}^{(k)}$ given by

$$\hat{y}_{ik}^{(k)} = \begin{cases} -1, & \text{if } \hat{g}_k^-(\mathbf{x}_i^{(k)}) > \hat{g}_k^+(\mathbf{x}_i^{(k)}); \\ +1, & \text{otherwise.} \end{cases} \quad (14)$$

The predicted output of the example shown in Fig. 2 is $\hat{y}_i^{(k)} = [+1, -1, -1]$.

Also, it is necessary to use the inverse process of Eq. (14) when training the model. When the target label is $+1$ or -1 , pair value $(\hat{g}_k^-, \hat{g}_k^+)$ should be set to $(0, 1)$ or $(1, 0)$. The value of six target outputs of the example shown in Fig. 2 are $[g_k^-, g_k^+, g_{c_k}^-, g_{c_k}^+, g_{f_k}^-, g_{f_k}^+] = [0, 1, 1, 0, 1, 0]$.

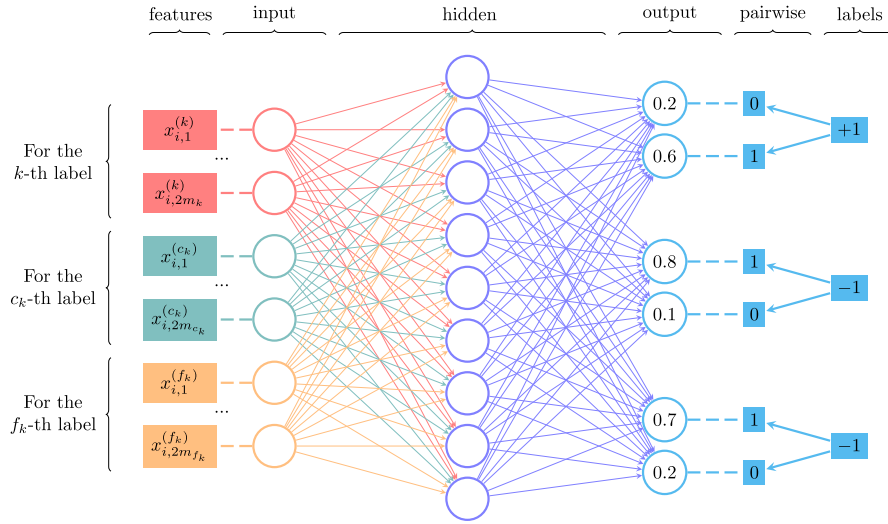


Fig. 2. A paired output network with a width of $2(m_k + m_{c_k} + m_{f_k})$ at the input and six output nodes applied to each label. The number of nodes in the hidden layer does not represent the actual scheme used and biases are omitted for brevity.

So it is possible to convert the F_k mapping to a new mapping G_k that better describes a particular k -th label output

$$G_k : \mathbb{R}^{2(m_k + m_{c_k} + m_{f_k})} \rightarrow [0, 1]^2, \text{ and} \quad (15)$$

$$F_k = [G_k, G_{c_k}, G_{f_k}].$$

As shown in the example of Fig. 2, the F_k of this network Θ_k is $[G_k, G_{c_k}, G_{f_k}] = [0.2, 0.6, 0.8, 0.1, 0.7, 0.2]$.

MSE (Mean-Square Error), a loss function commonly used for regression problems, is a simple and efficient strategy for calculating loss value in many network settings. Briefly, MSE computes the error between the predicted value and the target value. Moreover, it uses the magnitude of the error as a source of penalty during back propagation of the network. This scheme is exactly the loss function used by the network in LSTC, i.e.,

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n \left(g_k(y_i^{(k)}) - \hat{g}_k(x_i^{(k)}) \right)^2, \quad (16)$$

where $g_k(y_i^{(k)})$ can be defined by the target labels vector $y_i^{(k)}$ of the network, while $\hat{g}_k(x_i^{(k)})$ can be defined by the input node vector $x_i^{(k)}$.

When the loss function is determined, it is important to calculate the activation function backwards, layer by layer. In fact, LSTC employs more particular activation functions and layer settings for various data set properties. The configuration file of the source code contains this setting.

3.3. Prediction

The prediction approach used in Eq. (14) may be sufficient in training and computing accuracy. But it should be noted that because of general sparse for the label matrix, the idea of using Eq. (14) in actual prediction may be insufficient. Therefore, a new strategy will be adopted here to predict and evaluate the algorithm.

As shown in Fig. 3, we use softmax function to calculate the probability that label is positive

$$p_{ik} = \frac{\exp g_k^+(x_i^{(k)})}{\exp g_k^+(x_i^{(k)}) + \exp g_k^-(x_i^{(k)})}. \quad (17)$$

Consequently, we have a parametric prediction

$$\hat{y}_{ik}^{(k)} = \begin{cases} -1, & \text{if } p_{ik} < \theta; \\ +1, & \text{otherwise,} \end{cases} \quad (18)$$

where the threshold θ with domain $[0, 1]$ can be varied according to the requirement of the application. In fact, for ranking-based measures, we do not need to learn this parameter. However, if we need to predict logic labels, this parameter is very important.

Assuming a semi-supervised learning scenario, a model was previously built successfully from a set of n multi-label training examples \mathbf{X} and \mathbf{Y} . Now, if there is a new testing matrix \mathbf{U} with labels to be predicted, LSTC should be able to predict the label matrix $\hat{\mathbf{Y}}$ of the \mathbf{U} .

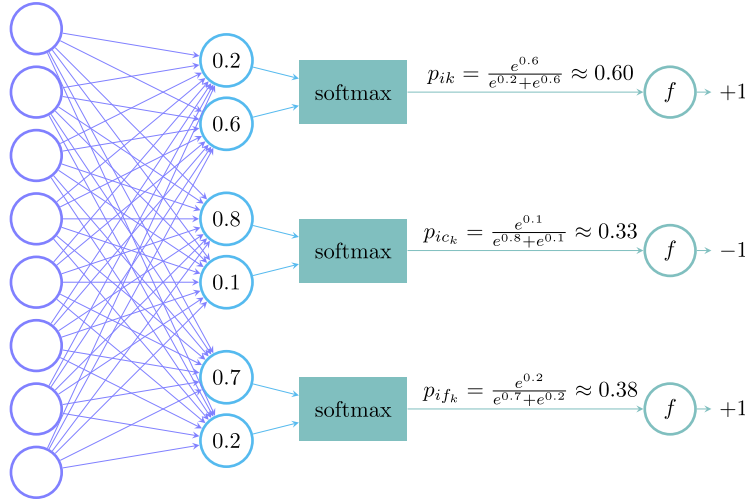


Fig. 3. The prediction process of a paired output network. Predicting from six nodes to three determined label classes requires calculating the probability that the label is positive through softmax function in Eq. (17), after then, outputs the label prediction through the threshold θ of the probability in Eq. (18) (θ in the figure is assumed to be 0.35).

If q well-established models are obtained, for any instance u_i in U , we can use $P_c^{(k)}$ and $N_c^{(k)}$ defined during training to compute the $U^{(k)}$ according to Eqs. (9) and (10). At the same time, it is natural to get the $U^{(c_k)}$ and $U^{(f_k)}$ based on two other auxiliary labels and their representative instances.

Hence $u_i^{(k)}$ is used as the input of the network Θ_k to obtain six real-valued output. Ultimately, the output of the model is parsed into the $\hat{y}_i^{(k)}$ by Eqs. (17) and (18).

Also, a note about the difference between training and prediction: we do not need to know the predictions of the three labels during the prediction process but only the results of the first label k . In fact, the training process only uses the auxiliary label c_k and f_k to assist in label correlation training. Consequently, when defining the model, this model already incorporates the correlation of auxiliary labels.

We can use a straightforward functional h that only concentrate on the relationship between the input and result to describe prediction process:

$$\hat{y}_i^{(k)} = h(u_i, P_c^j, N_c^j, \Theta_k), j \in \{k, f_k, c_k\}. \quad (19)$$

After training, only the first parameter in Eq. (19) is unknown. Hence given the determined independent variable u_i , the only dependent variable $\hat{y}_i^{(k)}$ must be determined. Simultaneously, if we apply Eq. (19) to each label and instance, we can get \hat{Y} of testing data matrix U .

4. Experiments

We conducted experiments to analyze the effectiveness of the LSTC algorithm by answering the following questions.

- 1) Is the time complexity analysis of Algorithm 1 valid?
- 2) Are the feature conversion and label prediction modules of LSTC better than the alternatives?
- 3) Does LSTC outperform state-of-the-art algorithms?
- 4) What kind of data is LSTC not suitable for?

The experiments were conducted on a computer using Pytorch and NVIDIA GTX 1660. The relevant parameters settings of different comparison algorithm were taken from the respective publications. In addition, some parameter settings were partially adjusted to ensure the good performance.

4.1. Datasets and evaluation measures

Table 3 presents datasets covering various domains such as music, medicine, image, biology, etc. They can be downloaded for free from Mulan,¹ the research group KDIS² at the University of Cordoba, and PALM³ at the Southeast University of China. Twelve

¹ <http://mulan.sourceforge.net/datasets-mlc.html>.

² <https://www.uco.es/kdis/mlresources/>.

³ <http://palm.seu.edu.cn/zhangml/files/Image.rar>.

Table 3

The description for each dataset. Cardinality stands for the average number of positive labels for each instance. Density stands for the proportion of non-zero elements in data matrix. GpositiveGO and WaterQuality are abbreviated as Gpositive and Water, respectively.

Dataset	#instance	#label	#feature	Cardinality	Density	Domain
Birds ¹	645	19	260	1.01	0.5850	Audio
Cal500 ¹	502	174	68	26.04	0.9960	Music
CHD49 ²	555	6	49	2.58	0.3947	Medicine
Emotion ¹	593	6	72	1.87	0.9893	Music
Flags ¹	194	7	19	3.39	0.4034	Image
Foodtruck ²	407	12	21	2.29	0.8378	Recommend
Image ³	2,000	5	294	1.24	0.9984	Image
Gpositive ²	519	4	912	1.24	0.0091	Biology
Scene ¹	2407	6	294	1.07	0.9885	Image
VirusGo ²	207	6	749	1.22	0.0152	Biology
Water ²	1,060	14	16	5.07	0.9455	Chemistry
Yeast ¹	2,417	14	103	4.24	0.9992	Biology
Art ¹	5,000	26	462	1.64	0.0754	Text
Business ¹	5,000	30	438	1.59	0.0775	Text
Enron ¹	1,702	53	1,001	3.38	0.0840	Text
Recreation ¹	5,000	22	606	1.42	0.0727	Text
Social ¹	5,000	39	1047	1.28	0.0613	Text

of them are from non-text domains and five are from the text domain. The reason for this division will be discussed experimentally.

MLL usually faces label sparsity, which corresponds to class-imbalance [24] for single-label learning. Therefore, ranking-based measures may be more suitable for this task. Then we will use six of them, including Peak F_1 -score, NDCG, One-Error, Coverage and Ranking Loss. Meanwhile, for the sake of objectivity, we also use label-based and sample-based measures including micro-AUC, macro-AUC and Hamming Loss [7].

