



# Synthetic Sample Generation for Label Distribution Learning

Manuel González<sup>a,\*</sup>, Julián Luengo<sup>a</sup>, José-Ramón Cano<sup>b</sup>, Salvador García<sup>a</sup>

<sup>a</sup> Department of Computer Science and Artificial Intelligence, University of Granada, Granada 18071, Spain

<sup>b</sup> Department of Computer Science, University of Jaén, EPS of Linares, Avenida de la Universidad S/N, Linares 23700, Jaén, Spain

## ARTICLE INFO

### Article history:

Received 2 November 2019

Received in revised form 23 July 2020

Accepted 25 July 2020

Available online 4 August 2020

### Keywords:

Label Distribution Learning

Data pre-processing

Oversampling

Machine learning

## ABSTRACT

Label Distribution Learning (LDL) is a general learning framework that assigns an instance to a distribution over a set of labels rather than a single label or multiple labels. Current LDL methods have proven their effectiveness in many machine learning applications. As of the first formulation of the LDL problem, numerous studies have been carried out that apply the LDL methodology to various real-life problem solving. Others have focused more specifically on the proposal of new algorithms. The purpose of this article is to start addressing the LDL problem as of the data pre-processing stage. The baseline hypothesis is that, due to the high dimensionality of existing LDL data sets, it is very likely that this data will be incomplete and/or that poor data quality will lead to poor performance once applied to the learning algorithms. In this paper, we propose an oversampling method, which creates a superset of the original dataset by creating new instances from existing ones. Then, we apply already existing algorithms to the pre-processed training set in order to validate the efficacy of our method. The effectiveness of the proposed SSG-LDL is verified on several LDL datasets, showing significant improvements to the state-of-the-art LDL methods.

© 2020 Elsevier Inc. All rights reserved.

## 1. Introduction

A supervised learning process essentially consists of assigning a label or set of labels to each of the samples. In existing learning paradigms, there are mainly two cases of label assignment: (1) a single label is assigned to an instance, and (2) multiple labels are assigned to each instance. Single-label learning (SLL) assumes that all instances in the training set are labeled according to the first case. Multi-Label Learning (MLL) [40,29] allows instances in the training set to be labeled using the second method. Thus, MLL can deal with ambiguity when an instance belongs to more than one class (label) [40].

However, in real-life problems, we can find cases for which MLL is not sufficient since the level of description of each label is not the same. The Label Distribution Learning (LDL) concept showed up for first time in 2013 [26] and was formally described in 2016 [23] in order to deal with label ambiguity in mapping when one instance is not necessarily mapped to one label. The aim of this paradigm is to answer the question “how much does each label describe the instance?” instead of “which label can describe the instance?” [23].

From the first formulation of the LDL problem, numerous studies have been carried out applying the LDL methodology to various real-life problem solving situations, e.g., sense of beauty recognition [34], facial age estimation [26], personality recognition in social media [44], image emotion classification [45], pre-release prediction of crowd opinion on movies

\* Corresponding author.

E-mail addresses: [manuel.gonzalez.es@gmail.com](mailto:manuel.gonzalez.es@gmail.com) (M. González), [julianlm@decsai.ugr.es](mailto:julianlm@decsai.ugr.es) (J. Luengo), [jrcano@ujaen.es](mailto:jrcano@ujaen.es) (J.-R. Cano), [salvagl@decsai.ugr.es](mailto:salvagl@decsai.ugr.es) (S. García).

[24], crowd counting in public video surveillance [49], ... Other studies have focused more particularly on the creation of new learners or the improvement of existing learners such as: AA-kNN, SA-IIS, SA-BFGS [23], LDL forests [36], Logistic boosting regression for LDL [43], Deep Label Distribution [20], LDL-SCL [50], Structured random forest for LDL [14].

The purpose of this article is to address the LDL problem from the data pre-processing stage [21,22]. Studying the existing LDL datasets, we observe that the number of samples they contain is significantly low as opposed to the high dimensionality of input characteristics as well as output labels. This leads us to the hypothesis that the datasets with sparse data will lead to poor performance once applied to the learning algorithms.

The pursued objective is to improve the performance of the underlying LDL learners without modifying them, by adding a preliminary pre-processing stage applied to the original dataset. It is an extra layer that can be added in a simple way to the machine learning process.

Over the last few years, many solutions have been proposed to deal with this problem. One of the most well-known and commonly used is the data sampling technique [4] in which the training instances are modified in such a way as to produce a more efficient class distribution that allow learners to perform in a similar manner to standard learning.

In the specialized literature, we can find several studies that show the effect of resampling the training set to treat incomplete datasets. Those studies have empirically proved that, applying a pre-processing step is usually a useful solution [12,35]. Furthermore, the main advantage of these techniques is that they are independent of the underlying learner.

Nowadays, researchers have proposed a few attempts to apply pre-processing in LDL, as the feature selection algorithm described in Ref. [42], but the oversampling approach has not yet been studied or applied to the LDL paradigm and that is why, in this paper, we have developed an oversampling method [33], which creates a superset of the original dataset by creating new instances from existing ones. The technique devised is based on one of the most renowned approaches in the oversampling area called: "Synthetic Minority Oversampling Technique" (SMOTE) [11]. Briefly, its main objective is to create new minority class examples to oversample the training set by interpolating several minority class instances that lie together.

The problem statement can thus be described as building a new optimized dataset, an extension of the original, which improves the performance of the underlying learners.

Synthetic Sample Generation for LDL, or from now on SSG-LDL, is our proposal for a pre-processing algorithm adapted to LDL restrictions. We will start from the basic structure of SMOTE, which we will extend to be able to deal with LDL datasets. In our case we have to take a step away from SMOTE because we do not want to deal with minority classes, rather we want to create synthetic samples of the most spatially distant points. To this end, we will also introduce a new technique to select remote points, inspired by the probabilistic selection used in genetic algorithms [27]. Another main challenge is how to generate synthetic samples from the ones available in the training set since, unlike SMOTE, we do not only have to deal with the dimension of the characteristics but we must also handle the dimension of the labels to get a synthetic sample that is compatible with LDL. It is noteworthy to mention that SSG-LDL is not a simple data augmentation such as those commonly used in deep learning to preprocess and to enlarge the set of training images, such as those used in Ref. [20].

In order to evaluate the proposal proficiency, we will compare three state-of-the-art LDL learners applying our SSG-LDL pre-processing step to the training set and measuring six aspects of their performance. We will repeat the experiment over 15 real-world datasets and validate the results of the empirical comparisons using Wilcoxon and Bayesian Sign tests.

The rest of the paper is organized as follows. First, a brief review and discussion of the related work on LDL and synthetic feature generation are given in Section 2. The proposed Synthetic Sample Generation LDL pre-processing method is described in Section 3. Then the details of the experiments are reported in Section 4. Finally, the results and conclusions are drawn in Section 5 and Section 6.

## 2. Preliminaries

In this section, the formulation of LDL (Section 2.1) and the LDL algorithms considered for our case study (Section 2.2) are presented. Furthermore, some basic concepts on synthetic sample generation are introduced (Section 2.3), providing the necessary background required to properly present the study carried out in this paper.

### 2.1. Formulation of Label Distribution Learning

Suppose that we are given a set of  $m$  samples  $S = \{(x_1, D_1), \dots, (x_m, D_m)\}$  to train, where  $x_i = \{x_{i1}, x_{i2}, \dots, x_{iq}\}$  is a  $q$ -dimensional vector. For each instance  $x_i$ , the label distribution is denoted by  $D_i = \{d_{x_i}^{y_1}, d_{x_i}^{y_2}, \dots, d_{x_i}^{y_c}\}$  where  $y \in Y$ ,  $Y = \{y_1, y_2, \dots, y_c\}$  denotes the complete set of labels. The constant  $c$  is the number of possible labels and  $d_{x_i}^{y_j}$  is the description degree of the particular  $j$ th label  $y_j$  to the particular  $i$ th instance  $x_i$ . According to the definition, each description degree should satisfy the constraints  $d_{x_i}^{y_j} \in [0, 1]$  and  $\sum_y d_x^y = 1$ .

### 2.2. Label Distribution Learning algorithms

The resolution of an LDL problem can be approached from different points of view. Depending on the approach chosen, the algorithm to be developed will vary substantially, being either a completely new algorithm developed specifically to deal

with LDL constraints, or an adaptation of existing classification algorithms, reformulated to work with such constraints. In all cases, the algorithms chosen for our study are “state of the art” algorithms.

The first one is the AA- $k$ NN adaptation described in [23], given a new instance  $x$ , its  $k$  nearest neighbors are first found in the training set. Then, the mean of the label distributions of all the  $k$  nearest neighbors is calculated as the label distribution of  $x$ . We have chosen this algorithm due to its simplicity and the ease with which it can be enhanced, it will provide a good starting point for a comparative analysis with other more complex methods.

The second is the specialized algorithm SA-BFGS that learns by optimizing an energy function based on the maximum entropy model using a quasi-Newton method. Of the six LDL algorithms proposed by Geng [23], the SA-BFGS stood out due to its great advantage in precision and efficiency.

Another state-of-the-art specialized algorithm is called Structured Random Forest (StructRF) [14]. It is a general LDL model that treats the distribution as an integral whole. In StructRF, all label distributions are mapped to a discrete space at each split node in a random forest. StructRF has proven to be able to train fast and it reaches higher accuracies and lower standard deviations among different measurements.

### 2.3. Motivation for using oversampling technology

Nowadays, one of the main challenges for supervised learning algorithms is still how to deal with datasets that are either incomplete or have significantly skewed class distributions. Both problems can be addressed using an initial data preprocessing that allows us to improve the quality obtained by the learning algorithm used.

The synthetic generation of samples allows us to approach this problem by generating new artificial samples from the originals, thus creating an expanded set of data that we can plug in as input to our learning algorithm.

An example of a widely used technique that achieves this task is the Synthetic Minority Oversampling Technique now widely known as SMOTE [11]. This preprocessing technique carries out an oversampling technique to rebalance the original training set. The basic concept of the SMOTE procedure is to accomplish an interpolation using neighboring minority class samples to improve the generalization capacity of the classifier. The SMOTE preprocessing technique has become a pioneer in the research community and has been widely applied in different scenarios such as face recognition, social media, software engineering, medical diagnosis or social networks. Due to its popularity and influence, SMOTE is considered to be one of the most influential data preprocessing/sampling algorithms in machine learning and data mining [22,19].

As a result, the synthetic generation of samples has been used successfully in multiple areas of study such as:

- Dynamic environments in which data arrives continuously, or data stream. An example of a pre-processing technique to deal with these data streams is Learn++.NSE-SMOTE [16].
- When a data stream is received over time and we have time information at our disposal, we refer to time series classification. The methods SPO and INOS [7] propose an integration of SMOTE in time series classification.
- In a perfect situation, we want to train classifiers using diverse labeled data that thoroughly represent all classes. However, in many real-life applications, there is a huge amount of unlabeled data and thus obtaining a representative subset is a complex process. We refer to semi-supervised classification that utilizes unlabeled data to improve the predictive performance. Several methods based on SMOTE have been developed for this learning paradigm, GS4 [32], SEG-SSC [38] and OCHS-SSC [17]. These methods generate synthetic examples to diminish the drawbacks produced by the absence of labeled examples.
- Several ideas, also based on SMOTE, have been proposed to tackle multi-instance learning, e.g. Instance-SMOTE and Bag-SMOTE [41]. The Instance-SMOTE algorithm creates synthetic minority instances in each bag, without creating new bags. Besides, Bag-SMOTE creates new synthetic minority bags with new instances.
- In multilabel classification [29] each data instance is associated with a vector of outputs, instead of only one value. MLSMOTE [10] is the most popular extension of SMOTE designed for multilabel classification. Its objective is to produce synthetic instances related to minority labels.
- Regression tasks consider the output variable as continuous and hence, the values are represented by real numbers. SMOTER is the SMOTE-based contribution of oversampling regression [37].

