

Multi-label classification via incremental clustering on an evolving data stream

Author

Tien, Thanh Nguyen, Manh, Truong Dang, Anh, Vu Luong, Liew, Alan Wee-Chung, Liang, Tiancai, McCall, John

Published

2019

Journal Title

Pattern Recognition

Version

Accepted Manuscript (AM)

DOI

[10.1016/j.patcog.2019.06.001](https://doi.org/10.1016/j.patcog.2019.06.001)

Rights statement

© 2019 Elsevier. Licensed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International Licence (<http://creativecommons.org/licenses/by-nc-nd/4.0/>) which permits unrestricted, non-commercial use, distribution and reproduction in any medium, providing that the work is properly cited.

Downloaded from

<http://hdl.handle.net/10072/386259>

Griffith Research Online

<https://research-repository.griffith.edu.au>

Highlights:

- An incremental clustering-based multi-label online classification algorithm for multi-label data stream is proposed.
- To handle concept drift, our algorithm evolves with time, giving higher attention to more recent samples than older samples through a weight decay mechanism.
- Our algorithm dynamically determines the number of predicted labels based on Hoeffding inequality and the label cardinality.
- Extensive comparative experiments with the state-of-the-art algorithms validated the superior performance of our algorithm in both the stationary and concept drift settings.

ACCEPTED MANUSCRIPT

Multi-Label Classification via Incremental Clustering on Evolving Data Stream

Tien Thanh Nguyen¹, Manh Truong Dang², Anh Vu Luong³, Alan Wee-Chung Liew³, Tiancai Liang⁴,
John McCall¹

¹School of Computing Science and Digital Media, Robert Gordon University, Aberdeen, Scotland, UK

²School of Information and Communication Technology, Hanoi University of Science and Technology,
Vietnam

³School of Information and Communication Technology, Griffith University, Australia

⁴GRGBanking Technology Co., Ltd, China

*Corresponding Author: Alan Wee-Chung Liew

Email: a.liew@griffith.edu.au

Abstract: With the advancement of storage and processing technology, an enormous amount of data is collected on a daily basis in many applications. Nowadays, advanced data analytics have been used to mine the collected data for useful information and make predictions, contributing to the competitive advantages of companies. The increasing data volume, however, has posed many problems to classical batch learning systems, such as the need to retrain the model completely with the newly arrived samples or the impracticality of storing and accessing a large volume of data. This has prompted interest on incremental learning that operates on data streams. In this study, we develop an incremental online multi-label classification (OMLC) method based on a weighted clustering model. The model is made to adapt to the change of data via the decay mechanism in which each sample's weight dwindles away over time. The clustering model therefore always focuses more on newly arrived samples. In the classification process, only clusters whose weights are greater than a threshold (called mature clusters) are employed to assign labels for the samples. In our method, not only is the clustering model incrementally maintained with the revealed ground truth labels of the arrived samples, the number of predicted labels in a sample are also adjusted based on the Hoeffding inequality and the label cardinality. The experimental results show that our method is competitive compared to several well-known benchmark algorithms on six performance measures in both the stationary and the concept drift settings.

Keywords: multi-label classification, incremental learning, online learning, clustering, data stream, concept drift

1. Introduction

Nowadays, an enormous amount of data is collected on a daily basis, from mobile devices to social networking sites. Data is growing at an astonishing rate according to IDC Research, at a compound annual growth rate of 42% through to 2020. This means that 90% of the data in the whole world has been created over the past two years. In recent years, ‘big data’ is one of the most popular terms mentioned in the media.

Realizing that data is a hidden resource, many companies have invested heavily in advanced data analytics to mine customer and sales data for useful information in order to gain a competitive edge. For example, about 35% of Amazon.com’s revenue is generated by its recommendation system. However, the tremendous growth in the volume of data has also posed many challenges to classical machine learning systems. First, batch learning performed on large volume of data is sometime impractical, resulting in more and more data not been processed [1]. In addition, learning models trained on existing data become outdated with the appearance of new data. Periodic re-training using the accumulated data can only be a temporary solution, and resulted in huge resource consumption [2, 3]. Learning models that are maintained incrementally by integrating new information acquired from newly arrived data are, therefore, more practical [3].

In this paper, we introduce an incremental online learning method for the multi-label classification (MLC) problem. The MLC problem arises from many real-world applications where an entity is described by multiple terms or having multiple semantic meanings. As a generalization of the multi-class classification problem, the task of MLC is to assign a set of labels to an object to express its semantics. In the literature, there are MLC algorithms for both the batch [4] and stream settings [5, 6]. In this study, we derive a clustering based MLC algorithm based on an online clustering algorithm adopted from [7] to solve the MLC problem. The proposed method can adapt to concept drift in the data stream by focusing more on the newly arrived samples using a decay mechanism [8]. In addition, the label distributions in the K closest clusters are used during the MLC to predict a set of labels for a newly arrived sample. In contrast to [5, 9] where an update of the number of labels assigned to each arrived sample is performed only after a fixed set of samples are received, our algorithm performs continuous update of the number of labels for each arrived sample.

The main contributions of this paper are:

- An online learning algorithm using the clustering model is proposed for the MLC problem.
- The clusters in our algorithm evolve with time, giving higher attention to more recent samples than older samples through a weight decay mechanism.
- A novel approach for learning the number of predicted labels for MLC based on the Hoeffding inequality and the label cardinality is proposed.

- An empirical demonstration that our method is competitive to several well-known benchmark algorithms in both the stationary and concept drift settings.

The paper is organized as follows. In Section 2, we briefly review several well-known learning algorithms for MLC in the supervised learning and stream learning settings. In section 3, we describe the proposed online MLC algorithm based on the online clustering model. The setting for experimental studies is described in section 4. Section 5 presents the detailed experimental results and discussion. Finally, Section 6 provides the conclusions and suggestions for further study.

TABLE.1. SUMMARY OF MAIN NOTATION

Notation	Description
\mathcal{D}	The stream of data
m_C	A mature cluster
im_C	An immature cluster
$m_C = \{m_C\}$	The set of all mature clusters
$\mathbf{c} = (c_j)_{j=1, \dots, d}$	The center of a cluster
r	The radius of a cluster
θ	The boundary of a cluster
λ	Decay control parameter
$\mathbf{x} = (x_j)_{j=1, \dots, d}$	A sample
$\mathbb{Y}_x = \{y\}$	The label set of \mathbf{x}
$\hat{\mathbb{Y}}_x$	The predicted label set for \mathbf{x}
\mathbf{y}	The label set
W	The mature weight of a cluster
W_0	The threshold of mature weight
$\mathbf{p} = (p(l))_{l=1, \dots, \mathbf{y} }$	The label distribution of a cluster
$\mathcal{K}(\mathbf{x})$	K -nearest mature clusters of \mathbf{x}
h	The number of predicted labels
z	Label cardinality

2. Background

2.1. Multi-label classification

Let \mathbb{X} denote the input space and $\mathbb{Y} = \{y_i | i = 1, \dots, M\}$ denote the label set. The purpose of a multi-label learning task is to search for a mapping function f from input space \mathbb{X} to output space $2^{\mathbb{Y}}$ so that each sample $\mathbf{x} \in \mathbb{X}$ is assigned with a subset of the output space. This is a generalization of the traditional multi-class classification problem in which each sample is associated with only a single label.

MLC algorithms can often be categorized into two approaches: problem transformation and algorithm adaptation [4]. In the first category, the MLC problem is transformed into some well-established learning problems such as binary classification. Two common approaches in this category are the Binary Relevance (BR) and Classifier Chains (CC) approaches, where a multi-label task is transformed into M binary classification tasks. The difference between BR and CC is that BR treats the labels independently in the learning process as each binary classification task is associated with a label in the label set. Meanwhile, CC creates the new training set for each binary problem by appending each instance with binary values that indicate which of the previous labels were assigned to that sample [4]. CC therefore has the advantage over BR of addressing label correlation. In practice, an ensemble of CC classifiers are generated via random orders over the label space instead of using a single CC to overcome the issue of label ordering in the chain. Several methods have also been introduced to improve CC's effectiveness, such as replacing binary values by probabilistic outputs [10] and using recurrent neural networks focusing only on positive labels as an extension of probabilistic CC [11]. Kumar et al. [12] improved PCC by using beam search, a classical heuristic search algorithm, so that instead of evaluating on 2^M possible labellings, only bM combinations need to be assessed (b is the beam width). The search also integrate to the learning algorithm to obtain the best order of labels. Ghamrawi and McCallum [13] modeled the dependencies between labels by constructing a graphical model to parameterize the pairwise relationships of feature-label, label-label, and feature-label-label triple. The Label Powerset (LP) is another popular method in this category which treats each different combination of labels as a single label [14]. Although LP can capture the label correlations in the learning model, it has high-complexity in training due to the exponential increase of the number of label combinations with the number of labels. LP is also unable to predict the label combinations that do not appear in the training set [6]. Read et al. [15] proposed the Pruned Set (PS) method which removes the samples belonging to the infrequent label sets to reduce the number of label combinations. Other MLC approaches considered the subsets of labels in a random way such as Random k -Labelsets [14], or in a deterministic way like in dependency network [16].

The algorithm adaptation approach is a group of methods that are adapted from multi-class classification algorithms to solve the MLC problem. Several methods can be mentioned, for example, k -Nearest Neighbor for MLC [17], Support Vector Machine for MLC [18], and Decision Tree for MLC [19].

In the era of big data, recent research on multi-label learning mainly focuses on dealing with large-scale multi-label data, especially on data with a large label set and data that come in the form of a stream. Ubaru and Mazumdar [20] used group testing and coding techniques to compress the label set to reduce the label dimension. Kapoor et al. [21] used compressed sensing in the Bayesian framework for label dimension reduction. SVD techniques was used to project the label vector onto a low dimensional space to reduce its dimension [22]. Besides, the performance of MLC systems can be enhanced by selecting an optimal subset of features to learn the MLC model. Some examples of feature selection method for MLC are feature rank-based stream feature selection method [23] and scalable relevance evaluation feature selection that measures feature dependency [24]. Several expansions of MLC for multi-dimensional classification (MDC, also known as multi-output classification) which is viewed as a general case of MLC where each label can take a number of discrete values. This setting increases the search space of the chain-sequences, resulting in costly testing time. Read et al. [25] proposed the classifier trellis to effectively capture the label correlation in MDC by sequentially placing the labels to the pre-defined trellis structure for the underlying graphical model. Read et al. [26] introduced double Monte Carlo optimization technique to search for the best classifier chain in the training phase and best label vector for the test sample for MDC.

2.2. Multi-label classification for data stream

The large volumes and the rapid growth of data have posed many challenges for traditional offline machine learning systems. First, it is often impractical or even infeasible for a learning algorithm to learn on the entire dataset at once. The newly arrived data also often make the model learned on the old data outdated, causing the degradation of the system performance. Although we can re-train the learning model with the arrival of new data, the re-training process on the continuously arriving data will consume much more time and resource. Incremental learning methods in which the learning model can be updated on-the-fly from the data stream are therefore highly desirable.

The characteristics of data stream present unique challenges to the design of learning algorithms. Bifet and Gavaladà [27] defined four characteristics of learning on data stream: (1) the model must be ready to make prediction on any sequentially arriving samples, (2) there are potentially infinitely many samples which need to be processed with finite resources (time and memory), (3) the samples need not be statistically stationary (the appearance of concept drift), and (4) samples can only be processed one at a time, and can only be inspected once before it is discarded. In this study, we aim to develop an incremental OMLC method in which the current learning model trained on the old data is used to predict unlabelled data. Only samples where true label can be revealed is used to update the learning model. As a

result, both prediction and training are considered in the proposed method. We also address the challenges of data stream's characteristics in designing the proposed method.

One of the earliest approaches that solves the MLC problem in the stream context is the batch-based incremental method [28]. In this method, an ensemble of BR-based classifiers is generated by learning the BR on each sequence of same-size-chunks. The outputs of these classifiers are concatenated to the original feature space as the meta-data which is learned by another BR to obtain the meta-classifier. Based on the dynamic classifier ensemble approach, the classifiers are weighted on each test sample via their performance on the test sample's neighbors obtained from the latest chunk. A similar approach was introduced by Wang et al. [29] in which the data stream is divided into many fixed numbers of chunks. Multi-label K Nearest Neighbor method [17] is then used to learn on these chunks to generate the ensemble of classifiers. These classifiers are also weighted and the weights are incrementally maintained based on the newly arrived chunks. Despite the ability to operate in both the stationary and concept drift settings, these methods face the memory-fill-up problem so they cannot satisfy the time and memory constraints of stream learning [5].

Read et al. [5] adapted the Hoeffding tree [30], a well-known member of the decision tree family, for MLC problem on data stream. In their method, the arrived samples are temporarily kept and the Hoeffding bound is used to determine how many samples are needed to achieve a certain level of confidence for tree splitting. The PS classifier [15] was used to prune the label combination at each leaf node when the buffer of arrived samples is full at the node. That framework was also combined with ADWIN [31] to form a new algorithm named EaHTps which can handle concept drift. The Pruned Set-based label combination module of EaHTps was improved by Shi et al. [32] in which the new frequent label combinations are dynamically recognized to update the set of label combinations. Xioufis et al. [9] used BR to solve the MLC by transforming the multi-label task into several binary classification tasks. Concept drift was handled by maintaining two variable-size windows per label for positive and negative samples. Compared to the single window approach for each label, that method can oversample the positive sample by adding the previous positive samples to the associated window and undersample the negative sample by keeping only the most recent ones. Shi et al. [33] created the super-label which is a set of class labels grouped based on their dependencies. The generated super-labels are treated as new class labels to annotate each arrived sample. To handle concept drift, the authors first measured the distribution of features and new class labels by a multi-label entropy approach. The change of the distribution was then monitored by the difference of the entropy measure between the two windows that keep the old and the recent samples. These windows are resized when the difference exceeds a pre-defined threshold. Osojnik et al. [6] applied multi-target regression to multi-label learning on data stream. In their approach,

the MLC problem and the multi-target regression problem are transformed back and forth as the solution is obtained by transforming the multi-target regression problem back to the MLC problem to obtain the predicted labels. The four multi-target tree-based methods introduced in [6], however, are only focused on learning stationary concept. A comparative study between several MLC methods for data stream such as BR with different learning algorithms and multi-label Hoeffding tree (with PS and Naïve Bayes at the leaves) was made by Karponi and Tsoumakas [34]. However, the experiments were conducted on just one reduced dataset, and therefore the conclusion is not convincing.