4.1.1. AUC

AUC (Area Under Curve) [20] is an evaluation measure obtained by calculating the area of the ROC (Receiver Operating Characteristic) curve [35]. Meanwhile, ROC curve is obtained based on the change process of TPR and FPR of confusion matrix [29]. AUC is measured on the whole, so if we care about the global performance of prediction, AUC can be more representative than F_1 -score [6]. In short, if the cost of our false positives is high, such as recommended system, we will choose AUC as the more appropriate measure. On the contrary, if the cost of missing the correct target is high, such as medical-domain, then it is better to use F_1 -score. AUC is usually divided into micro-AUC and macro-AUC due to differences in calculation details. The micro-AUC considers the entire label matrix as a whole, i.e.,

$$\text{micro-AUC} = \frac{|\{(y', y'') | \hat{y}' \geq \hat{y}''(y', y'') \in \mathcal{P} \times \mathcal{N}\}|}{|\mathcal{P}| |\mathcal{N}|}, \quad (20)$$

where \mathcal{P} and \mathcal{N} stand for the set of positive and negative label values in the label matrix \mathbf{Y} , respectively. This label pair (y', y'') is an element in the Cartesian product result of \mathcal{P}_k and \mathcal{N}_k . (\hat{y}', \hat{y}'') is the result of the label pair (y', y'') prediction at the original position. For macro-AUC, similar operations are limited to each label vector. In the end, we average the AUC values of each label, i.e.,

$$\text{macro-AUC} = \frac{1}{q} \sum_{k=1}^q \text{AUC}_k = \sum_{k=1}^q \frac{1}{q} \frac{|\{(y', y'') | \hat{y}' \geq \hat{y}''(y', y'') \in \mathcal{P}_k \times \mathcal{N}_k\}|}{|\mathcal{P}_k| |\mathcal{N}_k|}. \quad (21)$$

For an instance \mathbf{x}_k without positive labels, we give $\text{AUC}_k = 0$. Conversely, if this instance \mathbf{x}_k has no negative labels, given $\text{AUC}_k = 1$.

4.1.2. Peak F_1 -score

Peak F_1 -score is based on the F_1 -score [6], which in turn is an evaluation measure based on precision and recall of confusion matrix [29], i.e.,

$$F_1 = \frac{2PR}{P + R}, \quad (22)$$

where P represents precision and R represents recall. However, in order to calculate precision and recall, a threshold is needed to convert numerical predictions to binary classification. Similar to ROC curve, Peak F_1 -score [9,22] avoids this difficulty as follows. Labels are sorted according to their predicted probability in descendant order. By predicting the top- k labels as positive and the remaining ones as negative, we obtain respective F_1 -score. By enumerating k from 0 to n_q , the maximal value will be the Peak F_1 -score. In this way, Peak F_1 -score is threshold free.

NDCG [14] is commonly used in the evaluation of search and recommendation system. Also, it is applicable to evaluation of multi-label classification. First, DCG (Discounted Cumulative Gain) is a measure that emphasizes the influence of overall order on the prediction result. It can be expressed as the relationship between the gain g_i and the logarithm of the ordinal number, i.e.,

$$DCG = \sum_{i=1}^{nq} \frac{g_i}{\log_2(i+1)} \quad (23)$$

where g_i stands for the sorted label according to their predicted probability in descendant order. Second, IDCG can be considered as the DCG in the best sorted state. Therefore, we just sort the label according to the original label matrix \mathbf{Y} . Assuming that there are k positive labels in the \mathbf{Y} matrix, then the sorted array is all positive labels in domain $[1, k]$ and all negative labels in domain $[k+1, nq]$. IDCG can be calculated according to the idea of Eq. (23) and this array. Finally, we can describe the NDCG (Normalized Discounted Cumulative Gain) as normalized DCG, i.e.,

$$NDCG = \frac{DCG}{IDCG}. \quad (24)$$

The reason why we consider NDCG as the evaluation measure of this paper is to consider the following advantages of NDCG.

- 1) NDCG can judge whether a multi-label algorithm reasonably predicts those labels with a significant probability of being positive.
- 2) NDCG always increase regardless of whether individual labels are correctly predicted or not. So the penalty for wrong prediction is not very serious.
- 3) DCG uses logarithms to mitigate the loss, therefore compared with other ranking-based evaluation measures, the change of NDCG is relatively stable.

4.1.3. Hamming Loss

Hamming Loss can be considered as an inverse problem of accuracy. It evaluates the proportion of labels for which the predicted label matrix does not match the true label matrix, i.e.,

$$\text{Hamming Loss} = \sum_{1 < i < n, 1 < j < q} \frac{\lambda(h(\hat{y}_{ij}) \neq y_{ij})}{nq}, \quad (25)$$

where $h(\cdot)$ denotes

$$h(\hat{y}_{ij}) = \begin{cases} -1, & \text{if } \hat{y}_{ij} < 0.5; \\ +1, & \text{otherwise.} \end{cases} \quad (26)$$

For any predicate π , $\lambda(\pi)$ returns 1 if π holds, and 0 otherwise. Therefore, the lower the Hamming Loss is, the closer the predicted label is to the true label, then the more prominent the model prediction effect is.

4.1.4. One-Error

The One-Error evaluation measure reflects the proportion of instances with errors to the total number of instances. These erroneous instances give the label with the highest predicted probability the wrong prediction, i.e.,

$$\text{One-Error} = \frac{1}{n} \sum_{i=1}^n \lambda(y_{ik} = +1), k = \underset{1 < j < q}{\operatorname{argmax}} \hat{y}_{ij}. \quad (27)$$

Therefore, the lower the One-Error, the better the ability to classify the highest probability label of the model. Also, we need to note that One-Error does not count instances without positive labels.

4.1.5. Coverage

Coverage discusses how many steps the ranked label list of an instance traverses to cover all positive labels in the ground truth vector, i.e.,

$$\text{Coverage} = \sum_{i=1}^n \text{Coverage}_i = \frac{1}{n} \sum_{i=1}^n \max \{ \text{rank}(\hat{\mathbf{y}}_i, j) \mid j \in \mathcal{V}_i^+ \} - 1. \quad (28)$$

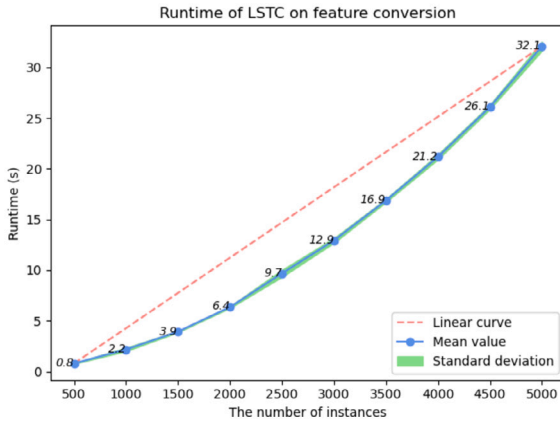
Any number j , as long as it satisfies $1 < j < q$ and belongs to the set \mathcal{V}_i^+ at the same time, then $y_{ij} = +1$ must exist. First, in the function $\text{rank}(\hat{\mathbf{y}}_i, j)$, the vector $\hat{\mathbf{y}}_i$ (the predicted label array of \mathbf{x}_i) will be sorted and the rank of each \hat{y}_{ij} will be obtained. Second, we select those ranks that satisfy the \mathcal{V}_i^+ constraints (the minimum rank is 1). In the end, we use the largest of the selected ranks as the Coverage_i . It should be noted that when $\mathcal{V}_i^+ = \emptyset$, the Coverage_i is 0.

4.1.6. Ranking Loss

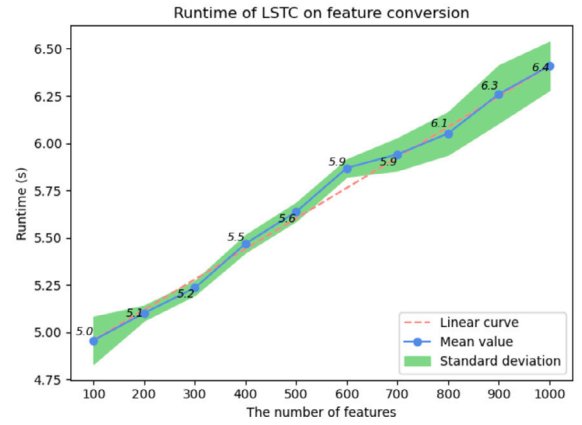
Ranking loss reflects the penalty average fraction of the model for misestimating the rank. Furthermore, an irrelevant label of an instance is ranked higher than its relevant one, i.e.,

$$\text{Ranking Loss} = \sum_{i=1}^n \text{Ranking Loss}_i = \sum_{i=1}^n \frac{1}{n} \frac{|(y', y'') \mid \hat{y}' \geq \hat{y}'', (y', y'') \in \mathcal{N}_i \times \mathcal{P}_i|}{|\mathcal{N}_i| |\mathcal{P}_i|}. \quad (29)$$

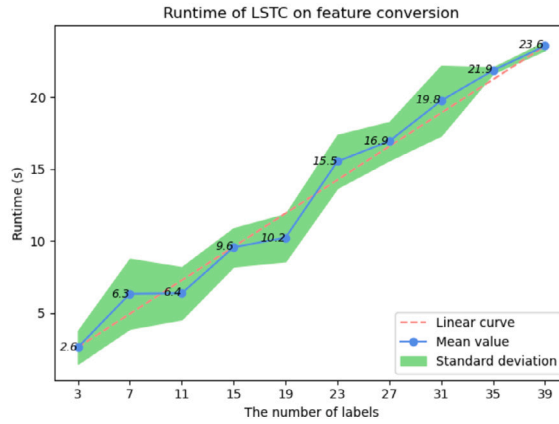
In fact, Ranking Loss can be seen as an inverse problem of macro-AUC. The actual difference between them is as follows.



(a) Runtime for different number of instances with $d = 500$ and $q = 10$.



(b) Runtime for different number of features with $n = 2000$ and $q = 10$.



(c) Runtime for different number of labels with $n = 2000$ and $d = 500$.

Fig. 4. The runtime change of feature conversion with the increase of the number of instances, features, and labels.

- 1) The order of the Cartesian product of sets is reversed.
- 2) Ranking Loss is the average evaluation of each instance, while macro-AUC is each label.

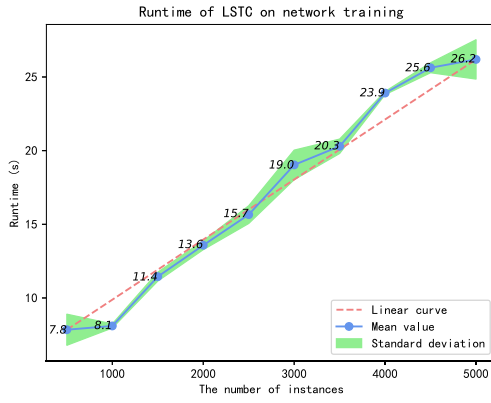
Thus the smaller the Ranking Loss, the less likely the model is to rank the labels incorrectly. Also we need to note that Ranking Loss_{*i*} is not discussed when any of $|\mathcal{N}_i|$ and $|\mathcal{P}_i|$ is 0.

4.2. Runtime analysis

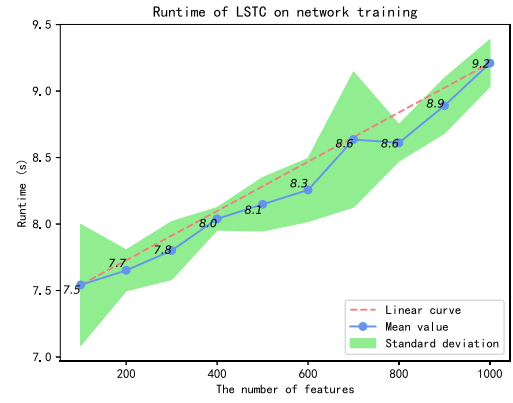
We would like to check the validity of the time complexity analysis of data conversion, as indicated by Table 2. Fig. 4 illustrates the runtime variation of the Algorithm 1 on different samples of the Social dataset. Specifically, Fig. 4(a) shows quadratic growth, while Figs. 4(b) and 4(c) show approximately linear growth. These trends of these curves are in line with the description of time complexity $O((d + q)n^2)$.

