

Multi-Dimensional Classification via k NN Feature Augmentation

Bin-Bin Jia^{a,b,c}, Min-Ling Zhang^{a,c,d,*}

^a*School of Computer Science and Engineering, Southeast University, Nanjing 210096, China*

^b*College of Electrical and Information Engineering, Lanzhou University of Technology, Lanzhou 730050, China*

^c*Key Laboratory of Computer Network and Information Integration (Southeast University), Ministry of Education, China*

^d*Collaborative Innovation Center of Wireless Communications Technology, China*

Abstract

In multi-dimensional classification (MDC), each training example is represented by a single instance (feature vector) while associated with multiple class variables, each of which specifies its class membership w.r.t. one specific class space. Most existing MDC approaches try to model dependencies among class variables in output space when inducing predictive functions, while the potential usefulness of manipulating feature space hasn't been investigated. As a first attempt towards feature manipulation in input space for MDC, a simple yet effective approach named KRAM is proposed which enriches the original feature space with augmented features based on k NN techniques. Specifically, simple counting statistics on the class membership of neighboring MDC examples as well as distance information between MDC examples and their k nearest neighbors are used to generate augmented feature vector. In this way, discriminative information from class space is expected to be brought into the feature space which would be helpful to the following MDC predictive model induction. To validate the effectiveness of the proposed feature augmentation techniques, comprehensive comparative studies are conducted over fifteen benchmark data sets. Compared to the original feature space, it is clearly shown that the k NN-

*Corresponding author

Email address: zhangml@seu.edu.cn (Min-Ling Zhang)

augmented features generated by the proposed KRAM approach can significantly improve generalization abilities of existing MDC approaches.

Keywords: machine learning, multi-dimensional classification, feature augmentation, class dependencies

1. Introduction

Multi-dimensional classification (MDC) aims to build learning models for real-world objects with a wealth of semantics, which assumes several class spaces to characterize the semantics of objects along different dimensions. Here, each example in MDC training set is represented by a single instance while associated with a number of class variables, and all these class variables respectively specify their class membership with regard to one specific class space. Specifically, there are many scenarios where we need to learn from MDC examples. For example, in image classification, the semantics of a natural scene image can be classified along the **landscape** dimension (with possible class labels *lake*, *grassland*, *mountain*, etc.), along the **time** dimension (with possible class labels *morning*, *afternoon*, *evening*, etc.), and along the **weather** dimension (with possible class labels *sunny*, *rainy*, *snowy*, etc.). For another example, in music classification, the semantics of a piece of song can be classified along the **language** dimension (with possible class labels *Chinese*, *English*, *Spanish*, etc.), along the **genre** dimension (with possible class labels *classical*, *popular*, *rock*, etc.), and along the **instrument** dimension (with possible class labels *guitar*, *violin*, *piano*, etc.). More applications of MDC techniques also include text classification [1, 2], bioinformatics [3, 4], web mining [5], etc.

Formally speaking, let $\mathcal{X} = \mathbb{R}^d$ be the d -dimensional input (feature) space and $\mathcal{Y} = C_1 \times C_2 \times \dots \times C_q$ be the output space which corresponds to the Cartesian product of q heterogeneous class spaces. Here, each class space C_j ($1 \leq j \leq q$) consists of K_j possible class labels, i.e., $C_j = \{c_1^j, c_2^j, \dots, c_{K_j}^j\}$. Furthermore, let $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i) \mid 1 \leq i \leq m\}$ be the MDC training set with m training examples, where $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{id}]^\top \in \mathcal{X}$ is a d -dimensional

feature vector and $\mathbf{y}_i = [y_{i1}, y_{i2}, \dots, y_{iq}]^\top \in \mathcal{Y}$ is the associated class vector, each of which is one possible value in j -th class space, i.e., $y_{ij} \in C_j$. Then, the learning task of multi-dimensional classification is to induce a mapping function $f : \mathcal{X} \mapsto \mathcal{Y}$ from \mathcal{D} which can predict a proper class vector $f(\mathbf{x}) \in \mathcal{Y}$ for the
30 unseen instance \mathbf{x} .

To accomplish MDC learning tasks, there are two intuitive solutions. The first one is to train multiple multi-class classifiers independently, one per class space. The second one is to train a single multi-class classifier by treating each distinct class combination in training set as a new class label, i.e., powerset
35 method. These two intuitive strategies are simple and sometimes effective, but both of them also have obvious disadvantages. As all class spaces share the same feature space, there must be some relationships among them. However, due to the limited training examples, not all possible combination of class variables will appear in training set. The first intuitive strategy above ignores all dependen-
40 cies among class spaces which would impact the generalization abilities of the induced MDC model, while the second strategy cannot make prediction for class combinations not appearing in the training set and also usually has prohibitively large number of combinations. In other words, both ignoring and overfitting class dependencies should be avoided. Therefore, most existing MDC approach-
45 es try to model dependencies among class variables from different dimensions in the output space in different ways, such as capturing dependencies between each pair of class spaces [6], specifying chaining order over class spaces [7, 8], assuming directed acyclic graph (DAG) structure over class spaces [9, 10], and partitioning class spaces into groups [11], etc.

50 Different from most existing works which directly model dependencies among different class variables in the output space, in this paper, we show the potential usefulness of manipulating input (feature) space for inducing MDC predictive models. Accordingly, a simple yet effective MDC approach named KRAM, i.e., *kNN featuRe Augmentation for Multi-dimensional classification*, is proposed.

55 Specifically, the main contributions of this paper correspond to:¹

- We propose a first attempt towards manipulating feature space for MD-
C, where the proposed KRAM approach works by enriching the original
features of MDC examples with their k NN-augmented features. In this
way, discriminative information from output space can be brought into
60 the input (feature) space to facilitate MDC predictive model induction.
- We design two versions of k NN-augmented features, i.e., discrete version
and continuous version. For the discrete version, it is based on standard
 k NN techniques which employs simple counting statistics according to the
class membership of k nearest neighbors in training set. For the continuous
65 version, it is based on weighted k NN techniques which combines the count-
ing statistics with extra bias terms by considering the distance between
current instance and its nearest neighbors.
- We conduct comprehensive comparative studies over fifteen benchmark
data sets to validate the effectiveness of KRAM in improving the gen-
70 eralization performance of MDC approaches. Accordingly, properties of
the proposed k NN augmented features have been analyzed based on the
empirical results.

The rest of the paper is structured as follows. Section 2 briefly discuss-
es existing works related to multi-dimensional classification. Section 3 presents
75 technical details of the proposed KRAM approach. Section 4 reports comprehen-
sive experimental results of comparative studies over a wide range of benchmark
data sets. Finally, Section 5 concludes this paper.

¹This paper is an extension of our preliminary work [12]. The main differences include:
(1) The introduction and conclusion parts have been updated with further discussions on
the proposed approaches; (2) Another strategy for generating k NN-augmented features (i.e.
the continuous version) has been developed based on weighted k NN techniques; (3) More
comprehensive comparative studies have been conducted in terms of newly added benchmark
data sets, enriched parameter sensitivity analyses, and one recently proposed compared ap-
proach [13].

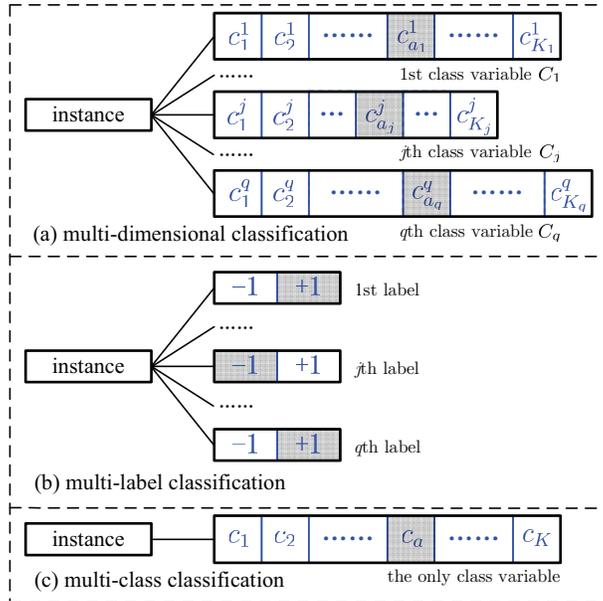


Figure 1: An intuitive comparison among multi-dimensional classification, multi-label classification, and multi-class classification.