3. Proposed method

3.1. The online clustering model

The change in data often makes the learning model that is built on old data inconsistent with the new data, therefore it is crucial to adapt a learning model to concept drift. In the literature, there are two main strategies to deal with concept drift, namely, using a sliding window or a decay function (also called weighted samples [8]) [31]. In the first strategy, the windows are maintained that keep the most recent samples inside while discarding the old samples outside. Several methods use a single window of a fixed size or of variable size. For example, in ADWIN, a well-known concept drift handling method, the window will grow or shrink depending on whether data changes or not. The improved version, ADWIN2, is more efficient than the first one in time and memory consumption [31]. Other research such as [32, 33] use two adjustable windows to represent old and new samples. Xioufis et al. [9] used two windows per label to capture the positive and negative groups of samples. The decay function approach, meanwhile, use weights to mitigate or strengthen the importance of samples based on their age, i.e., the older sample is less important than the new one [35]. The learning model, therefore, can adapt to the changes that appear in the new data.

In this study, we follow the decay function approach to develop an incremental online MLC method that adapts to the changes of data in the data stream. Our approach aggregates the information from incoming samples into clusters based on their proximity with each other as well as their time of arrival. We require that the weight of each sample decreases gradually over time. The incremental MLC, therefore, focuses more on the new samples than the old ones. Here the weight of each sample is decayed *exponentially* with time t via the fading function $f(t) = 2^{-\lambda t}$ where $\lambda > 0$ is the parameter that controls the decay rate [7]. From these data points and their weights, we build the online clustering model for the incremental online MLC method. First, we define the online clustering model:

Definition 1 (Mature cluster) [7]: A mature cluster C (m_C) at time t for the group of close points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, $\mathbf{x}_i = (x_{i1}, \dots, x_{id})$ with time stamps T_1, T_2, \dots, T_n is defined as the 3-tuple $\{\mathbf{CF}^1, \mathbf{CF}^2, W\}$ in which the weight $W = \sum_{i=1}^n f(t - T_i) > W_0$, W_0 is a pre-defined mature threshold, $\mathbf{CF}^1 = (CF_j^1)$, $CF_j^1 = \sum_{i=1}^n f(t - T_i) x_{ij}$ is the weighted linear sum of the j^{th} feature, $\mathbf{CF}^2 = (CF_j^2)$, $CF_j^2 = \sum_{i=1}^n f(t - T_i) x_{ij}^2$ is the weighted squared sum of the j^{th} feature. The center of m_C is $\mathbf{c} = (c_j)$, $c_j =$

$$\frac{CF_j^1}{W}, \text{ the radius of } m_C \text{ is } r = \max_{j=1, \dots, d} \left\{ \sqrt{\frac{CF_j^2}{W} - \left(\frac{CF_j^1}{W}\right)^2} \right\} \leq \theta$$

Definition 2 (Immature cluster) [7]: A immature cluster C (im_C) at time t for the group of close points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, $\mathbf{x}_i = (x_{i1}, \dots, x_{id})$ with time stamps T_1, T_2, \dots, T_n is defined as $\{\mathbf{CF}^1, \mathbf{CF}^2, W\}$. The definitions of $\mathbf{CF}^1, \mathbf{CF}^2, W, \mathbf{c}$, and r are the same as those of m_C , however $W \leq W_0$.

From the definition 1 and 2, a cluster (mature or immature) is a set of data point bounded by the pre-defined threshold θ . The center of a cluster is the weighted average of all data points inside the cluster whereas the radius is the maximum value among all standard deviations of the d features. We illustrate an example of the cluster evolution from immature to mature cluster on the EURON dataset in Fig 1. In this example we set $\lambda = 0.25, W_0 = 2$, and $\theta = 0.495$. At time $t = 0$, a cluster is generated with a newly arrived sample \mathbf{x}_{47} with time stamp 0. This is an immature cluster because its weight is 1 (equal to the weight of \mathbf{x}_{47}). At $t = 2$, a new sample \mathbf{x}_{298} arrives to the cluster with the weight 1 and time stamp 2. Meanwhile, the weight of \mathbf{x}_{47} is $1 \times 2^{-0.25 \times 2} = 0.7071$. As the total weight of this cluster is smaller than W_0 , the cluster is still an immature cluster. At $t = 3$, we have a new sample \mathbf{x}_{351} which arrives with time stamp 3. The weights of \mathbf{x}_{47} and \mathbf{x}_{298} are 0.5946 and 0.8409, respectively. That makes the weight of the cluster greater than W_0 and the cluster becomes a mature cluster. It is noted that in all cases, the radius of the cluster must be smaller than the pre-defined maximum threshold of radius $\theta = 0.495$.

In this model, the mature cluster is trusted more than the immature cluster so that only the mature clusters will be used during the classification process. It is noted that definitions 1 and 2 is for unsupervised learning, i.e. the samples do not have labels. For MLC, each sample in a cluster is associated with a set of labels. We extend the mature and immature cluster for MLC as:

Definition 3 (Cluster for MLC): A cluster C (mature or immature) at time t for the group of close points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, $\mathbf{x}_i = (x_{i1}, \dots, x_{id})$ with time stamps T_1, T_2, \dots, T_n is defined as the 4-tuple $\{\mathbf{CF}^1, \mathbf{CF}^2, W, \mathbf{p}\}$

in which $\mathbf{CF}^1, \mathbf{CF}^2, W, \mathbf{c}$, and r are defined as in Definition 1 and $\mathbf{p} = (p(l))_{l=1, \dots, |\mathbf{y}|}$ is the label distribution of the points.

In this definition, along with the components $\mathbf{CF}^1, \mathbf{CF}^2$, and W , we store the label distribution as the fourth component of each cluster. The distribution is approximated by the label frequency of all data points inside:

$$p(l) = \frac{\sum_{\mathbf{x} \in \mathcal{C}} \llbracket l \in Y_{\mathbf{x}} \rrbracket}{|\mathcal{C}|} \quad (1)$$

in which $\llbracket \cdot \rrbracket$ returns 1 if the predicate holds and 0 otherwise.

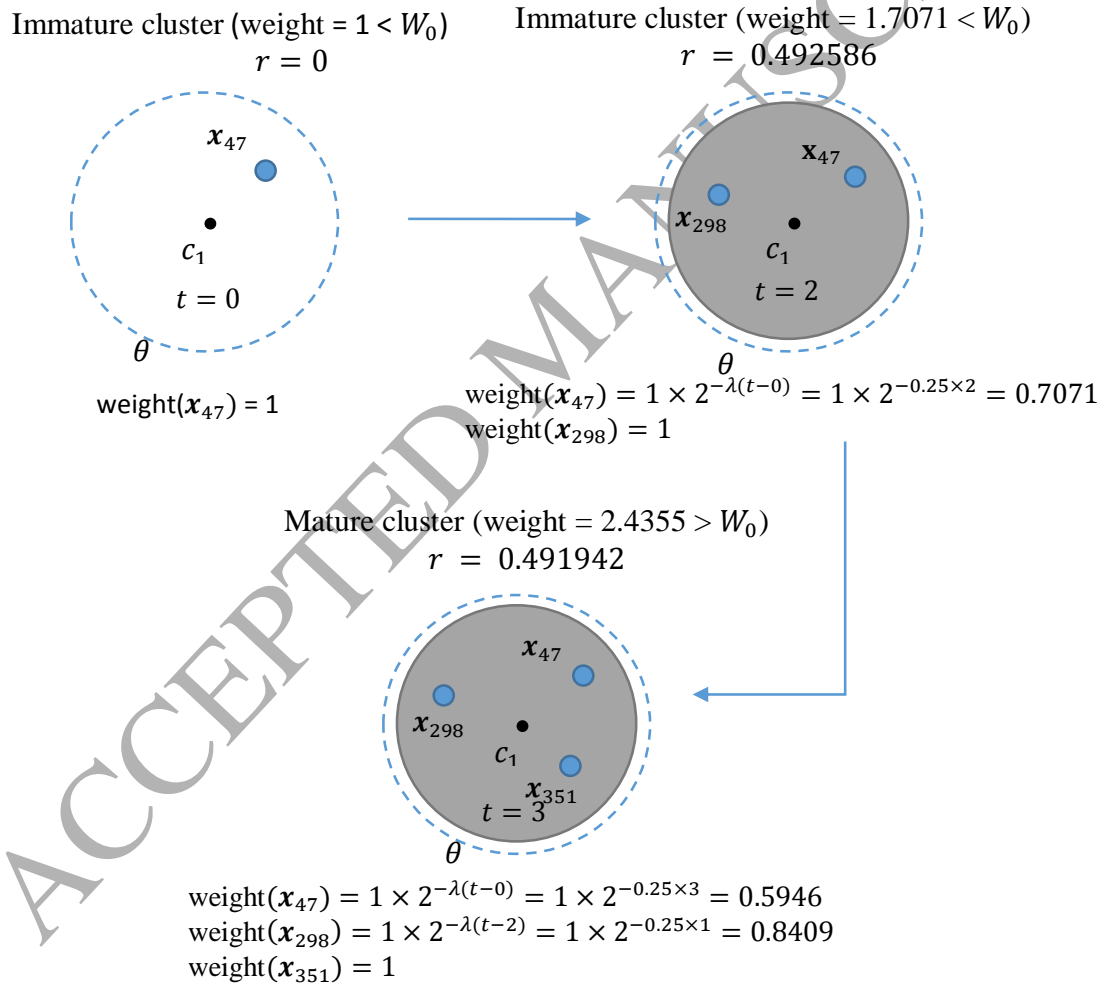


Fig.1. An illustration of cluster evolution on the EURON dataset

3.2. The model update

From the arrived data points, we incrementally learned the clustering model to solve the MLC problem. In this section, we introduce a method to construct the clustering structure for the MLC problem. The incremental process that maintains the clustering structure include:

- **Update the cluster's components:** Due to the decay mechanism over the data points, the components of a cluster are updated over time, even if it does not receive any new data. In this work, we follow the update equations introduced in [7]. First, when the new data point $\mathbf{x} = (x_j)$ (its weight is 1) is merged into a cluster C , the new value of the first three components are computed as:

$$\{(CF_j^1)^{new}, (CF_j^2)^{new}, W^{new}\} = \{(CF_j^1)^{old} + x_j, (CF_j^2)^{old} + x_j^2, W^{old} + 1\} \quad (2)$$

Whereas, after the time interval δt , if the cluster C does not receive any new points, its components will be decayed due to the fact that the weight of all data points decreases by a factor of $2^{-\lambda\delta t}$ in the interval:

$$\{(CF_j^1)^{new}, (CF_j^2)^{new}, W^{new}\} = \{2^{-\lambda\delta t}(CF_j^1)^{old}, 2^{-\lambda\delta t}(CF_j^2)^{old}, 2^{-\lambda\delta t}W^{old}\} \quad (3)$$

Based on the new values of these components, the new center and radius are updated respectively:

$$(c_j)^{new} = \frac{(CF_j^1)^{new}}{W^{new}} \quad (4)$$

$$r^{new} = \max_{j=1, \dots, d} \left\{ \sqrt{\frac{(CF_j^2)^{new}}{W^{new}} - \left(\frac{(CF_j^1)^{new}}{W^{new}}\right)^2} \right\} \quad (5)$$

The label frequency of the cluster that the data point is merged to is updated by using the ground truth labels if available. The incremental update of the label frequency is given by:

$$p^{new}(l) = \frac{p^{old}(l) + \mathbb{I}\{l \in \mathbb{Y}_x\}}{|C| + 1} \quad (6)$$

- **An immature cluster can become a mature cluster:** When an immature cluster receives a new data point, its weight is increased by 1 (see Eqn. (2)). When the weight of the immature cluster is greater than the pre-defined mature threshold W_0 , this cluster will become a mature cluster.
- **A new immature cluster can appear:** When a new data point cannot be merged into any cluster as it makes the radius of the cluster exceed the boundary threshold θ , it becomes the first data point of a new immature cluster. In this case, the new cluster only has one data point which is also the center of this cluster. The radius of this cluster is 0 and the weight is 1.
- **A mature cluster can become an immature cluster:** Over time, if a mature cluster does not get any new samples, its weight will gradually decrease as a consequence of the weight decay of all data points inside the cluster. The mature cluster will become an immature cluster when its weight is smaller than or equal to the threshold W_0 . In this case, we need to check the weight of all mature clusters to detect

the change. Assume that after the time span T_s , the weight of the mature cluster is $2^{-\lambda T_s}W$ ($W > W_0$) (see Eqn. 2), and a mature cluster become an immature cluster, we obtain the inequality $2^{-\lambda T_s}W_0 < 2^{-\lambda T_s}W \leq W_0$. Thus, the minimal time span for a mature cluster to become an immature cluster is $T_s = \left\lceil \frac{1}{\lambda} \left(\log \left(\frac{W_0}{W_0 - 1} \right) \right) \right\rceil$ which is computed from the equation $2^{-\lambda T_s}W_0 + 1 = W_0$. Hence, we only periodically check the mature clusters for every time period T_s [7].