Fig. 5 presents the runtime curve of network training. Among them, the curves of Figs. 5(a) and 5(c) are more stable than curves of 5(b). Simultaneously, both of the more stable curves have larger increments than the Fig. 5(b). In fact, this is due to the deliberate design of Algorithm 1. Data conversion makes the number of features actually participating in network training have no obvious relationship with the original number. Therefore, the change in Fig. 5(b) may only be caused by the vector distance calculation with small changes and poor stability.

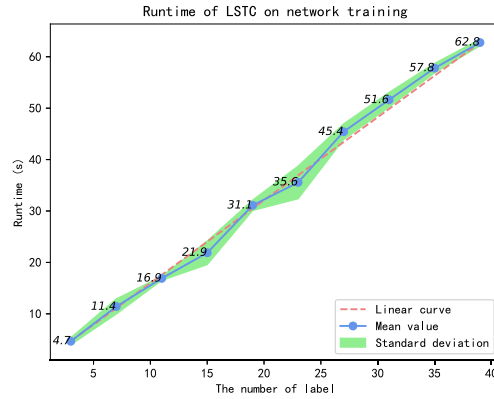
Obviously, the network training is not quite time-consuming. For example, on the biggest dataset Social we experienced, the training time is only 102 seconds, which is quite acceptable in applications.



(a) Runtime for different number of instances with $d = 100$ and $q = 5$.



(b) Runtime for different number of features with $n = 500$ and $q = 5$.



(c) Runtime for different number of labels with $n = 500$ and $d = 100$.

Fig. 5. The runtime change of network training with the increase of the number of instances, features, and labels.

4.3. Ablation studies

We undertake ablation studies to show the validity of LSTC in comparison with some alternatives.

Fig. 6 shows the ablation study between LSTC and LSTC with 2NN. The purpose of this ablation study is to explore whether the selection strategy of auxiliary labels of LSTC is effective.

The selection strategy of LSTC is to select two similar and dissimilar labels based on Eqs. (11) and (12) to achieve label correlation and alleviate the issue of getting in local optima. On the contrary, the common approach k NN only implements the former but not the latter. Hence, we will choose 2NN (2-nearest neighbors) as ablation comparison to demonstrate that mitigating local optima is effective.

Fig. 7 shows the ablation study between LSTC and LSTC with random representative. In feature conversion, we select the most representative m_k positive/negative instances by density peak as conversion reference. In this ablation study, in order to verify the effectiveness of the density peak is effective, we will use random representative instances selection as the ablation comparison of the density peak strategy. Specifically, after obtaining \mathbf{P} and \mathbf{N} of the data matrix \mathbf{X} , m_k representative instances are randomly selected to form $\mathbf{P}_c^{(k)}$ and $\mathbf{N}_c^{(k)}$.

Fig. 8 shows the ablation study between third-order label and first-order label. The purpose of this ablation study is to explore whether it is meaningful to consider label correlations.

LSTC deals with third-order label strategy through two auxiliary labels and special neural network with $2(m_k + m_{c_k} + m_{f_k})$ input nodes and six output nodes. Here, if we need to use the first-order label strategy, we can make the following changes in the model training. On the one hand, we can train the model with only a single label. On the other hand, we can continue to adopt the paired output prediction network, but the input and output ports will be set to single port.

From the results in Figs. 6, 7 and 8, we observe that the ablation study that consider label correlation has a more significant boosts than the other two strategies. For example, if the first order label strategy is adopted, the mean rank of micro-AUC and NDCG

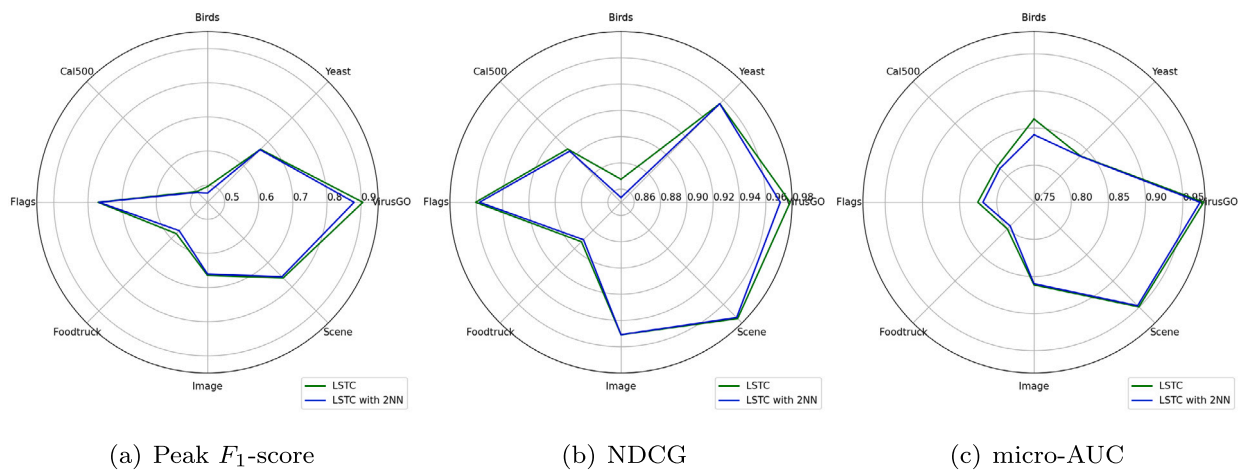


Fig. 6. Ablation study on the auxiliary label selection strategy. For brevity, only ten datasets are selected for comparison.

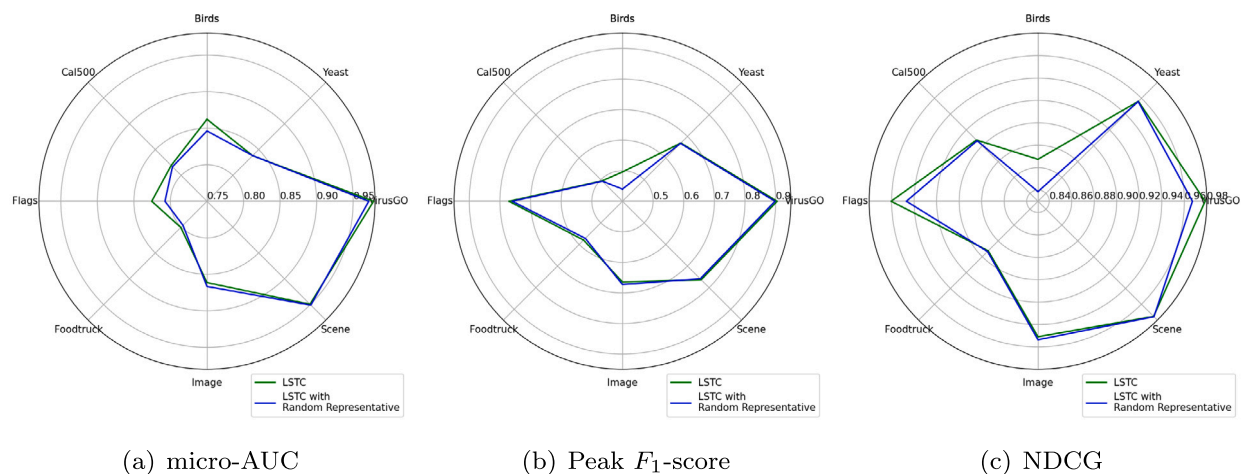


Fig. 7. Ablation study on representative instances selection in feature conversion.

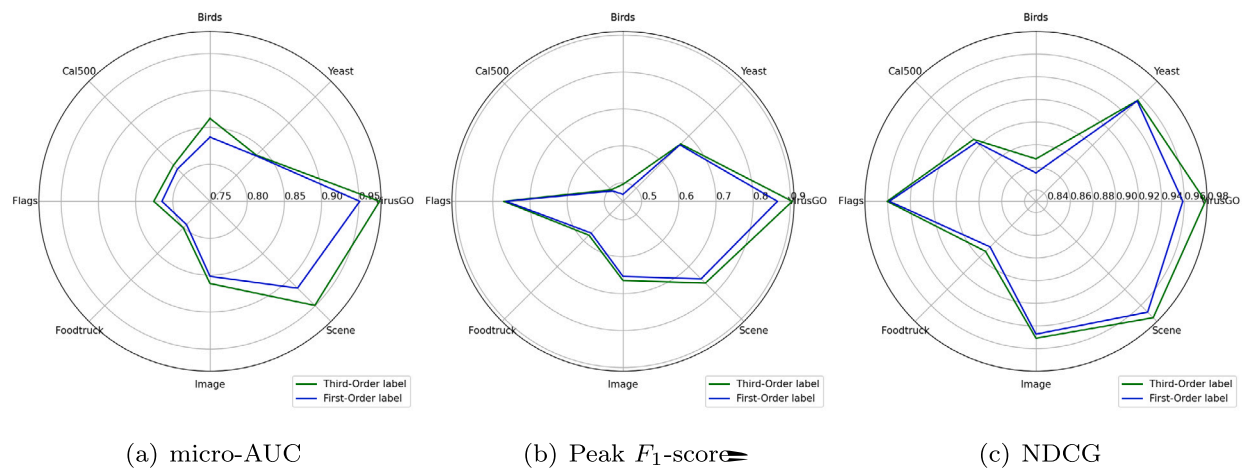


Fig. 8. Ablation study on label correlations.

will increase from 1.50 to 3.75 and 1.58 to 3.58. Therefore, we can consider that the algorithm is significantly and reliably improved after considering the label correlation.

Although the changes in Figs. 6 and 7 are slight compared to Fig. 8, such increment is also meaningful. For example, for NDCG, a small improvement may lead to a large change in the mean rank. So if the random representative instances selection strategy is adopted, the mean rank of NDCG will still increase from 1.58 to 3.08.

4.4. Effectiveness comparison

The brief introduction and parameter settings of the comparison algorithms are as follows.

- 1) **BP-MLL**⁴ [46] is the first multi-label neural network. It proposes a loss function for multi-label and considers label correlation exploitation. The settings are as follows: The size of the hidden layer is 20% of the input layer. epochs = 100, $\alpha = 0.05$.
- 2) **ML-KNN**⁵ [47] is a multi-label lazy learning approach based on k NN and maximum a posteriori (MAP) principle. The settings are as follows: $k = 15$ and smooth term is set to 1.
- 3) **IMLLA** [44] is an improved multi-label lazy learning approach through extending label correlation and modifying nearest neighbor selection scheme based on ML-KNN. The settings are as follows: $k = 15$ and use Minkowski distance with order 1.
- 4) **LIFT**⁶ [45] is multi-label embedding classification algorithms, also is the first algorithm to propose label-specific features through k -Means clustering. The settings are as follows: $r = 0.1$ and the SVM adopts linear scheme.
- 5) **LLSF**⁷ [13] is a second-order MLL algorithm combining shared features of two strongly correlated labels with an improved binary classifier. The settings are as follows: $\alpha = 2^{-5}$, $\beta = 2^{-3}$, $\gamma = 0.1$.
- 6) **GLOCAL**⁸ [48] has the ability to handle missing labels through learning a low-rank model with global and local label correlations. The settings are as follows: We cancel the setting of missing labels. $g = 3$ (the number of clusters), $k = 20$ (the latent matrix dimension), the regularization parameters $\lambda_1, \lambda_2, \lambda_3$ and λ_4 are set to 1, 10^{-2} , 0.125 and 0.125, respectively.
- 7) **LSML**⁹ [10] is augmented from the incomplete label matrix by learning a high-order label correlation to obtain label-specific features for multi-label classification with missing labels. The settings are as follows: $\lambda_1 = 10^2$, $\lambda_2 = 10^{-5}$, $\lambda_3 = 10^{-3}$, $\lambda_4 = 10^{-5}$.
- 8) **PML-NI**¹⁰ [42] has the advantages of multi-label processing for noisy environments by introducing a noise matrix. The settings are as follows: We minimize noise here to ensure fair comparisons. $\lambda = 100$, $\beta = 0.5$, $\gamma = 0.5$.
- 9) **MASP**¹¹ [22] is a multi-label active learning through serial-parallel neural network, it can address exploitation of label correlation, missing labels and selection of query labels. The settings are as follows: We unset the missing labels and perform semi-supervised learning. The serial and parallel parts adopt the Sigmoid activation function, while the layer between them adopt the Relu activation function.
- 10) **LRMML**¹² [15] handles missing labels and discriminatory label correlation by using an auxiliary label matrix and low-rank label subspace constraints. The settings are as follows: We set the missing rate to 0. Meantime, we set $\lambda_T = 10^{-1}$, $\lambda_R = 2$, and $\lambda_L = 10^{-1}$, which is consistent with the optimal setting of the released code.