2. Related Work

The most related learning frameworks to multi-dimensional classification
80 include the traditional multi-class classification (MCC) and the popular multi-
label classification (MLC) [14, 15, 16]. As shown in Figure 1, the MDC problem
corresponds to several coupled MCC problems while the MLC problem corre-
sponds to several coupled binary classification problems, i.e., both MDC and
MCC can be regarded as one possible instantiation of multi-output learning [17]
85 where each object is associated with multiple output variables. However, it is
worth emphasizing that we should not simply consider that the difference be-
tween these two learning frameworks is whether the type of output variables is
multi-class or binary class. Specifically, the key difference between MDC and
MLC is whether semantic spaces are *heterogenous* or *homogeneous*, where each
90 class variable in MDC corresponds to one specific class space, while each label
in MLC specifies whether one concept is relevant or not in the only class space.

Therefore, it is unreasonable and will get suboptimal solutions to directly align class labels from different dimensions when trying to design MDC approaches. Anyway, the MLC problem can be regarded as a degenerated case of MDC by
95 restricting all class variables to be binary-valued ones [11, 18, 19]. Additionally, recently proposed learning frameworks, such as duel set multi-label learning [20] and multiple ordinal output classification [21], can also be regarded as special cases of MDC when the number of dimensions equals two or class labels in each class space have ordinal relationship.

100 Intuitively, MDC can be solved by decomposing the original problem into a number of MCC problems, i.e., training a multi-class classifier independently according to each class variable. However, this natural strategy neglects the dependencies among class variables which may exist in real-world MDC tasks, and thus leads to suboptimal MDC model. Therefore, when trying to induce
105 better MDC models, one of the core ways is to model dependencies among class spaces.

Dependencies between each pair of class spaces could be modeled by a set of base classifiers, and the final multi-dimensional inference is accomplished by combining predictive outputs from base classifiers via Markov random field [6].
110 Similar to classifier chains in MLC [22], the MDC problem can be converted into a chain of MCC problems, where the outputs of preceding multi-class classifiers in the chain are treated as extra input features when building subsequent ones. Obviously, its effectiveness is largely affected by the chaining order over class spaces which can be specified in deterministic manner [7] or random manner [8].

115 In addition, lots of existing works explicitly model dependencies among class variables with different families of directed acyclic graph (DAG) structures [9, 10], which form a family of probabilistic graphical models for MDC called multi-dimensional Bayesian network classifiers. Recent works focus on more efficient structure learning strategies which is still challenging [23, 24, 25].
120 Class powerset (CP) approach models class dependencies by converting the original MDC problem into only one MCC problem, where each distinct class combination existing in training set is regarded as a new class label in the new

MCC problem. Considering the possible huge number of class labels (with at most $\prod_{j=1}^q K_j$ class labels after powerset transformation), it will be helpful to group MDC class spaces into super-classes so as to facilitate the subsequent model induction for MDC [11].

3. The KRAM Approach

It has been widely acknowledged that modeling class dependencies in output space plays a crucial role when attempting to induce better MDC models, the importance of manipulating input (feature) space, however, hasn't been well investigated for MDC studies. In this section, we will present the technical details of our proposed KRAM approach, which aims at improving the generalization performance of existing MDC approaches via augmenting MDC examples' original features with the help of k NN techniques.

Following the same notations as used in previous sections, let $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i) \mid 1 \leq i \leq m\}$ be the MDC training set where $\mathbf{y}_i = [y_{i1}, y_{i2}, \dots, y_{iq}]^\top \in \mathcal{Y}$ is the corresponding class vector associated with \mathbf{x}_i . For each instance \mathbf{x} , let $\mathcal{N}(\mathbf{x}) = \{i_r \mid 1 \leq r \leq k\}$ be the set of indices for \mathbf{x} 's k nearest neighbors which are identified in \mathcal{D} . Here, Euclidean distance is utilized to measure the similarities between two instances, and for convenience, we assume that distance between \mathbf{x} and \mathbf{x}_{i_r} is in ascending order, which means that the smaller the value of r , the closer the distance between \mathbf{x} and \mathbf{x}_{i_r} . Then, an indicating vector $\mathbf{v}_{ja}^{\mathbf{x}} = [v_{ja}^{\mathbf{x}}(1), v_{ja}^{\mathbf{x}}(2), \dots, v_{ja}^{\mathbf{x}}(k)]^\top \in \{0, 1\}^k$ is defined as follows:

$$v_{ja}^{\mathbf{x}}(r) = \llbracket y_{i_r j} = c_a^j \rrbracket \quad (1 \leq r \leq k, i_r \in \mathcal{N}(\mathbf{x})) \quad (1)$$

Here, $1 \leq a \leq K_j, 1 \leq j \leq q$. $\mathbf{y}_{i_r} = [y_{i_r 1}, y_{i_r 2}, \dots, y_{i_r q}]^\top$ is the corresponding class vector associated with \mathbf{x}_{i_r} . The predicate $\llbracket \pi \rrbracket$ returns 1 if π holds and 0 otherwise. Therefore, $v_{ja}^{\mathbf{x}}(r)$ records whether the \mathbf{x} 's r -th nearest neighbor \mathbf{x}_{i_r} has class label c_a^j in the j -th class space or not.

Based on $\mathbf{v}_{ja}^{\mathbf{x}}$, the following discrete version of statistics $\delta_j^{\mathbf{x}} = [\delta_{j1}^{\mathbf{x}}, \delta_{j2}^{\mathbf{x}}, \dots, \delta_{jK_j}^{\mathbf{x}}]^\top$ can be defined w.r.t. the j -th class space:

$$\delta_{ja}^{\mathbf{x}} = \langle \mathbf{1}_k, \mathbf{v}_{ja}^{\mathbf{x}} \rangle \quad (1 \leq a \leq K_j) \quad (2)$$

Here, $\mathbf{1}_k$ is a column vector of all ones with length k , and $\langle \cdot, \cdot \rangle$ returns inner
140 product of two vectors. Therefore, $\delta_{ja}^{\mathbf{x}}$ records the number of examples in \mathbf{x} 's k
nearest neighbors whose class label equals c_a^j in the j -th class space. According
to the definition in Eq.(2), it is easy to know that $\sum_{a=1}^{K_j} \delta_{ja}^{\mathbf{x}} = k$ holds.

After traversing each class space one by one, a total of q different counting
statistics $\delta_j^{\mathbf{x}}$ ($1 \leq j \leq q$), each of which contains K_j elements, can be generated.
By concatenating all these q counting statistics, we can define the following
augmented feature vector $\Delta_{\mathbf{x}}$ for \mathbf{x} :

$$\Delta_{\mathbf{x}} = [\delta_1^{\mathbf{x}}, \delta_2^{\mathbf{x}}, \dots, \delta_q^{\mathbf{x}}] \quad (3)$$

Based on the above feature vector, the original MDC training set \mathcal{D} can be
transformed into:

$$\tilde{\mathcal{D}} = \{(\tilde{\mathbf{x}}_i, \mathbf{y}_i) \mid 1 \leq i \leq m\} \quad (4)$$

145 Here, $\tilde{\mathbf{x}}_i = [\mathbf{x}_i, \Delta_{\mathbf{x}_i}] \in \tilde{\mathcal{X}}$ which means $\tilde{\mathbf{x}}_i$ is comprised of a concatenation of
 \mathbf{x}_i and $\Delta_{\mathbf{x}_i}$. $\tilde{\mathcal{X}}$ represents a synthetic feature space which corresponds to the
Cartesian product between the original feature space (i.e., \mathcal{X}) and a $(\sum_{j=1}^q K_j)$ -
dimensional augmented one. Thereafter, an MDC predictive model $f: \tilde{\mathcal{X}} \mapsto \mathcal{Y}$
can be trained over the new constructed data set $\tilde{\mathcal{D}}$ by applying any off-the-
150 shelf MDC training algorithm \mathcal{L} , i.e., $f \leftarrow \mathcal{L}(\tilde{\mathcal{D}})$. Given an unseen instance
 \mathbf{x}^* , its predicted class vector \mathbf{y}^* can be assigned by feeding its corresponding
augmented instance $\tilde{\mathbf{x}}^*$ into f .

Obviously, discrete version of statistics in Eq.(2) utilizes standard k NN tech-
niques. Standard k NN approach simply assigns a test instance the class of the
155 majority of its k nearest neighbors, i.e., the class of test instance is voted u-
niformly by its k nearest neighbors. As a popular variation, weighted k NN
predicts the class of a test instance via non-uniformly voting by its k near-
est neighbors. The closer distance between test instance and neighbor is, the
greater weight the neighbor has. Following the idea of weighted k NN, another
160 continuous version of augmented feature vector is designed based on the discrete
one defined in Eq.(1)~Eq.(3).