### 2.4. Principles of the proposed method

Up until now, all existing LDL studies and algorithms have dealt with raw datasets [23,36,43,20,50,14] and have not applied any kind of pre-processing that might later facilitate the learning of the algorithm. The pursued objective of this study is to improve the performance of the underlying LDL learners without modifying them. To achieve this goal we have created an oversampling method adapted to the LDL restrictions in order to extend the initial datasets by applying synthetic sample generation and use them as input for the state-of-the-art LDL algorithms. Later we will analyze the results and see if they can improve the results obtained by learners. The mechanism of the proposed method can be depicted as in Fig. 1.

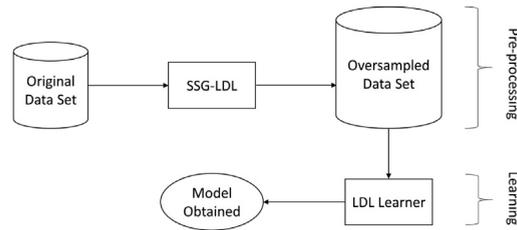


Fig. 1. Principle of the proposed method.

### 3. Synthetic Sample Generation for Label Distribution Learning

SSG-LDL is the proposed pre-processing algorithm that has been adapted to LDL restrictions.

The main idea of this method is to obtain an optimized dataset that will be made up of the original samples plus the artificially generated ones. We define a synthetic sample as a sample generated based on a set of real samples, on which we apply interpolation techniques between remote points that are within a defined neighborhood. The new samples obtained using this procedure are thus added to the raw dataset. Throughout the following subsections we will explain the entire process in detail, covering from how to select the candidate samples to how to obtain a new synthetic sample specially designed to meet the constraints of an LDL problem.

The formal procedure works as follows. First, the total percentage of oversampling  $N$  (an integer value) is set up, then, an iterative process is carried out, composed of several steps. In the initial step, a remote sample is selected at random from the training set (this process is described in Section 3.2). Next, its  $k$  nearest neighbors are obtained. Finally, one of these  $k$  instances is randomly chosen to compute the new sample by interpolation. This synthetic sample generation process is described in detail in Section 3.3 where we will explain how to interpolate the features vector and also how the label dimension is created using the  $k$  nearest neighbor label vectors. This process is repeated iteratively until  $N \times m/100$  synthetic samples are obtained. The whole process is summarized in Algorithm 1.

---

#### Algorithm 1. SSG-LDL

---

Function SSG-LDL( $S[]$ ,  $N$ ,  $k$ ):

```

1  input :  $S[] \leftarrow$  original training set ;
    $N \leftarrow$  % of oversampling ;
    $k \leftarrow$  considered neighbors ;
2  output:  $S'[] \leftarrow$  the oversampled training set
3   $S' = S$  ;
4   $T = m * N/100 \leftarrow$  n° synthetic samples to create;
5  for  $i=1$  to  $T$  do
6      $j = \text{SelectSample}(S)$  ;
7      $ss = \text{CreateSyntheticSample}(S, j, k)$  ;
8      $S' \leftarrow ss$  ;
9      $i = i + 1$ ;
10 |
    end
End Function
  
```

---

#### 3.1. Distances

As previously mentioned, the first step of the algorithm consists in selecting a random sample located far away from the others. To measure the distance between two samples we have chosen the Euclidean distance, adapted to LDL, that we will

apply not only to the dimension of the features ( $x$ ) but also to the dimension of the labels ( $D$ ). The two independent parameters  $f_x$  and  $f_y$  allow us to have two degrees of freedom when calculating the distance between two samples. Keeping this in mind is of interest if a wider space must be covered when creating synthetic LDL samples. Therefore, the final distance between two samples is the sum of the feature's distance multiplied by the factor  $f_x$  and the label's distance multiplied by the factor  $f_y$ . What we calculate at this step, and what will later help us to select remote samples, is the average distance between a sample and all the others as summarized in Eq. 1.

$$\overline{DIST}[x_i] = f_x \frac{\sum_{l=1}^m \text{euclidean}(x_i, x_l)}{m} + f_y \frac{\sum_{l=1}^m \text{euclidean}(D_i, D_l)}{m}. \quad (1)$$

### 3.2. Divergence cumulative selection

A part of the SSG-LDL algorithm is based on choosing isolated samples from which synthetic samples are generated in order to fill in the data set in less populated areas. To do so we will rely on the distances calculated in the previous section, adding some randomness to the process. The implemented method is inspired by the selection mechanism called “roulette wheel selection”, widely used in genetic algorithms [27].

In divergence cumulative selection, the probability of being selected will increase as the distance increases. This has been achieved by dividing the mean distance of a sample to other samples by the total distance, thereby normalizing them to 1 (cumulative probability vector). Then a random selection is made, similarly to how a roulette wheel is spun. Therefore, the probability  $P$  that sample  $x_i$  will be selected is defined according to Eq. (2).

$$P_{x_i} = \frac{\overline{DIST}[x_i]}{\sum_{j=1}^m \overline{DIST}[x_j]}. \quad (2)$$

The whole process is summarized in Algorithm 2. The algorithm generates a random number between 0 and the sum of the frequencies. Then, it goes through the array of frequencies, producing a running total. At some point the running total exceeds the generated threshold. The index at that point corresponds to the selected instance.

---

#### Algorithm 2. Divergence Cumulative Selection Algorithm

---

```

Function SelectSample(S[]):
1  input  : S[] ← original training set
2  output: i ← the selected sample index
      #Calculate the total distance
3  total_dist = 0 ;
4  for i=1 to m do
5  | total_dist = total_dist +  $\overline{DIST}[i]$  ;
   end
      #Return a random position according to the cumulative
      probabilities vector
6  r = random(0, 1) * total_dist ;
7  for i=1 to m do
8  | r = r -  $\overline{DIST}[i]$  ;
9  | if r < 0 then
10 | | return i ;
11 | |
   end
12 |
   end
End Function

```

---

### 3.3. Synthetic features generation

Remember that in the formulation of an LDL problem we have, on the one hand, a vector of features that quantifies the sample and, on the other hand, a vector of labels that qualifies the sample. It is an important constraint to take into account

when generating a synthetic sample as we will have to create both dimensions. The algorithm we have developed here is therefore divided as follows:

To generate the synthetic vector of characteristics we randomly choose one of the closest neighbors (the definition of close is made by taking into account the distance vector calculated in Section 3.1). From these two samples a new one is interpolated that will be located in an intermediate point between the two. The calculation of this intermediate point also introduces a randomness factor in order to carry out the most uniform possible sampling of the data space. Fig. 2 could be used to help understand this process.

Regarding the label synthetic creation, the value assigned to each position of the label is the mean of the  $k$  nearest neighbors.

The formal procedure works as follows. First, the  $k$  nearest neighbors of the sample selected by the selection algorithm are calculated and their indexes are stored. The second step is to create the different elements of the feature array, and to do so we randomly choose one of the neighbors obtained in the previous step and calculate a random intermediate value between the sample and the selected neighbor. This process is repeated iteratively for each of the  $q$  features of the sample. Finally, the vector of labels is generated by taking into account all the nearest  $k$  neighbors and obtaining the value for the new sample as the arithmetic mean of each label vector. The whole process is summarized in Algorithm 3.

---

### Algorithm 3. Create a Synthetic Sample

---

```

Function CreateSyntheticSample( $S[]$ ,  $i$ ,  $k$ ):
1  input :  $S[]$  ← original training set
            $i$  ← selected sample index
            $k$  ← nearest neighbors to be considered
2  output:  $ss$  ← the synthetic sample created
3  #Step 1 - Compute  $k$  nearest neighbors for  $i$ , and save the indexes
           in the  $mnarray$ 
4  #Step 2 - Generate the synthetic features
5   $nm = \text{random}(1, k)$ ;
6  for  $attr = 1$  to  $q$  do
7     Compute:  $dif = S[mnarray[nm]].x_{attr} - S[i].x_{attr}$ ;
8     Compute:  $gap = \text{random}(0, 1)$ ;
9      $ss.x_{attr} = S[i].x_{attr} + gap \times dif$ ;
10
11 end
12 #Step 3 - Generate the synthetic labels
13 for  $lbl = 1$  to  $c$  do
14     for  $nm = 1$  to  $k$  do
15          $ss.d_{lbl} += S[mnarray[nm]].d_{lbl}$ ;
16     end
17      $ss.d_{lbl} /= k$ ;
18 end
19 return  $ss$ ;
End Function

```

---

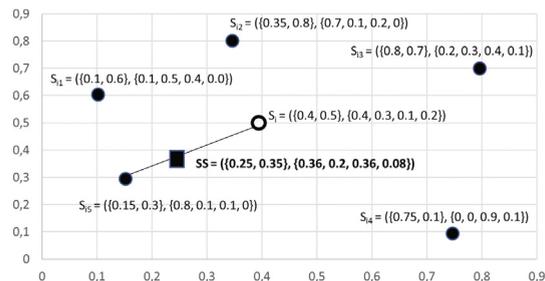


Fig. 2. An illustration of how to create the synthetic data points in the SSG-LDL algorithm.

A simple example of SSG-LDL is illustrated in Fig. 2. Here, an  $S_i$  remote sample is selected as the basis for the creation of a new synthetic data point. Based on a distance metric and using a value of  $k = 5$ , five nearest neighbors (samples  $S_{i1}$  to  $S_{i5}$ ) are chosen from the training set. Finally, the synthetic sample generation is carried out in order to obtain a new synthetic sample  $SS$ . Specifically in this example the value obtained for the first characteristic of the synthetic sample is equal to  $0.4 + (0.15 - 0.4) \times \text{random}(0, 1)$ . For a random value of 0.6 we get a result of 0.25. The rest of the positions of the characteristic vector are calculated in the same way. As for the label vector, the synthetic value obtained corresponds to the average of the values of the 5 neighbors. In our case, the value of the first label is equal to  $(0.1 + 0.7 + 0.2 + 0.0 + 0.8)/5 = 0.36$ . The other positions of the label vector are calculated in the same way.

#### 4. Experimental framework

This section is devoted to introducing the experimental framework used in the different empirical studies of the paper. In our experiments, we have included 15 datasets of a good variety of real-world problems. These are referred to in the following Section 4.1.

In order to evaluate the learners' proficiency, we have employed six measures (described in Section 4.2) for the different aspects of their performance.

For each dataset and learner, these measures were computed over a set merged from the test predictions of a 10-fold cross validation set (10-fcv). In the first pre-processing step we apply our synthetic sample generation algorithm described in Section 3 to each training set and then we use them to train the learners.

Finally, the Wilcoxon test and Bayesian Sign test [15,5] are used to validate the results of the empirical comparisons. In the Bayesian Sign test, a distribution of the differences of the results achieved using methods  $L$  (original learner without applying SSG-LDL) and  $R$  (learner using the SSG-LDL pre-processing method) is computed into a graphical space divided into 3 regions: left, rope and right. The location of most of the distribution in these sectors indicates the final decision: the superiority of algorithm  $L$ , statistical equivalence and the superiority of algorithm  $R$ , respectively. The KEEL package [39] has been used to compute the Wilcoxon test and the R package rNPBST [8] was used to extract the graphical representations of the Bayesian Sign tests analyzed in the following empirical studies. The Rope limit parameter used to represent the Bayesian Sign test is 0.0001.