Table 4 lists the comparison of LSTC with different MLL algorithms under the four evaluation measures micro-AUC, macro-AUC [20], Peak F_1 -score [9] and NDCG [14]. The results are obtained through 5-fold cross-validation. The higher the value, the better the effect of the algorithm.

Table 5 lists the comparison of LSTC with different MLL algorithms under the four evaluation measures Hamming Loss, One-Error, Coverage and Ranking Loss. The results are obtained through 5-fold cross-validation. The lower the value, the better the algorithm.

From Tables 4 and 5 we observe the following:

- 1) Overall, it can be found that the performance of LSTC generally outperforms the counterparts.
- 2) LSTC generally performs well in micro-AUC, macro-AUC, Peak F_1 -score, NDCG and Ranking Loss measures. But relatively speaking, it performs slightly worse on the Hamming Loss, One-Error, and Coverage measures.
- 3) There are six datasets with outstanding performance in more than five measures. Among them, CHD49 performs well on seven evaluation measures.
- 4) Under certain evaluation measures, LSTC performs very well on some datasets but mediocly on others. For example, on macro-AUC, Peak F_1 -score, and NDCG, the performance of the Birds dataset of LSTC is 1.26, 1.75, and 2.57 percentage points higher than the second place, respectively. Relatively, the One-Error and Coverage measures are at a disadvantage on the Yeast dataset. Although performs very well on this dataset.

⁴ http://www.lamda.nju.edu.cn/code_BPMLL.ashx.

⁵ <http://www.lamda.nju.edu.cn/files/MLkNN.rar>.

⁶ <http://palm.seu.edu.cn/zhangml/Resources.htm#ijcai11>.

⁷ <https://github.com/jiunhwang/LLSF>.

⁸ http://www.lamda.nju.edu.cn/code_GLOCAL.ashx.

⁹ <http://www.escience.cn/people/huangjun/index.html>.

¹⁰ <http://www.xiemk.pro/code/PMLNIcode.zip>.

¹¹ <https://gitee.com/fansmale/masp>.

¹² <https://github.com/sanjaysau/lrmml-2>.

Table 4

The performance of LSTC on twelve datasets under the four evaluation measures of micro-AUC, macro-AUC, Peak F_1 -score, and NDCG. **Bold** indicates the best performance, and *italics* indicates the second best.

Dataset	micro-AUC										
	BP-MLL'06	ML-KNN'07	IMLLA'12	LIFT'14	LLSF'15	GLOCAL'18	LSML'19	PML-NI'21	MASP'22	LRMML'22	LSTC
Birds	.5581 _{±.0966}	.7698 _{±.0357}	.4367 _{±.0235}	.7011 _{±.0288}	.7860 _{±.0247}	.5462 _{±.0328}	.8204 _{±.0141}	.8537 _{±.0223}	.7592 _{±.0208}	.8225 _{±.0156}	.8622 _{±.0131}
Cal500	.8009 _{±.0110}	.8149 _{±.0057}	.6793 _{±.0086}	.8143 _{±.0039}	.7716 _{±.0114}	.8053 _{±.0079}	.8142 _{±.0051}	.8009 _{±.0078}	.8001 _{±.0243}	.7837 _{±.0275}	.8193 _{±.0057}
CHD49	.6193 _{±.1197}	.7730 _{±.0195}	.7721 _{±.0189}	.7882 _{±.0175}	.7687 _{±.0147}	.7725 _{±.0118}	.7879 _{±.0202}	.7812 _{±.0154}	.6598 _{±.0243}	.7798 _{±.0243}	.7886 _{±.0207}
Emotion	.5750 _{±.0514}	.8419 _{±.0326}	.8535 _{±.0109}	.8683 _{±.0068}	.8433 _{±.0195}	.8464 _{±.0080}	.8573 _{±.0048}	.8562 _{±.0116}	.8341 _{±.0146}	.8555 _{±.0166}	.8697 _{±.0063}
Flags	.5397 _{±.1036}	.7236 _{±.0402}	.7339 _{±.0095}	.7402 _{±.0304}	.7712 _{±.0430}	.726 _{±.03590}	.7737 _{±.0175}	.7983 _{±.0285}	.7872 _{±.0421}	.7757 _{±.0308}	.8256 _{±.0093}
Foodtruck	.6897 _{±.0456}	.7913 _{±.0305}	.7868 _{±.0189}	.7917 _{±.0241}	.7980 _{±.0269}	.7907 _{±.0216}	.8001 _{±.0241}	.7923 _{±.0229}	.7486 _{±.0189}	.8000 _{±.0301}	.8007 _{±.0211}
Image	.5152 _{±.0240}	.8395 _{±.0090}	.8418 _{±.0084}	.8593 _{±.0072}	.8321 _{±.0076}	.8279 _{±.0164}	.8270 _{±.0088}	.8308 _{±.0113}	.8267 _{±.0148}	.8277 _{±.0113}	.8613 _{±.0146}
GpositiveGO	.7724 _{±.0911}	.9879 _{±.0042}	.9888 _{±.0058}	.9896 _{±.0074}	.9911 _{±.0058}	.9853 _{±.0090}	.9908 _{±.0029}	.9901 _{±.0069}	.8327 _{±.0158}	.9913 _{±.0046}	.9934 _{±.0039}
Scene	.5264 _{±.0171}	.9434 _{±.0047}	.9435 _{±.0079}	.9504 _{±.0053}	.9139 _{±.0038}	.9188 _{±.0080}	.9323 _{±.0068}	.9151 _{±.0109}	.9394 _{±.0033}	.9327 _{±.0041}	.9493 _{±.0061}
VirusGO	.8050 _{±.0571}	.9492 _{±.0192}	.9403 _{±.0199}	.9666 _{±.0270}	.9593 _{±.0197}	.9331 _{±.0279}	.9846 _{±.0063}	.9796 _{±.0103}	.9082 _{±.0248}	.9851 _{±.0081}	.9780 _{±.0177}
WaterQuality	.5736 _{±.0492}	.7415 _{±.0102}	.7507 _{±.0085}	.7484 _{±.0128}	.7291 _{±.0073}	.7284 _{±.0102}	.7166 _{±.0124}	.7253 _{±.0084}	.6709 _{±.0104}	.7149 _{±.0080}	.7529 _{±.0086}
Yeast	.7417 _{±.0461}	.8416 _{±.0113}	.8406 _{±.0083}	.8351 _{±.0071}	.8303 _{±.0018}	.8304 _{±.0114}	.8345 _{±.0050}	.8312 _{±.0046}	.8339 _{±.0085}	.8351 _{±.0061}	.8383 _{±.0074}
Mean rank	10.50	5.67	6.50	4.17	6.75	7.83	4.75	5.17	8.33	5.00	1.50
macro-AUC											
Birds	.4807 _{±.0354}	.6050 _{±.0417}	.3776 _{±.0700}	.5528 _{±.0842}	.6790 _{±.0354}	.4766 _{±.0552}	.7184 _{±.0256}	.7932 _{±.0641}	.7264 _{±.0538}	.7749 _{±.0355}	.8058 _{±.0228}
Cal500	.4807 _{±.0202}	.4944 _{±.0238}	.4891 _{±.0103}	.5165 _{±.0141}	.5213 _{±.0184}	.5415 _{±.0252}	.5385 _{±.0160}	.5201 _{±.0147}	.5280 _{±.0244}	.5316 _{±.0100}	.5457 _{±.0098}
CHD49	.5236 _{±.0567}	.5139 _{±.0248}	.5283 _{±.0417}	.6034 _{±.0430}	.5719 _{±.0369}	.5798 _{±.0183}	.5906 _{±.0135}	.5742 _{±.0307}	.5561 _{±.0313}	.5842 _{±.0306}	.6126 _{±.0306}
Emotion	.5861 _{±.0698}	.8257 _{±.0212}	.8336 _{±.0192}	.8487 _{±.0195}	.8372 _{±.0151}	.8360 _{±.0140}	.8456 _{±.0180}	.8333 _{±.0134}	.8321 _{±.0162}	.8406 _{±.0176}	.8615 _{±.0168}