According to the motivation of weighted k NN, in this paper, the weight vector for k nearest neighbors is simply set as $\mathbf{w} = [1, 1/\sqrt{2}, \dots, 1/\sqrt{k}]^\top$. Then, a bias $\zeta_{ja}^{\mathbf{x}}$ for $\delta_{ja}^{\mathbf{x}}$ in Eq.(2) is defined according to $\mathbf{v}_{ja}^{\mathbf{x}}$ in Eq.(1) as follows:

$$\zeta_{ja}^{\mathbf{x}} = \frac{\langle \mathbf{w}, \mathbf{v}_{ja}^{\mathbf{x}} \rangle - \min(\mathbf{v}_{ja}^{\mathbf{x}})}{\max(\mathbf{v}_{ja}^{\mathbf{x}}) - \min(\mathbf{v}_{ja}^{\mathbf{x}})} (\zeta_{max} - \zeta_{min}) + \zeta_{min} \quad (5)$$

Here, $\max(\mathbf{v}_{ja}^{\mathbf{x}}) = \sum_{r=1}^{\langle \mathbf{1}_k, \mathbf{v}_{ja}^{\mathbf{x}} \rangle} w(r)$, $\min(\mathbf{v}_{ja}^{\mathbf{x}}) = \sum_{r=k-\langle \mathbf{1}_k, \mathbf{v}_{ja}^{\mathbf{x}} \rangle+1}^k w(r)$ represent the possible maximum and minimum of $\langle \mathbf{w}, \mathbf{v}_{ja}^{\mathbf{x}} \rangle$ respectively, where $w(r)$ denotes the r -th element of weight vector \mathbf{w} . ζ_{max} and ζ_{min} are two hyper-
165 parameters, and $\zeta_{max} - \zeta_{min} < 1$ holds. In this paper, we set ζ_{max} as 0.5 and ζ_{min} as 0. We can easily know that $\zeta_{min} \leq \zeta_{ja}^{\mathbf{x}} \leq \zeta_{max}$.

Then, compared with the discrete version of statistics defined in Eq.(2), we can define another continuous version of statistics $\boldsymbol{\delta}_j^{\mathbf{x}} = [\delta_{j1}^{\mathbf{x}}, \delta_{j2}^{\mathbf{x}}, \dots, \delta_{jK_j}^{\mathbf{x}}]^\top$ as follows:

$$\delta_{ja}^{\mathbf{x}} = \langle \mathbf{1}_k, \mathbf{v}_{ja}^{\mathbf{x}} \rangle + \zeta_{ja}^{\mathbf{x}} \quad (1 \leq a \leq K_j) \quad (6)$$

Based on Eq.(6), we can have the continuous version of $\boldsymbol{\Delta}_{\mathbf{x}}$ in Eq.(3) and then transformed data set $\tilde{\mathcal{D}}$ in Eq.(4).

Algorithm 1 summarizes the complete procedure of the proposed KRAM ap-
170 proach. Firstly, the original feature vector of each training example is enriched by two different k NN-augmented features: for discrete type, simple counting statistics derived from neighboring MDC examples is used, and for continuous type, extra bias is added into the simple counting statistics, and then a transformed MDC training set is gradually constructed (steps 1-17). After that, a
175 MDC model is trained over the transformed MDC training set $\tilde{\mathcal{D}}$ (step 18). Finally, for unseen instance, its class vector can be predicted based on the original features combined with augmented features (steps 19-21). In the remaining parts of this paper, we denote the discrete version as KRAM_d and the continuous version as KRAM_c respectively.

180 As shown in Algorithm 1, the k NN-augmented features generation procedure (steps 1-17) and the MDC predictive model induction procedure (step 18) are detached. In other words, the proposed KRAM approach is actually a meta-

Algorithm 1 The proposed KRAM approach.

Input: \mathcal{D} : MDC training set $\{(\mathbf{x}_i, \mathbf{y}_i) \mid 1 \leq i \leq m\}$

k : number of nearest neighbors considered

\mathcal{L} : MDC training algorithm

\mathbf{x}^* : unseen instance

Output: \mathbf{y}_* : predicted class vector for \mathbf{x}^*

```

1:  $\tilde{\mathcal{D}} = \emptyset$ ;
2: for  $i = 1$  to  $m$  do
3:   Identify  $k$  nearest neighbors of  $\mathbf{x}_i$  in  $\mathcal{D}$  and store their indices in  $\mathcal{N}(\mathbf{x}_i)$ ;
4:   for  $j = 1$  to  $q$  do
5:     for  $a = 1$  to  $K_j$  do
6:       switch type do
7:         case discrete:
8:           Calculate  $\delta_{ja}^{\mathbf{x}_i}$  according to Eq.(2);
9:         case continuous:
10:          Calculate  $\delta_{ja}^{\mathbf{x}_i}$  according to Eq.(6);
11:        end switch
12:      end for
13:      Set  $\boldsymbol{\delta}_j^{\mathbf{x}_i} = [\delta_{j1}^{\mathbf{x}_i}, \delta_{j2}^{\mathbf{x}_i}, \dots, \delta_{jK_j}^{\mathbf{x}_i}]^\top$ ;
14:    end for
15:    Set  $\boldsymbol{\Delta}_{\mathbf{x}_i} = [\boldsymbol{\delta}_1^{\mathbf{x}_i}, \boldsymbol{\delta}_2^{\mathbf{x}_i}, \dots, \boldsymbol{\delta}_q^{\mathbf{x}_i}]$ ;
16:     $\tilde{\mathcal{D}} = \tilde{\mathcal{D}} \cup (\tilde{\mathbf{x}}_i, \mathbf{y}_i)$ , where  $\tilde{\mathbf{x}}_i = [\mathbf{x}_i, \boldsymbol{\Delta}_{\mathbf{x}_i}]$ ;
17:  end for
18: Train MDC model  $f$  over  $\tilde{\mathcal{D}}$ , i.e.,  $f \leftarrow \mathcal{L}(\tilde{\mathcal{D}})$ ;
19: Identify  $k$  nearest neighbors of  $\mathbf{x}^*$  in  $\mathcal{D}$  and store their indices in  $\mathcal{N}(\mathbf{x}^*)$ ;
20: Augment  $\mathbf{x}^*$  with  $\boldsymbol{\Delta}_{\mathbf{x}^*}$  being calculated the same as training set, i.e.,  $\tilde{\mathbf{x}}^* = [\mathbf{x}^*, \boldsymbol{\Delta}_{\mathbf{x}^*}]$ ;
21: Return  $\mathbf{y}^* = f(\tilde{\mathbf{x}}^*)$ .

```

strategy for MDC model induction, where any existing MDC training algorithm (i.e., \mathcal{L} in step 18) can be employed to instantiate KRAM. Besides, both the two versions of k NN-augmented features designed in this paper should only be considered as a first attempt towards feature manipulation techniques for MDC and are not meant to be the best possible implementation among other feasible choices in future.

Generally speaking, KRAM embodies two major merits: 1) *Simplicity*: As shown in Algorithm 1, the KRAM is very succinct and can be implemented

easily. Specifically, the most time-consuming operation of KRAM is the k nearest neighbors identification process which has been well studied in k NN researches. The number of augmented features’ dimension equals $\sum_{j=1}^q K_j$ which is not large usually, so there will not be too much extra computation. 2) *Effectiveness*:
 195 experimental studies reported in Section 4 clearly validate the fact that KRAM can improve the generalization abilities of any off-the-shelf MDC approaches.

Table 1: The characteristics of the employed benchmark data sets.

Data Set	#Exam.	#Dim.	#Labels/Dim.	#Features [†]
Edm	154	2	3	16 <i>n</i>
Flare1	323	3	3,4,2	10 <i>x</i>
Song	785	3	3	98 <i>n</i>
WQplants	1060	7	4	16 <i>n</i>
WQanimals	1060	7	4	16 <i>n</i>
WaterQuality	1060	14	4	16 <i>n</i>
Voice	3136	2	4,2	19 <i>n</i>
Thyroid	9172	7	5,5,3,2,4,4,3	7 <i>n</i> , 20 <i>b</i> , 2 <i>x</i>
Flickr	12198	5	3,4,3,4,4	1536 <i>n</i>
Music	591	6	2	71 <i>n</i>
Enron	1677	10	2	1001 <i>b</i>
Image	2000	5	2	294 <i>n</i>
Scene	2407	6	2	294 <i>n</i>
Yeast	2417	14	2	103 <i>n</i>
Mediamill	41583	11	2	120 <i>n</i>

[†] n , x and b denote numeric, nominal and binary features respectively.

4. Experiments

4.1. Experimental Setup

4.1.1. Benchmark data sets

200 To validate the effectiveness of our proposed KRAM approach in improving the predictive abilities of existing MDC approaches, a total of 15 data sets have been used for performance comparison. Table 1 summarizes the characteristics of all benchmark data sets, including *the number of examples* (#Exam.),

the number of class spaces ($\#Dim.$), the number of class labels per class space
205 ($\#Labels/Dim.$),² and the number of features ($\#Features$).