##### 4.1. Datasets

There are 15 real-world datasets employed in the experiments in total. This is a compilation of the datasets most used in the referenced LDL articles [23,24,36,43,50,14]. The statistics are shown in Table 1.

The first to the tenth datasets were collected from ten biological experiments on the budding yeast *Saccharomyces cerevisiae* [18]. Each dataset includes 2465 yeast genes, and an associated phylogenetic profile vector with a length of 24 is utilized to represent each gene. In a biological experiment, the gene expression level is usually disparate at each discrete time point, so the labels correspond to the time point. Data from different time courses of gene expression in yeast *S. cerevisiae* were combined and clustered. The data were obtained from the time courses during the following processes: the cell division cycle after synchronization by alpha factor arrest (Yeast\_alpha, 18 time points); centrifugal elutriation (Yeast\_elu, 14 time points), and with a temperature-sensitive *cdc15* mutant (Yeast\_cdc, 15 time points); sporulation (Yeast\_spo/spo5/spoem, 7 time points plus four additional samples); shock by high temperature (Yeast\_heat, 6 time points); reducing agents (Yeast\_dtt, 4 time points) and low temperature (Yeast\_cold, 4 time points); and the diauxic shift (Yeast\_diau, 7 time points). All data were gathered using DNA microarrays with elements that represented almost all the ORFs (Open Reading Frame) of the fully sequenced *S. cerevisiae* genome; all measurements were made against a reference sample of time 0 except for the

**Table 1**  
Datasets used in experiments.

No.	Datasets	Examples ( $n$ )	Features ( $q$ )	Labels ( $l$ )
1	Yeast_alpha	2465	24	18
2	Yeast_cdc	2465	24	15
3	Yeast_cold	2465	24	4
4	Yeast_diau	2465	24	7
5	Yeast_dtt	2465	24	4
6	Yeast_elu	2465	24	14
7	Yeast_heat	2465	24	6
8	Yeast_spo	2465	24	6
9	Yeast_spo5	2465	24	3
10	Yeast_spoem	2465	24	2
11	SJAFFE	213	243	6
12	SBU_3DFE	2500	243	6
13	Movie	7755	1869	5
14	Natural_Scene	2000	294	9
15	Human_Gene	30542	36	68

cell-cycle experiments, in which an unsynchronized sample was used. The contribution of each sample to the gene similarity score from a given process was weighted by the inverse of the square root of the number of samples analyzed from that process.

Datasets JAFFE [31] and BU\_3DFE [47] are two widely used facial expression image datasets. In JAFFE, ten expressors posed 3 or 4 examples of each of the six basic facial expressions (fear, disgust, happiness, anger, sadness, surprise) and a neutral face for a total of 213 images of facial expressions. For the simplicity of the experimental design, only Japanese women participated. Each woman took pictures of herself while looking through a semi-reflective plastic sheet into the camera. Tungsten lights were placed to create a uniform illumination of the face. The images were then printed in monochrome and digitized using a flatbed scanner. Finally, images were labelled by 92 people with a 5 or 6 component vector with ratings averaged over all subjects. The similarities between these semantic vectors were computed using Euclidean distance. There were 100 subjects who participated in BU\_3DFE face scans. The resulting database consists of 60% women and 40% men with a variety of ethnic/racial backgrounds. Each subject performed the same basic facial expressions as in JAFFE. With the exception of the neutral expression, each of the six prototypic expressions includes four levels of intensity. Therefore, there are 25 instant 3D expression models for each subject, resulting in a total of 2500 3D facial expression models in the database. Each of these images were subsequently scored by 23 people using the same scale as used in JAFFE. Each dataset is represented by a 243-dimensional feature vector extracted using the Local Binary Patterns method (LBP) [3]. The score for each emotion is regarded as the description degree, and the description degrees (normalized gene expression level) of all the six emotions constitute a label distribution for a particular facial expression image.

Dataset Movie includes 7755 movies. There is a total of 54,243,292 ratings from 478,656 different users on a scale from 1 to 5 integral stars from <sup>®</sup>Netflix. The percentage of each rating level is regarded as the label distribution. There are numeric and categorical attributes in the dataset such as genre, director, country, year, budget and so on. After transforming the categorical attributes into binary vectors, the final feature vector of each movie is 1,869-dimensional.

The Natural Scene dataset [6] is collected from 2,000 natural scene images. Each image is divided into 49 blocks using a 7×7 grid, then the first and second moments (mean and variance) of each band are computed, corresponding to a low-resolution image and to computationally inexpensive texture features, respectively. The end result is a 49×2×3 = 294-dimension feature vector per image. The images were later labeled by three human observers with multiple labels selected from 9 possible values, i.e., sun, sky, water, cloud, mountain, snow, desert, building, and plant. Rankers were required to rank the relevant labels in descending order of relevance and, as expected, the rankings for the same image from different rankers are highly inconsistent. So, a nonlinear programming process [25] is applied to achieve the label distribution.

The Human Gene dataset is much larger than the other datasets used in this experiment. This dataset is collected from biological research on the relationship between human genes and diseases. Each of the 30,542 human genes is represented by the 36 numerical descriptors for a gene sequence proposed in [48]. This 36-D vector was deduced using a modified method based on I-TN curve to display the specific features of protein-coding genes in *Aeropyrum pernix* K1 genome. There are 68 different disease labels in total, and the normalized gene expression level for each disease is considered to be the description degree of the corresponding disease label. The gene expression level of different diseases provides a measure of the description degree of the label for every human gene.

#### 4.2. Evaluation measure selection

We propose using a set of six measures when comparing different LDL algorithms: Chebyshev Distance, Clark Distance, Canberra Metric, Kullback–Leibler Divergence, Cosine Coefficient and Intersection Similarity as shown in Table 2. Each of these measures come from a different syntactic family summarized in [9] and are relatively widely used in the related areas. Thus, they can represent the different aspects of the algorithms well. Another reason that has led us to choose these mea-

**Table 2**  
Evaluation measure for LDL learners. ↓ means that the lowest value is the best and ↑ means the opposite.

Name	Formula
Chebyshev (Cheby)↓	$Dis(D, \hat{D}) = \max_j  d_j - \hat{d}_j $
Clark↓	$Dis(D, \hat{D}) = \sqrt{\sum_{j=1}^c \frac{(d_j - \hat{d}_j)^2}{(d_j + \hat{d}_j)^2}}$
Canberra (Can)↓	$Dis(D, \hat{D}) = \sum_{j=1}^c \frac{ d_j - \hat{d}_j ^2}{d_j + \hat{d}_j}$
Kullback–Leibler (KL)↓	$Dis(D, \hat{D}) = \sum_{j=1}^c d_j \ln \frac{d_j}{\hat{d}_j}$
Cosine (Cos)↑	$Sim(D, \hat{D}) = \frac{\sum_{j=1}^c d_j \hat{d}_j}{\sqrt{\sum_{j=1}^c d_j^2} \sqrt{\sum_{j=1}^c \hat{d}_j^2}}$
Intersection (Inter)↑	$Sim(D, \hat{D}) = \sum_{j=1}^c \min(d_j, \hat{d}_j)$

tures over others is that they are the most widely used in the referenced LDL studies [23,45,24,36,43,50,14], thus allowing us to make a more accurate comparison.

### 4.3. Parameters

Fig. 3 shows the influence of the different values of the parameter  $N$  on the results and the percentage of oversampling applied to the initial dataset. We can observe that when the oversampling percentage is over 300%, the value obtained by the learning algorithm gets stuck. This is the value selected for our synthetic sample generation algorithm in all experiments. Increasing this parameter to over 300% would negatively affect performance without achieving any significant improvements. The figure above represents how the Kullback–Leibler divergence on the *yeast\_spoem* dataset varies; these variations are similar for all other measures and datasets. Other parameter values in synthetic sample generation are:  $k = 5$ , where  $k$  is the number of neighbors considered. With regard to the calculation of distances between points, we have taken 50% of the distance between features and 50% of the distance between labels into account.

There were 4 neighbors selected for the  $k$ -NN algorithm. For the BFGS algorithm, the convergence criterion  $\epsilon$  has been setup to  $10^{-5}$ . Regarding the parameters used for the StructRF algorithm, we have respected the same values as those used in the original proposal papers: numbers of trees = 50, sampling ratio = 0.8, max. depth = 20, min. size of leaf = 5.

An overview of all these parameters can be found in Table 3.

The choice of parameters has been made according to the standards and recommendations given by the authors in the original proposal papers. Due to the fact that the experimental evaluation comprises a high number of datasets, having to adjust each parameter individually for each dataset would be unreasonable. In fact, our idea is just the opposite; we aim to compare the datasets in the most general scenario possible. We have not performed any tuning to adapt these parameters,

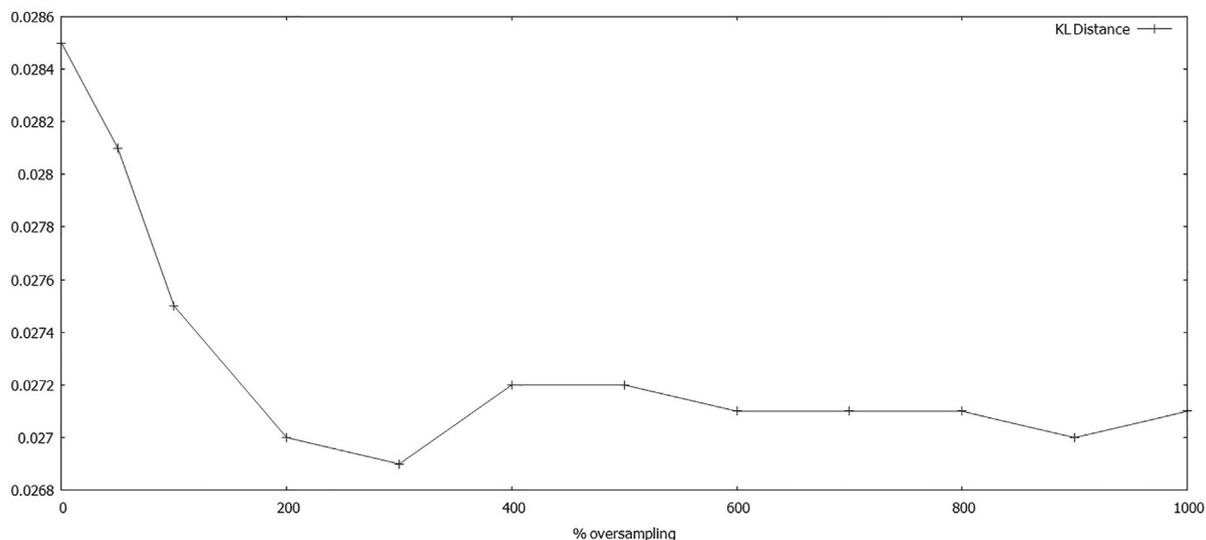


Fig. 3. Influence of experiment parameters on the results using the yeast\_spoem dataset.

Table 3

Summary of the parameters.

Algorithm	Parameter	Description	Value
$k$ -NN	$k$	Number of selected neighbors	4
BFGS	$\epsilon$	Convergence criterion: must be less than $\epsilon$ before successful termination	$10^{-5}$
StructRF	trees	Number of trees	50
	sampling	Sampling ratio of data	0.8
	max. depth	Maximum depth of the tree	20
	min. leaf	Minimum size of the leaf	5
SSG-LDL	$N$	% of oversampling	300
	$k$	Number of selected neighbors	5
	$f_x$	Feature distance factor	0.5
	$f_y$	Label distance factor	0.5

because our objective is not to maximize the accuracy or any other performance measure, but to fairly compare the algorithms and their robustness in a common environment on different datasets.