Flags	.4844 _{±.0164}	.5113 _{±.0550}	.5360 _{±.0539}	.5093 _{±.0629}	.6191 _{±.0400}	.5669 _{±.0757}	.6442 _{±.0144}	.7207 _{±.0538}	.7270 _{±.0383}	.6467 _{±.0405}	.7355 _{±.0186}
Foodtruck	.5163 _{±.0431}	.5356 _{±.0225}	.5685 _{±.0390}	.5593 _{±.0352}	.5593 _{±.0274}	.5791 _{±.0395}	.5963 _{±.0160}	.5837 _{±.0214}	.5581 _{±.0342}	.5808 _{±.0295}	.5921 _{±.0345}
Image	.5216 _{±.0660}	.8316 _{±.0117}	.8395 _{±.0060}	.8565 _{±.0066}	.8261 _{±.0191}	.8243 _{±.0111}	.8208 _{±.0143}	.8315 _{±.0049}	.8323 _{±.0150}	.8213 _{±.0111}	.8567 _{±.0124}
GpositiveGO	.6993 _{±.0848}	.9775 _{±.0206}	.9779 _{±.0259}	.9881 _{±.0069}	.9727 _{±.0241}	.9745 _{±.0185}	.9895 _{±.0057}	.9911 _{±.0025}	.9471 _{±.0288}	.9857 _{±.0102}	.9918 _{±.0043}
Scene	.6160 _{±.0710}	.9385 _{±.0072}	.9374 _{±.0031}	.9477 _{±.0017}	.9113 _{±.0100}	.9130 _{±.0105}	.9284 _{±.0079}	.9110 _{±.0061}	.9342 _{±.0090}	.9276 _{±.0047}	.9436 _{±.0033}
VirusGO	.6780 _{±.1193}	.8649 _{±.1014}	.9128 _{±.0917}	.9381 _{±.0795}	.8710 _{±.0708}	.9105 _{±.0691}	.9497 _{±.0644}	.9195 _{±.0951}	.8822 _{±.0512}	.9796 _{±.0165}	.9412 _{±.0665}
WaterQuality	.4946 _{±.0483}	.7022 _{±.0077}	.7156 _{±.0125}	.7158 _{±.0103}	.6918 _{±.0148}	.6914 _{±.0118}	.6843 _{±.0062}	.6894 _{±.0125}	.7207 _{±.0078}	.6833 _{±.0129}	.7229 _{±.0098}
Yeast	.5308 _{±.0204}	.6879 _{±.0118}	.6835 _{±.0102}	.6820 _{±.0098}	.6872 _{±.0137}	.6915 _{±.0091}	.6956 _{±.0122}	.6818 _{±.0173}	.7036 _{±.0180}	.6970 _{±.0115}	.6826 _{±.0151}
Mean rank	10.75	7.58	6.75	5.08	7.00	6.42	4.42	5.83	5.58	4.75	1.92
Peak F ₁ -score											
Birds	.1420 _{±.0294}	.2886 _{±.0326}	.1038 _{±.0024}	.2312 _{±.0269}	.3118 _{±.0194}	.1971 _{±.0078}	.3140 _{±.0180}	.4785 _{±.0295}	.2312 _{±.0007}	.3195 _{±.0354}	.4960 _{±.0414}
Cal500	.4641 _{±.0083}	.4810 _{±.0116}	.3844 _{±.0075}	.4859 _{±.0030}	.4626 _{±.0109}	.4856 _{±.0100}	.4907 _{±.0107}	.4846 _{±.0062}	.4661 _{±.0261}	.4666 _{±.0247}	.4943 _{±.0059}
CHD49	.6357 _{±.0605}	.7046 _{±.0110}	.7012 _{±.0265}	.7078 _{±.0111}	.7024 _{±.0302}	.7074 _{±.0191}	.7087 _{±.0163}	.7056 _{±.0119}	.6688 _{±.0192}	.7097 _{±.0241}	.7157 _{±.0239}
Emotion	.5106 _{±.0500}	.6882 _{±.0368}	.7047 _{±.0118}	.7103 _{±.0065}	.6884 _{±.0294}	.6896 _{±.0143}	.7029 _{±.0266}	.7004 _{±.0107}	.6885 _{±.0221}	.7003 _{±.0185}	.7162 _{±.0168}
Flags	.6585 _{±.0178}	.7107 _{±.0272}	.7122 _{±.0292}	.7202 _{±.0210}	.7372 _{±.0307}	.6967 _{±.0290}	.7343 _{±.0091}	.7514 _{±.0368}	.7447 _{±.0400}	.7331 _{±.0249}	.7708 _{±.0118}
Foodtruck	.4528 _{±.0240}	.5572 _{±.0370}	.5623 _{±.0260}	.5618 _{±.0415}	.5673 _{±.0352}	.5580 _{±.0302}	.5660 _{±.0197}	.5608 _{±.0359}	.4958 _{±.0185}	.5698 _{±.0267}	.5796 _{±.0214}
Image	.4080 _{±.0045}	.6305 _{±.0127}	.6327 _{±.0129}	.6581 _{±.0143}	.6319 _{±.0076}	.6173 _{±.0224}	.6156 _{±.0135}	.6282 _{±.0167}	.6364 _{±.0251}	.6150 _{±.0152}	.6642 _{±.0125}
GpositiveGO	.6924 _{±.0730}	.9321 _{±.0203}	.9351 _{±.0131}	.9387 _{±.0165}	.9366 _{±.0255}	.9432 _{±.0179}	.9469 _{±.0101}	.9537 _{±.0163}	.6367 _{±.0204}	.9506 _{±.0066}	.9457 _{±.0165}
Scene	.3190 _{±.0066}	.7523 _{±.0177}	.7533 _{±.0136}	.7816 _{±.0154}	.6891 _{±.0112}	.6953 _{±.0197}	.7240 _{±.0113}	.6919 _{±.0176}	.7487 _{±.0086}	.7245 _{±.0116}	.7633 _{±.0209}
VirusGO	.6285 _{±.0483}	.8212 _{±.0378}	.8483 _{±.0368}	.8743 _{±.0461}	.8668 _{±.0402}	.8249 _{±.0659}	.9035 _{±.0128}	.9139 _{±.0366}	.8021 _{±.0348}	.9165 _{±.0237}	.9051 _{±.0498}
WaterQuality	.5349 _{±.0097}	.6162 _{±.0106}	.6219 _{±.0124}	.6196 _{±.0225}	.6075 _{±.0048}	.6055 _{±.0139}	.5964 _{±.0205}	.6039 _{±.0103}	.3884 _{±.0073}	.5964 _{±.0059}	.6241 _{±.0137}
Yeast	.5844 _{±.0259}	.6714 _{±.0115}	.6715 _{±.0135}	.6734 _{±.0095}	.6593 _{±.0048}	.6646 _{±.0103}	.6628 _{±.0134}	.6628 _{±.0073}	.6625 _{±.0108}	.6643 _{±.0093}	.6695 _{±.0102}
Mean rank	10.58	6.83	6.08	4.00	6.83	6.92	5.17	5.08	8.00	5.00	1.75
NDCG											
Birds	.5422 _{±.0239}	.7231 _{±.0139}	.5696 _{±.0504}	.6501 _{±.0472}	.7122 _{±.0126}	.6676 _{±.0192}	.6000 _{±.4899}	.8418 _{±.0334}	.8037 _{±.0255}	.7331 _{±.0367}	.8675 _{±.0169}
Cal500	.8884 _{±.0054}	.9018 _{±.0021}	.8521 _{±.0036}	.9011 _{±.0086}	.8836 _{±.0043}	.8992 _{±.0045}	.9044 _{±.0011}	.8981 _{±.0051}	.8932 _{±.0139}	.8916 _{±.0101}	.9073 _{±.0024}
CHD49	.8998 _{±.0133}	.9171 _{±.0174}	.9093 _{±.0160}	.9273 _{±.0236}	.9207 _{±.0101}	.9194 _{±.0211}	.9187 _{±.0119}	.9261 _{±.0086}	.7007 _{±.0334}	.9203 _{±.0157}	.9245 _{±.0177}
Emotion	.7887 _{±.0681}	.9363 _{±.0068}	.9444 _{±.0084}	.9478 _{±.0063}	.9239 _{±.0123}	.9338 _{±.0158}	.9331 _{±.0052}	.9365 _{±.0081}	.9302 _{±.0085}	.9342 _{±.0123}	.9445 _{±.0151}
Flags	.8480 _{±.0491}	.9186 _{±.0261}	.9189 _{±.0083}	.9027 _{±.0240}	.9348 _{±.0270}	.9053 _{±.0175}	.9396 _{±.0106}	.9552 _{±.0079}	.9524 _{±.0110}	.9362 _{±.0227}	.9605 _{±.0070}
Foodtruck	.7965 _{±.0515}	.8703 _{±.0164}	.8708 _{±.0282}	.8909 _{±.0110}	.8837 _{±.0263}	.8852 _{±.0160}	.8853 _{±.0191}	.8896 _{±.0147}	.8401 _{±.0231}	.8862 _{±.0191}	.8925 _{±.0193}
Image	.7717 _{±.0169}	.9439 _{±.0030}	.9449 _{±.0049}	.9526 _{±.0064}	.9307 _{±.0056}	.9281 _{±.0052}	.9283 _{±.0061}	.9266 _{±.0067}	.9351 _{±.0076}	.9297 _{±.0067}	.9509 _{±.0057}
GpositiveGO	.7332 _{±.0700}	.9941 _{±.0040}	.9960 _{±.0015}	.9965 _{±.0026}	.9940 _{±.0068}	.9899 _{±.0092}	.9959 _{±.0018}	.9897 _{±.0167}	.9303 _{±.0079}	.9956 _{±.0036}	.9974 _{±.0015}
Scene	.7423 _{±.0143}	.9719 _{±.0047}	.9717 _{±.0032}	.9781 _{±.0008}	.9531 _{±.0018}						