Algorithm 1 summarized the update process to build the incremental online clustering model for the MLC problem. For each arrived sample, we try to merge it into the nearest mature cluster. If the new radius is still smaller than the boundary of a cluster, the sample is successfully merged into this mature cluster (step 2-3). The three components, as well as the cluster center and radius, are updated using the features and weight of the new member (step 4). In contrast, if the new sample makes the mature cluster exceed the boundary, we try to merge this sample to the nearest immature cluster (step 7-9). In this case, we check the weight of the immature cluster that the sample merged into. If the new weight is larger than the pre-defined mature threshold, the immature cluster will become a mature cluster (step 10). If the sample cannot be merged into any cluster as it makes the cluster's radius exceeds the boundary threshold, we build a new immature cluster for this sample (step 14-15). To complete the update process, we incrementally update the label frequency component of the cluster that the sample is merged into (step 18-20). Finally, we periodically check all clusters in the mature list. If there exist any clusters with weights smaller than or equal to the mature threshold, these clusters will become immature clusters (step 23-25).

Algorithm 1: Incremental clustering model for MLC

Input: Arrived sample \mathbf{x} at the current time t

Output: The updated incremental learning model

- 1 Try to merge \mathbf{x} into the nearest m_C
- 2 If (the new r (of m_C) $\leq \theta$)
- 3 Merge \mathbf{x} into m_C
- 4 Update r , \mathbf{c} , and W_{m_C} of m_C
- 5 Else
- 6 Try to merge \mathbf{x} into the nearest im_C
- 7 If (the new r (of im_C) $\leq \theta$)
- 8 Merge \mathbf{x} into im_C
- 9 Update new r , \mathbf{c} , and W_{im_C} of im_C
- 10 If ($W_{im_C} > W_0$)
- 11 $im_C \rightarrow mC$
- 12 End
- 13 Else
- 14 Create a new immature cluster using \mathbf{x}
- 15 Compute r , \mathbf{c} , and W for the new cluster
- 16 End
- 17 End
- 18 For each label i in \mathcal{Y}
- 19 Update label frequency by (6)
- 20 End
- 21 Set $T_x = \left\lceil \frac{1}{\lambda} \left(\log \left(\frac{W_0}{W_0 - 1} \right) \right) \right\rceil$
- 22 If $(t \bmod T_x) = 0$
- 23 For each m_C in m_C
- 24 If ($W_{m_C} \leq W_0$)
- 25 $m_C \rightarrow im_C$
- 26 End
- 27 End
- 28 End

3.3. The classification process

Based on the constructed clustering model, we predict the labels for each newly arrived sample. As mentioned above, we trust the mature cluster more than the immature cluster so we only use the information of the mature clusters to assign labels for each newly arrived sample. In this study, we develop a classification method based on the similarity between the sample and the mature clusters. First, we measure the distances between the sample and all mature clusters. The top K mature clusters with the shortest distance are selected as the K -nearest neighbors of that sample. We then use the label frequencies of these K nearest neighbors to compute the posterior probability that the sample belongs to a class label. The top h labels associated with the largest posterior probabilities are assigned to the sample as the prediction result.

In detail, for sample \mathbf{x} , let $\mathcal{K}(\mathbf{x})$ represent the set of its K -nearest neighbors in $m_{\mathcal{C}}$. Generally, the similarity between the sample and the mature cluster $m_{\mathcal{C}}$ is measured with the Euclidean distance between the samples and the clusters' center.

$$\mathcal{K}(\mathbf{x}) = \{m_{\mathcal{C}_i}, 1 \leq i \leq K | i \in \arg \text{smallest}_k\{d(\mathbf{x}, \mathbf{c}_i)\}\} \quad (7)$$

Here $\text{smallest}_k(\cdot)$ returns the K shortest distances between \mathbf{x} and the mature clusters in $m_{\mathcal{C}}$ and $d(\mathbf{x}, \mathbf{c}_i)$ is the distance between \mathbf{x} and the cluster center \mathbf{c}_i . We calculate the posterior probability that the sample belongs to a label l as the sum of the products of the weight of $m_{\mathcal{C}_i}$ and label frequencies of label l among all mature clusters in $\mathcal{K}(\mathbf{x})$:

$$P(l|\mathbf{x}) \sim \frac{1}{k} \sum_{m_{\mathcal{C}_i} \in \mathcal{K}(\mathbf{x})} p_{m_{\mathcal{C}_i}}(l) W_{m_{\mathcal{C}_i}} \quad (8)$$

Let $\hat{\mathcal{Y}}_{\mathbf{x}}$ denote the predicted label set for \mathbf{x} . We obtain $\hat{\mathcal{Y}}_{\mathbf{x}}$ by getting the labels associated with the h -largest values of these posterior probabilities:

$$\hat{\mathcal{Y}}_{\mathbf{x}} = \{l \in \mathcal{Y} | P(l|\mathbf{x}) \text{ in top } h\} \quad (9)$$

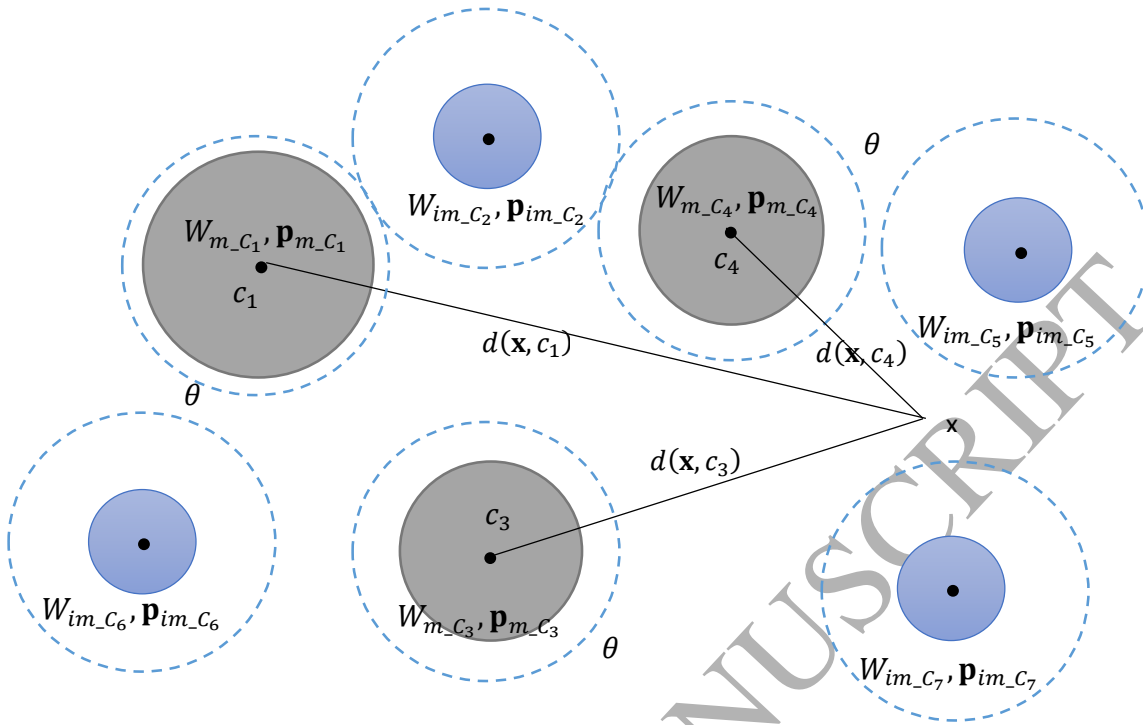


Fig.2. The classification process based on the clustering model

The proposed MLC algorithm is different to the K -nearest neighbor approach [17] and weighted K -nearest neighbor [36] approach in label assignment, although they all used the K -nearest neighbors for label assignment. In weighted K -nearest neighbor [36], each class label is assigned to the sample via the decision function computed from the posterior probability that the sample belongs to the considered label (positive label) or not (negative label) in the K neighboring samples. In our method, clustering is performed, then each mature cluster is treated as a single point having its mature degree and the distribution of label frequencies. Instead of selecting a label based on the decision function computed from the label's posterior probability of positive or negative label, we select h labels associated with the largest posterior probability, where h is learned from the data sequence. In our method, the number of predicted labels is updated adaptively based on the revealed ground truth labels of the arrived samples. The update procedure for h will be introduced in the next section.

The proposed classification process is summarized in Algorithm 2.

Algorithm 2: Predict labels based on weights of clusters

Input:	Sample \mathbf{x} , the set of mature cluster $m_{\mathcal{C}}, h, K$
Output	Predicted labels for \mathbf{x}
1	For each $m_{\mathcal{C}}$ in $m_{\mathcal{C}}$
2	Compute $d(\mathbf{x}, \mathbf{c})$
3	End
4	Select $\mathcal{K}(\mathbf{x})$ by (7)
5	Initialize $P(l \mathbf{x}) = 0 \forall l \in \mathbf{y}$
6	For each $l \in \mathbf{y}$
7	For each $m_{\mathcal{C}}$ in $\mathcal{K}(\mathbf{x})$
8	Compute $P(l \mathbf{x})$ by (8)
9	End
10	End
11	Return $\hat{\mathbf{Y}}_{\mathbf{x}}$ by (9)

3.4. Label set learning

We predict the h labels associated with the top h -posterior probabilities for the arrived sample. Very often, the number of labels to be learned is fixed beforehand. In general, however, the number of predicted labels for each sample should be flexible and depends on the sample itself. In this study, we propose a method to adjust the number of predicted labels by using the Hoeffding inequality [37] and the label cardinality.

The label cardinality of a dataset is the average number of labels per sample in the dataset. It is independent of the number of labels M , and naturally can be used to quantify the number of predicted labels h for each sample. In this study, not only the learning model but also the number of predicted labels are learned with the arrival of data samples. Here, we introduce an approach to incrementally learn the value of h in which h is adjusted if it is different from the label cardinality by more than a certain amount as determined by the Hoeffding inequality.

In probability theory, the Hoeffding inequality provides a bound on the probability that the sum of independent random variables will deviate from its expected value by more than a certain amount [37]. We restate the result of the Hoeffding inequality as the theoretical basis of our approach: If $X_i, i = 1, \dots, N$ are independent random variables and if $a_i \leq X_i \leq b_i, S = \sum_{i=1}^N X_i, \bar{X} = S/N, \mu = \mathbb{E}[\bar{X}]$ then for $\varepsilon > 0$

$$\begin{aligned} \text{a) } & P\{\bar{X} - \mu \geq \varepsilon\} \leq \exp\left\{\frac{-2N^2\varepsilon^2}{\sum_{i=1}^N(a_i-b_i)^2}\right\} \text{ and } P\{\mu - \bar{X} \geq \varepsilon\} \leq \exp\left\{\frac{-2N^2\varepsilon^2}{\sum_{i=1}^N(a_i-b_i)^2}\right\} \\ \text{b) } & P\{|\bar{X} - \mu| \geq \varepsilon\} \leq 2\exp\left\{\frac{-2N^2\varepsilon^2}{\sum_{i=1}^N(a_i-b_i)^2}\right\} \end{aligned}$$

Assume that we have a stream with N arrived samples (\mathbf{x}_i, Y_i) $i = 1, \dots, N$, in which the i^{th} sample has $|Y_i|$ labels. Applying the Hoeffding inequality with the note that $0 \leq |Y_i| \leq L$ for all $i = 1, \dots, N$, and L is the number of distinct labels of the data stream, we have:

$$P\{|\bar{z} - \mathbb{E}[\bar{z}]| \geq \varepsilon\} \leq 2\exp\left\{\frac{-2N^2\varepsilon^2}{NL^2}\right\} = 2\exp\left\{\frac{-2N\varepsilon^2}{L^2}\right\} \quad (10)$$

in which the label cardinality as the average number of labels of the stream is given by:

$$\bar{z} = \frac{1}{N} \sum_{i=1}^N |Y_i| \quad (11)$$

Denoting the right hand of Eqn. (10) by δ , ε is computed as:

$$\varepsilon = \sqrt{\frac{L^2 \ln(2/\delta)}{2N}} \quad (12)$$

Eqn. (10) becomes:

$$P\{|\bar{z} - \mathbb{E}[\bar{z}]| \leq \varepsilon\} \geq 1 - \delta \quad (13)$$

If $|h - \bar{z}| > \varepsilon$, based on the inequality $|h - \mathbb{E}[\bar{z}]| = |(h - \bar{z}) - (\mathbb{E}[\bar{z}] - \bar{z})| \geq |h - \bar{z}| - |\mathbb{E}[\bar{z}] - \bar{z}|$ and combine with Eqn. (12), we have:

$$P\{|h - \mathbb{E}[\bar{z}]| > 0\} \geq 1 - \delta \quad (14)$$

That means h is different from $\mathbb{E}[\bar{z}]$ with a probability of at least $1 - \delta$ if $|h - \bar{z}| > \varepsilon$. The key idea of our approach is that we only update h if there is a certain difference between h and the current label cardinality \bar{z} . In this case, we update h by $h = \text{round}(\bar{z})$.

Fig 3 shows the proposed approach to determine the size of the predicted label set. We first initialize a value for h . After receiving $N - 1$ samples in the sequence in which i^{th} sample has $|Y_i|$ true labels, the label cardinality of the stream at the $(N - 1)^{\text{th}}$ sample, denoted by \bar{z}_{N-1} is given by:

$$\bar{z}_{N-1} = \frac{1}{N-1} \sum_{i=1}^{N-1} |Y_i| \quad (15)$$

After predicting for the new N^{th} sample, the label cardinality at the N^{th} sample is updated by:

$$\bar{z} = \frac{1}{N} \sum_{i=1}^N |Y_i| = \frac{1}{N} (\sum_{i=1}^{N-1} |Y_i| + |Y_N|) = \frac{1}{N} ((N-1)\bar{z}_{N-1} + |Y_N|) \quad (16)$$

For each arrived sample, we compute the current label cardinality \bar{z} based on the size of the true label $|Y_N|$ and the previous label cardinality given by (15). We then check whether the difference between \bar{z} and h is greater than the threshold ε , and the value of h used for the next sample will be re-calculated by the rounded value of current \bar{z} . If the update condition is met, we reset $\bar{z} = 0$, $L = 0$, and $N = 0$ to begin a period with the new value of h .