**Table 4**  
Experimental Results (mean  $\pm$  std) measured by Chebyshev Distance  $\downarrow$ .

Dataset	<i>k</i> -NN	SSG-LDL+ <i>k</i> -NN	BFGS	SSG-LDL+BFGS	StructRF	SSG-LDL+StructRF
Yeast_alpha	0.0148 $\pm$ 0.0007	<b>0.0144 <math>\pm</math> 0.0007</b>	<b>0.0135 <math>\pm</math> 0.0008</b>	<b>0.0135 <math>\pm</math> 0.0008</b>	<b>0.0134 <math>\pm</math> 0.0008</b>	<b>0.0134 <math>\pm</math> 0.0008</b>
Yeast_cdc	0.0177 $\pm$ 0.001	<b>0.0172 <math>\pm</math> 0.001</b>	0.0163 $\pm$ 0.0009	<b>0.0162 <math>\pm</math> 0.0009</b>	0.0162 $\pm$ 0.0009	<b>0.0161 <math>\pm</math> 0.0009</b>
Yeast_cold	0.0554 $\pm$ 0.0021	<b>0.0542 <math>\pm</math> 0.002</b>	0.0512 $\pm$ 0.0018	<b>0.051 <math>\pm</math> 0.0019</b>	0.0501 $\pm$ 0.0016	<b>0.05 <math>\pm</math> 0.0016</b>
Yeast_diau	0.0393 $\pm$ 0.0009	<b>0.0381 <math>\pm</math> 0.0012</b>	<b>0.037 <math>\pm</math> 0.0014</b>	<b>0.037 <math>\pm</math> 0.0015</b>	<b>0.0359 <math>\pm</math> 0.0013</b>	0.036 $\pm$ 0.0016
Yeast_dtt	0.0393 $\pm$ 0.0016	<b>0.038 <math>\pm</math> 0.0016</b>	0.0361 $\pm$ 0.0013	<b>0.0359 <math>\pm</math> 0.0013</b>	0.0353 $\pm$ 0.0013	<b>0.035 <math>\pm</math> 0.0012</b>
Yeast_elu	0.0177 $\pm$ 0.0005	<b>0.0172 <math>\pm</math> 0.0005</b>	<b>0.0163 <math>\pm</math> 0.0006</b>	<b>0.0163 <math>\pm</math> 0.0006</b>	<b>0.0161 <math>\pm</math> 0.0005</b>	<b>0.0161 <math>\pm</math> 0.0006</b>
Yeast_heat	0.0451 $\pm$ 0.0012	<b>0.0438 <math>\pm</math> 0.001</b>	0.0423 $\pm$ 0.0008	<b>0.0422 <math>\pm</math> 0.0008</b>	0.0407 $\pm$ 0.0009	<b>0.0406 <math>\pm</math> 0.0009</b>
Yeast_spo	0.0643 $\pm$ 0.0024	<b>0.0621 <math>\pm</math> 0.0024</b>	0.0584 $\pm$ 0.0037	<b>0.0582 <math>\pm</math> 0.0036</b>	0.0578 $\pm$ 0.003	<b>0.0574 <math>\pm</math> 0.0031</b>
Yeast_spo5	0.0962 $\pm$ 0.0043	<b>0.0949 <math>\pm</math> 0.0043</b>	0.0914 $\pm$ 0.005	<b>0.0913 <math>\pm</math> 0.0051</b>	<b>0.0874 <math>\pm</math> 0.0046</b>	0.0882 $\pm$ 0.0044
Yeast_spoem	0.0924 $\pm$ 0.0036	<b>0.0905 <math>\pm</math> 0.0043</b>	<b>0.0868 <math>\pm</math> 0.0049</b>	0.0869 $\pm$ 0.0052	<b>0.0836 <math>\pm</math> 0.0045</b>	0.0841 $\pm$ 0.0037
SJAFFE	0.1155 $\pm$ 0.018	<b>0.1153 <math>\pm</math> 0.0183</b>	0.1603 $\pm$ 0.016	<b>0.1142 <math>\pm</math> 0.0147</b>	<b>0.1094 <math>\pm</math> 0.0108</b>	0.1141 $\pm$ 0.0136
SBU_3DFE	0.135 $\pm$ 0.0048	<b>0.1327 <math>\pm</math> 0.0052</b>	<b>0.11 <math>\pm</math> 0.0039</b>	0.1231 $\pm$ 0.0048	<b>0.1183 <math>\pm</math> 0.0058</b>	0.1225 $\pm$ 0.0053
Movie	<b>0.1284 <math>\pm</math> 0.0066</b>	0.1317 $\pm$ 0.0081	0.1398 $\pm$ 0.0161	<b>0.1276 <math>\pm</math> 0.0095</b>	<b>0.1104 <math>\pm</math> 0.0048</b>	0.1195 $\pm$ 0.0071
Natural_Scene	<b>0.3168 <math>\pm</math> 0.0081</b>	0.318 $\pm$ 0.0084	0.3549 $\pm$ 0.0159	<b>0.3217 <math>\pm</math> 0.0142</b>	<b>0.2738 <math>\pm</math> 0.0116</b>	0.2863 $\pm$ 0.0117
Human_Gene	0.0652 $\pm$ 0.0048	<b>0.0618 <math>\pm</math> 0.0042</b>	<b>0.0533 <math>\pm</math> 0.0036</b>	<b>0.0533 <math>\pm</math> 0.0036</b>	0.0553 $\pm$ 0.0048	<b>0.0538 <math>\pm</math> 0.004</b>
Average	0.0829 $\pm$ 0.0040	<b>0.0820 <math>\pm</math> 0.0042</b>	0.0845 $\pm$ 0.0051	<b>0.0792 <math>\pm</math> 0.0046</b>	<b>0.0736 <math>\pm</math> 0.0038</b>	0.0755 $\pm$ 0.0040

**Table 5**  
Experimental Results (mean  $\pm$  std) measured by Clark Distance  $\downarrow$ .

Dataset	<i>k</i> -NN	SSG-LDL+ <i>k</i> -NN	BFGS	SSG-LDL+BFGS	StructRF	SSG-LDL+StructRF
Yeast_alpha	0.2321 $\pm$ 0.0112	<b>0.2256 <math>\pm</math> 0.0119</b>	0.2107 $\pm$ 0.0126	<b>0.2103 <math>\pm</math> 0.0126</b>	0.2095 $\pm$ 0.0125	<b>0.2094 <math>\pm</math> 0.0125</b>
Yeast_cdc	0.2375 $\pm$ 0.0134	<b>0.2314 <math>\pm</math> 0.0139</b>	0.2165 $\pm$ 0.0129	<b>0.2161 <math>\pm</math> 0.013</b>	0.2158 $\pm$ 0.0134	<b>0.2152 <math>\pm</math> 0.0132</b>
Yeast_cold	0.1509 $\pm$ 0.007	<b>0.1475 <math>\pm</math> 0.0066</b>	0.1398 $\pm$ 0.0055	<b>0.1394 <math>\pm</math> 0.0058</b>	0.1367 $\pm$ 0.0053	<b>0.1366 <math>\pm</math> 0.0054</b>
Yeast_diau	0.2122 $\pm$ 0.0042	<b>0.2064 <math>\pm</math> 0.006</b>	0.2008 $\pm$ 0.0079	<b>0.2007 <math>\pm</math> 0.0082</b>	<b>0.1947 <math>\pm</math> 0.0076</b>	0.195 $\pm$ 0.0078
Yeast_dtt	0.1068 $\pm$ 0.0045	<b>0.1035 <math>\pm</math> 0.0045</b>	0.0984 $\pm$ 0.0039	<b>0.0979 <math>\pm</math> 0.0038</b>	0.0961 $\pm$ 0.0034	<b>0.0956 <math>\pm</math> 0.0035</b>
Yeast_elu	0.2182 $\pm$ 0.0048	<b>0.2118 <math>\pm</math> 0.0053</b>	0.1992 $\pm$ 0.0055	<b>0.1988 <math>\pm</math> 0.0056</b>	0.1974 $\pm$ 0.0054	<b>0.1969 <math>\pm</math> 0.0057</b>
Yeast_heat	0.1955 $\pm$ 0.0044	<b>0.1901 <math>\pm</math> 0.0034</b>	0.1828 $\pm$ 0.0027	<b>0.1824 <math>\pm</math> 0.0028</b>	0.177 $\pm$ 0.0028	<b>0.1766 <math>\pm</math> 0.0028</b>
Yeast_spo	0.2715 $\pm$ 0.0112	<b>0.2633 <math>\pm</math> 0.0112</b>	0.25 $\pm$ 0.0164	<b>0.2494 <math>\pm</math> 0.0159</b>	0.2467 $\pm$ 0.0139	<b>0.2458 <math>\pm</math> 0.014</b>
Yeast_spo5	0.1933 $\pm$ 0.0099	<b>0.1906 <math>\pm</math> 0.0099</b>	0.1842 $\pm$ 0.0112	<b>0.1839 <math>\pm</math> 0.0113</b>	<b>0.1767 <math>\pm</math> 0.0104</b>	0.178 $\pm$ 0.0096
Yeast_spoem	0.1374 $\pm$ 0.0053	<b>0.1345 <math>\pm</math> 0.0063</b>	<b>0.1293 <math>\pm</math> 0.0072</b>	<b>0.1293 <math>\pm</math> 0.0077</b>	<b>0.1247 <math>\pm</math> 0.0066</b>	0.1254 $\pm$ 0.0054
SJAFFE	0.4137 $\pm$ 0.0489	<b>0.4119 <math>\pm</math> 0.0502</b>	0.6466 $\pm$ 0.0462	<b>0.4254 <math>\pm</math> 0.0427</b>	<b>0.39 <math>\pm</math> 0.035</b>	0.405 $\pm$ 0.0421
SBU_3DFE	0.4255 $\pm$ 0.0134	<b>0.4166 <math>\pm</math> 0.0139</b>	<b>0.3784 <math>\pm</math> 0.0122</b>	0.3852 $\pm$ 0.0144	<b>0.368 <math>\pm</math> 0.0173</b>	0.3786 $\pm$ 0.0167
Movie	<b>0.5686 <math>\pm</math> 0.0393</b>	0.5849 $\pm$ 0.0466	0.6035 $\pm$ 0.056	<b>0.5675 <math>\pm</math> 0.0405</b>	<b>0.5039 <math>\pm</math> 0.0233</b>	0.5467 $\pm$ 0.0371
Natural_Scene	<b>1.8253 <math>\pm</math> 0.0343</b>	1.9844 $\pm$ 0.0382	<b>2.3817 <math>\pm</math> 0.0239</b>	2.4503 $\pm$ 0.0203	<b>2.3946 <math>\pm</math> 0.0227</b>	2.4306 $\pm$ 0.0239
Human_Gene	2.3872 $\pm$ 0.1012	<b>2.3171 <math>\pm</math> 0.0845</b>	2.1111 $\pm$ 0.082	<b>2.1109 <math>\pm</math> 0.0794</b>	2.1776 $\pm$ 0.1232	<b>2.1291 <math>\pm</math> 0.0905</b>
Average	<b>0.505 <math>\pm</math> 0.0209</b>	0.508 $\pm$ 0.0208	0.5289 $\pm$ 0.0204	<b>0.5165 <math>\pm</math> 0.0189</b>	<b>0.5073 <math>\pm</math> 0.0202</b>	0.511 $\pm$ 0.0193