Table 5

The performance of LSTC on twelve datasets under the four evaluation measures of Hamming Loss, One-Error, Coverage, and Ranking Loss. **Bold** indicates the best performance, and *italics* indicates the second best.

Dataset	Hamming Loss										
	BP-MLL'06	ML-KNN'07	IMLLA'12	LIFT'14	LLSF'15	GLOCAL'18	LSML'19	PML-NI'21	MASP'22	LRMML'22	LSTC
Birds	.0650 _{±.0205}	.0529 _{±.0062}	.1213 _{±.0174}	.0539 _{±.0059}	.0535 _{±.0057}	.0979 _{±.0098}	.0533 _{±.0046}	.0551 _{±.0052}	.0550 _{±.0044}	.0534 _{±.0036}	.0480_{±.0054}
Cal500	.1540 _{±.0078}	.1382 _{±.0017}	.1787 _{±.0028}	.1368 _{±.0034}	.1435 _{±.0022}	.1388 _{±.0013}	.1373 _{±.0023}	.5683 _{±.3862}	.1734 _{±.0162}	.1413 _{±.0026}	.1365_{±.0018}
CHD49	.3578 _{±.0431}	.3103 _{±.0276}	.3053 _{±.0214}	.2841 _{±.0103}	.3583 _{±.0164}	.2974 _{±.0188}	.2847 _{±.0125}	.3171 _{±.0121}	.3339 _{±.0080}	.2877 _{±.0225}	.2823_{±.0181}
Emotion	.3506 _{±.0352}	.2097 _{±.0102}	.2007 _{±.0146}	.1902_{±.0148}	.2628 _{±.0050}	.2002 _{±.0115}	.1995 _{±.0144}	.2010 _{±.0082}	.2114 _{±.0097}	.1984 _{±.0142}	<i>.1917_{±.0105}</i>
Flags	.3416 _{±.0264}	.3404 _{±.0505}	.3327 _{±.0495}	.3460 _{±.0130}	.3821 _{±.0219}	.3049 _{±.0535}	.2959 _{±.0264}	.2810 _{±.0194}	.2783 _{±.0280}	.2969 _{±.0281}	.2775_{±.0141}
Foodtruck	.2029 _{±.0650}	.1563 _{±.0194}	.7521 _{±.0157}	.1572 _{±.0180}	.1554 _{±.0190}	.7286 _{±.0300}	.1627 _{±.0156}	.1780 _{±.0192}	.1994 _{±.0068}	<i>.1533_{±.0160}</i>	.1532_{±.0105}
Image	.3327 _{±.0898}	.1754 _{±.0107}	.1756 _{±.0135}	.1594 _{±.0099}	.2264 _{±.0053}	.1841 _{±.0054}	.1989 _{±.0031}	.1825 _{±.0060}	.1754 _{±.0088}	.2001 _{±.0063}	.1590_{±.0071}
GpositiveGO	.2039 _{±.0296}	.0377 _{±.0116}	.0381 _{±.0097}	.0347 _{±.0150}	.0525 _{±.0071}	.0315 _{±.0110}	.0318 _{±.0055}	.0271_{±.0082}	.0617 _{±.0045}	<i>.0313_{±.0076}</i>	.0323 _{±.0073}
Scene	.2244 _{±.0608}	.0870 _{±.0062}	.0879 _{±.0060}	.0815_{±.0069}	.1511 _{±.0028}	.1111 _{±.0045}	.1048 _{±.0042}	.1077 _{±.0077}	.0906 _{±.0058}	.1072 _{±.0054}	<i>.0848_{±.0056}</i>
VirusGO	.2223 _{±.0546}	.0739 _{±.0258}	.0652 _{±.0314}	.0549 _{±.0181}	.0888 _{±.0209}	.0820 _{±.0260}	.0578 _{±.0239}	.0420_{±.0148}	.0854 _{±.0178}	<i>.0498_{±.0129}</i>	.0572 _{±.0151}
WaterQuality	.4025 _{±.0368}	.2893 _{±.0079}	.6112 _{±.0092}	.2813_{±.0193}	.3057 _{±.0049}	.6288 _{±.0116}	.3846 _{±.0164}	.3572 _{±.0385}	.2832 _{±.0091}	.3119 _{±.0074}	<i>.2842_{±.0107}</i>
Yeast	.2790 _{±.0387}	.1974 _{±.0046}	.1972 _{±.0039}	.1928_{±.0063}	.2443 _{±.0058}	.2000 _{±.0048}	.2004 _{±.0023}	.2009 _{±.0084}	.2067 _{±.0073}	.1999 _{±.0035}	.1984 _{±.0058}
Mean rank	10.00	5.08	7.17	3.33	8.42	7.08	5.08	6.25	6.92	4.58	2.17
Dataset	One-Error										
	BP-MLL'06	ML-KNN'07	IMLLA'12	LIFT'14	LLSF'15	GLOCAL'18	LSML'19	PML-NI'21	MASP'22	LRMML'22	LSTC
Birds	.8990 _{±.1070}	.7324 _{±.0444}	.7312 _{±.0604}	.7386 _{±.0611}	.6570 _{±.0729}	.6919 _{±.0359}	.6538 _{±.0776}	.6514 _{±.0556}	.2853 _{±.0570}	.6793 _{±.0630}	.2589_{±.0324}
Cal500	.2835 _{±.2082}	.1162_{±.0357}	.4074 _{±.0202}	.1509 _{±.0589}	.2312 _{±.0425}	.1526 _{±.0572}	<i>.1194_{±.0427}</i>	.1622 _{±.0232}	.1594 _{±.0508}	.1989 _{±.0671}	.1334 _{±.0235}
CHD49	.3306 _{±.1530}	.2417 _{±.0682}	.2440 _{±.0613}	.2426 _{±.0338}	.2825 _{±.0436}	.2702 _{±.0565}	.2553 _{±.0307}	.2707 _{±.0812}	.3369 _{±.0591}	.2699 _{±.0463}	.2360_{±.0396}
Emotion	.6533 _{±.1208}	.3019 _{±.0110}	.2959 _{±.0598}	.2383 _{±.0598}	.2751 _{±.0507}	.2571 _{±.0314}	.2496 _{±.0159}	.2817 _{±.0446}	.2817 _{±.0245}	.2681 _{±.0282}	.2344_{±.0404}
Flags	.3548 _{±.2555}	.2708 _{±.0644}	.2050 _{±.0869}	.2398 _{±.0570}	.2618 _{±.0972}	.2243 _{±.0791}	.2495 _{±.0589}	.2224 _{±.0751}	.2883 _{±.0643}	.2342 _{±.0349}	.2015_{±.0469}
Foodtruck	.5428 _{±.2794}	.2853 _{±.0223}	.2977 _{±.0597}	.2901 _{±.0647}	.2848 _{±.0387}	.2875 _{±.0560}	.2754 _{±.0490}	.2945 _{±.0494}	.4054 _{±.0383}	<i>.2746_{±.0466}</i>	.2675_{±.0611}
Image	.7531 _{±.0324}	.3275 _{±.0431}	.3264 _{±.0226}	.2732_{±.0261}	.3190 _{±.0145}	.3309 _{±.0212}	.3415 _{±.0136}	.3268 _{±.0151}	.3190 _{±.0214}	.3415 _{±.0146}	.2855 _{±.0180}
GpositiveGO	.5114 _{±.1449}	.0715 _{±.0369}	.0751 _{±.0244}	.0636 _{±.0224}	.0674 _{±.0227}	.0637 _{±.0282}	<i>.0580_{±.0138}</i>	.0637 _{±.0113}	.1021 _{±.0324}	.0636 _{±.0179}	.0539_{±.0167}
Scene	.8196 _{±.0256}	.2272 _{±.0113}	.2314 _{±.0300}	.1925_{±.0303}	.2078 _{±.0125}	.2636 _{±.0152}	.2301 _{±.0272}	.2656 _{±.0189}	.2264 _{±.0211}	.2326 _{±.0192}	<i>.2123_{±.0056}</i>
VirusGO	.6904 _{±.2085}	.1468 _{±.0328}	.1345 _{±.0684}	.1159 _{±.0201}	.1011 _{±.0402}	.0961 _{±.0566}	.0722_{±.0296}	.0947 _{±.0273}	.2175 _{±.0412}	.0821 _{±.0320}	.0966 _{±.0487}
WaterQuality	.5447 _{±.1713}	.2781 _{±.0114}	.2825 _{±.0234}	.2685_{±.0331}	.3235 _{±.0284}	.3118 _{±.0444}	.3303 _{±.0211}	.3146 _{±.0263}	.2745 _{±.0150}	.3273 _{±.0254}	<i>.2689_{±.0358}</i>
Yeast	.5625 _{±.3048}	.2345 _{±.0061}	.2214_{±.0148}	.2223 _{±.0205}	.2279 _{±.0229}	.2227 _{±.0231}	.2274 _{±.0143}	.2272 _{±.0156}	.2685 _{±.0264}	.2292 _{±.0076}	.2271 _{±.0291}
Mean rank	10.83	6.42	6.42	4.00	7.00	5.75	4.83	5.75	7.25	6.08	2.00
Dataset	Coverage										
	BP-MLL'06	ML-KNN'07	IMLLA'12	LIFT'14	LLSF'15	GLOCAL'18	LSML'19	PML-NI'21	MASP'22	LRMML'22	LSTC
Birds	5.5472 _{±.3871}	4.0621 _{±.3885}	4.5363 _{±.7044}	4.1560 _{±.8223}	3.7023 _{±.3407}	7.4132 _{±.4880}	3.8527 _{±.2839}	2.5980_{±.1689}	3.8729 _{±.5765}	3.7535 _{±.4129}	<i>2.8124_{±.2288}</i>
Cal500	137.1812 _{±.5428}	130.5605_{±.0816}	162.9398 _{±.6650}	133.3209 _{±.1576}	148.0021 _{±.3322}	133.7434 _{±.4446}	134.2643 _{±.3957}	142.3020 _{±.0200}	141.2599 _{±.4585}	144.6494 _{±.6256}	<i>133.1583_{±.6702}</i>
CHD49	3.7588 _{±.7227}	2.8464 _{±.0761}	2.8359 _{±.0902}	2.6829 _{±.0834}	2.7405 _{±.0832}	2.7461 _{±.1416}	2.7243 _{±.1212}	2.7565 _{±.1633}	2.9928 _{±.0946}	2.7045 _{±.0959}	2.6649_{±.0733}
Emotion	2.7393 _{±.6446}	1.8163 _{±.1202}	1.7552 _{±.0844}	<i>1.6963_{±.1824}</i>	1.8504 _{±.0766}	1.8065 _{±.1495}	1.7791 _{±.1175}	1.8045 _{±.1297}	1.8026 _{±.0321}	1.7783 _{±.1521}	1.6951_{±.0436}
Flags	4.1351 _{±.4898}	4.0209 _{±.1942}	3.9025 _{±.4274}	3.9080 _{±.1149}	3.8919 _{±.0783}	3.8804 _{±.3755}	3.8800 _{±.2228}	3.9002 _{±.1856}	3.9163 _{±.2753}	3.8888 _{±.2817}	3.8308_{±.1429}
Foodtruck	4.1526 _{±.7362}	3.6950 _{±.3479}	3.8176 _{±.1905}	3.6620 _{±.7342}	3.6525 _{±.2674}	3.7588 _{±.6372}	<i>3.5320_{±.2847}</i>	3.6805 _{±.4746}	4.5958 _{±.2690}	3.5306_{±.412}	3.6456 _{±.4342}
Image	2.1799 _{±.0642}	0.9686 _{±.0745}	0.9675 _{±.0771}	<i>0.8704_{±.0763}</i>	0.9585 _{±.0366}	0.9729 _{±.0488}	0.9975 _{±.0744}	0.9476 _{±.0215}	0.9880 _{±.0418}	0.9925 _{±.0550}	0.8640_{±.0757}
GpositiveGO	0.8537 _{±.3514}	0.1024 _{±.0442}	0.1176 _{±.0581}	0.0982 _{±.0392}	0.0924 _{±.0414}	0.0954 _{±.0438}	0.0847 _{±.0323}	0.0753_{±.0253}	0.1773 _{±.0130}	0.0829 _{±.0180}	<i>0.0825_{±.0211}</i>
Scene	2.5753 _{±.1478}	0.4601 _{±.0441}	0.4772 _{±.0244}	0.3989_{±.0403}	0.5929 _{±.0291}	0.5637 _{±.0639}	0.4836 _{±.0593}	0.5899 _{±.0450}	0.4873 _{±.0538}	0.4919 _{±.0379}	<i>0.4400_{±.0367}</i>
VirusGO	1.5607 _{±.6952}	0.5435 _{±.2693}	0.5094 _{±.2999}	0.4065 _{±.0850}	0.4451 _{±.1501}	0.3736 _{±.1241}	0.3491_{±.0748}	<i>0.3529_{±.1231}</i>	0.6573 _{±.1605}	0.3714 _{±.0886}	0.3823 _{±.1422}
WaterQuality	10.7404 _{±.3190}	8.7302 _{±.2293}	8.6854 _{±.2502}	8.7352 _{±.0747}	8.9632 _{±.3357}	8.9705 _{±.1453}	9.1009 _{±.0905}	8.9931 _{±.2503}	<i>8.6642_{±.1443}</i>	9.1179 _{±.1371}	8.6292_{±.2445}
Yeast	7.9741 _{±.5760}	<i>6.3198_{±.1826}</i>	6.2823_{±.0773}	6.3298 _{±.1620}	6.4693 _{±.0892}	6.3798 _{±.1238}	6.3508 _{±.2888}	6.4436 _{±.1366}	6.3847 _{±.2280}	6.3473 _{±.1218}	6.3244 _{±.0737}
Mean rank	10.42	6.25	6.42	4.42	6.58	6.58	4.83	5.58	7.75	5.17	2.00
Dataset	Ranking Loss										
	BP-MLL'06	ML-KNN'07	IMLLA'12	LIFT'14	LLSF'15	GLOCAL'18	LSML'19	PML-NI'21	MASP'22	LRMML'22	LSTC
Birds	.3952 _{±.0279}	.2985 _{±.0230}	.3655 _{±.0432}	.3038 _{±.0131}	.2609 _{±.0233}	.2855 _{±.0171}	.2664 _{±.0260}	.1608_{±.0211}	.3283 _{±.0303}	.2604 _{±.0262}	<i>.1876_{±.0305}</i>
Cal500	.1935 _{±.0075}	.1817 _{±.0091}	.3175 _{±.0084}	.1826 _{±.0032}	.2232 _{±.0116}	.1877 _{±.0012}	.1825 _{±.0045}	.2004 _{±.0054}	<i>.1811_{±.0057}</i>	.2142 _{±.0272}	.1792_{±.0034}
CHD49	.3123 _{±.1042}	.2226 _{±.0324}	.2213 _{±.0317}	.2009_{±.0083}	.2204 _{±.0156}	.2222 _{±.0330}	<i>.2006_{±.0229}</i>	.2171 _{±.0278}	.2818 _{±.0231}	.2075 _{±.0278}	.1987_{±.0174}
Emotion	.3807 _{±.1417}	.1717 _{±.0054}	.1624 _{±.0269}	<i>.1426_{±.0298}</i>	.1694 _{±.0190}	.1635 _{±.0231}	.1556 _{±.0099}	.1571 _{±.0153}	.1880 _{±.0230}	.1585 _{±.0205}	.1422_{±.0182}
Flags	.2577 _{±.0598}	.2454 _{±.0328}	.2419 _{±.0551}	.2375 _{±.0195}	.2298 _{±.0347}	.2324 _{±.0524}	<i>.2232_{±.0235}</i>	.2370 _{±.0336}	.2520 _{±.0092}	.2239 _{±.0349}	.2165_{±.0301}
Foodtruck	.2039 _{±.0619}	.1549 _{±.0216}	.1681 _{±.0091}	.1556 _{±.0197}	.1539 _{±.0220}	.1593 _{±.0167}	.1501_{±.0236}	.1533 _{±.0232}	.2249 _{±.0244}	<i>.1510_{±.0248}</i>	.1532 _{±.0213}
Image	.4888 _{±.0200}	.1744 _{±.0200}	.1753 _{±.0130}	<i>.1502_{±.0149}</i>	.1720 _{±.0093}	.1749 _{±.0123}	.1777 _{±.0070}	.1755 _{±.0149}	.1898 _{±.0163}	.1794 _{±.0096}	.1493_{±.0052}
GpositiveGO	.2826 _{±.1164}	.0314 _{±.0160}	.0360 _{±.0176}	.0302 _{±.0121}	<i>.0244_{±.0117}</i>	.0294 _{±.0136}	.0255 _{±.0047}	.0250 _{±.0056}	.0523 _{±.0141}	.0249 _{±.0071}	.0236