The first nine benchmark data sets in Table 1 are collected from different
MDC tasks in real world:

- **Edm** aims at predicting control operations during electrical discharge ma-
chining process [26]. This data set includes totally 2 class spaces which
210 correspond to two parameters of controlling gap and flow respectively.
- **Flare1** aims at predicting the number of times certain types of solar flare
occurred within 24 hours period [27]. This data set includes totally 3
class spaces which correspond to common, moderate, and severe solar
flares respectively.
- 215 • **Song** aims at predicting different characteristics of songs. This data set
includes totally 3 class spaces which correspond to the scenarios, genre
and emotion of one song respectively. Besides, all songs are collected and
annotated by ourselves.
- **Water Quality** aims at predicting the amount of different species in Slove-
220 nian rivers [28]. This data set includes totally 14 class spaces which cor-
respond to 7 plants and 7 animals species respectively. By focusing on 7
plants species or the 7 animals species, then we have data sets **WQplants**
and **WQanimals** [29].
- **Voice** aims at predicting some characteristics of a piece of human voice.³
225 This data set includes totally 2 class spaces which correspond to mean
frequency range and speaker’s gender respectively [20].

²If all class spaces have the same number of class labels, then only this number is recorded;
If different class spaces have different number of class labels, the number of class labels for
each class space is recorded in turn.

³ <http://www.primaryobjects.com/2016/06/22/identifying-the-gender-of-a-voice-using-machine-learning>

- **Thyroid** aims at estimating types of thyroid problems based on personal information of patients [27]. This data set includes totally 7 class spaces which correspond to seven different diagnosis respectively.
- 230 • **Flickr** aims at predicting objects in MIRFLICKR25000 [30] which are re-annotated by ourselves according to MDC framework and just part of pictures are reserved. This data set includes totally 5 class spaces which correspond to sky, people, night, plant, indoor respectively.

The last six data sets in Table 1 are selected from multi-label classification tasks including audio classification: **Music** [11], text classification: **Enron**,⁴ image classification: **Image** [31], **Scene** [32], gene functional analysis: **Yeast** [33], and video classification: **mediamill** [34]. For these data sets, each class variable is binary-valued which is widely known as *label*.⁵

To the best of our knowledge, this paper employs more real-world MDC data sets than most state-of-the-art works on multi-dimensional classification [9, 11, 240 13]. Moreover, as shown in Table 1, the characteristics of all benchmark data sets are very diversified, e.g., the number of examples ranges from 154 to 41583, the number of features ranges from 10 to 1536, and the ratio of sum of all the numbers of class labels in each class space to the number of features ranges from 245 0.01 to 3.50 (i.e., the ratio of the number of augmented features to the number of original features). Therefore, the reported experimental results in this paper are quite comprehensive which can validate the effectiveness of KRAM more thoroughly.

4.1.2. Evaluation metrics

250 Given the test data set $\mathcal{S} = \{(\mathbf{x}_i, \mathbf{y}_i) \mid 1 \leq i \leq p\}$ with p MDC examples, where $\mathbf{y}_i = [y_{i1}, y_{i2}, \dots, y_{iq}]^\top \in \mathcal{Y}$ is the ground-truth class vector of \mathbf{x}_i . More-

⁴<http://mulan.sourceforge.net/datasets-mlc.html>

⁵For Enron, we just use 10 out of all 53 labels with most positive instances. For Mediamill, we just use 11 out of all 120 labels similar to Enron. And instances without relevant labels are removed.

over, let $f : \mathcal{X} \mapsto \mathcal{Y}$ denote the induced MDC predictive model which is to be evaluated, and let $\hat{\mathbf{y}}_i = f(\mathbf{x}_i) = [\hat{y}_{i1}, \hat{y}_{i2}, \dots, \hat{y}_{iq}]^\top$ be the predicted class vector for \mathbf{x}_i by function f . Then, $r^{(i)} = \sum_{j=1}^q \mathbb{1}[y_{ij} = \hat{y}_{ij}]$ denotes the number of class labels correctly classified by f . Based on these notations, the definitions of the three metrics employed in this paper are given as follows:

- *Hamming Score:*

$$\text{HScore}_{\mathcal{S}}(f) = \frac{1}{p} \sum_{i=1}^p \frac{1}{q} \cdot r^{(i)}$$

The value of this metric can be regarded as the probability of that any class label of one test example is correctly predicted by the induced MDC model.

- *Exact Match:*

$$\text{EMatch}_{\mathcal{S}}(f) = \frac{1}{p} \sum_{i=1}^p \mathbb{1}[r^{(i)} = q]$$

The value of this metric can be regarded as the probability of that all q class labels of one test example are correctly predicted simultaneously by the induced MDC model. Generally, the value of *exact match* might be rather low when the number of class spaces is large.

- *Sub-Exact Match:*

$$\text{SEMatch}_{\mathcal{S}}(f) = \frac{1}{p} \sum_{i=1}^p \mathbb{1}[r^{(i)} \geq q - 1]$$

The value of this metric can be regarded as the probability of that at least $q-1$ class labels of one test example are correctly predicted simultaneously by the induced MDC model. Obviously, it is a relaxed version of *exact match*.

4.1.3. Compared approaches

As stated before, the proposed KRAM approach is a meta-strategy to induce MDC predictive models. This means that any off-the-shelf MDC approaches can be employed to instantiate KRAM to further improve their generalization

performance. Here, a total of four well-established MDC approaches [11] serve this purpose:

- 275 • *Binary Relevance* (BR): This approach solves the MDC problem by training a number of multi-class classifiers independently, one per class space. In other words, BR completely ignores class dependencies.
- 280 • *Ensembles of Classifier Chains* (ECC): This approach solves the MDC problem by training a chain of multi-class classifiers, one per class space. Specifically, the subsequent multi-class classifiers in the chain are built by augmenting the feature space with the predictions of preceding multi-class classifiers. In other words, ECC models class dependencies via assuming chaining structures over class spaces. Besides, different base models in an ensemble of classifier chains consider different random chaining orders.
- 285 • *Ensembles of Class Powerset* (ECP): This approach solves the MDC problem by training a single multi-class classifier. Specifically, each distinct class combination in training set is treated as a new class label. In other words, ECP models class dependencies via powerset transformation. Besides, different base models in an ensemble of class powerset classifiers are built over different sub-training sets which are randomly sampled from the original training set.
- 290 • *Ensembles of Super Class classifiers* (ESC): This approach solves the MDC problem by partitioning class spaces into super-classes, where the partition process is fulfilled according to conditional dependencies among class spaces. In other words, ESC models class dependencies via the new generated super-classes. Specifically, different base models in an ensemble of super-class classifiers are built over different sub-training sets which are randomly sampled from the original training set.

For each ensemble approach (i.e., ECC, ECP and ESC), its base MDC model is induced over a sub-training set which contains 67% examples sampled from the original one randomly, and a total of 10 base MDC models are employed

in this paper [11]. Moreover, majority voting strategy is used to combine the predictions of all base MDC models for each example.

As implementing each MDC approach also necessitate a base multi-class classifier, both support vector machine (SVM) and Naïve Bayes (NB) are investigated in this paper. Specifically, SVM is implemented by LIBSVM [35] where
305 the type of kernel function is linear and the regularization parameter C is set to 1. NB takes the common implementation where Gaussian pdf is used for continuous features and frequency counting with Laplacian correction is used for discrete features. Besides, the only parameter k for KRAM, which denotes
310 the number of nearest neighbors considered, is set to 8 when conducting comparative experiments.

Let KRAM- \mathcal{A} be the instantiation of KRAM with \mathcal{A} , where \mathcal{A} denotes one of the compared approaches, i.e., $\mathcal{A} \in \{\text{BR}, \text{ECC}, \text{ECP}, \text{ESC}\}$. To show whether KRAM could improve the generalization abilities of MDC approaches, our aim
315 is to compare predictive performance of KRAM- \mathcal{A} against \mathcal{A} . On each configuration (each approach in terms of each metric on each data set), ten-fold cross-validation is conducted where both the mean metric value and corresponding standard deviation are recorded for performance comparison.

4.2. Experimental Results

The detailed experimental results are reported in Tables 2 to 5 for all MD-
320 C approaches and their KRAM counterparts in terms of *hamming score*, *exact match*, and *sub-exact match* respectively. Based on the ten-fold cross-validation results of each data set, pairwise t -test at 0.05 significance level is further conducted to show whether the performance of each KRAM counterpart is signif-
325 icantly different to its corresponding MDC approach with different multi-class classifiers (SVM or NB). Accordingly, the resulting win/tie/loss counts are summarized in Table 6 in terms of each evaluation metric over the 15 benchmark data sets.