**Table 6**  
Experimental Results (mean  $\pm$  std) measured by Canberra Metric  $\downarrow$ .

Dataset	<i>k</i> -NN	SSG-LDL+ <i>k</i> -NN	BFGS	SSG-LDL+BFGS	StructRF	SSG-LDL+StructRF
Yeast_alpha	0.7577 $\pm$ 0.0372	<b>0.7343 <math>\pm</math> 0.0404</b>	0.6847 $\pm$ 0.0432	<b>0.6831 <math>\pm</math> 0.043</b>	0.6812 $\pm$ 0.0434	<b>0.6801 <math>\pm</math> 0.0428</b>
Yeast_cdc	0.7179 $\pm$ 0.0395	<b>0.6981 <math>\pm</math> 0.04</b>	0.6493 $\pm$ 0.0394	<b>0.6479 <math>\pm</math> 0.0401</b>	0.6472 $\pm$ 0.041	<b>0.6448 <math>\pm</math> 0.0401</b>
Yeast_cold	0.2611 $\pm$ 0.0129	<b>0.2552 <math>\pm</math> 0.0118</b>	0.2407 $\pm$ 0.0097	<b>0.2402 <math>\pm</math> 0.0099</b>	0.2359 $\pm$ 0.0092	<b>0.2356 <math>\pm</math> 0.0097</b>
Yeast_diau	0.456 $\pm$ 0.0109	<b>0.4443 <math>\pm</math> 0.014</b>	<b>0.431 <math>\pm</math> 0.0186</b>	0.4313 $\pm$ 0.0195	<b>0.4177 <math>\pm</math> 0.0178</b>	0.4184 $\pm$ 0.0187
Yeast_dtt	0.1836 $\pm$ 0.0075	<b>0.1779 <math>\pm</math> 0.007</b>	0.1693 $\pm$ 0.0061	<b>0.1685 <math>\pm</math> 0.0058</b>	0.1649 $\pm$ 0.005	<b>0.1643 <math>\pm</math> 0.0052</b>
Yeast_elu	0.6444 $\pm$ 0.0155	<b>0.6242 <math>\pm</math> 0.0164</b>	0.5838 $\pm$ 0.0159	<b>0.5824 <math>\pm</math> 0.0162</b>	0.578 $\pm$ 0.0156	<b>0.5769 <math>\pm</math> 0.0159</b>
Yeast_heat	0.3924 $\pm$ 0.0096	<b>0.3814 <math>\pm</math> 0.0076</b>	0.3647 $\pm$ 0.0058	<b>0.3637 <math>\pm</math> 0.0059</b>	0.3541 $\pm$ 0.0063	<b>0.3528 <math>\pm</math> 0.006</b>
Yeast_spo	0.5594 $\pm$ 0.0232	<b>0.543 <math>\pm</math> 0.0234</b>	0.5137 $\pm$ 0.0335	<b>0.513 <math>\pm</math> 0.0321</b>	0.5066 $\pm$ 0.0274	<b>0.5046 <math>\pm</math> 0.0275</b>
Yeast_spo5	0.2972 $\pm$ 0.0145	<b>0.2934 <math>\pm</math> 0.0143</b>	0.2829 $\pm$ 0.0163	<b>0.2826 <math>\pm</math> 0.0165</b>	<b>0.2713 <math>\pm</math> 0.0151</b>	0.2735 $\pm$ 0.0143
Yeast_spoem	0.1913 $\pm$ 0.0074	<b>0.1872 <math>\pm</math> 0.0088</b>	<b>0.1798 <math>\pm</math> 0.01</b>	0.1799 $\pm$ 0.0107	<b>0.1733 <math>\pm</math> 0.0092</b>	0.1744 $\pm$ 0.0076
SJAFFE	0.8527 $\pm$ 0.0962	<b>0.8471 <math>\pm</math> 0.0998</b>	1.3595 $\pm$ 0.11	<b>0.8783 <math>\pm</math> 0.0992</b>	<b>0.8089 <math>\pm</math> 0.0769</b>	0.84 $\pm$ 0.0938
SBU_3DFE	0.8746 $\pm$ 0.029	<b>0.8611 <math>\pm</math> 0.0302</b>	<b>0.7831 <math>\pm</math> 0.0246</b>	0.8265 $\pm$ 0.0338	<b>0.7817 <math>\pm</math> 0.0354</b>	0.8087 $\pm$ 0.0361
Movie	<b>1.0946 <math>\pm</math> 0.0766</b>	1.1255 $\pm$ 0.0927	1.1679 $\pm$ 0.1163	<b>1.0905 <math>\pm</math> 0.0833</b>	<b>0.9595 <math>\pm</math> 0.0441</b>	1.0412 $\pm$ 0.0737
Natural_Scene	<b>4.2609 <math>\pm</math> 0.0866</b>	4.8189 $\pm$ 0.1051	<b>6.5546 <math>\pm</math> 0.0914</b>	6.7306 $\pm$ 0.0831	<b>6.4559 <math>\pm</math> 0.095</b>	6.6115 $\pm$ 0.1026
Human_Gene	16.2737 $\pm$ 0.7555	<b>15.7894 <math>\pm</math> 0.6271</b>	14.4531 $\pm$ 0.6124	<b>14.4487 <math>\pm</math> 0.5935</b>	14.8633 $\pm$ 0.8857	<b>14.5496 <math>\pm</math> 0.6531</b>
Average	1.8545 $\pm$ 0.0815	<b>1.8521 <math>\pm</math> 0.0759</b>	1.8945 $\pm$ 0.0769	<b>1.8711 <math>\pm</math> 0.0728</b>	1.86 $\pm$ 0.0885	<b>1.8584 <math>\pm</math> 0.0765</b>

**Table 7**Experimental Results (mean  $\pm$  std) measured by Kullback–Leibler Divergence  $\downarrow$ .

Dataset	k-NN	SSG-LDL+k-NN	BFGS	SSG-LDL+BFGS	StructRF	SSG-LDL+StructRF
Yeast_alpha	0.0066 $\pm$ 0.0006	<b>0.0063 <math>\pm</math> 0.0006</b>	<b>0.0055 <math>\pm</math> 0.0006</b>	<b>0.0055 <math>\pm</math> 0.0006</b>	<b>0.0055 <math>\pm</math> 0.0006</b>	<b>0.0055 <math>\pm</math> 0.0006</b>
Yeast_cdc	0.0083 $\pm$ 0.0009	<b>0.0079 <math>\pm</math> 0.0009</b>	<b>0.0070 <math>\pm</math> 0.0008</b>	<b>0.0070 <math>\pm</math> 0.0008</b>	0.0070 $\pm$ 0.0009	<b>0.0069 <math>\pm</math> 0.0008</b>
Yeast_cold	0.0142 $\pm$ 0.0014	<b>0.0136 <math>\pm</math> 0.0014</b>	<b>0.0122 <math>\pm</math> 0.0011</b>	<b>0.0122 <math>\pm</math> 0.0012</b>	<b>0.0118 <math>\pm</math> 0.0011</b>	<b>0.0118 <math>\pm</math> 0.0011</b>
Yeast_diau	0.0150 $\pm$ 0.0008	<b>0.0142 <math>\pm</math> 0.0009</b>	<b>0.0131 <math>\pm</math> 0.0011</b>	<b>0.0131 <math>\pm</math> 0.0011</b>	<b>0.0125 <math>\pm</math> 0.001</b>	<b>0.0125 <math>\pm</math> 0.0011</b>
Yeast_dtt	0.0073 $\pm$ 0.0007	<b>0.0069 <math>\pm</math> 0.0007</b>	0.0063 $\pm$ 0.0006	<b>0.0062 <math>\pm</math> 0.0006</b>	0.0061 $\pm$ 0.0006	<b>0.0060 <math>\pm</math> 0.0006</b>
Yeast_elu	0.0074 $\pm$ 0.0003	<b>0.007 <math>\pm</math> 0.0004</b>	<b>0.0062 <math>\pm</math> 0.0004</b>	<b>0.0062 <math>\pm</math> 0.0004</b>	<b>0.0061 <math>\pm</math> 0.0004</b>	<b>0.0061 <math>\pm</math> 0.0004</b>
Yeast_heat	0.0146 $\pm$ 0.0007	<b>0.0138 <math>\pm</math> 0.0005</b>	<b>0.0126 <math>\pm</math> 0.0004</b>	<b>0.0126 <math>\pm</math> 0.0004</b>	0.0120 $\pm$ 0.0004	<b>0.0119 <math>\pm</math> 0.0004</b>
Yeast_spo	0.0302 $\pm$ 0.0024	<b>0.0283 <math>\pm</math> 0.0024</b>	0.0246 $\pm$ 0.0029	<b>0.0245 <math>\pm</math> 0.0028</b>	0.0243 $\pm$ 0.0024	<b>0.0241 <math>\pm</math> 0.0025</b>
Yeast_spo5	0.0333 $\pm$ 0.0033	<b>0.0322 <math>\pm</math> 0.0036</b>	<b>0.0293 <math>\pm</math> 0.0028</b>	<b>0.0293 <math>\pm</math> 0.0028</b>	<b>0.027 <math>\pm</math> 0.0025</b>	0.0276 $\pm$ 0.0022
Yeast_spoem	0.0285 $\pm$ 0.0023	<b>0.0272 <math>\pm</math> 0.0026</b>	<b>0.0245 <math>\pm</math> 0.0023</b>	0.0246 $\pm$ 0.0024	<b>0.0228 <math>\pm</math> 0.002</b>	0.0233 $\pm$ 0.0018
SJAFFE	0.0730 $\pm$ 0.0162	<b>0.0724 <math>\pm</math> 0.0165</b>	0.1639 $\pm$ 0.0245	<b>0.0754 <math>\pm</math> 0.0166</b>	<b>0.0603 <math>\pm</math> 0.0089</b>	0.0655 $\pm$ 0.0116
SBU_3DFE	0.0907 $\pm$ 0.0048	<b>0.0872 <math>\pm</math> 0.0048</b>	<b>0.0634 <math>\pm</math> 0.0038</b>	0.0689 $\pm$ 0.0038	<b>0.0659 <math>\pm</math> 0.0054</b>	0.0703 $\pm$ 0.0046
Movie	<b>0.1255 <math>\pm</math> 0.0137</b>	0.1313 $\pm$ 0.0167	0.1543 $\pm$ 0.0405	<b>0.1211 <math>\pm</math> 0.0177</b>	<b>0.0901 <math>\pm</math> 0.0061</b>	0.1046 $\pm$ 0.0114
Natural_Scene	1.1924 $\pm$ 0.0765	<b>1.1011 <math>\pm</math> 0.063</b>	1.112 $\pm$ 0.0999	<b>0.8061 <math>\pm</math> 0.0292</b>	<b>0.6436 <math>\pm</math> 0.0211</b>	0.696 $\pm$ 0.0259
Human_Gene	0.3006 $\pm$ 0.0248	<b>0.2815 <math>\pm</math> 0.0206</b>	0.2365 $\pm$ 0.0184	<b>0.236 <math>\pm</math> 0.018</b>	0.248 $\pm$ 0.0265	<b>0.239 <math>\pm</math> 0.0197</b>
Average	0.1298 $\pm$ 0.01	<b>0.1221 <math>\pm</math> 0.009</b>	0.1248 $\pm$ 0.0133	<b>0.0966 <math>\pm</math> 0.0066</b>	<b>0.0829 <math>\pm</math> 0.0053</b>	0.0874 $\pm$ 0.0056