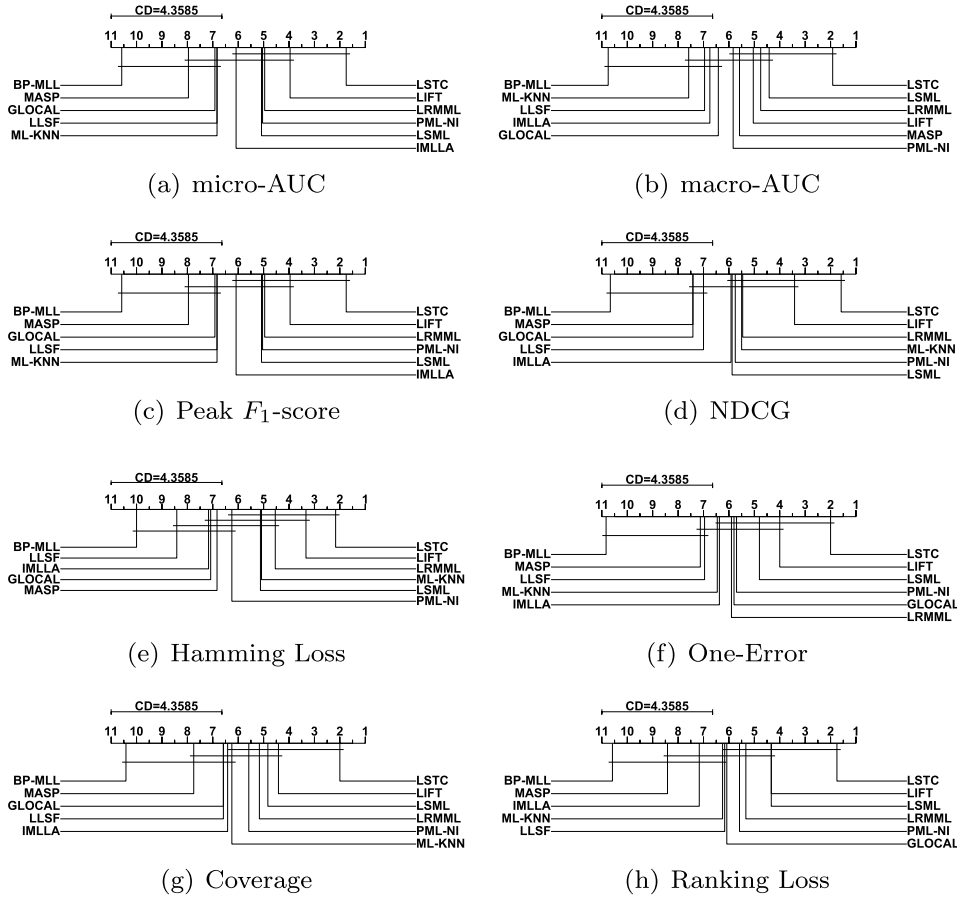


Fig. 9. Comparison of LSTC against ten comparing algorithms with the Bonferroni-Dunn test. CD = 3.4681 at significance level $\alpha = 0.05$.

7) LSTC ranks second across all evaluation measures on Scene datasets, while LIFT consistently ranks first. This reflects the superior performance of the feature embedding idea on this dataset. Meanwhile, it also reflects that the improvement based on LIFT in this paper is not sensitive to this dataset.

Fig. 9 illustrates the CD diagrams on each evaluation measure (# comparing algorithms $k = 11$, # dataset $N = 12$, $q_\alpha = 3.102$). We can observe that LSTC always maintains the best average rank. Meanwhile, in the Fig. 9(a), there is no evidence of a statistical difference among LSTC, LIFT, LRMMML, PML-NI, LSML and IMLLA. However, they outperform the other five algorithms. Similarly, this performance also exists in other evaluation measures, but the algorithms that perform better are different.

4.5. Limitation

Although LSTC performs well on non-text-domain datasets, but faced with text-domain datasets, LSTC may no longer have a significant advantage according to Table 6. Because it is only to understand the limitations of LSTC on text-domain datasets, Table 6 only compares three of the eight evaluation measures.

The reasons are as follows.

- 1) The data matrix of the datasets on text-domain are generally sparse, LSTC does not perform well on some sparse datasets such as VirusGO in Tables 4 and 5.
- 2) Meanwhile, the label correlation of datasets on text-domain is not obvious. LSTC lacks to deal with the misjudgment of some independent labels caused by over-consideration of label correlation [36].
- 3) Considering label correlation through neural network may be prone to overfitting and getting stuck local optima. From Table 6, it can be found that the BP-MLL and MASP algorithms using neural network generally do not perform well.

At the same time, from Table 6, we also found that LSML and LIFT have excellent results when dealing with datasets on text-domain.

Table 6

The performance of LSTC on text-domain datasets under the three evaluation measures of micro-AUC, Peak F_1 -score and NDCG. **Bold** indicates the best performance, and *italics* indicates the second best.

Dataset	micro-AUC										
	BP-MLL'06	ML-KNN'07	IMLLA'12	LIFT'14	LLSF'15	GLOCAL'18	LSML'19	PML-NI'21	MASP'22	LRMML'22	LSTC
Art	.8263 \pm .0062	.8271 \pm .0030	.8017 \pm .0023	.8620\pm.0034	.7798 \pm .0092	.7859 \pm .0071	.8254 \pm .0032	.8362 \pm .0058	.7941 \pm .0241	<i>.8401\pm.0065</i>	.8266 \pm .0046
Business	.9243 \pm .0095	<i>.9508\pm.0042</i>	.9284 \pm .0051	.9598\pm.0039	.9109 \pm .0055	.9118 \pm .0052	.9341 \pm .0056	.9438 \pm .0054	.8797 \pm .0186	.9456 \pm .0037	.9376 \pm .0075
Enron	.8884 \pm .0076	.8994 \pm .0024	.8781 \pm .0045	.9137\pm.0084	.8345 \pm .0080	.7955 \pm .0097	.9024 \pm .0079	.8356 \pm .0049	.8181 \pm .0101	<i>.9083\pm.0044</i>	.8907 \pm .0030
Recreation	.7795 \pm .0366	.7912 \pm .0069	.7813 \pm .0081	.8618\pm.0034	.7857 \pm .0079	.7927 \pm .0097	.8176 \pm .0032	.8351 \pm .0082	.7351 \pm .0111	<i>.8373\pm.0046</i>	.7922 \pm .0071
Social	.9052 \pm .0049	<i>.9367\pm.0016</i>	.8928 \pm .0088	.9437\pm.0035	.8725 \pm .0070	.8722 \pm .0050	.9051 \pm .0059	.9111 \pm .0040	.7608 \pm .0100	.9163 \pm .0039	.9065 \pm .0069
Mean rank	7.20	3.80	7.80	1.00	9.40	9.00	5.40	4.40	10.40	2.40	5.20
Dataset	Peak F_1 -score										
	BP-MLL'06	ML-KNN'07	IMLLA'12	LIFT'14	LLSF'15	GLOCAL'18	LSML'19	PML-NI'21	MASP'22	LRMML'22	LSTC
Art	.3394 \pm .0198	.3374 \pm .0075	.3490 \pm .0052	.4476 \pm .0183	.3948 \pm .0106	.3965 \pm .0145	<i>.4609\pm.0036</i>	.4600 \pm .0097	.3498 \pm .0292	.4657\pm.0111	.3603 \pm .0108
Business	.6353 \pm .0585	.7255 \pm .0084	.7231 \pm .0056	<i>.7452\pm.0128</i>	.7153 \pm .0096	.7257 \pm .0054	<i>.7459\pm.0085</i>	<i>.7461\pm.0099</i>	.7093 \pm .0114	.7469\pm.0050	.7129 \pm .0101
Enron	.5149 \pm .0215	.5350 \pm .0106	.5575 \pm .0082	.6122 \pm .0128	.5119 \pm .0161	.4774 \pm .0208	.6225\pm.0170	.5273 \pm .0150	.5276 \pm .0183	<i>.6196\pm.0049</i>	.5423 \pm .0101
Recreation	.2853 \pm .0399	.2929 \pm .0102	.3107 \pm .0089	.4766 \pm .0051	.3992 \pm .0118	.3990 \pm .0134	.4791\pm.0074	.4711 \pm .0177	.4020 \pm .0175	<i>.4777\pm.0083</i>	.2989 \pm .0096
Social	.4329 \pm .0173	.6212 \pm .0087	.6199 \pm .0198	.6508 \pm .0141	.6096 \pm .0211	.6195 \pm .0114	.6525\pm.0115	.6468 \pm .0055	.5690 \pm .0153	<i>.6513\pm.0053</i>	.5131 \pm .0110
Mean rank	10.40	7.60	6.80	3.40	7.60	7.00	1.60	4.20	7.80	1.60	8.00
Dataset	NDCG										
	BP-MLL'06	ML-KNN'07	IMLLA'12	LIFT'14	LLSF'15	GLOCAL'18	LSML'19	PML-NI'21	MASP'22	LRMML'22	LSTC
Art	.8099 \pm .0203	.8467 \pm .0047	.8458 \pm .0066	.8906 \pm .0041	.8622 \pm .0047	.8653 \pm .0026	.8912\pm.0027	.8893 \pm .0050	.8488 \pm .0079	<i>.8911\pm.0044</i>	.8508 \pm .0020
Business	.9244 \pm .0048	.9585 \pm .0026	.9559 \pm .0021	.9650\pm.0006	.9474 \pm .0035	.9498 \pm .0026	.9641 \pm .0013	.9641 \pm .0018	.9400 \pm .0046	<i>.9643\pm.0013</i>	.9560 \pm .0031
Enron	.8881 \pm .0135	.8992 \pm .0058	.9004 \pm .0069	.9126 \pm .0029	.8822 \pm .0050	.8654 \pm .0109	.9236\pm.0068	.8802 \pm .0044	.8815 \pm .0078	<i>.9235\pm.0027</i>	.9032 \pm .0015
Recreation	.7882 \pm .0351	.8324 \pm .0051	.8369 \pm .0044	.8971\pm.0037	.8651 \pm .0034	.8685 \pm .0056	.8941 \pm .0040	.8922 \pm .0050	.8520 \pm .0067	<i>.8946\pm.0036</i>	.8179 \pm .0061
Social	.8345 \pm .0083	.9346 \pm .0020	.9267 \pm .0029	.9415\pm.0042	.9121 \pm .0049	.9139 \pm .0034	<i>.9378\pm.0028</i>	.9348 \pm .0029	.8780 \pm .0048	.9373 \pm .0021	.8995 \pm .0093
Mean rank	10.20	6.80	7.20	1.80	7.40	7.20	2.20	5.20	8.80	2.20	7.20

LSML introduces a high-order label correlation control matrix C and manifold regularizer [2] to adjust the correlation, so that it can be stabilized within a reasonable range.

LIFT is the first algorithm that adopts the idea of label-specific features. Also, it is the source of the embedding idea of LSTC. But to reflect label correlation, we use neural network to train the model instead of LIFT using SVM classifier. As for why LIFT works well, on the one hand, the first-order multi-label algorithm obviously does not fall into overfitting due to label correlation. On the other hand, compared to neural network, SVM classifiers focus on global optima rather than local optima. However, the sparsity and weak label correlation of text-domain datasets are prone to getting stuck in local optima. Therefore, although LIFT is a first-order multi-label algorithm, it still outperforms most algorithms that consider label correlation on text-domain datasets.