Based on the experimental results mentioned above, it is interesting to obtain
330 the following observations:

Table 2: Predictive performance of BR and its two version KRAM counterparts (mean \pm std. deviation) in terms of all three evaluation metrics. In addition, \bullet/\circ indicates whether KRAM_d-BR or KRAM_c-BR is statistically superior/inferior to BR on each data set, and results in bold face of KRAM_d-BR or KRAM_c-BR indicate that current KRAM version is statistically superior to the other version on each data set (pairwise t -test at 0.05 significance level).

(a) Multi-class classifier: SVM

Data Set	Hamming Score			Exact Match			Sub-Exact Match		
	BR	KRAM _d -BR	KRAM _c -BR	BR	KRAM _d -BR	KRAM _c -BR	BR	KRAM _d -BR	KRAM _c -BR
Edm	.689 \pm .070	.734 \pm .083 \bullet	.760\pm.095\bullet	.442 \pm .125	.521 \pm .141 \bullet	.560\pm.150\bullet	.935 \pm .061	.947 \pm .076	.960 \pm .072
Flare1	.922 \pm .034	.922 \pm .033	.921 \pm .034	.821 \pm .073	.818 \pm .072	.814 \pm .077	.947 \pm .039	.951 \pm .036	.951 \pm .036
Song	.793 \pm .023	.787 \pm .023 \circ	.786 \pm .024	.479 \pm .059	.476 \pm .050	.476 \pm .049	.903 \pm .033	.888 \pm .046	.885 \pm .054
WQpla.	.657 \pm .016	.664 \pm .013	.667 \pm .014 \bullet	.097 \pm .033	.099 \pm .034	.104\pm.032	.287 \pm .055	.300 \pm .042	.299 \pm .045
WQani.	.630 \pm .014	.635 \pm .012 \bullet	.637 \pm .012 \bullet	.058 \pm .022	.063 \pm .014	.065 \pm .016	.229 \pm .034	.232 \pm .030	.233 \pm .031
WQ	.644 \pm .013	.646 \pm .010	.648 \pm .011	.007 \pm .008	.008 \pm .007	.008 \pm .007	.051 \pm .024	.053 \pm .017	.058 \pm .020
Voice	.964 \pm .007	.957 \pm .008 \circ	.957 \pm .007 \circ	.929 \pm .014	.915 \pm .016 \circ	.915 \pm .014 \circ	.999 \pm .002	.999 \pm .002	.998 \pm .002
Thyroid	.965 \pm .002	.969 \pm .003 \bullet	.969 \pm .003 \bullet	.773 \pm .015	.800 \pm .018 \bullet	.801 \pm .016 \bullet	.982 \pm .004	.983 \pm .004	.983 \pm .004
Flickr	.791 \pm .005	.790 \pm .006	.790 \pm .006	.313 \pm .014	.310 \pm .018	.310 \pm .020	.720 \pm .014	.719 \pm .014	.718 \pm .015
Music	.808 \pm .023	.818 \pm .022 \bullet	.818 \pm .023	.272 \pm .075	.331 \pm .082 \bullet	.321 \pm .079 \bullet	.674 \pm .067	.682 \pm .054	.686 \pm .061
Enron	.808 \pm .010	.807 \pm .010	.805 \pm .010	.179 \pm .036	.182 \pm .031	.181 \pm .030	.416 \pm .034	.397 \pm .041	.398 \pm .032 \circ
Image	.828 \pm .010	.841 \pm .011 \bullet	.842 \pm .012 \bullet	.394 \pm .028	.459 \pm .033 \bullet	.464 \pm .039 \bullet	.782 \pm .031	.783 \pm .027	.785 \pm .029
Scene	.895 \pm .009	.918 \pm .008 \bullet	.918 \pm .008 \bullet	.530 \pm .035	.651 \pm .038 \bullet	.649 \pm .035 \bullet	.855 \pm .018	.867 \pm .020	.867 \pm .018
Yeast	.801 \pm .006	.811 \pm .007 \bullet	.812 \pm .006 \bullet	.151 \pm .017	.199 \pm .015 \bullet	.199 \pm .015 \bullet	.269 \pm .029	.307 \pm .020 \bullet	.310 \pm .011 \bullet
Mediamill	.830 \pm .001	.850 \pm .001 \bullet	.851\pm.001\bullet	.148 \pm .005	.239 \pm .009 \bullet	.247\pm.007\bullet	.411 \pm .008	.500 \pm .006 \bullet	.507\pm.006\bullet

(b) Multi-class classifier: NB

Data Set	Hamming Score			Exact Match			Sub-Exact Match		
	BR	KRAM _d -BR	KRAM _c -BR	BR	KRAM _d -BR	KRAM _c -BR	BR	KRAM _d -BR	KRAM _c -BR
Edm	.677 \pm .096	.680 \pm .088	.693 \pm .073	.432 \pm .166	.445 \pm .153	.444 \pm .151	.922 \pm .074	.916 \pm .060	.941 \pm .038
Flare1	.886 \pm .061	.872\pm.051	.838 \pm .059 \circ	.774 \pm .099	.756\pm.095	.722 \pm .093 \circ	.910 \pm .066	.895\pm.055	.858 \pm .068 \circ
Song	.626 \pm .038	.629 \pm .034	.623 \pm .033	.238 \pm .054	.224 \pm .050	.213 \pm .050 \circ	.678 \pm .071	.695 \pm .068	.689 \pm .065
WQpla.	.397 \pm .028	.506\pm.033\bullet	.475 \pm .037 \bullet	.001 \pm .003	.036\pm.026\bullet	.016 \pm .018 \bullet	.018 \pm .012	.113\pm.040\bullet	.082 \pm .035 \bullet
WQani.	.381 \pm .021	.419\pm.019\bullet	.400 \pm .022 \bullet	.004 \pm .009	.008\pm.010	.003 \pm .006	.041 \pm .016	.049 \pm .019	.040 \pm .020
WQ	.389 \pm .017	.488\pm.022\bullet	.443 \pm .016 \bullet	.000 \pm .000	.000 \pm .000	.000 \pm .000	.000 \pm .000	.003 \pm .005	.001 \pm .003
Voice	.882 \pm .008	.921 \pm .008 \bullet	.939\pm.010\bullet	.782 \pm .015	.847 \pm .016 \bullet	.880\pm.020\bullet	.982 \pm .006	.996 \pm .003 \bullet	.998\pm.003\bullet
Thyroid	.926 \pm .005	.925\pm.003	.704 \pm .030 \circ	.580 \pm .027	.575\pm.015	.047 \pm .013 \circ	.916 \pm .011	.912\pm.009	.307 \pm .053 \circ
Flickr	.648 \pm .007	.654 \pm .007 \bullet	.659\pm.007\bullet	.139 \pm .011	.143 \pm .011 \bullet	.147\pm.011\bullet	.436 \pm .013	.444 \pm .012 \bullet	.453\pm.014\bullet
Music	.743 \pm .018	.761 \pm .023 \bullet	.760 \pm .027 \bullet	.206 \pm .043	.218 \pm .058	.208 \pm .060	.552 \pm .057	.591 \pm .050 \bullet	.601 \pm .057 \bullet
Enron	.551 \pm .012	.576 \pm .015 \bullet	.574 \pm .014 \bullet	.025 \pm .017	.027 \pm .017	.029 \pm .017	.109 \pm .022	.119 \pm .026	.114 \pm .029
Image	.573 \pm .016	.586 \pm .018 \bullet	.595\pm.016\bullet	.069 \pm .016	.074 \pm .021 \bullet	.080\pm.019\bullet	.255 \pm .028	.279 \pm .034 \bullet	.297\pm.035\bullet
Scene	.763 \pm .009	.777 \pm .009 \bullet	.791\pm.011\bullet	.177 \pm .023	.198 \pm .022 \bullet	.238\pm.034\bullet	.561 \pm .021	.591 \pm .026 \bullet	.605\pm.024\bullet
Yeast	.699 \pm .010	.695 \pm .014	.704\pm.013\bullet	.095 \pm .018	.115 \pm .018 \bullet	.115 \pm .022 \bullet	.149 \pm .020	.182 \pm .027 \bullet	.175 \pm .014 \bullet
Mediamill	.620 \pm .001	.659 \pm .002 \bullet	.673\pm.001\bullet	.012 \pm .001	.048 \pm .003 \bullet	.067\pm.005\bullet	.086 \pm .005	.135 \pm .005 \bullet	.161\pm.007\bullet

Table 3: Predictive performance of ECC and its two version KRAM counterparts (mean \pm std. deviation) in terms of all three evaluation metrics. In addition, \bullet / \circ indicates whether KRAM_d-ECC or KRAM_c-ECC is statistically superior/inferior to ECC on each data set, and results in bold face of KRAM_d-ECC or KRAM_c-ECC indicate that current KRAM version is statistically superior to the other version on each data set (pairwise t -test at 0.05 significance level).