**Table 8**Experimental Results (mean  $\pm$  std) measured by Cosine Coefficient  $\uparrow$ .

Dataset	k-NN	SSG-LDL+k-NN	BFGS	SSG-LDL+BFGS	StructRF	SSG-LDL+StructRF
Yeast_alpha	0.9935 $\pm$ 0.0006	<b>0.9938 <math>\pm</math> 0.0006</b>	<b>0.9946 <math>\pm</math> 0.0006</b>	<b>0.9946 <math>\pm</math> 0.0006</b>	<b>0.9946 <math>\pm</math> 0.0006</b>	<b>0.9946 <math>\pm</math> 0.0006</b>
Yeast_cdc	0.992 $\pm$ 0.0008	<b>0.9924 <math>\pm</math> 0.0008</b>	<b>0.9933 <math>\pm</math> 0.0007</b>	<b>0.9933 <math>\pm</math> 0.0007</b>	<b>0.9933 <math>\pm</math> 0.0008</b>	<b>0.9933 <math>\pm</math> 0.0007</b>
Yeast_cold	0.9866 $\pm$ 0.0012	<b>0.9872 <math>\pm</math> 0.0011</b>	0.9885 $\pm$ 0.0009	<b>0.9886 <math>\pm</math> 0.0010</b>	<b>0.9889 <math>\pm</math> 0.0009</b>	<b>0.9889 <math>\pm</math> 0.0009</b>
Yeast_diau	0.9862 $\pm$ 0.0008	<b>0.987 <math>\pm</math> 0.0008</b>	<b>0.9879 <math>\pm</math> 0.0010</b>	<b>0.9879 <math>\pm</math> 0.0010</b>	<b>0.9885 <math>\pm</math> 0.0010</b>	<b>0.9885 <math>\pm</math> 0.0011</b>
Yeast_dtt	0.9931 $\pm$ 0.0005	<b>0.9934 <math>\pm</math> 0.0005</b>	<b>0.9941 <math>\pm</math> 0.0004</b>	<b>0.9941 <math>\pm</math> 0.0004</b>	<b>0.9943 <math>\pm</math> 0.0004</b>	<b>0.9943 <math>\pm</math> 0.0004</b>
Yeast_elu	0.9928 $\pm$ 0.0003	<b>0.9932 <math>\pm</math> 0.0003</b>	0.994 $\pm$ 0.0004	<b>0.9941 <math>\pm</math> 0.0004</b>	<b>0.9941 <math>\pm</math> 0.0003</b>	<b>0.9941 <math>\pm</math> 0.0004</b>
Yeast_heat	0.9861 $\pm$ 0.0007	<b>0.9869 <math>\pm</math> 0.0006</b>	<b>0.9880 <math>\pm</math> 0.0004</b>	<b>0.9880 <math>\pm</math> 0.0004</b>	0.9886 $\pm$ 0.0004	<b>0.9887 <math>\pm</math> 0.0005</b>
Yeast_spo	0.9716 $\pm$ 0.002	<b>0.9734 <math>\pm</math> 0.002</b>	0.9769 $\pm$ 0.0026	<b>0.977 <math>\pm</math> 0.0025</b>	0.9773 $\pm$ 0.0021	<b>0.9775 <math>\pm</math> 0.0021</b>
Yeast_spo5	0.9705 $\pm$ 0.0025	<b>0.9714 <math>\pm</math> 0.0028</b>	<b>0.9741 <math>\pm</math> 0.0023</b>	<b>0.9741 <math>\pm</math> 0.0023</b>	<b>0.9762 <math>\pm</math> 0.002</b>	0.9756 $\pm$ 0.0018
Yeast_spoem	0.9754 $\pm$ 0.0019	<b>0.9764 <math>\pm</math> 0.0023</b>	<b>0.979 <math>\pm</math> 0.0019</b>	0.9789 $\pm$ 0.002	<b>0.9804 <math>\pm</math> 0.0018</b>	0.9799 $\pm$ 0.0016
SJAFFE	0.9308 $\pm$ 0.0163	<b>0.9313 <math>\pm</math> 0.0166</b>	0.8739 $\pm$ 0.0179	<b>0.9293 <math>\pm</math> 0.0154</b>	<b>0.9429 <math>\pm</math> 0.0083</b>	0.9381 $\pm$ 0.0107
SBU_3DFE	0.9118 $\pm$ 0.0047	<b>0.915 <math>\pm</math> 0.0047</b>	<b>0.9389 <math>\pm</math> 0.0034</b>	0.9324 $\pm$ 0.0037	<b>0.9348 <math>\pm</math> 0.0055</b>	0.9309 $\pm$ 0.0044
Movie	<b>0.9174 <math>\pm</math> 0.0074</b>	0.9135 $\pm$ 0.0093	0.9072 $\pm$ 0.0188	<b>0.9209 <math>\pm</math> 0.0103</b>	<b>0.9407 <math>\pm</math> 0.0039</b>	0.9319 $\pm$ 0.0066
Natural_Scene	<b>0.7043 <math>\pm</math> 0.0137</b>	0.704 $\pm$ 0.0108	0.6724 $\pm$ 0.0149	<b>0.719 <math>\pm</math> 0.0104</b>	<b>0.7895 <math>\pm</math> 0.0091</b>	0.767 $\pm$ 0.0129
Human_Gene	0.7663 $\pm$ 0.0184	<b>0.7847 <math>\pm</math> 0.0136</b>	0.8343 $\pm$ 0.0102	<b>0.8347 <math>\pm</math> 0.0099</b>	0.8202 $\pm$ 0.0205	<b>0.8306 <math>\pm</math> 0.0126</b>
Average	0.9386 $\pm$ 0.0048	<b>0.9402 <math>\pm</math> 0.0045</b>	0.9398 $\pm$ 0.0051	<b>0.9471 <math>\pm</math> 0.0041</b>	<b>0.9536 <math>\pm</math> 0.0038</b>	0.9516 $\pm$ 0.0038

**Table 9**Experimental Results (mean  $\pm$  std) measured by Intersection Similarity  $\uparrow$ .

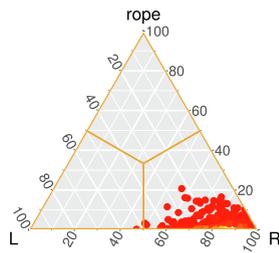
Dataset	k-NN	SSG-LDL+k-NN	BFGS	SSG-LDL+BFGS	StructRF	SSG-LDL+StructRF
Yeast_alpha	0.9581 $\pm$ 0.0020	<b>0.9594 <math>\pm</math> 0.0022</b>	0.9622 $\pm$ 0.0024	<b>0.9623 <math>\pm</math> 0.0024</b>	<b>0.9624 <math>\pm</math> 0.0024</b>	<b>0.9624 <math>\pm</math> 0.0023</b>
Yeast_cdc	0.9527 $\pm$ 0.0025	<b>0.954 <math>\pm</math> 0.0026</b>	0.9573 $\pm$ 0.0026	<b>0.9574 <math>\pm</math> 0.0026</b>	0.9574 $\pm$ 0.0027	<b>0.9576 <math>\pm</math> 0.0026</b>
Yeast_cold	0.9356 $\pm$ 0.0031	<b>0.937 <math>\pm</math> 0.0028</b>	0.9407 $\pm$ 0.0023	<b>0.9408 <math>\pm</math> 0.0023</b>	<b>0.9419 <math>\pm</math> 0.0022</b>	<b>0.9419 <math>\pm</math> 0.0023</b>
Yeast_diau	0.9367 $\pm$ 0.0017	<b>0.9383 <math>\pm</math> 0.002</b>	<b>0.9403 <math>\pm</math> 0.0027</b>	0.9402 $\pm$ 0.0028	<b>0.9421 <math>\pm</math> 0.0026</b>	0.942 $\pm$ 0.0027
Yeast_dtt	0.9547 $\pm$ 0.0017	<b>0.9561 <math>\pm</math> 0.0016</b>	0.9582 $\pm$ 0.0013	<b>0.9584 <math>\pm</math> 0.0013</b>	0.9593 $\pm$ 0.0011	<b>0.9595 <math>\pm</math> 0.0011</b>
Yeast_elu	0.9545 $\pm$ 0.0011	<b>0.956 <math>\pm</math> 0.0012</b>	0.9588 $\pm$ 0.0011	<b>0.9589 <math>\pm</math> 0.0012</b>	0.9592 $\pm$ 0.0011	<b>0.9593 <math>\pm</math> 0.0011</b>
Yeast_heat	0.9356 $\pm$ 0.0017	<b>0.9374 <math>\pm</math> 0.0014</b>	0.9401 $\pm$ 0.001	<b>0.9403 <math>\pm</math> 0.001</b>	0.9419 $\pm$ 0.0011	<b>0.9421 <math>\pm</math> 0.001</b>
Yeast_spo	0.9076 $\pm$ 0.0036	<b>0.9104 <math>\pm</math> 0.0037</b>	0.9154 $\pm$ 0.0053	<b>0.9156 <math>\pm</math> 0.0051</b>	0.9167 $\pm$ 0.0042	<b>0.917 <math>\pm</math> 0.0043</b>
Yeast_spo5	0.9038 $\pm$ 0.0043	<b>0.9051 <math>\pm</math> 0.0043</b>	0.9086 $\pm$ 0.005	<b>0.9087 <math>\pm</math> 0.0051</b>	<b>0.9126 <math>\pm</math> 0.0046</b>	0.9118 $\pm$ 0.0044
Yeast_spoem	0.9076 $\pm$ 0.0036	<b>0.9095 <math>\pm</math> 0.0043</b>	<b>0.9132 <math>\pm</math> 0.0049</b>	0.9131 $\pm$ 0.0052	<b>0.9164 <math>\pm</math> 0.0045</b>	0.9159 $\pm$ 0.0037
SJAFFE	0.8529 $\pm$ 0.0176	<b>0.8539 <math>\pm</math> 0.0182</b>	0.7788 $\pm$ 0.0162	<b>0.8496 <math>\pm</math> 0.0185</b>	<b>0.8622 <math>\pm</math> 0.0132</b>	0.8568 $\pm$ 0.0164
SBU_3DFE	0.8394 $\pm$ 0.0053	<b>0.8425 <math>\pm</math> 0.0053</b>	<b>0.8616 <math>\pm</math> 0.0041</b>	0.8521 $\pm$ 0.0056	<b>0.8592 <math>\pm</math> 0.0065</b>	0.8546 $\pm$ 0.0063
Movie	<b>0.8153 <math>\pm</math> 0.0109</b>	0.8095 $\pm$ 0.014	0.804 $\pm$ 0.0199	<b>0.8171 <math>\pm</math> 0.0137</b>	<b>0.8432 <math>\pm</math> 0.006</b>	0.8285 $\pm$ 0.0107
Natural_Scene	<b>0.5634 <math>\pm</math> 0.0098</b>	0.5588 $\pm$ 0.0079	<b>0.5248 <math>\pm</math> 0.0145</b>	0.5226 $\pm$ 0.0082	<b>0.5919 <math>\pm</math> 0.0097</b>	0.5633 $\pm$ 0.0098
Human_Gene	0.7417 $\pm$ 0.0133	<b>0.7525 <math>\pm</math> 0.0104</b>	0.7842 $\pm$ 0.0092	<b>0.7845 <math>\pm</math> 0.0089</b>	0.7739 $\pm$ 0.0159	<b>0.7814 <math>\pm</math> 0.0105</b>
Average	0.8773 $\pm$ 0.0055	<b>0.8787 <math>\pm</math> 0.0055</b>	0.8765 $\pm$ 0.0062	<b>0.8814 <math>\pm</math> 0.0056</b>	<b>0.8894 <math>\pm</math> 0.0052</b>	0.8863 $\pm$ 0.0053