At the same time, we realize that label correlation is unsuitable for all multi-label applications. In the field of personality psychology research, correlations between personality traits are low [23]. Almost all personality trait identification methods assume no dependencies among labels. Therefore, MLL algorithms that exploit label correlation like LSTC may not perform well on such datasets. These flaws may be a possible commonality among such MLL algorithms.

4.6. Discussion

We can now answer the questions proposed at the beginning of this section.

- 1) According to Fig. 4, the runtime of Algorithm 1 is in line with our time complexity analysis. In this way, the feature conversion is not the bottleneck of the LSTC algorithm. Meanwhile, Fig. 5 also verifies that the network training is not quite time-consuming.
- 2) According to Figs. 6, 8, and 7, all components of LSTC are appropriate.
- 3) According to Tables 4, 5 and Fig. 9, LSTC outperforms state-of-the-art algorithms on non-text-domain data. However, this does not indicate that LSTC is really superior than these algorithms. The reason is that LSTC only handles the basic MLL problem, while these algorithms may be able to handle missing data, missing labels, active learning, etc.
- 4) According to Table 6, LSTC is inappropriate for text-domain data. This is an important limitation.

5. Conclusions and further works

In this study, we proposed the LSTC algorithm with new feature conversion and label correlation learning techniques. The algorithm outperformed state-of-the-art algorithms on datasets from various domains other than text. Ablation studies also validated the superiority of our detailed techniques. Naturally, there are still some topics that require further investigation.

- 1) Adjustment for text-domain data. We have analyzed some possible reasons why LSTC is not suitable for this type of data. Consequently, we should adjust the algorithm, most likely the distance calculation strategy, to handle them more effectively.
- 2) Adaptation to missing labels, missing values, and active learning scenarios. This study is quite fundamental without considering these important issues in real world applications. Therefore, we will extend the algorithm to deal with these complicated situations.
- 3) Dealing with extreme multi-label learning (XML) [17]. XML is characterized by millions of instances, thousands to millions of features and labels, very high proportion (> 99.9%) of missing values, and very low (< 0.1%) label sparsity. It has attracted much research interests.¹³ Hence we are also planning to enter this challenging and promising field.

CRedit authorship contribution statement

Xing-Yi Zhang: Methodology, Investigation, Software, Writing - original draft, Experimentation. **Fan Min:** Supervision, Methodology, Writing - review & editing. **Guojie Song:** Supervision, Writing - review & editing. **Hong Yu:** Supervision, Methodology.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

Data will be made available on request.

Acknowledgements

This work is supported by Central Government Funds of Guiding Local Scientific and Technological Development (No. 2021ZYD0003), the National Natural Science Foundation of China (Nos. 62136002, 41674141), and Nanchong Municipal Government-Universities Scientific Cooperation Project (No. SXHZ045, 19XHZ0030).

References

- [1] R. Agrawal, R. Srikant, Fast algorithms for mining association rules, in: VLDB, 1994, pp. 487–499.
- [2] M. Belkin, P. Niyogi, V. Sindhwani, Manifold regularization: a geometric framework for learning from labeled and unlabeled examples, *J. Mach. Learn. Res.* 7 (2006) 2399–2434.
- [3] M.R. Boutell, J. Luo, X. Shen, C.M. Brown, Learning multi-label scene classification, *Pattern Recognit.* 37 (2004) 1757–1771.
- [4] W.-C. Chang, H.-F. Yu, K. Zhong, Y. Yang, I.S. Dhillon, Taming pretrained transformers for extreme multi-label text classification, in: SIGKDD, 2020, pp. 3163–3171.
- [5] Z. Cheng, Z. Zeng, Joint label-specific features and label correlation for multi-label learning with missing label, *Appl. Intell.* 50 (2020) 4029–4049.
- [6] N. Chinchor, MUC-4 evaluation metrics, in: MUC, 1992.
- [7] E. Gibaja, S. Ventura, A tutorial on multilabel learning, *ACM Comput. Surv.* 47 (2015) 1–38.
- [8] H. Haripriya, C. Prathibhamol, Y.R. Pai, M.S. Sandeep, A.M. Sankar, S.N. Veerla, P. Nedungadi, Multi label prediction using association rule generation and simple k-means, in: ICCTICT, 2016, pp. 159–163.
- [9] H. Huang, H. Xu, X. Wang, W. Silamu, Maximum f1-score discriminative training criterion for automatic mispronunciation detection, *IEEE/ACM Trans. Audio Speech Lang. Process.* 23 (2015) 787–797.
- [10] J. Huang, F. Qin, X. Zheng, Z.-K. Cheng, Z.-X. Yuan, W.-G. Zhang, Q.-M. Huang, Improving multi-label classification with missing labels by learning label-specific features, *Inf. Sci.* 492 (2019) 124–146.
- [11] S.-J. Huang, Z.-H. Zhou, Multi-label learning by exploiting label correlations locally, in: AAAI, 2012, pp. 949–955.
- [12] X.-Y. Jia, S.-S. Zhu, W.-W. Li, Joint label-specific features and correlation information for multi-label learning, *J. Comput. Sci. Technol.* 35 (2020) 247–258.
- [13] H. Jun, G. Li, Q. Huang, X. Wu, Learning label specific features for multi-label classification, in: ICDM, 2015, pp. 181–190.
- [14] J. Kekäläinen, K. Järvelin, Using graded relevance assessments in IR evaluation, *J. Am. Soc. Inf. Sci. Technol.* 53 (2002) 1120–1129.
- [15] S. Kumar, R. Rastogi, Low rank label subspace transformation for multi-label learning with missing labels, *Inf. Sci.* 596 (2022) 53–72.
- [16] Y. Li, Y. Song, J. Luo, Improving pairwise ranking for multi-label image classification, in: CVPR, 2017, pp. 3617–3625.
- [17] J. Liu, W.-C. Chang, Y. Wu, Y. Yang, Deep learning for extreme multi-label text classification, in: SIGIR, 2017, pp. 115–124.
- [18] Y. Liu, K. Wen, Q. Gao, X. Gao, F. Nie, SVM based multi-label learning with missing labels for image annotation, *Pattern Recognit.* 78 (2018) 307–317.
- [19] Z. Ma, S. Chen, Expand globally, shrink locally: discriminant multi-label learning with missing labels, *Pattern Recognit.* 111 (2021) 107675.
- [20] S.J. Mason, N.E. Graham, Areas beneath the relative operating characteristics (ROC) and relative operating levels (ROL) curves: statistical significance and interpretation, *Q. J. R. Meteorol. Soc.* 128 (2002) 2145–2166.
- [21] A.K. McCallum, Multi-label text classification with a mixture model trained by EM, in: AAAI Workshop on Text Learning, 1999.
- [22] X.-Y. Min, K. Qian, B.-W. Zhang, G. Song, F. Min, Multi-label active learning through serial-parallel neural networks, *Knowl.-Based Syst.* 251 (2022) 109226.
- [23] N.K. Mishra, A. Singh, P.K. Singh, Multi-label personality trait identification from text, *Multimed. Tools Appl.* 81 (2022) 1–17.
- [24] N.K. Mishra, P.K. Singh, Feature construction and smote-based imbalance handling for multi-label learning, *Inf. Sci.* 563 (2021) 342–357.
- [25] J.M. Moyano, E.L. Gibaja, K.J. Cios, S. Ventura, Review of ensembles of multi-label classifiers: models, experimental study and prospects, *Inf. Fusion* 44 (2018) 33–45.
- [26] J. Nam, J. Kim, E.L. Mencia, I. Gurevych, J. Fürnkranz, Large-scale multi-label text classification revisiting neural networks, in: ECML, 2014, pp. 437–452.

¹³ <http://manikvarma.org/downloads/XC/XMLRepository.html>.

- [27] G. Nasierding, G. Tsoumakas, A.Z. Kouzani, Clustering based multi-label classification for image annotation and retrieval, in: IEEE SMC, 2009, pp. 4514–4519.
- [28] T.T. Nguyen, M.T. Dang, A.V. Luong, A.W.-C. Liew, T. Liang, J. McCall, Multi-label classification via incremental clustering on an evolving data stream, *Pattern Recognit.* 95 (2019) 96–113.
- [29] F. Provost, R. Kohavi, Guest editors' introduction: on applied research in machine learning, *Mach. Learn.* 30 (1998) 127–132.
- [30] G.-J. Qi, X.-S. Hua, Y. Rui, J. Tang, T. Mei, H.-J. Zhang, Correlative multi-label video annotation, in: ACM MM, 2007, pp. 17–26.
- [31] K. Qian, X.-Y. Min, Y. Cheng, F. Min, Weight matrix sharing for multi-label learning, *Pattern Recognit.* 136 (2023) 109156.
- [32] W. Qian, J. Huang, Y. Wang, W. Shu, Mutual information-based label distribution feature selection for multi-label learning, *Knowl.-Based Syst.* 195 (2020) 105684.
- [33] A. Rodriguez, A. Laio, Clustering by fast search and find of density peaks, *Science* 344 (2014) 1492–1496.
- [34] S. Shu, F. Lv, Y. Yan, L. Li, S. He, J. He, Incorporating multiple cluster centers for multi-label learning, *Inf. Sci.* 590 (2022) 60–73.
- [35] K.A. Spackman, Signal detection theory: valuable tools for evaluating inductive learning, in: IWML, 1989, pp. 160–163.
- [36] F. Sun, J. Tang, H. Li, G.-J. Qi, T.S. Huang, Multi-label image categorization with sparse factor representation, *IEEE Trans. Image Process.* 23 (2014) 1028–1037.
- [37] L. Sun, T. Yin, W. Ding, Y. Qian, J. Xu, Feature selection with missing labels using multilabel fuzzy neighborhood rough sets and maximum relevance minimum redundancy, *IEEE Trans. Fuzzy Syst.* 30 (2021) 1197–1211.
- [38] A. Tan, J. Liang, W.-Z. Wu, J. Zhang, Semi-supervised partial multi-label classification via consistency learning, *Pattern Recognit.* 131 (2022) 108839.
- [39] G. Tsoumakas, I. Katakis, I. Vlahavas, Random k-labelsets for multilabel classification, *IEEE Trans. Knowl. Data Eng.* 23 (2010) 1079–1089.
- [40] M. Wang, F. Min, Z.-H. Zhang, Y.-X. Wu, Active learning through density clustering, *Expert Syst. Appl.* 85 (2017) 305–317.
- [41] R. Wang, S. Kwong, X. Wang, Y. Jia, Active k-labelsets ensemble for multi-label classification, *Pattern Recognit.* 109 (2021) 107583.
- [42] M.-K. Xie, S.-J. Huang, Partial multi-label learning with noisy label identification, *IEEE Trans. Pattern Anal. Mach. Intell.* 44 (2022) 1–12.
- [43] S. Xu, X. Yang, H. Yu, D.-J. Yu, J. Yang, E.C. Tsang, Multi-label learning with label-specific feature reduction, *Knowl.-Based Syst.* 104 (2016) 52–61.
- [44] M.-L. Zhang, An improved multi-label lazy learning approach, *J. Comput. Res. Dev.* 49 (2012) 2271–2282.
- [45] M.-L. Zhang, L. Wu, LIFT: multi-label learning with label-specific features, *IEEE Trans. Pattern Anal. Mach. Intell.* 37 (2014) 107–120.
- [46] M.-L. Zhang, Z.-H. Zhou, Multi-label neural networks with applications to functional genomics and text categorization, *IEEE Trans. Knowl. Data Eng.* 18 (2006) 1338–1351.
- [47] M.-L. Zhang, Z.-H. Zhou, ML-KNN: a lazy learning approach to multi-label learning, *Pattern Recognit.* 40 (2007) 2038–2048.
- [48] Y. Zhu, J.T. Kwok, Z.-H. Zhou, Multi-label learning with global and local label correlation, *IEEE Trans. Knowl. Data Eng.* 30 (2018) 1081–1094.