In the literature, we found only three incremental MLC thresholding methods to determine the size of the predicted label set. In all of them, a label is selected if its associated confident score (e.g., posterior probability) is higher than a threshold. The threshold is initialized, and is employed to obtain the predicted labels, and is then re-calculated via sample adjustment [38] or batch adjustment [5, 9]. In [38], Read et al. performed small adjustment on the threshold on a per-sample basis depending on the predicted and the actual label cardinality. Xioufis et al. [9] by contrast incrementally adjusted a threshold for each label on each fixed window of samples so that the frequency of the predicted label approximates that of the ground truth label. In [5], Read et al. incrementally adjusted a threshold for all labels on each batch to ensure that the predicted label cardinality approximates the true label cardinality. Compared to the above methods, our approach is significantly different and is more flexible since the number of samples used to adjust h is not fixed but is based on using a condition derived from the Hoeffding bound.

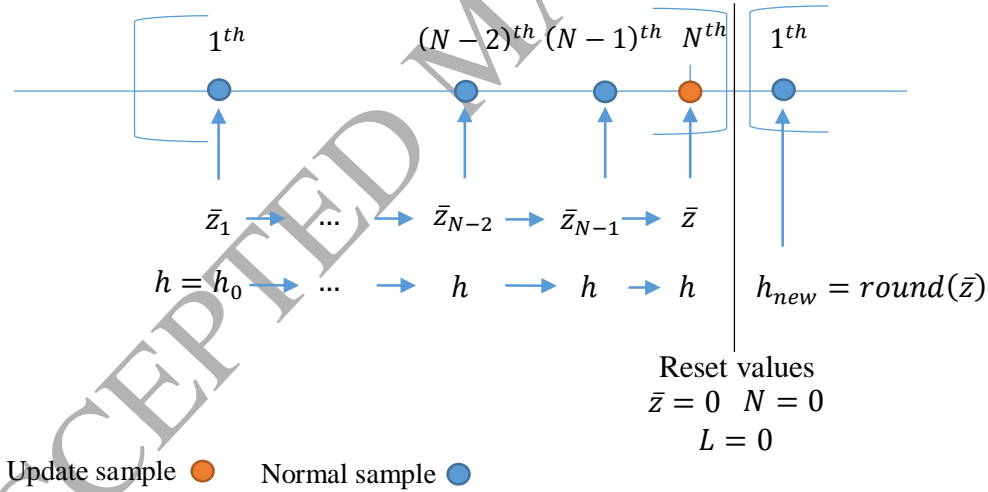


Fig.3. Illustration of adjusting the number of predicted labels

As mentioned before, we developed the incremental MLC method via the online learning paradigm. By this way, we performed the following three steps on each arrived sample \mathbf{x} :

- *Predict labels:* The current learning model is used to predict the label set \hat{Y}_x for \mathbf{x} .
- *Compute the update condition:* We update the learning model if the set of true labels Y_x of \mathbf{x} can be revealed for the environment.

- *Update learning model*: If the update condition is satisfied, the new classification model is obtained by updating from the previous one based on sample \mathbf{x} and \mathbb{Y}_x .

The incremental learning algorithm based on the clustering model for OMLC is summarized as:

Algorithm 3: Incremental online multi-label classification based on a clustering approach

Input:	Data sequence $\mathcal{D}, K, \varepsilon, W_0, \lambda$
Output	Predict labels for each arrived samples
1	For each arrived sample \mathbf{x} from \mathcal{D}
2	$h = h_0$;
3	Obtain $\hat{\mathbb{Y}}_x$ by Algorithm 2
4	If (\mathbb{Y}_x can be revealed from the environment)
5	Update clustering structure by Algorithm 1
6	End
7	Compute \bar{z} by (16)
8	Compute ε by (12)
9	If ($ \bar{z} - h > \varepsilon$)
10	Update $h = \text{round}(\bar{z})$
11	$\bar{z} = 0, N = 0, L = 0$
12	End
13	End

4. Experimental Studies

4.1. Datasets

To evaluate the performance of the benchmark algorithms and the proposed method, we selected 5 popular multi-label datasets [5, 6] for the stationary setting and generated 12 datasets for the concept drift setting, respectively. All these datasets were treated as sequential datasets by processing them in the order they were collected.

[1] We generated a synthetic dataset named *SynRTG* using the Random Tree Generator (RTG) in MOA library (<http://moa.cms.waikato.ac.nz>). The RTG constructs a decision tree by randomly choosing attributes to split and assigning a random class label to each leaf. Once the tree is built, new data points

are generated by assigning uniformly distributed random values to attributes which are then used to determine the class label via the tree. We generated a 5-level tree to create *SynRTG*. To generate the concept drift for *SynRTG* (named *SynRTG-drift*), all generation schemes used in our work are initialized as binary generators with parameters as in Read et al. [5]. For a dataset with N generated samples, the drifts are centred at the $(N/4)^{th}$, $(N/2)^{th}$, and $(3N/4)^{th}$ sample, extending over $N/1000$, $N/100$, and $N/10$ samples, respectively. In the first drift, only 10% of label dependencies are changed. In the second drift, the underlying concepts are changed and more labels are associated on average with each sample (a higher label cardinality). In the third drift, 20% of label dependencies are changed. We also generated the ‘gradual’ drift version for the five multi-label datasets (named *ENRON-drift*, *IMDB-drift*, *OHSUMED-drift*, *SLASHDOT-drift*, and *TMC2007-drift*) by concatenating half of the samples from the original data and the other half generated by RTG. The gradual concept drift is formed by randomly choosing S samples around the joint from the original data and the synthetic data to transition the concept drift. We also generated the ‘break’ drift for the five multi-label datasets (named *ENRON-break*, *IMDB-break*, *OHSUMED-break*, *SLASHDOT-break*, *TMC2007-break*) using the similar scheme mentioned above except with $S = 0$ (Fig 4). The details of the experimental datasets are described in Table 2.

TABLE 2. THE EXPERIMENTAL DATASETS

Dataset	# of samples	# of features	# of labels	Label cardinality	S
ENRON	1702	1001 binary	53	3.4	-
IMDB	120919	1001 binary	28	2.0	-
OHSUMED	13929	1002 binary	23	1.7	-
SLASHDOT	3782	1079 binary	22	1.2	-
TMC2007	28596	500 binary	22	2.2	-
ENRON-drift	3000	1001 binary	53	3.4 -> 5.0	150
IMDB-drift	200000	1001 binary	28	2.0 -> 4.0	5000
OHSUMED-drift	26000	1002 binary	23	1.7 -> 4.0	1300
SLASHDOT-drift	7000	1079 binary	22	1.2 -> 3.0	350
TMC2007-drift	56000	500 binary	22	2.2 -> 5.0	2800
SynRTG-drift	1000000	30 binary	8	1.8 -> 3.0	-
ENRON-break	3000	1001 binary	53	3.4 -> 5.0	0
IMDB-break	200000	1001 binary	28	2.0 -> 4.0	0
OHSUMED-break	26000	1002 binary	23	1.7 -> 4.0	0

SLASHDOT-break	7000	1079 binary	22	1.2 -> 3.0	0
TMC2007-break	56000	500 binary	22	2.2 -> 5.0	0
SynRTG-break	1000000	30 binary	8	1.8 -> 3.0	-

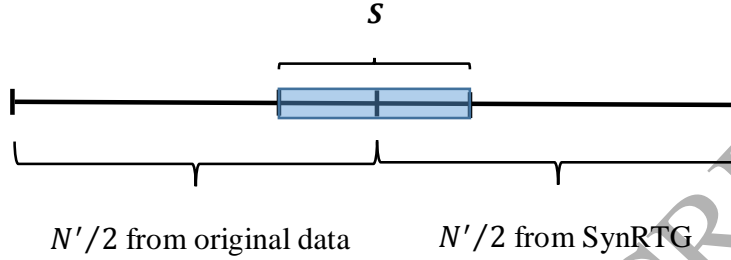


Fig.4. Structure of datasets with concept drift

4.2. Experimental Setup

For each arrived sample, we first predict the labels for each arrived sample and then use this sample with its true labels to update the learning model. For the benchmark algorithms, we use the parameters set by the MOA library. For the proposed method, there are four parameters, i.e., K , θ , W_0 , and λ , that need to be considered. In the next section, we evaluate the influence of K and λ on each of the performance measures. Meanwhile, the mature weight W_0 was set to 3 in the experiment as in [7].

We set the cluster bound θ for binary data since all datasets used in the experiment have only binary features. Assume that on the feature x_j , s samples have value 1 and $n - s$ samples have value 0. The standard variation computed on n samples on this feature is given by:

$$\sigma(x_j) = \sqrt{\mathbb{E}(x_j^2) - \mathbb{E}(x_j)^2} = \sqrt{\frac{s}{n} - \left(\frac{s}{n}\right)^2} = \sqrt{\frac{1}{4} - \left(\frac{1}{2} - \frac{s}{n}\right)^2} \leq \frac{1}{2} \quad (17)$$

Since the inequality (17) holds on all features, we have the upper bound of the radius of the cluster given by $r = \max_{j=1, \dots, d} \{\sigma(x_j)\} \leq 1/2$. If the cluster bound θ is set with a value greater than or equal to 0.5, step 2 in Algorithm 1 will always be satisfied and all samples will be merged into a single mature cluster. Therefore, θ was set to 0.495.

4.3. Performance Measures and Benchmark Algorithms

In multi-label learning, each sample is associated with a set of labels. Zhang and Zhou [4] stated that the performance of MLC algorithms should be tested on a range of measures instead of only the one being optimized to ensure a fair evaluation. In our experiments, we compute six performance measures based on

the predicted label \hat{Y}_x and the ground truth Y_x on each arrived sample x from the data sequence \mathcal{D} . These measures are grouped into three groups: sample-based measures (accuracy and F1), label-based measures (micro F1 and macro F1), and ranking-based measures (average precision and ranking loss). The details of the performance measures can be found in the Appendix. The running time including the time for prediction and update are also reported for all methods.

Since our method is applicable for both the stationary and concept drift settings, we compared the proposed method with both well-known non-adaptive and adaptive incremental MLC algorithms. For the stationary setting, the proposed method was compared with the following well-known incremental MLC algorithms: incremental Classifier Chain (denoted by iCC), incremental Pruned Set (denoted by iPS), and Majority Label Set (denoted by MLS). The base classifier for these methods is the Hoeffding tree. All these benchmark algorithms were run with the default parameters as given in the MEKA library (<http://waikato.github.io/meka>). We also implemented incremental Binary Relevance with SVM as base classifier using the scikit-learn library (<https://scikit-learn.org>) (denoted by iBR(SVM)). The proposed method was also compared with four recent incremental multi-target tree based MLC algorithms introduced in [6] including the multi-target model trees (denoted by iSOUP-MT), multi-target regression tree (denoted by iSOUP-RT), online Bagging for iSOUP-MT (denoted by iSOUP-EBMT), and online Bagging for iSOUP-RT (denoted by iSOUP-EBRT). The parameters for these methods were set as recommended in [6]. The algorithms iCC, iPS, MLS, iBR(SVM) were combined with the state-of-the-art adaptive method named ADWIN2 [31] to adapt to concept drift in the data. However, the four incremental multi-target-tree based MLC algorithms [6] are only introduced for the stationary concept. Therefore, we omitted them from the comparison under the concept drift setting.

4.4. Statistical Test

To assess the statistical significance of the experimental results of incremental MLC methods, Read et al. [5] and Osojnik et al. [6] used the Friedman test [39] to test the difference between the performances of multiple methods on multiple datasets. The Friedman test is preferred over the ANOVA test for the following reasons. First, the Friedman test does not assume normal distribution as in the ANOVA test. Second, an important assumption of repeated-measures ANOVA is sphericity (similar to the requirement that the random variables have equal variance), which cannot be taken for granted because of the nature of learning algorithms and datasets [40]. Here the Friedman test is used to test the null hypothesis that “all methods perform equally”. If the null hypothesis is rejected, a post-hoc test is then conducted. Read et al. [5], and Osojnik et al. [6] used the Nemenyi test for all pairwise comparisons based on the rankings of algorithms on all datasets. The difference in performance of two methods is treated as statistically

significant if the p – value computed from the post-hoc test statistic is smaller than an adjust value of confident level computed from Nemenyi’s procedure. The confident level of the test was set to 0.05.

5. Experimental Results

5.1. The influence of parameters

In this section, we examine the influence of two parameters, i.e., the number of nearest neighbors K and the decay rate λ on the 6 performance measures. The value of λ was set in the range of $\{0, 0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5\}$ to examine its influence on the 6 performance measures. Fig. 5 and Fig. 6 present the relationship between λ and the 6 measures on the 5 datasets at $K = 3$ and $K = 10$, respectively. The up or down arrow beside each measure in Fig. 5 and 6 indicates whether a higher or lower value is better, respectively. It can be seen that the performance measures at $K = 3$ and $K = 10$ with the different values of λ are nearly similar. Therefore, we only analyze the influence of λ on the 6 measures at $K = 3$.

We can observe that the 4 measures, i.e., sample-based F1, sample-based accuracy, micro-average F1, and average precision, remain nearly constant for different values of λ for the datasets OHSUMED, IMDB, and SLASHDOT. However, for the TMC2007 dataset, the value of these 4 measures fluctuates and decreases to a minimum at $\lambda = 0.4$ before slightly increases. For the ENRON dataset, these 4 measures show a fluctuating but generally still increasing trend with the increase of λ . For the TMC2007 dataset, the macro-average F1 decreases with the increase of λ like the 4 measures mentioned above but the decrease here is more significant (decreases from 0.3 to 0.1). The ranking loss for TMC2007 increases slightly when λ is between 0.05 and 0.5. The ranking loss for the SLASHDOT and IMDB datasets also remains nearly constant for different values of λ . For the OHSUMED dataset, the ranking loss decreases to a minimum at $\lambda = 0.15$ and then fluctuates. For the ENRON dataset, the ranking loss becomes better in general with the increase of λ , reaching the best value at $\lambda = 0.5$.