(a) Multi-class classifier: SVM

Data Set	Hamming Score			Exact Match			Sub-Exact Match		
	ECC	KRAM _d -ECC	KRAM _c -ECC	ECC	KRAM _d -ECC	KRAM _c -ECC	ECC	KRAM _d -ECC	KRAM _c -ECC
Edm	.695 \pm .065	.769 \pm .087 \bullet	.756 \pm .090 \bullet	.454 \pm .123	.598 \pm .169 \bullet	.585 \pm .147 \bullet	.935 \pm .069	.940 \pm .058	.928 \pm .066
Flare1	.922 \pm .034	.922 \pm .034	.922 \pm .035	.817 \pm .078	.818 \pm .073	.818 \pm .078	.951 \pm .036	.951 \pm .036	.951 \pm .036
Song	.790 \pm .024	.788 \pm .026	.788 \pm .027	.481 \pm .057	.476 \pm .051	.479 \pm .048	.891 \pm .036	.891 \pm .047	.888 \pm .051
WQpla.	.654 \pm .016	.663 \pm .014 \bullet	.666\pm.014\bullet	.093 \pm .037	.105 \pm .037	.105 \pm .038	.283 \pm .049	.295 \pm .044	.297 \pm .043
WQani.	.630 \pm .014	.637 \pm .014 \bullet	.637 \pm .012 \bullet	.061 \pm .023	.064 \pm .010	.067 \pm .015	.229 \pm .032	.241 \pm .040	.238 \pm .037
WQ	.643 \pm .013	.644 \pm .013	.646 \pm .013	.006 \pm .008	.009 \pm .006	.009 \pm .004	.050 \pm .023	.048 \pm .018	.050 \pm .023
Voice	.961 \pm .008	.953 \pm .009 \circ	.953 \pm .009 \circ	.923 \pm .016	.908 \pm .017 \circ	.908 \pm .019	.998 \pm .002	.998 \pm .003	.998 \pm .003
Thyroid	.965 \pm .002	.969 \pm .003 \bullet	.969\pm.003\bullet	.772 \pm .014	.800 \pm .016 \bullet	.802 \pm .017 \bullet	.981 \pm .004	.982 \pm .004	.983 \pm .004
Flickr	.797 \pm .004	.797 \pm .005	.797 \pm .005	.328 \pm .013	.325 \pm .014	.325 \pm .017	.730 \pm .015	.732 \pm .014	.732 \pm .015
Music	.814 \pm .025	.810 \pm .022	.813 \pm .021	.346 \pm .079	.343 \pm .078	.352 \pm .086	.676 \pm .064	.677 \pm .051	.676 \pm .053
Enron	.824 \pm .010	.822\pm.011	.819 \pm .012 \circ	.215 \pm .030	.203 \pm .037	.199 \pm .027 \circ	.461 \pm .037	.455 \pm .037	.443 \pm .040
Image	.831 \pm .012	.844 \pm .012 \bullet	.846 \pm .012 \bullet	.479 \pm .033	.522 \pm .036 \bullet	.525 \pm .036 \bullet	.730 \pm .033	.745 \pm .031	.749 \pm .036
Scene	.905 \pm .011	.921 \pm .008 \bullet	.921 \pm .007 \bullet	.649 \pm .035	.708 \pm .026 \bullet	.710 \pm .025 \bullet	.796 \pm .030	.825 \pm .020 \bullet	.825 \pm .018 \bullet
Yeast	.797 \pm .007	.808 \pm .007 \bullet	.808 \pm .008 \bullet	.207 \pm .014	.252 \pm .014 \bullet	.253 \pm .018 \bullet	.288 \pm .023	.316 \pm .022 \bullet	.319 \pm .025 \bullet
Mediamill	.830 \pm .002	.849 \pm .001 \bullet	.851\pm.001\bullet	.188 \pm .008	.267 \pm .006 \bullet	.276\pm.006\bullet	.434 \pm .006	.513 \pm .005 \bullet	.521\pm.005\bullet

(b) Multi-class classifier: NB

Data Set	Hamming Score			Exact Match			Sub-Exact Match		
	ECC	KRAM _d -ECC	KRAM _c -ECC	ECC	KRAM _d -ECC	KRAM _c -ECC	ECC	KRAM _d -ECC	KRAM _c -ECC
Edm	.690 \pm .084	.674 \pm .097	.689 \pm .070	.451 \pm .145	.438 \pm .162	.444 \pm .148	.929 \pm .064	.909 \pm .062	.935 \pm .044
Flare1	.883 \pm .059	.875\pm.053	.838 \pm .062 \circ	.774 \pm .087	.771\pm.088	.737 \pm .095	.904 \pm .073	.889\pm.060	.852 \pm .064 \circ
Song	.621 \pm .036	.623\pm.034	.613 \pm .037	.228 \pm .036	.219\pm.043	.191 \pm .047 \circ	.671 \pm .068	.683 \pm .066	.680 \pm .069
WQpla.	.353 \pm .033	.494\pm.038\bullet	.444 \pm .038 \bullet	.001 \pm .003	.035\pm.018\bullet	.024 \pm .019 \bullet	.013 \pm .010	.123\pm.037\bullet	.073 \pm .030 \bullet
WQani.	.377 \pm .024	.416\pm.020\bullet	.395 \pm .021 \bullet	.007 \pm .008	.006 \pm .007	.004 \pm .007	.039 \pm .016	.049 \pm .015 \bullet	.042 \pm .018
WQ	.360 \pm .020	.487\pm.021\bullet	.431 \pm .018 \bullet	.000 \pm .000	.000 \pm .000	.000 \pm .000	.000 \pm .000	.001 \pm .003	.000 \pm .000
Voice	.880 \pm .009	.921 \pm .008 \bullet	.939\pm.010\bullet	.780 \pm .015	.847 \pm .015 \bullet	.879\pm.020\bullet	.980 \pm .007	.996 \pm .003 \bullet	.998\pm.003\bullet
Thyroid	.926 \pm .007	.929\pm.004	.758 \pm .016 \circ	.593 \pm .026	.592\pm.022	.058 \pm .013 \circ	.906 \pm .020	.922\pm.008\bullet	.392 \pm .063 \circ
Flickr	.649 \pm .007	.655 \pm .007 \bullet	.660\pm.007\bullet	.140 \pm .012	.143 \pm .011 \bullet	.147\pm.011\bullet	.438 \pm .014	.447 \pm .014 \bullet	.455\pm.015\bullet
Music	.745 \pm .020	.761 \pm .023 \bullet	.763 \pm .026 \bullet	.230 \pm .058	.221 \pm .065	.221 \pm .064	.557 \pm .051	.603 \pm .048 \bullet	.608 \pm .049 \bullet
Enron	.551 \pm .011	.573 \pm .015 \bullet	.573 \pm .013 \bullet	.027 \pm .016	.029 \pm .016	.029 \pm .017	.105 \pm .026	.115 \pm .029	.115 \pm .027 \bullet
Image	.576 \pm .014	.587 \pm .014 \bullet	.596\pm.017\bullet	.069 \pm .019	.074 \pm .020 \bullet	.078\pm.021\bullet	.261 \pm .028	.283 \pm .033 \bullet	.296\pm.036\bullet
Scene	.767 \pm .010	.780 \pm .010 \bullet	.793\pm.010\bullet	.181 \pm .024	.200 \pm .021 \bullet	.238\pm.030\bullet	.569 \pm .027	.595 \pm .031 \bullet	.612\pm.023\bullet
Yeast	.696 \pm .009	.698 \pm .013	.704\pm.013\bullet	.102 \pm .016	.125 \pm .024 \bullet	.124 \pm .022 \bullet	.163 \pm .020	.193 \pm .027 \bullet	.188 \pm .016 \bullet
Mediamill	.628 \pm .001	.667 \pm .002 \bullet	.679\pm.002\bullet	.013 \pm .001	.050 \pm .003 \bullet	.069\pm.005\bullet	.095 \pm .004	.142 \pm .005 \bullet	.167\pm.007\bullet

Table 4: Predictive performance of ECP and its two version KRAM counterparts (mean \pm std. deviation) in terms of all three evaluation metrics. In addition, \bullet/\circ indicates whether KRAM_d-ECP or KRAM_c-ECP is statistically superior/inferior to ECP on each data set, and results in bold face of KRAM_d-ECP or KRAM_c-ECP indicate that current KRAM version is statistically superior to the other version on each data set (pairwise t -test at 0.05 significance level).