### 5. Results and analysis

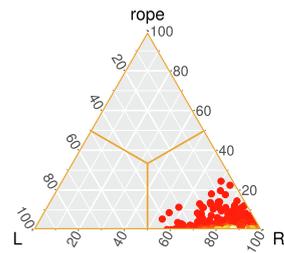
This section presents the results of the empirical studies and their analyses. For each learner we will compare the version with and without pre-processing, showing their different strengths. As those algorithms have been tested using 10-fcv, the performance is represented using “mean ± standard deviation”. Each of the tables shows the experimental results obtained for each learner, comparing them with the results obtained from same learner but using SSG-LDL pre-processing. In each

**Table 10**  
Wilcoxon Signed Ranks test: SyntheticLDL+k-NN vs. k-NN.

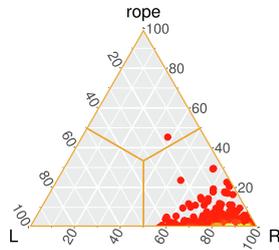
Measure	$R^+$	$R^-$	$p$ -value
Cheby	101	19	0.018066
Clark	92	28	0.073
Can	92	28	0.073
KL	107	13	0.005372
Cos	103	17	0.012452
Inter	93	27	0.06372



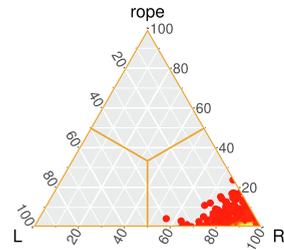
(a) Chebyshev Distance



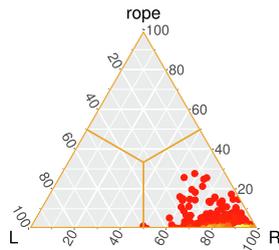
(b) Clark Distance



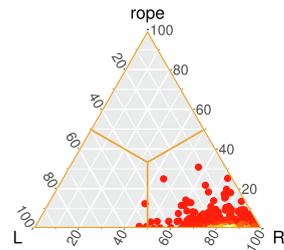
(c) Canberra Metric



(d) Kullback-Leibler Divergence



(e) Cosine Coefficient



(f) Intersection Similarity

**Fig. 4.** Bayesian Sign test comparing k-NN (L) vs. SSG-LDL+k-NN (R).

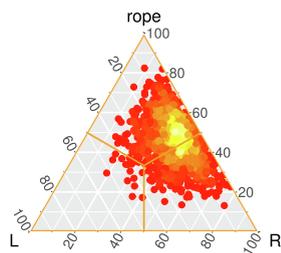
case the best result is highlighted in bold. The last row is the mean of each column. The best average is also highlighted in bold. The results of the different measures are shown in Tables 4–9.

5.1. Evaluation of SSG-LDL+k-NN vs. k-NN

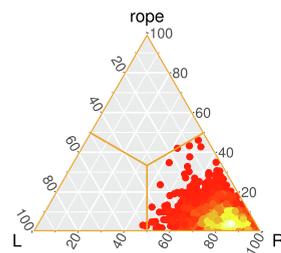
A comparison between k-NN and k-NN using a previously pre-processed dataset using SSG-LDL highlights the best ranking of the learner when using a preprocessed dataset.

**Table 11**  
Wilcoxon Signed Ranks test: SyntheticLDL+BFGS vs. BFGS.

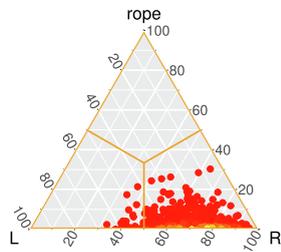
Measure	$R^+$	$R^-$	$p$ -value
Cheby	94	26	0.05536
Clark	81	24	0.0785
Can	90	30	0.0946
KL	76	29	0.15308
Cos	89	31	0.107
Inter	87	33	0.13538



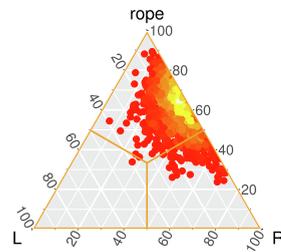
(a) Chebyshev Distance



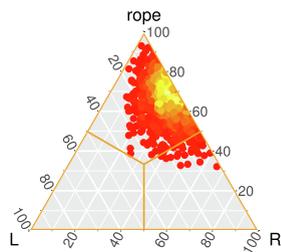
(b) Clark Distance



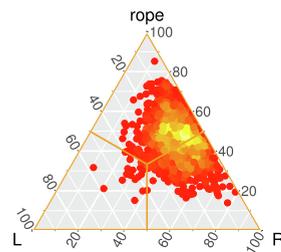
(c) Canberra Metric



(d) Kullback-Leibler Divergence



(e) Cosine Coefficient



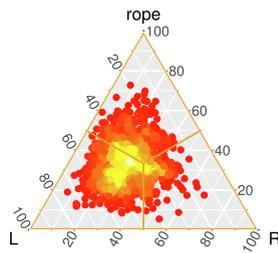
(f) Intersection Similarity

**Fig. 5.** Bayesian Sign test comparing BFGS (L) vs. SSG-LDL+BFGS (R).

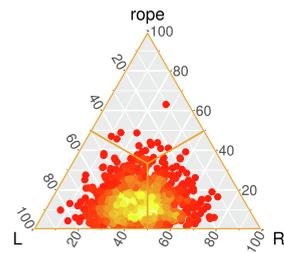
As previously mentioned, we have used the Wilcoxon Signed Ranks test and the Bayesian Sign test to corroborate the significance of the differences between our approach and the selected methods. Table 10 includes the outcome of the Wilcoxon test comparing both learners. All the hypotheses of equivalence are rejected with small p-values. Fig. 4 graphically represents the difference between using data pre-processing or not and its statistical significance in terms of accuracy. The fol-

**Table 12**  
Wilcoxon Signed Ranks test: SyntheticLDL+StructRF vs. StructRF.

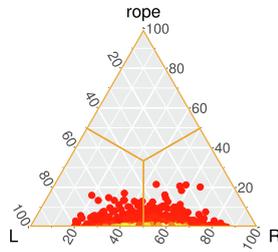
Measure	$R^+$	$R^-$	p-value
Cheby	42	78	$\geq 0.2$
Clark	49	71	$\geq 0.2$
Can	53	67	$\geq 0.2$
KL	44	76	$\geq 0.2$
Cos	40	81	$\geq 0.2$
Inter	44	77	$\geq 0.2$



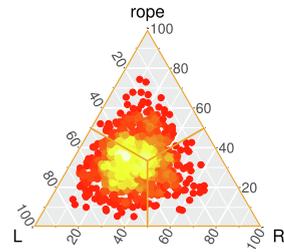
(a) Chebyshev Distance



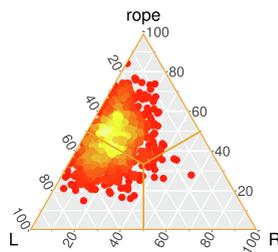
(b) Clark Distance



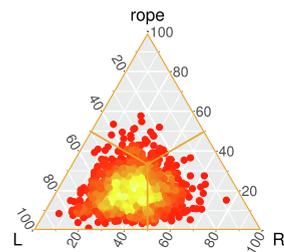
(c) Canberra Metric



(d) Kullback-Leibler Divergence



(e) Cosine Coefficient



(f) Intersection Similarity

**Fig. 6.** Bayesian Sign test comparing StructRF (L) vs. SSG-LDL+StructRF (R).

lowing heat-maps clearly indicate the significant superiority of SSG-LDL, as the computed distributions are always located in the right region.

### 5.2. Evaluation of SSG-LDL+BFGS vs. BFGS

A comparison between BFGS and BFGS using a previously pre-processed dataset with synthetic feature generation also highlights the best ranking (for all distance measures) of the learner when using a preprocessed dataset. Table 11 includes the outcome of the Wilcoxon test comparing both learners. All the hypotheses of equivalence are rejected with small p-values. Fig. 5 graphically represents the difference between using data pre-processing or not and its statistical significance in terms of accuracy. The heat-maps clearly indicate significant superiority, for almost all measures, when applying SSG-LDL, since the computed distributions are always located in the right region. In view of the results of this statistical analysis, there is a high probability that the results obtained will improve after applying this pre-processing method.

### 5.3. Evaluation of SSG-LDL+StructRF vs. StructRF

Let's see how the SSG-LDL algorithm behaves when it is combined with a learning method as robust as StructRF. As for the general classification, we see that applying pre-processing improves the results in 40% of the cases and ties in 15% of the cases. If we go into further detail we observe that the behavior varies according to the dataset: for the *Human\_Gene* for example we notice that applying synthetic sample generation does significantly improve the original results although it does not substantially alter the results for the *Yeast* datasets and causes slightly worse outputs for the sets related to facial expression images, movies or natural scenes. Table 12 includes the outcome of the Wilcoxon test comparing both learners. Fig. 6 also corroborates these results statistically speaking, showing that applying pre-processing on the StructRF algorithm can be advantageous in a high percentage of cases. In summary, by applying a data pre-processing method such as the one presented here, SSG-LDL, significantly improves the results obtained by the learning algorithm. In the case of *k*-NN, previously applying an SSG-LDL treatment makes the learner behave more efficiently in practically all datasets. This also happens when the same type of pre-processing is carried out on a BFGS learning method, obtaining notable improvement in the results. Even when using a method as robust as StructRF, it is also advantageous to apply the SSG-LDL as it can improve results in a high percentage of experiments.

## 6. Conclusion

LDL datasets suffer high dispersion that can negatively affect learners' performance.

In this paper, we proposed a novel method of synthetic data generation that adapts to LDL constraints. The SSG-LDL proposal can be applied, in a previous pre-processing stage, to any training set and then can be plugged in as an input to any learning algorithm.

We have based the design of this algorithm on techniques that have already demonstrated their potential in the area of data pre-processing. We have then extended and adapted them to be compatible with a distribution of labels as the output.

In order to verify the effectiveness of the solution designed, it has been applied to three state of the art learners in the LDL scope and we have verified that the results obtained are very encouraging.

The SSG-LDL method is the first proposal of oversampling adapted to LDL restrictions. In future research we want to make improvements on the current proposal to make it even stronger. Some ideas that we can anticipate are as follows:

- To select the instances to be oversampled, SSG-LDL uses the Euclidean distance. Some recent methods of oversampling such as MDO [2] build synthetic examples that have the same Mahalanobis distance from each examined class mean as the other minority examples. Thus, an alternative procedure that we would consider integrating in future studies is to improve learning in the region of minority instances by preserving the covariance and the probability during the generation of synthetic examples. The Hellinger distance metric, based on probability distributions, is more tolerant to skewed class distributions. This metric has been applied in the context of data reduction [46] and it also be of interest to introduce it in an oversampling method such as SSG-LDL.
- ProLSFEO-LDL [28] is a novel method that simultaneously addresses the prototype selection and the label-specific feature selection pre-processing techniques, specifically developed to deal with LDL restrictions. Devising a method that intelligently combines SSG-LDL as an oversampling method and ProLSFEO-LDL as a data reduction strategy could yield a very effective pre-processing method.
- Combining SMOTE and a boosting procedure such as SMOTEBoost [13] or MDOBoost [1] usually improves the prediction performance. Extending either of these two methods to use SSG-LDL and making them compatible with LDL problems is an interesting approach to follow.
- The filtering of artificial samples is a frequent operation that supports the success of SMOTE on real data. Adding a rebalancing stage to SSG-LDL similar to AMSCO [30] could mitigate the potential drawback of generating overlapping and noisy examples.