Recall that the decay speed λ determines how fast the model forgets the old samples as reflected by the reduction of the sample’s weight. Hence, this parameter significantly affects the clustering model. A large value of λ would cause a significant reduction of the weights of old samples, reducing the number of mature clusters that would appear. Based on the observations from Fig 5 and 6, there does not exist a common value of λ in which the proposed method would perform well on all the experimental datasets for all 6 measures. Therefore, we set $\lambda = 0.25$ when comparing the performance of the proposed method to the benchmark algorithms.

Fig. 7 presents the relationship between K and the six measures on the 5 non-drift datasets with $\lambda = 0.25$. In this study, the value of K was set to be in the range of $\{3, 5, 7, 10\}$. Generally, there is not a common trend of the performance values with the different values of K among the 5 datasets.

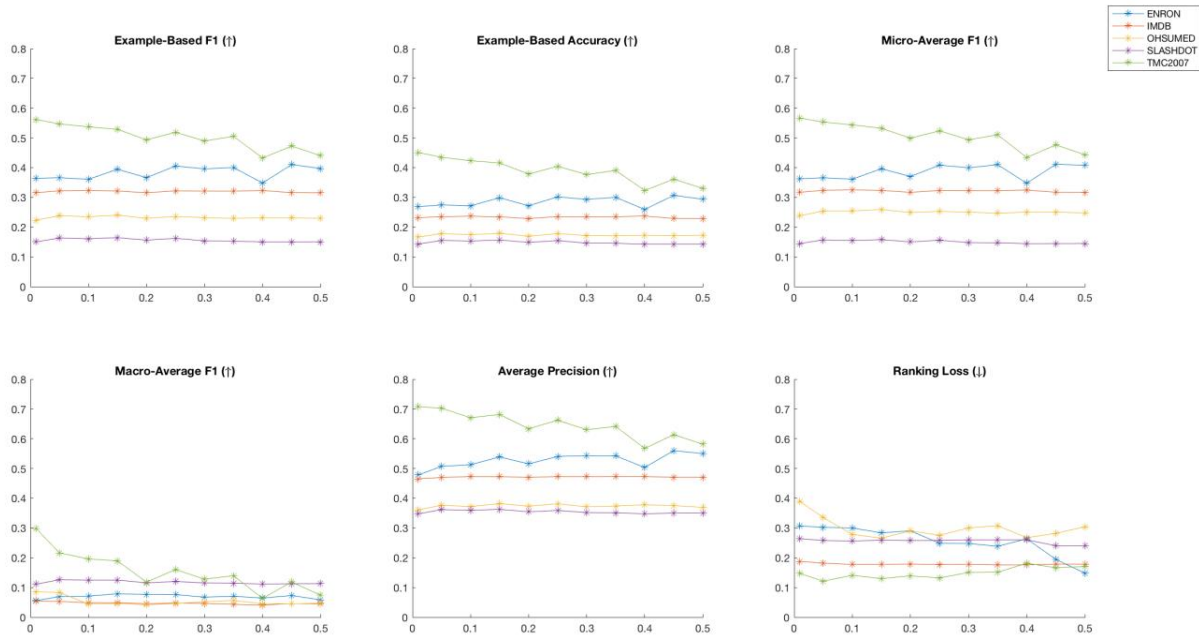


Fig.5. Performance measures (y-axis) for different values of λ (x-axis) at $K = 3$

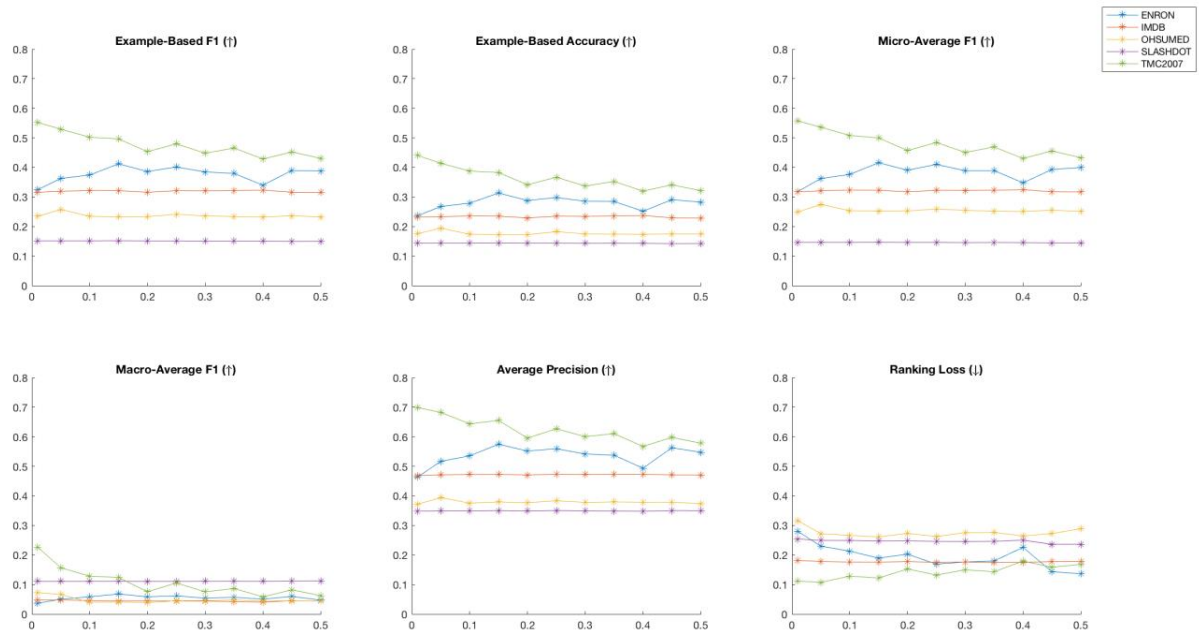


Fig.6. Performance measures (y-axis) for different values of λ (x-axis) at $K = 10$

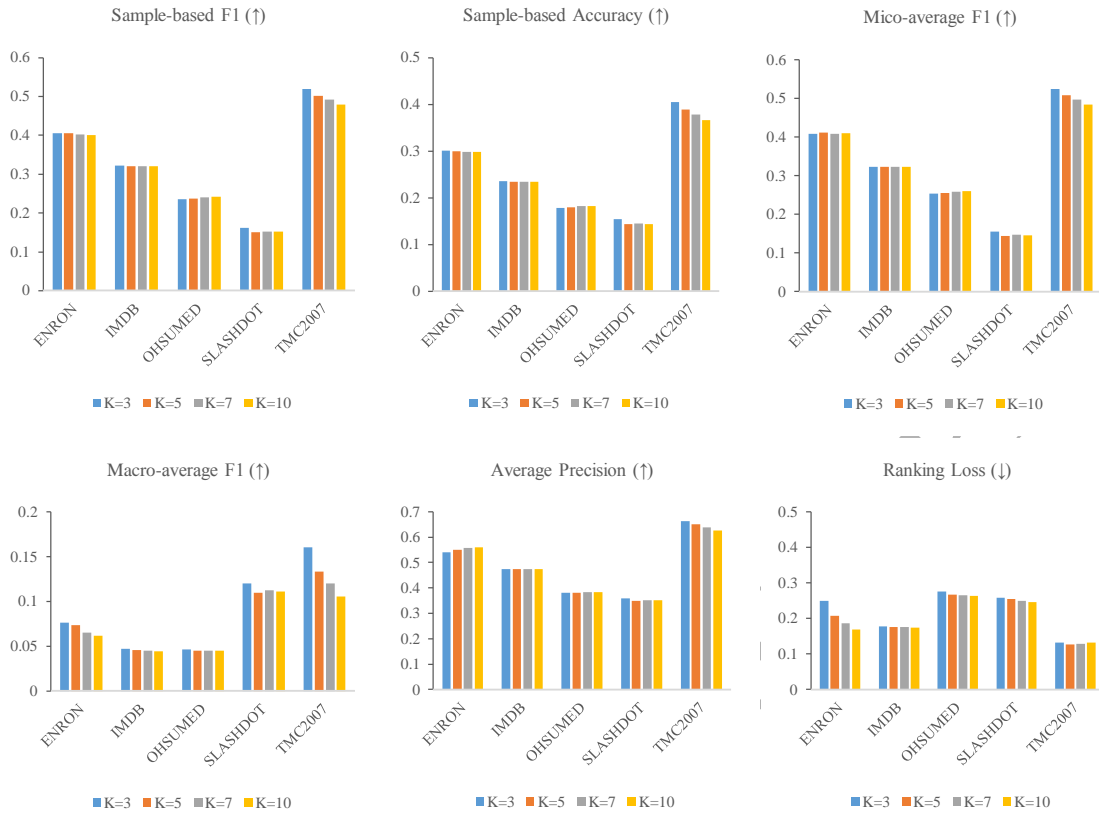
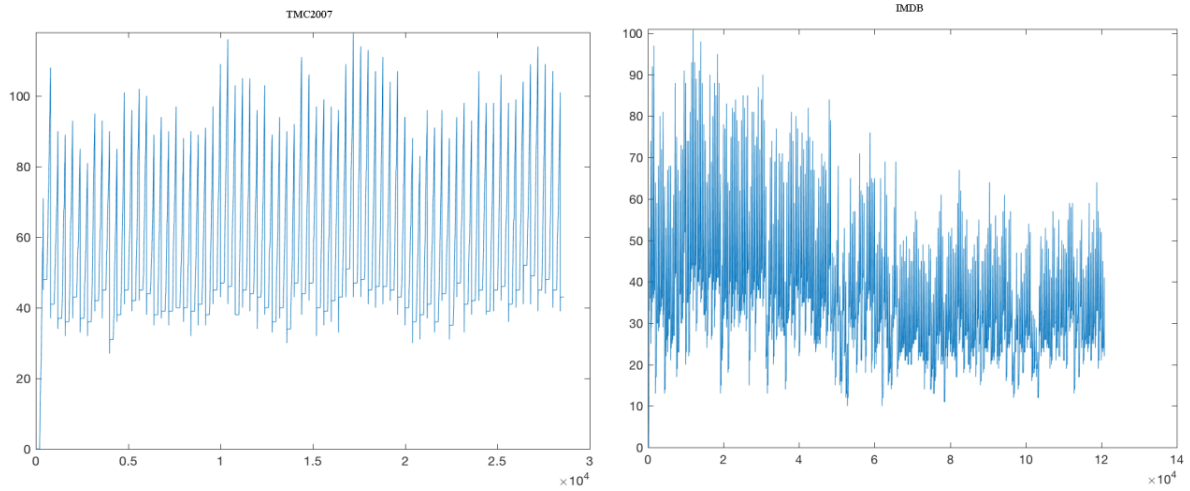


Fig.7. Performance measures for different values of K nearest neighbors



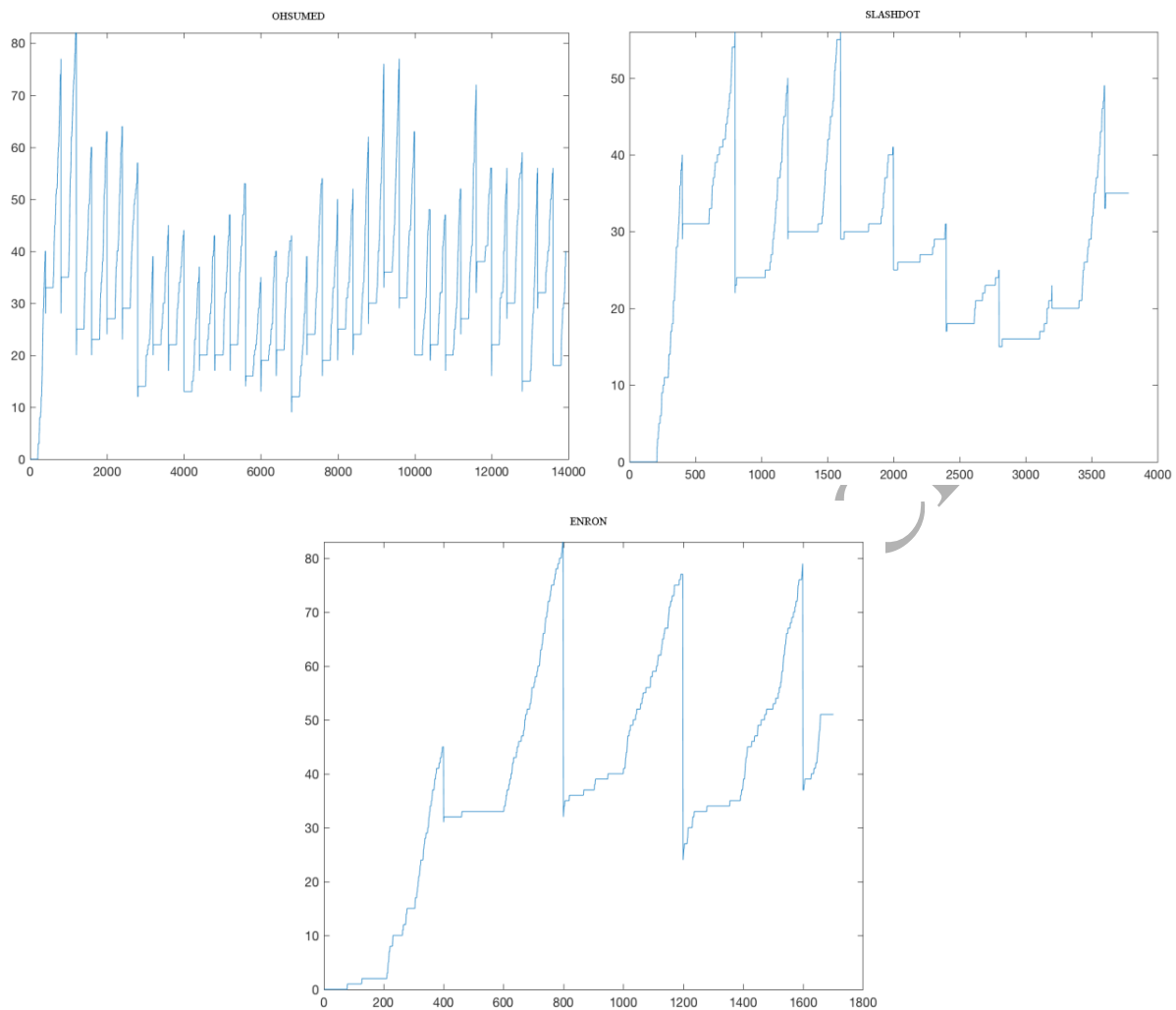


Fig.8. Number of mature clusters (y-axis) generated on the sequence of samples (x-axis) under the stationary setting