(a) Multi-class classifier: SVM

Data Set	Hamming Score			Exact Match			Sub-Exact Match		
	ECP	KRAM _d -ECP	KRAM _c -ECP	ECP	KRAM _d -ECP	KRAM _c -ECP	ECP	KRAM _d -ECP	KRAM _c -ECP
Edm	.721 \pm .082	.763 \pm .107 \bullet	.766 \pm .094 \bullet	.559 \pm .136	.612 \pm .170	.618 \pm .140 \bullet	.883 \pm .074	.915 \pm .075	.915 \pm .075
Flare1	.921 \pm .036	.922 \pm .034	.923 \pm .035	.817 \pm .078	.821 \pm .073	.821 \pm .073	.947 \pm .039	.947 \pm .039	.951 \pm .042
Song	.786 \pm .029	.781 \pm .028	.781 \pm .035	.484 \pm .054	.467 \pm .059	.470 \pm .071	.878 \pm .040	.877 \pm .040	.874 \pm .043
WQpla.	.647 \pm .015	.585 \pm .027 \circ	.586 \pm .026 \circ	.093 \pm .028	.067 \pm .029 \circ	.065 \pm .029 \circ	.281 \pm .049	.187 \pm .040 \circ	.191 \pm .042 \circ
WQani.	.629 \pm .013	.556 \pm .014 \circ	.550 \pm .015 \circ	.065 \pm .018	.029 \pm .011 \circ	.024 \pm .012 \circ	.230 \pm .032	.151 \pm .030 \circ	.136 \pm .022 \circ
WQ	.628 \pm .015	.557 \pm .010 \circ	.555 \pm .012 \circ	.001 \pm .003	.004 \pm .005	.004 \pm .005	.035 \pm .018	.019 \pm .016	.016 \pm .015 \circ
Voice	.955 \pm .013	.950 \pm .010	.950 \pm .010	.912 \pm .025	.903 \pm .020	.903 \pm .019	.998 \pm .003	.998 \pm .003	.997 \pm .004
Thyroid	.965 \pm .002	.968 \pm .002 \bullet	.969 \pm .003 \bullet	.773 \pm .014	.802 \pm .015 \bullet	.804\pm.015\bullet	.981 \pm .005	.979 \pm .003 \circ	.979 \pm .003 \circ
Flickr	.772 \pm .004	.760 \pm .006 \circ	.761 \pm .006 \circ	.297 \pm .012	.281 \pm .009 \circ	.281 \pm .010 \circ	.680 \pm .011	.658 \pm .016 \circ	.658 \pm .018 \circ
Music	.799 \pm .032	.802 \pm .025	.800 \pm .030	.343 \pm .076	.341 \pm .073	.343 \pm .083	.640 \pm .064	.659 \pm .066	.650 \pm .072
Enron	.830 \pm .008	.824 \pm .009 \circ	.822 \pm .010 \circ	.235 \pm .029	.224 \pm .026	.224 \pm .033	.482 \pm .021	.459 \pm .043	.462 \pm .034
Image	.832 \pm .012	.842 \pm .009 \bullet	.841 \pm .011 \bullet	.513 \pm .024	.540 \pm .024 \bullet	.535 \pm .029 \bullet	.710 \pm .036	.727 \pm .029 \bullet	.725 \pm .032
Scene	.914 \pm .009	.925 \pm .008 \bullet	.925 \pm .007 \bullet	.700 \pm .029	.731 \pm .029 \bullet	.729 \pm .029 \bullet	.796 \pm .028	.825 \pm .018 \bullet	.825 \pm .021 \bullet
Yeast	.795 \pm .007	.795 \pm .007	.794 \pm .010	.252 \pm .012	.262 \pm .018	.264 \pm .018 \bullet	.304 \pm .020	.317 \pm .018	.317 \pm .022 \bullet
Mediamill	.818 \pm .002	.841 \pm .002 \bullet	.843\pm.002\bullet	.218 \pm .008	.286 \pm .004 \bullet	.293\pm.003\bullet	.430 \pm .008	.506 \pm .007 \bullet	.513\pm.008\bullet

(b) Multi-class classifier: NB

Data Set	Hamming Score			Exact Match			Sub-Exact Match		
	ECP	KRAM _d -ECP	KRAM _c -ECP	ECP	KRAM _d -ECP	KRAM _c -ECP	ECP	KRAM _d -ECP	KRAM _c -ECP
Edm	.731 \pm .062	.722 \pm .089	.699 \pm .076	.554 \pm .112	.548 \pm .120	.528 \pm .118	.909 \pm .047	.896 \pm .081	.870 \pm .069
Flare1	.908 \pm .045	.903\pm.046	.867 \pm .048 \circ	.790 \pm .081	.777\pm.084	.734 \pm .094 \circ	.941 \pm .057	.938\pm.057	.883 \pm .059 \circ
Song	.674 \pm .044	.684\pm.042	.674 \pm .044	.311 \pm .053	.317 \pm .051	.308 \pm .057	.733 \pm .079	.749 \pm .080	.733 \pm .077
WQpla.	.607 \pm .015	.647\pm.019\bullet	.609 \pm .028	.034 \pm .021	.067\pm.038\bullet	.043 \pm .028	.175 \pm .043	.258\pm.056\bullet	.208 \pm .049 \bullet
WQani.	.590 \pm .020	.625\pm.017\bullet	.598 \pm .016	.020 \pm .014	.042\pm.016\bullet	.030 \pm .011 \bullet	.143 \pm .054	.221\pm.049\bullet	.188 \pm .031 \bullet
WQ	.599 \pm .018	.597\pm.018	.562 \pm .019 \circ	.008 \pm .009	.004 \pm .007 \circ	.007 \pm .008	.032 \pm .024	.033 \pm .020	.027 \pm .018
Voice	.903 \pm .010	.927 \pm .011 \bullet	.937\pm.008\bullet	.811 \pm .019	.857 \pm .022 \bullet	.877\pm.018\bullet	.995 \pm .004	.997 \pm .003	.997 \pm .004
Thyroid	.966 \pm .003	.963\pm.003\circ	.948 \pm .002 \circ	.793 \pm .017	.768\pm.015\circ	.678 \pm .018 \circ	.974 \pm .005	.973\pm.005	.962 \pm .006 \circ
Flickr	.714 \pm .007	.717 \pm .008 \bullet	.719\pm.007\bullet	.197 \pm .006	.201 \pm .006 \bullet	.200 \pm .006	.563 \pm .015	.568 \pm .016 \bullet	.571\pm.014\bullet
Music	.770 \pm .029	.784 \pm .019 \bullet	.789 \pm .018 \bullet	.249 \pm .078	.281 \pm .073 \bullet	.277 \pm .077 \bullet	.591 \pm .071	.617 \pm .053	.635 \pm .036 \bullet
Enron	.777 \pm .011	.784 \pm .009 \bullet	.788 \pm .009 \bullet	.166 \pm .034	.177 \pm .033	.176 \pm .036	.324 \pm .035	.343 \pm .034 \bullet	.358\pm.031\bullet
Image	.746 \pm .012	.754 \pm .011 \bullet	.759 \pm .008 \bullet	.285 \pm .022	.302 \pm .022 \bullet	.315\pm.020\bullet	.597 \pm .034	.612 \pm .033 \bullet	.622 \pm .028 \bullet
Scene	.867 \pm .011	.875 \pm .013 \bullet	.883\pm.012\bullet	.550 \pm .030	.575 \pm .040 \bullet	.591\pm.034\bullet	.693 \pm .031	.713 \pm .036 \bullet	.736\pm.033\bullet
Yeast	.773 \pm .011	.787 \pm .008 \bullet	.784 \pm .009 \bullet	.203 \pm .018	.240 \pm .024 \bullet	.234 \pm .022 \bullet	.258 \pm .022	.293 \pm .022 \bullet	.292 \pm .029 \bullet
Mediamill	.656 \pm .005	.689 \pm .004 \bullet	.692\pm.005\bullet	.031 \pm .003	.056\pm.003\bullet	.053 \pm .004 \bullet	.097 \pm .008	.147 \pm .008 \bullet	.145 \pm .008 \bullet

Table 5: Predictive performance of ESC and its two version KRAM counterparts (mean \pm std. deviation) in terms of all three evaluation metrics. In addition, \bullet/\circ indicates whether KRAM_d-ESC or KRAM_c-ESC is statistically superior/inferior to ESC on each data set, and results in bold face of KRAM_d-ESC or KRAM_c-ESC indicate that current KRAM version is statistically superior to the other version on each data set (pairwise t -test at 0.05 significance level).