## CRedit authorship contribution statement

**Manuel González:** Conceptualization, Methodology, Software, Investigation, Writing - review & editing. **Julián Luengo:** Validation, Supervision. **José-Ramón Cano:** Supervision, Resources. **Salvador García:** Funding acquisition, Project administration, Supervision.

## 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.

## Acknowledgements

This work is supported by the Spanish National Research Project TIN2017-89517-P.

## References

- [1] L. Abdi, S. Hashemi, To combat multi-class imbalanced problems by means of over-sampling and boosting techniques, *Soft Computing* 19 (12) (2015) 3369–3385.
- [2] L. Abdi, S. Hashemi, To combat multi-class imbalanced problems by means of over-sampling techniques, *IEEE Transactions on Knowledge and Data Engineering* 28 (1) (2016) 238–251.
- [3] T. Ahonen, A. Hadid, M. Pietikainen, Face description with local binary patterns: application to face recognition, *IEEE Transactions on Pattern Analysis & Machine Intelligence* 12 (2006) 2037–2041.
- [4] R. Barandela, J.S. Sánchez, V. García, E. Rangel, Strategies for learning in class imbalance problems, *Pattern Recognition* 36 (3) (2003) 849–851.
- [5] A. Benavoli, G. Corani, J. Demšar, M. Zaffalon, Time for a change: a tutorial for comparing multiple classifiers through bayesian analysis, *The Journal of Machine Learning Research* 18 (1) (2017) 2653–2688.
- [6] M.R. Boutell, J. Luo, X. Shen, C.M. Brown, Learning multi-label scene classification, *Pattern Recognition* 37 (9) (2004) 1757–1771.
- [7] H. Cao, X.-L. Li, D.Y. Woon, S.-K. Ng, Integrated oversampling for imbalanced time series classification, *IEEE Transactions on Knowledge and Data Engineering* 25 (12) (2013) 2809–2822.
- [8] J. Carrasco, S. García, M. del Mar Rueda, F. Herrera, rnpbst: An r package covering non-parametric and bayesian statistical tests, in: *International Conference on Hybrid Artificial Intelligence Systems*, Springer, 2017, pp. 281–292.
- [9] S.-H. Cha, Comprehensive survey on distance/similarity measures between probability density functions, *City* 1 (2) (2007) 1.
- [10] F. Charte, A.J. Rivera, M.J. del Jesus, F. Herrera, Mlsmote: approaching imbalanced multilabel learning through synthetic instance generation, *Knowledge-Based Systems* 89 (2015) 385–397.
- [11] N.V. Chawla, K.W. Bowyer, L.O. Hall, W.P. Kegelmeyer, Smote: synthetic minority over-sampling technique, *Journal of Artificial Intelligence Research* 16 (2002) 321–357.
- [12] N.V. Chawla, D.A. Cieslak, L.O. Hall, A. Joshi, Automatically countering imbalance and its empirical relationship to cost, *Data Mining and Knowledge Discovery* 17 (2) (2008) 225–252.
- [13] N.V. Chawla, A. Lazarevic, L.O. Hall, K.W. Bowyer, Smoteboost: improving prediction of the minority class in boosting, in: *European Conference on Principles of Data Mining and Knowledge Discovery*, Springer, 2003, pp. 107–119.
- [14] M. Chen, X. Wang, B. Feng, W. Liu, Structured random forest for label distribution learning, *Neurocomputing* 320 (2018) 171–182.
- [15] J. Derrac, S. García, D. Molina, F. Herrera, A practical tutorial on the use of nonparametric statistical tests as a methodology for comparing evolutionary and swarm intelligence algorithms, *Swarm and Evolutionary Computation* 1 (1) (2011) 3–18.
- [16] G. Ditzler, R. Polikar, Incremental learning of concept drift from streaming imbalanced data, *IEEE Transactions on Knowledge and Data Engineering* 25 (10) (2013) 2283–2301.
- [17] A. Dong, F.-L. Chung, S. Wang, Semi-supervised classification method through oversampling and common hidden space, *Information Sciences* 349 (2016) 216–228.
- [18] M.B. Eisen, P.T. Spellman, P.O. Brown, D. Botstein, Cluster analysis and display of genome-wide expression patterns, *Proceedings of the National Academy of Sciences* 95 (25) (1998) 14863–14868.
- [19] A. Fernández, S. García, F. Herrera, N.V. Chawla, Smote for learning from imbalanced data: progress and challenges, marking the 15-year anniversary, *Journal of Artificial Intelligence Research* 61 (2018) 863–905.
- [20] B. Gao, C. Xing, C. Xie, J. Wu, X. Geng, Deep label distribution learning with label ambiguity, *IEEE Transactions on Image Processing* 26 (6) (2017) 2825–2838.
- [21] S. García, J. Luengo, F. Herrera, *Data Preprocessing in Data Mining* vol. 72 (2015).
- [22] S. García, J. Luengo, F. Herrera, Tutorial on practical tips of the most influential data preprocessing algorithms in data mining, *Knowledge-Based Systems* 98 (2016) 1–29.
- [23] X. Geng, Label distribution learning, *IEEE Transactions on Knowledge and Data Engineering* 28 (7) (2016) 1734–1748.
- [24] X. Geng, P. Hou, Pre-release prediction of crowd opinion on movies by label distribution learning, in: *Twenty-Fourth International Joint Conference on Artificial Intelligence*, vol. 2015-January, Affiliation: School of Computer Science and Engineering, Southeast University, Nanjing, China, 2015, pp. 3511–3517..
- [25] X. Geng, L. Luo, Multilabel ranking with inconsistent rankers, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 3742–3747.
- [26] X. Geng, C. Yin, Z. Zhou, Facial age estimation by learning from label distributions, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35 (10) (2013) 2401–2412.
- [27] D.E. Goldberg, *Genetic algorithms in search, Optimization, and Machine Learning* (1989).
- [28] M. González, J.-R. Cano, S. García, Prolsfeo-ldl: Prototype selection and label-specific feature evolutionary optimization for label distribution learning, *Applied Sciences* 10 (9) (2020) 3089.
- [29] F. Herrera, F. Charte, A.J. Rivera, M.J. Del Jesus, *Multilabel Classification*, Springer, 2016, pp. 17–31..
- [30] J. Li, S. Fong, R.K. Wong, V.W. Chu, Adaptive multi-objective swarm fusion for imbalanced data classification, *Information Fusion* 39 (2018) 1–24.
- [31] M. Lyons, S. Akamatsu, M. Kamachi, J. Gyoba, Coding facial expressions with gabor wavelets, in: *Proceedings Third IEEE International Conference on Automatic Face and Gesture Recognition*, IEEE, 1998, pp. 200–205.
- [32] P. Moutafis, I.A. Kakadiaris, Gs4: Generating synthetic samples for semi-supervised nearest neighbor classification, in: *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, Springer, 2014, pp. 393–403.

- [33] R.C. Prati, G.E. Batista, D.F. Silva, Class imbalance revisited: a new experimental setup to assess the performance of treatment methods, *Knowledge and Information Systems* 45 (1) (2015) 247–270.
- [34] Y. Ren, X. Geng, Sense beauty by label distribution learning, in: *IJCAI*. Affiliation: MOE Key Laboratory of Computer Network and Information Integration, School of Computer Science and Engineering, Southeast University, Nanjing, China, 2017, pp. 2648–2654.
- [35] A. Roy, R.M. Cruz, R. Sabourin, G.D. Cavalcanti, A study on combining dynamic selection and data preprocessing for imbalance learning, *Neurocomputing* 286 (2018) 179–192.
- [36] W. Shen, K. Zhao, Y. Guo, A. Yuille, Label distribution learning forests, in: *Advances in Neural Information Processing Systems*, vol. 2017–December. Affiliation: Key Laboratory of Specialty Fiber Optics and Optical Access Networks, Shanghai Institute for Advanced Communication and Data Science, School of Communication and Information Engineering, Shanghai University, China; Affiliation: Department of Computer Science, Johns Hopkins University, United States, 2017, pp. 835–844.
- [37] L. Torgo, P. Branco, R.P. Ribeiro, B. Pfahringer, Resampling strategies for regression, *Expert Systems* 32 (3) (2015) 465–476.
- [38] I. Triguero, S. García, F. Herrera, Seg-ssc: A framework based on synthetic examples generation for self-labeled semi-supervised classification, *IEEE Transactions on Cybernetics* 45 (4) (2015) 622–634.
- [39] I. Triguero, S. Gonzalez, J.M. Moyano, S. García, J. Alcalá-Fdez, J. Luengo, A. Fernandez, M. Jose del Jesus, L. Sanchez, F. Herrera, Keel 3.0: An open source software for multi-stage analysis in data mining, *International Journal of Computational Intelligence Systems* 10 (1) (2017) 1238–1249.
- [40] G. Tsoumakas, I. Katakis, Multi-label classification: an overview, *International Journal of Data Warehousing and Mining (IJDDWM)* 3 (3) (2007) 1–13.
- [41] J. Wang, B. Yun, P. Huang, Y.-A. Liu, Applying threshold smote algorithm with attribute bagging to imbalanced datasets, in: *International Conference on Rough Sets and Knowledge Technology*, Springer, 2013, pp. 221–228.
- [42] Y. Wang, J. Dai, Label distribution feature selection based on mutual information in fuzzy rough set theory, 2019, Vol. 2019–July.
- [43] C. Xing, X. Geng, H. Xue, Logistic boosting regression for label distribution learning, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, vol. 2016–January. Affiliation: Key Lab of Computer Network and Information Integration, Ministry of Education, School of Computer Science and Engineering, Southeast University, Nanjing, 211189, China, 2016, pp. 4489–4497.
- [44] D. Xue, Z. Hong, S. Guo, L. Gao, L. Wu, J. Zheng, N. Zhao, Personality recognition on social media with label distribution learning, *IEEE Access* 5 (2017) 13478–13488.
- [45] J. Yang, D. She, M. Sun, Joint image emotion classification and distribution learning via deep convolutional neural network, in: *IJCAI*. Affiliation: College of Computer and Control Engineering, Nankai University, Tianjin, China, 2017, pp. 3266–3272.
- [46] L. Yin, Y. Ge, K. Xiao, X. Wang, X. Quan, Feature selection for high-dimensional imbalanced data, *Neurocomputing* 105 (2013) 3–11.
- [47] L. Yin, X. Wei, Y. Sun, J. Wang, M.J. Rosato, A 3d facial expression database for facial behavior research, in: *7th International Conference on Automatic Face and Gesture Recognition (FG06)*, IEEE, 2006, pp. 211–216.
- [48] J.-F. Yu, D.-K. Jiang, K. Xiao, Y. Jin, J.-H. Wang, X. Sun, Discriminate the falsely predicted protein-coding genes in *aeropyrum pernix* k1 genome based on graphical representation, *Match-Communications in Mathematical and Computer Chemistry* 67 (3) (2012) 845.
- [49] Z. Zhang, M. Wang, X. Geng, Crowd counting in public video surveillance by label distribution learning, *Neurocomputing* 166 (2015) 151–163.
- [50] X. Zheng, X. Jia, W. Li, Label distribution learning by exploiting sample correlations locally, in: *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.