In detail, on the ENRON dataset, the different values of K almost do not affect the first three measures, and Macro-average F1 obtains the best value at $K = 3$ whereas average precision and ranking loss obtain the best value at $K = 10$. On the IMDB dataset, all six measures almost remain unchanged with different values of K . On the OHSUMED dataset, K has little influence on 5 measures except for ranking loss in which it obtains the best value at $K = 10$. On the SLASHDOT dataset, the proposed method obtains the best value at $K = 3$ for the first 5 performance measures. In contrast, it performs the poorest on ranking loss at $K = 3$. Finally, on the TMC2007 dataset, on the first 5 measures, the performance of the proposed method reduces with the increase of K . Meanwhile, it performs the best at $K = 5$ for the ranking loss. In the next section, we will compare the performance of the proposed method using $K = 3$.

Fig 8 shows the number of mature clusters generated versus the sequence of arriving samples. It can be seen that the number of mature clusters fluctuates significantly because of the appearance of new mature clusters as well as the change of mature clusters to immature clusters.

5.2. Comparison with weighted KNN

In this section, we compare the effectiveness of the proposed method to the weighted KNN method [36], where we learn the number of labels h for a sample from the data stream and [36] uses a decision function to decide whether a label should be included in the label set of a sample. The experimental results on 5 non-drift datasets are shown in Fig 9. Comparison is only done on the non-drift datasets as the weighted KNN method is a batch learning method and is not designed to handle concept drift in data. Clearly, our method is better than weighted KNN on all datasets for the 4 performance measures: sample-based F1, sample-based accuracy, macro-average F1 and micro-average F1. In particular, our approach is significantly better on 3 datasets IMDB, OHSUMED, and SLASHDOT. For instance, on the IMDB dataset, the sample-based F1 of our method is 0.3220 while that of weighted KNN is only 0.0055. For the average-prediction and ranking loss, our method is only very slightly worse than the weighted KNN method. The result is particularly noteworthy since, in general, batch learning algorithms have better performance than online learning algorithms due to having information available all at once.

The success of the proposed classification method compared to the weighted KNN method originates from the effectiveness of our approach in predicting the number of labels h for each sample. The number of predicted labels h is learned from the true labels of the arrived samples. Here we assigned the top h labels associated with the top h values of the posterior probability to the sample. In our approach, the value of h is adaptive since the number of samples used to adjust h is not fixed, i.e., h is adjusted if the adjust condition based on Hoeffding inequality is satisfied.

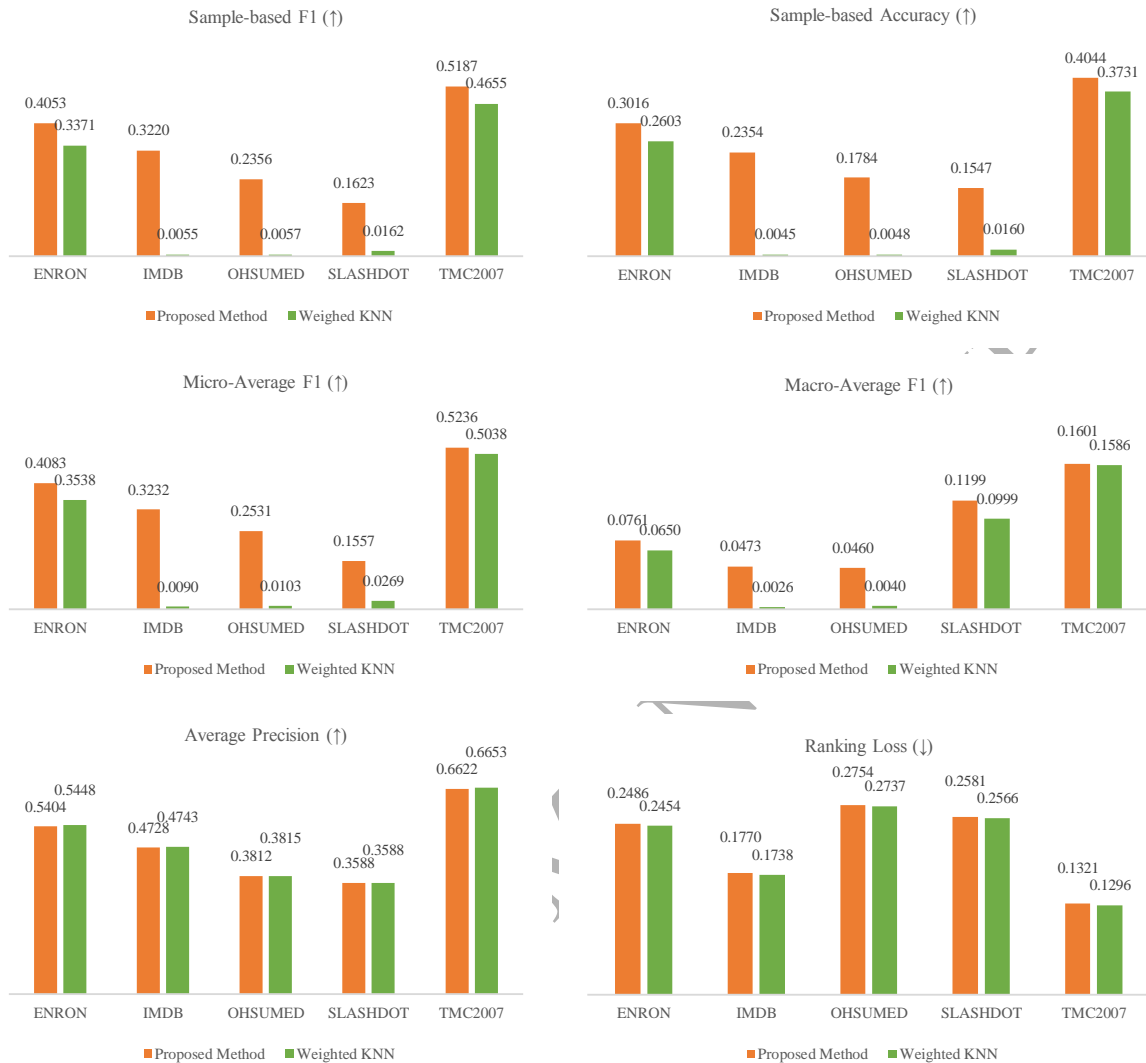


Fig.9. Performance of the proposed classification method and weighted KNN

5.3. Comparative study under the stationarity setting

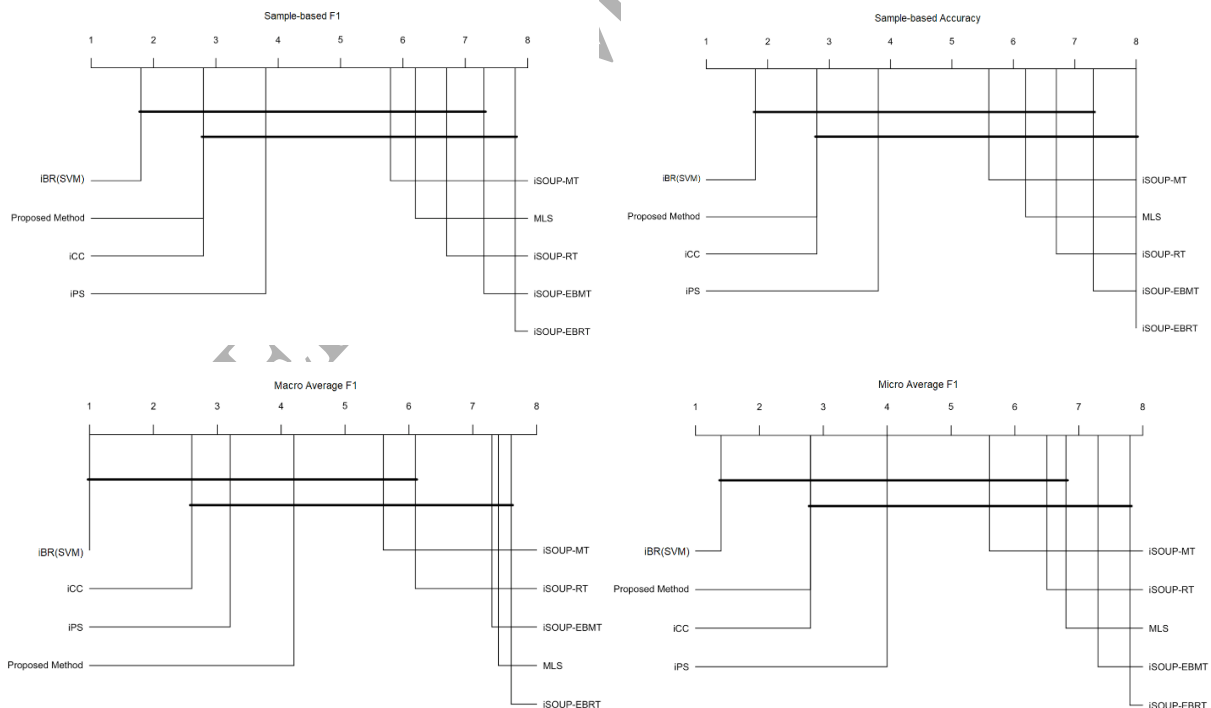
5.3.1. Results on sample-based measures

Table S3 in the supplement material shows the sample-based F1 and sample-based accuracy of the benchmark algorithms and the proposed method. The P-values computed based on the rankings with the Friedman test are $9.3060E-4$ and $7.3760E-4$ for F1 and accuracy, respectively, therefore we rejected the null hypotheses that the performances of all methods are equal. From the Nemenyi significance test results shown in Fig 10, there is a statistical difference in the pairwise comparison between iBR(SVM) and iSOUP-EBRT.

In detail, iBR(SVM) ranks first for both measures (1.8 for both F1 and accuracy), followed by the proposed method and iCC (2.8 for both F1 and accuracy). On dataset IMDB, the proposed method obtains the best results for both F1 and accuracy. iBR(SVM) ranks first on three datasets SLASHDOT, EURON, and OHSUMED. Meanwhile, the four multi-target tree-based methods and MLS obtain poor results on the experimental datasets, ranking at the bottom positions. Especially on two datasets SLASHDOT and IMDB, the four multi-target tree-based methods are significantly poorer than the other methods.

5.3.2. Results on label-based measures

For the two label-based measures, i.e., micro-average F1 and macro-average F1, the P-values for the label-based measures computed by the Friedman test are $4.0171E-4$ and $2.1023E-4$, respectively. Hence, we rejected the null hypotheses and conducted the post-hoc test for all pairwise comparisons among all the methods. On the micro F1 measure, iBR(SVM) ranks first (rank value 1.4), our method and iCC rank second (rank value 2.8), followed by iPS (rank value 4) and iSOUP-MT (rank value 5.6) (see Table S4 in the supplement material). Although on the macro F1 measure, the proposed method is worse than iCC and iPS (rank value 4.2 compared to 2.6 of iCC and 3.2 of iPS), our method continues to obtain better results on ENRON and IMDB. The four multi-target tree-based methods continue to perform poorly on IMDB and SLASHDOT.



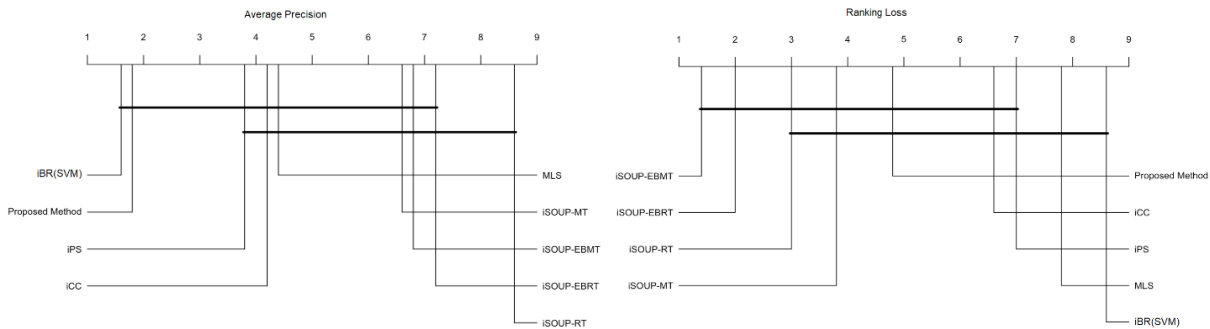


Fig.10. Nemenyi test for the six measures under the stationary setting

5.3.3. Results on ranking-based measures

The performance of the benchmark algorithms and the proposed method for the average precision and ranking loss are shown in Table S5 in the supplement material. We again conducted the Nemenyi post-hoc test and reported the results in Fig 10 for the pairwise comparison. The iSOUP-RT method performs the worst among all methods for the average precision measure, and Nemenyi test shows that it is worse than the proposed method and iBR(SVM) for the average precision measure. For the average precision measure, the proposed method is ranked the second (rank value 1.8) after iBR(SVM) (rank value 1.6), while iPS and iCC are ranked the third and fourth, respectively. The four multi-target tree-based methods continue to perform poorly for the average precision measure, ranking at the four bottom positions.

Surprisingly, the four multi-target tree-based methods perform well on all datasets for the ranking loss measure. The ranking loss of the proposed method is in the middle while iCC, iPS, and MLS obtain poor results (rank values are 6.6, 7, and 7.4 respectively). iBR(SVM) is the poorest among all methods (rank value 8.6). Based on the Nemenyi test, MLS and iBR(SVM) are worse than iSOUP-EBMT and iSOUP-EBRT.

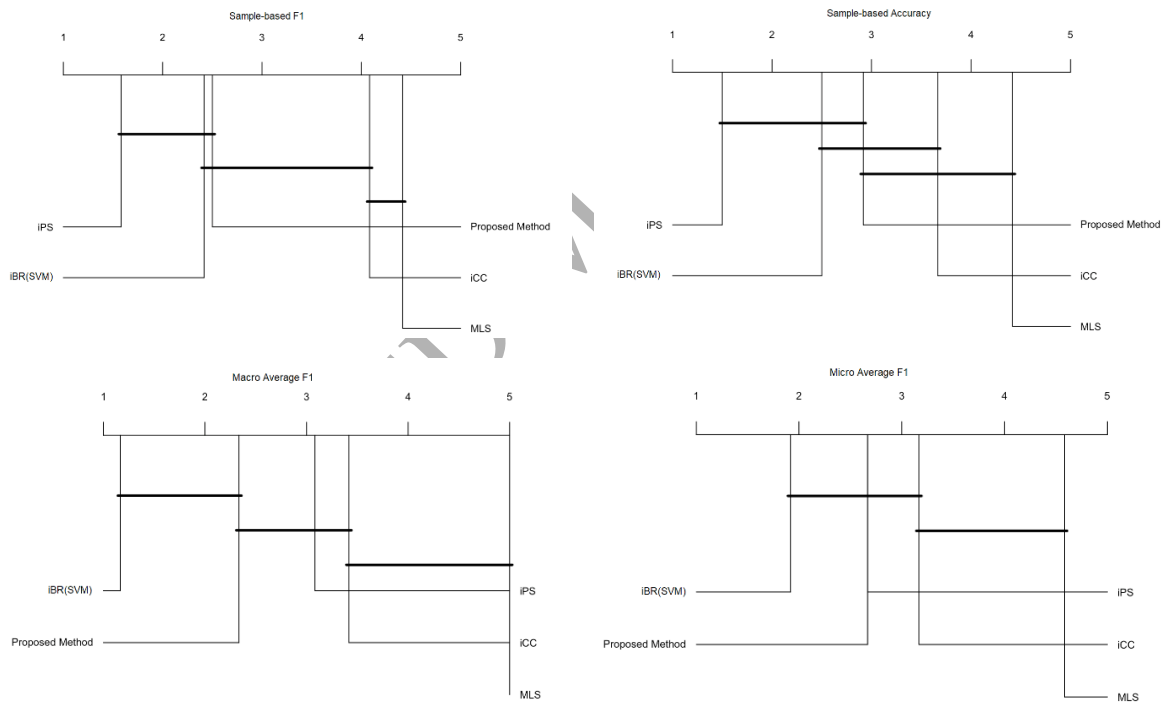
5.4. Comparative study under the concept drift setting

Table S6, S7, and S8 in the supplement material show the experimental results of all methods for the six measures under the concept drift setting. Based on the Nemenyi test, the proposed method is better than MLS on 5 measures except for the sample-based accuracy (Fig 11). The proposed method is also better than iCC on the average precision and better than iBR(SVM) on the ranking loss, while there is no statistical difference between ours and iPS.

In detail, the proposed method ranks first on 2 measures, i.e., average precision and ranking loss, and ranks second on two label-based measures. iPS is a competitive MLC method in handling concept drift

which ranks first on 2 measures, i.e., two sample-based measures and ranks second on 3 measures, i.e., micro-average F1, average precision and ranking loss. Although iBR(SVM) ranks first on two label-based measure, it performs the poorest for the ranking loss. MLS meanwhile is the weakest among all methods on 5 measures.

Unsurprisingly, the benchmark algorithms iPS and iCC obtain better performance on SynRTG-drift and SynRTG-break than the proposed method for the two sample-based measures, i.e., micro-average F1, and average precision. The SynRTG-drift and SynRTG-break are generated based on the decision tree generator, and in the experiment, we used the Hoeffding tree (an incremental decision tree) as the base classifier of these two methods. The performance difference between our algorithm and iCC and iPS for these two synthetic datasets, however, are not statistically significant for the average precision and ranking loss whereas the proposed method outperforms all benchmark algorithms for the macro-average F1.



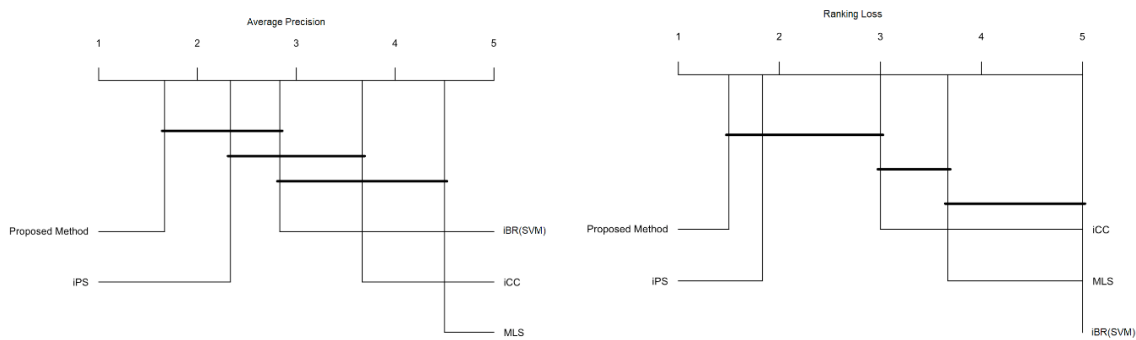


Fig.11. Nemenyi test for the six measures under the concept drift setting

5.5. Discussion

The four multi-target tree-based methods and MLS perform poorly on the experimental datasets. While MLS ranks at the bottom on 5 measures except for the average precision, the four multi-target tree-based methods only perform well for the ranking loss measure. In fact, multi-target tree-based methods are highly threshold-dependent. In these methods, the threshold is used twice, i.e., once in converting multi-target prediction to multi-label prediction and then in updating the regressor. It is noted that the sample-based measures and label-based measures are threshold-based measures. Therefore, multi-target tree-based methods could underperform on these measures if a sub-optimal threshold is used. MLS method, meanwhile, is the simplest multi-label classifier which assigns the most common label set from the training data for all test samples. It, therefore, is an uncompetitive incremental classifier for most of the performance measures.

iCC and iPS, especially iPS, are competitive MLC methods for all measures except the ranking loss. Because of taking into account the label correlation in the learning model, these methods can choose suitable labels to assign to each arrived sample, resulting in good MLC performance. In fact, iCC and iPS compare the prediction scores to a threshold to select the predicted labels. In this paper, we followed (Read et al., 2012) by adjusting the threshold via a batch-based approach on the predicted and true label cardinality. The dependence on choosing the appropriate threshold and the less flexible adjustment for the number of predicted labels per sample could have caused the under-performance of these methods on some datasets. iBR(SVM) meanwhile is a very high-performing MLC methods in both settings because of the fact that SVM is a state-of-the-art supervised learning algorithm. This method ranks first on 5 measures except for the ranking loss in the non-concept drift setting and ranks first on the 2 label-based measures and in the concept drift setting. However, it under-performs compared to the proposed method for the ranking loss.

The proposed method is competitive to the other benchmark algorithms like iBR(SVM) and iPS. In the concept drift setting, it ranks first on the two ranking based measures and ranks second on the two label-based measures. Although iBR(SVM) has higher ranking than the proposed method on 5 measures on the non-concept drift setting, the proposed method is a better choice as it always attain high ranks and does not rank last on any measures. In the proposed method, the clustering model is maintained properly with the arrival of each sample based on its ground truth label set. The classification stage, therefore, can benefit when computing the posterior probability that a sample belongs to a class label based on the weight of the mature clusters. The weight of each cluster is decayed exponentially with time so that the model can be adaptive to the concept drift. Moreover, the number of predicted labels is learned based on the Hoeffding inequality and the label cardinality to ensure that each sample is assigned with an adequate set of labels. The proposed method, therefore, performed well in our experiments.

6. Conclusions

In this paper, we have introduced an incremental online learning method to solve the online MLC problem. In detail, we aggregate information in the arriving samples through a clustering process that takes into account the sample's time of arrival to compute the sample's weight. Each cluster's weight is then computed from the weights of the samples inside. The clustering process also builds up a distribution of labels in each cluster that would be later used for MLC. To handle concept drift, we proposed a decay mechanism on the sample's weight so that the influence of old samples is reduced over time while attention is put on the new ones. The classification process on each sample is conducted by using the mature clusters and their weights to compute the posterior probability that the sample belongs to a class. For MLC, the h labels associated with the top h classes of posterior probabilities (computed from the cluster's label distributions) are selected to assign labels for the sample. The number of predicted labels h is determined adaptively in our algorithm. Specifically, given a confidence level, we conduct the adjustment if there is a certain difference between h and the current label cardinality \bar{z} . The difference needed for the update is computed based on the Hoeffding inequality. The clustering model and the number of predicted labels are updated on-the-fly with the arrived samples and their ground truth labels. Due to the incremental learning nature of our algorithm, the incoming samples do not need to be stored once they are processed.

The proposed method and the benchmark algorithms were evaluated on five popular multi-label datasets in the stationary setting, and twelve multi-label datasets in the concept drift setting. The experimental results showed that our method is highly competitive compared to several benchmark algorithms, especially under the concept drift setting. The proposed method is high desirable in practice as it always maintains the high ranks and does not rank last on any measures.

The performance of the proposed method can be enhanced by integrating label correlation [4] in the learning model. For MLC problems, label correlation is an important factor that can enhance the prediction quality. For example, if a sample has been assigned with the label ‘indoor’, labels like ‘table’ and ‘chair’ should have more chance of been assigned to the sample than labels like ‘car’ and ‘grass’. By considering the semantic relationship between the labels, we can choose the proper set of labels in the prediction process. In the streaming context, moreover, the label correlation will need to be updated with the incoming of each sample. Several methods capturing the label correlation such as classifier trellis [25] and graphical model for feature-label-label relationship triple [13] can be combined with the proposed method to enhance the performance of MLC task. Moreover, the proposed method can be expanded to handle the more general learning paradigm like MDC [25]. These will be our future works.

Acknowledgements

This work is partially funded by the Science and Technology Foundation of Guangdong Province, China (Project No. 2017A050501002)

References

- [1] T.T.T. Nguyen, A.W.C. Liew, T.T. Nguyen, S.L. Wang, A novel Bayesian framework for Online Imbalanced Learning, in *Proceeding of Digital Image Computing: Techniques and Applications*, 2017.
- [2] T.T.T. Nguyen, T.T. Nguyen, A.W.C. Liew, S.L. Wang, Variational inference based Bayes online classifiers with concept drift adaptation, *Pattern Recognition*.81 (2018), 280-293.
- [3] X.C. Pham, M.T. Dang, S.V. Dinh, S. Hoang, T.T. Nguyen, A.W.C. Liew, Learning from data stream based on Random Projection and Hoeffding Tree Classifier, in *Proceeding of Digital Image Computing: Techniques and Applications*, 2017.
- [4] M.L. Zhang, Z.H. Zhou, A review on Multi-Label Learning Algorithms, *IEEE Transactions on Knowledge and Data Engineering*. 26(8) (2014), 1819-1837.
- [5] J. Read, A. Bifet, G. Holmes, B. Pfahringer, Scalable and efficient multi-label classification for evolving data stream, *Machine Learning*. 88, 2012, 243-272.
- [6] A. Osojnik, P. Panov, S. Dzeroski, Multi-label classification via multi-target regression on data streams, *Machine Learning*. 106, 2017, 745-770.
- [7] F. Cao, M. Ester, W. Qian, A. Zhou, Density-Based Clustering over an Evolving Data Stream with Noise, in *Proceedings of the SIAM International Conference on Data Mining*, 2006.

- [8] J. Gama, P. Medas, G. Castillo, Pedro Rodrigues, Learning with Drift Detection, in *Advances in Artificial Intelligence*, Springer Berlin Heidelberg, 2004, vol 3171, pp. 268 – 295.
- [9] E.S. Xioufis, M. Spiliopoulou, G. Tsoumakas, I. Vlahavas, Dealing with Concept Drift and Class Imbalance in Multi-Label Stream Classification, in *Proceeding of 22nd International Conference on Artificial Intelligence*, 2011.
- [10] K. Dembczýnski, W. Cheng, E. Hullermeier, Bayes optimal multilabel classification via probabilistic classifier chains, in *Proceedings of the 27th International Conference on Machine Learning*, Haifa, Israel, 2010, pp. 279–286.
- [11] J. Nam, E.L. Mencia, H.J. Kim, J. Furnkranz, Maximizing Subset Accuracy with Recurrent Neural Networks in Multi-Label Classification, in *Proceeding of NIPS*, 2017.
- [12] A. Kumar, S. Vembu, A. Menon, C. Elkan, Beam search algorithms for multi-label learning, *Machine Learning* 92 (2013), 65-89.
- [13] N. Ghamrawi, A. McCallum, Collective multi-label classification, in: *CIKM '05: 14th ACM International Conference on Information and Knowledge Management*, ACM Press, New York, NY, USA, 2005, pp. 195-200.
- [14] G. Tsoumakas, I. Katakis, and I. Vlahavas, Random k-labelsets for multi-label classification, *IEEE Transactions on Knowledge and Data Engineering*, 23(7), 2011, 1079–1089.
- [15] J. Read, B. Pfahringer, G. Holmes, Multi-label Classification using Ensembles of Pruned Sets, in *Proc. of IEEE International Conference on Data Mining*, 2008, pp. 995–1000.
- [16] Y. Guo, S. Gu, Multi-label classification using conditional dependency networks, in *Proceeding of the International Joint Conference on Artificial Intelligence*, 2011, pp 1300-1305.
- [17] M.L. Zhang, Z.H. Zhou, ML-KNN: A lazy learning approach to multi-label learning, *Pattern Recognition*. 40 (7) (2007), 2038–2048.
- [18] X. Li, Y. Guo, Active Learning with Multi-Label SVM Classification, in *Proceeding of the International Joint Conference on Artificial Intelligence*, 2013, pp. 1479-1485.
- [19] C. Vens, J. Struyf, L. Schietgat, S. Džeroski, H. Blockeel, Decision trees for hierarchical multilabel classification. *Machine Learning*, 73(2) (2008), 185–214.
- [20] S. Ubaru, A. Mazumdar, Multilabel Classification with Group Testing and Codes, in *Proceeding of ICML*, 2017.
- [21] A. Kapoor, R. Viswanathan, P. Jain, Multilabel Classification using Bayesian Compressed Sensing, in *Proceeding of NIPS*, 2012.
- [22] Y.N. Chen, H.T. Lin, Feature-aware label space dimension reduction for multi-label classification, In *Advances in Neural Information Processing Systems*, 2012, pp. 1529-1537.

- [23] Y. Lin, Q. Hu, J. Zhang, Z. Wu, Multi-label feature selection with stream labels, *Information Sciences*. 372 (2016), 256-275.
- [24] J. Lee, D.W. Kim, SCLS: Multi-label feature selection based on scalable criterion for large dataset, *Pattern Recognition*. 66 (2017), 342-352.
- [25] J. Read, L. Martino, P. Olmos, D. Luengo, Scalable Multi-Output Label Prediction: From Classifier Chains to Classifier Trellises, *Pattern Recognition*. 48(6) (2015), 2096-2109.
- [26] J. Read, L. Martino, D. Luengo, Efficient Monte Carlo Methods for Multi-Dimensional Learning with Classifier Chains, *Pattern Recognition*. 47(3) (2014), 1535-1546.
- [27] A. Bifet, R. Gavaldà, Adaptive learning from evolving data streams, In *Advances in Intelligent Data Analysis VIII* (pp. 249–260). Springer, 2009.
- [28] W. Qu, Y. Zhang, J. Zhu, Q. Qiu, Mining multi-label concept-drifting data streams using dynamic classifier ensemble. In *Advances in machine learning*, pp. 308–321, Springer, 2009.
- [29] L. Wang, H. Shen, H. Tian, Weighted Ensemble Classification of Multi-label Data Streams, in J. Kim et al. (Eds.), *PAKDD Part II, Lecture Note in Artificial Intelligence*, pp. 551-562, 2017.
- [30] P. Domingos, G. Hulten. Mining high-speed data streams. In *Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining*, 2000, pp. 71–80.
- [31] A. Bifet, R. Gavaldà, Learning from Time-Changing Data with Adaptive Windowing, in *Proceeding of ICDM*, 2007.
- [32] Z. Shi, Y. Xue, Y. Wen, G. Cai, Efficient Class Incremental Learning for Multi-label classification of Evolving Data Streams, *International Joint Conference on Neural Networks 2014*.
- [33] Z. Shi, C. Feng, Y. Wen, H. Zhao, Drift Detection for Multi-label Data Streams Based on Label Grouping and Entropy, In *IEEE International Conference on Data Mining Workshop*, 2014.
- [34] K. Karponi, G. Tsoumakas, An Empirical Comparison of Methods for Multi-Label data Stream Classification, in P. Angelov et al. (eds.), *Advances in Big Data, Advantages in Intelligent Systems and Computing* 529, 2017.
- [35] E. Cohen, M. Strauss, Maintaining Time-Decaying Stream Aggregates, *Journal of Algorithms*. 59 (1) (2006), 19-36.
- [36] J. Xu, Multi-Label Weighted k-Nearest Neighbor Classifier with Adaptive Weight Estimation, *International Conference on Neural Information Processing*, 2011, pp 79-88
- [37] W. Hoeffding, Probability Inequalities for Sums of Bounded Random Variables, *Journal of the American Statistical Association*. 58 (301) (1963), pp.13-30.
- [38] J. Read, Scalable multi-label classification, PhD Thesis, University of Waikato, 2010.
- [39] S. Garcia, F. Herrera, An extension on statistical comparisons of classifiers over multiple data sets for all pairwise comparisons, *Journal of Machine Learning Research* 9 (2008), 2579–2596.

[40] T.T. Nguyen, M.P. Nguyen, X.C. Pham, A.W.C. Liew, Heterogeneous Classifier Ensemble with Fuzzy Rule-based Meta Learner, *Information Sciences*. 422 (2018), 144-160.

Appendix

The performance measures

In this paper, the comparisons of the proposed method and the benchmark algorithms are based on six well-known performance measures: sample-based F1/accuracy, label-based micro-averaged F1/macro-averaged F1, and ranking-based average precision/ranking loss. We briefly describe each measure supposing that N examples $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ are received with ground truth label sets $\mathbb{Y}_1, \mathbb{Y}_2, \mathbb{Y}_N$ and prediction label sets $\hat{\mathbb{Y}}_1, \hat{\mathbb{Y}}_2, \dots, \hat{\mathbb{Y}}_N$, respectively.

Sample-based measures

Sample-based measures evaluate the performance of a MLC algorithm on a per-sample basis. They are calculated for each sample and then averaged over all of them.

Sample-based accuracy is the average proportion of label values correctly classified out of the total number (predicted and true) of labels:

$$\text{Sample - based accuracy} = \frac{1}{N} \sum_{t=1}^N \frac{|\mathbb{Y}_t \cap \hat{\mathbb{Y}}_t|}{|\mathbb{Y}_t \cup \hat{\mathbb{Y}}_t|} \quad (\text{A1})$$

The sample-based F1 is the harmonic mean of sample-based precision and recall:

$$\text{Sample - based F1} = 2 / \left(\frac{1}{\text{Ex.Precision}} + \frac{1}{\text{Ex.Recall}} \right) \quad (\text{A2})$$

Here,

$$\text{Ex. Precision} = \frac{1}{N} \sum_{t=1}^N \frac{|\mathbb{Y}_t \cap \hat{\mathbb{Y}}_t|}{|\hat{\mathbb{Y}}_t|} \quad (\text{A3})$$

$$\text{Ex. Recall} = \frac{1}{N} \sum_{t=1}^N \frac{|\mathbb{Y}_t \cap \hat{\mathbb{Y}}_t|}{|\mathbb{Y}_t|} \quad (\text{A4})$$

The greater the sample-based accuracy and F1 of an MLC algorithm (with an optimal value of 1), the better its classification performance over different samples.

Label-based measures

Label-based measures evaluate the performance of a MLC algorithm on a per-label basis. They are calculated for each label and then averaged over all of them. Definitions of many of the label-based measures are based on four basic quantities named true positive (TP), true negative (TN), false positives (FP) and false negative (FN), which are calculated as follows for label $l \in \mathcal{Y}$:

$$TP_l = |\{\mathbf{x}_t | l \in \mathbb{Y}_t \wedge l \in \hat{\mathbb{Y}}_t, 1 \leq t \leq N\}| \quad (\text{A5})$$

$$TN_l = |\{\mathbf{x}_t | l \notin \mathbb{Y}_t \wedge l \in \widehat{\mathbb{Y}}_t, 1 \leq t \leq N\}| \quad (\text{A6})$$

$$FP_l = |\{\mathbf{x}_t | l \notin \mathbb{Y}_t \wedge l \in \widehat{\mathbb{Y}}_t, 1 \leq t \leq N\}| \quad (\text{A7})$$

$$FN_l = |\{\mathbf{x}_t | l \in \mathbb{Y}_t \wedge l \notin \widehat{\mathbb{Y}}_t, 1 \leq t \leq N\}| \quad (\text{A8})$$

The value of F1 can be obtained in the form of macro-averaged or micro-averaged:

$$\text{Macro - average F1} = \frac{1}{M} \sum_l F1(TP_l, TN_l, FP_l, FN_l) \quad (\text{A9})$$

$$\text{Micro - average F1} = F1(\sum_l TP_l, \sum_l TN_l, \sum_l FP_l, \sum_l FN_l) \quad (\text{A10})$$

Here,

$$F1(TP, TN, FP, FN) = \frac{2TP}{2TP+FP+FN} \quad (\text{A11})$$

Clearly, the greater the micro F1 and macro F1 (with an optimal value of 1), the better the predictive performance over different labels obtained by the learner.

Ranking-based measures

Ranking-based measures analyze the confidence outputs $f(\mathbf{x}_t, l) \in \mathbb{R}, l \in \mathbf{y}$ of a MLC methods directly, i.e. independent of the thresholding procedure. For \mathbf{x}_t , $rank_f(\mathbf{x}_t, l)$ returns the rank of l in \mathbf{y} based on the descending order induced from $f(\mathbf{x}_t, \cdot)$. That means label l is considered to be ranked higher than l' , i.e. $rank_f(\mathbf{x}_t, l) \leq rank_f(\mathbf{x}_t, l')$ if $f(\mathbf{x}_t, l) > f(\mathbf{x}_t, l')$.

Ranking loss evaluates the fraction of reversely ordered label pairs when an irrelevant label is ranked higher than a relevant label:

$$\text{Ranking Loss} = \frac{1}{N} \sum_{t=1}^N \frac{1}{|\mathbb{Y}_t| |\overline{\mathbb{Y}}_t|} |\{(l, l') | f(\mathbf{x}_t, l) \leq f(\mathbf{x}_t, l'), (l, l') \in \mathbb{Y}_t \times \overline{\mathbb{Y}}_t\}| \quad (\text{A12})$$

where $\overline{\mathbb{Y}}_t$ is the complementary set of \mathbb{Y}_t in \mathbf{y} . Small values of ranking loss are desired.

Average precision evaluates the average fraction of labels ranked above a particular label $l \in \mathbb{Y}_t$ which actually are in \mathbb{Y}_t .

$$\text{Average Precision} = \frac{1}{N} \sum_{t=1}^N \frac{1}{|\mathbb{Y}_t|} \sum_{l \in \mathbb{Y}_t} \frac{|\{l' | rank_f(\mathbf{x}_t, l') \leq rank_f(\mathbf{x}_t, l), l' \in \mathbb{Y}_t\}|}{rank_f(\mathbf{x}_t, l)} \quad (\text{A13})$$

Average precision reaches the maximum value of 1 when f ranks the labels for all samples perfectly so that there is no sample \mathbf{x}_t for which a label not in \mathbb{Y}_t has higher rank than a label in \mathbb{Y}_t .

Author Biographies

Tien Thanh Nguyen received his PhD degree in computer science from the School of Information & Communication Technology, Griffith University, Australia in 2017. He is currently a Research Fellow at the School of Computing Science and Digital Media, Robert Gordon University, Aberdeen, Scotland, UK. His research interest is in the field of machine learning, pattern recognition, and evolutionary computation. He is a member of the IEEE since 2014.

Manh Truong Dang is currently an advanced undergraduate student at the School of Information & Communication Technology, Hanoi University of Science and Technology, Vietnam. His research interest is in the field of machine learning and pattern recognition.

Anh Vu Luong is currently a PhD student at the School of Information & Communication Technology, Griffith University, Australia. His research interest is in the field of machine learning and pattern recognition.

Alan Wee-Chung Liew is currently an Associate Professor at the School of Information & Communication Technology, Griffith University, Australia. His research interest is in the field of machine learning, pattern recognition, computer vision, medical imaging, and bioinformatics. He has served on the technical program committee of many international conferences and is on the editorial board of several journals, including the IEEE Transactions on Fuzzy Systems. He is a senior member of the IEEE since 2005.

Tiancai Liang is currently the Chief Technology Officer and the Director of the Intelligent Security Institute of GRGBanking Technology Co., Ltd, China. He obtained his PhD from the South China University of Technology, majoring in Pattern Recognition and Intelligent Systems. His research interest is in the field of artificial intelligence and video image analysis.

John McCall is currently a Professor at the School of Computing Science and Digital Media, Robert Gordon University, Aberdeen, Scotland, UK. His research interest is in the area of naturally-inspired computing and their application to real-world problems arising in complex engineering and medical/biological systems.

Accepted Manuscript

Multi-Label Classification via Incremental Clustering on Evolving Data Stream

Tien Thanh Nguyen , Manh Truong Dang , Anh Vu Luong ,
Alan Wee-Chung Liew , Tiancai Liang , John McCall

PII: S0031-3203(19)30232-8
DOI: <https://doi.org/10.1016/j.patcog.2019.06.001>
Reference: PR 6935



To appear in: *Pattern Recognition*

Received date: 31 October 2018
Revised date: 5 March 2019
Accepted date: 1 June 2019

Please cite this article as: Tien Thanh Nguyen , Manh Truong Dang , Anh Vu Luong , Alan Wee-Chung Liew , Tiancai Liang , John McCall , Multi-Label Classification via Incremental Clustering on Evolving Data Stream, *Pattern Recognition* (2019), doi: <https://doi.org/10.1016/j.patcog.2019.06.001>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.