(a) Multi-class classifier: SVM

Data Set	Hamming Score			Exact Match			Sub-Exact Match		
	ESC	KRAM _d -ESC	KRAM _c -ESC	ESC	KRAM _d -ESC	KRAM _c -ESC	ESC	KRAM _d -ESC	KRAM _c -ESC
Edm	.701 \pm .079	.751 \pm .102 \bullet	.766 \pm .105 \bullet	.513 \pm .122	.592 \pm .165 \bullet	.624 \pm .164 \bullet	.890 \pm .076	.909 \pm .070	.909 \pm .070
Flare1	.923 \pm .033	.923 \pm .036	.923 \pm .035	.821 \pm .073	.821 \pm .073	.824 \pm .073	.951 \pm .036	.951 \pm .042	.947 \pm .039
Song	.790 \pm .030	.789 \pm .029	.787 \pm .029	.480 \pm .067	.481 \pm .058	.480 \pm .057	.893 \pm .038	.888 \pm .047	.884 \pm .047
WQpla.	.651 \pm .016	.664 \pm .015 \bullet	.665 \pm .016 \bullet	.094 \pm .038	.104 \pm .039	.101 \pm .037	.284 \pm .050	.283 \pm .049	.299\pm.049
WQani.	.630 \pm .014	.636 \pm .016	.633 \pm .013	.062 \pm .021	.065 \pm .018	.058 \pm .017	.232 \pm .033	.239 \pm .041	.230 \pm .039
WQ	.641 \pm .013	.638 \pm .017	.635 \pm .014	.006 \pm .008	.012 \pm .010	.009 \pm .006	.046 \pm .022	.050 \pm .022	.048 \pm .017
Voice	.961 \pm .008	.953 \pm .009 \circ	.953 \pm .010 \circ	.924 \pm .016	.908 \pm .018 \circ	.908 \pm .020 \circ	.998 \pm .002	.998 \pm .003	.998 \pm .003
Thyroid	.965 \pm .002	.969 \pm .002 \bullet	.969 \pm .003 \bullet	.771 \pm .014	.801 \pm .015 \bullet	.802 \pm .016 \bullet	.982 \pm .004	.981 \pm .004	.981 \pm .004
Flickr	.791 \pm .004	.784 \pm .005 \circ	.783 \pm .004 \circ	.320 \pm .011	.314 \pm .011 \circ	.313 \pm .012	.718 \pm .009	.704 \pm .013 \circ	.700 \pm .010 \circ
Music	.809 \pm .022	.810 \pm .028	.810 \pm .025	.330 \pm .069	.352 \pm .093	.335 \pm .086	.669 \pm .062	.671 \pm .065	.676 \pm .057
Enron	.833 \pm .009	.829 \pm .009	.830 \pm .011	.224 \pm .039	.216 \pm .032	.211 \pm .028	.485 \pm .038	.470 \pm .041	.479 \pm .040
Image	.833 \pm .008	.845 \pm .013 \bullet	.843 \pm .015 \bullet	.494 \pm .025	.529 \pm .034 \bullet	.528 \pm .040 \bullet	.719 \pm .028	.745 \pm .036 \bullet	.738 \pm .038 \bullet
Scene	.910 \pm .013	.923 \pm .008 \bullet	.924 \pm .008 \bullet	.668 \pm .045	.720 \pm .030 \bullet	.722 \pm .028 \bullet	.799 \pm .032	.825 \pm .022 \bullet	.827 \pm .027 \bullet
Yeast	.800 \pm .006	.807 \pm .007 \bullet	.807 \pm .007 \bullet	.236 \pm .019	.258 \pm .018 \bullet	.260 \pm .022 \bullet	.309 \pm .028	.320 \pm .025 \bullet	.320 \pm .024 \bullet
Mediamill	.824 \pm .002	.845 \pm .002 \bullet	.848\pm.001\bullet	.207 \pm .006	.277 \pm .005 \bullet	.285\pm.004\bullet	.434 \pm .007	.508 \pm .008 \bullet	.517\pm.008\bullet

(b) Multi-class classifier: NB

Data Set	Hamming Score			Exact Match			Sub-Exact Match		
	ESC	KRAM _d -ESC	KRAM _c -ESC	ESC	KRAM _d -ESC	KRAM _c -ESC	ESC	KRAM _d -ESC	KRAM _c -ESC
Edm	.674 \pm .095	.674 \pm .101	.696 \pm .064	.432 \pm .166	.438 \pm .162	.450 \pm .140	.915 \pm .063	.909 \pm .062	.941 \pm .049
Flare1	.896 \pm .059	.892\pm.053	.857 \pm .058 \circ	.780 \pm .093	.768\pm.086	.728 \pm .094 \circ	.929 \pm .064	.926\pm.060	.870 \pm .069 \circ
Song	.646 \pm .031	.666 \pm .037 \bullet	.662 \pm .037 \bullet	.274 \pm .047	.304 \pm .054 \bullet	.290 \pm .063	.692 \pm .066	.719 \pm .067 \bullet	.722 \pm .067 \bullet
WQpla.	.442 \pm .034	.549\pm.031\bullet	.487 \pm .030 \bullet	.001 \pm .003	.040\pm.025\bullet	.017 \pm .017 \bullet	.042 \pm .019	.133\pm.031\bullet	.082 \pm .026 \bullet
WQani.	.577 \pm .022	.598\pm.013\bullet	.574 \pm .022	.024 \pm .018	.026 \pm .023	.026 \pm .019	.139 \pm .050	.167 \pm .045	.136 \pm .043
WQ	.609 \pm .017	.609\pm.017	.553 \pm .023 \circ	.002 \pm .004	.002 \pm .004	.001 \pm .003	.023 \pm .013	.025 \pm .017	.011 \pm .012
Voice	.881 \pm .012	.922 \pm .008 \bullet	.939\pm.010\bullet	.782 \pm .020	.847 \pm .015 \bullet	.880\pm.020\bullet	.981 \pm .008	.996 \pm .003 \bullet	.998\pm.003\bullet
Thyroid	.958 \pm .004	.952\pm.006\circ	.945 \pm .003 \circ	.738 \pm .022	.703\pm.036\circ	.660 \pm .015 \circ	.970 \pm .007	.966\pm.006	.960 \pm .007 \circ
Flickr	.703 \pm .008	.705 \pm .009	.710\pm.010\bullet	.178 \pm .013	.181 \pm .014	.186 \pm .016 \bullet	.533 \pm .017	.537 \pm .015	.546\pm.020\bullet
Music	.738 \pm .023	.764 \pm .030 \bullet	.762 \pm .026 \bullet	.210 \pm .070	.242 \pm .089 \bullet	.218 \pm .063	.524 \pm .039	.581 \pm .082 \bullet	.603 \pm .051 \bullet
Enron	.768 \pm .016	.776 \pm .018	.766 \pm .022	.098 \pm .029	.106 \pm .039	.080 \pm .045	.302 \pm .049	.337\pm.051	.288 \pm .078
Image	.593 \pm .017	.608 \pm .015 \bullet	.618\pm.020\bullet	.069 \pm .021	.074 \pm .021	.084\pm.021\bullet	.289 \pm .032	.315 \pm .029 \bullet	.332\pm.035\bullet
Scene	.866 \pm .010	.868 \pm .013	.878\pm.015\bullet	.541 \pm .024	.528 \pm .046	.525 \pm .072	.703 \pm .031	.733 \pm .038 \bullet	.793\pm.021\bullet
Yeast	.716 \pm .006	.743\pm.006\bullet	.731 \pm .008 \bullet	.110 \pm .014	.154\pm.015\bullet	.135 \pm .022 \bullet	.167 \pm .019	.217\pm.022\bullet	.195 \pm .021 \bullet
Mediamill	.655 \pm .003	.693 \pm .008 \bullet	.695\pm.009\bullet	.038 \pm .004	.071 \pm .005 \bullet	.076\pm.004\bullet	.114 \pm .006	.177 \pm .009 \bullet	.184\pm.007\bullet

Table 6: Win/tie/loss counts of pairwise t -test (at 0.05 significance level) between each MD-C approach and its KRAM counterpart in terms of *hamming score* (HScore), *exact match* (EMatch), and *sub-exact match* (SEMatch).

	multi-class classifier: SVM			multi-class classifier: NB			In Total
	HScore	EMatch	SEMatch	HScore	EMatch	SEMatch	
KRAM _d -BR against BR	8/5/2	7/7/1	2/13/0	10/5/0	7/8/0	8/7/0	42/45/3
KRAM _d -ECC against ECC	8/6/1	6/8/1	3/12/0	10/5/0	7/8/0	10/5/0	44/44/2
KRAM _d -ECP against ECP	5/5/5	4/8/3	3/8/4	10/4/1	9/4/2	8/7/0	39/36/15
KRAM _d -ESC against ESC	7/6/2	6/7/2	4/10/1	8/6/1	6/8/1	8/7/0	38/44/7
KRAM _c -BR against BR	8/6/1	7/7/1	2/12/1	11/2/2	7/5/3	8/5/2	43/37/10
KRAM _c -ECC against ECC	8/5/2	6/8/1	3/12/0	11/2/2	7/6/2	9/4/2	44/37/9
KRAM _c -ECP against ECP	5/5/5	6/6/3	3/7/5	8/4/3	7/6/2	9/4/2	38/32/20
KRAM _c -ESC against ESC	7/6/2	6/8/1	4/10/1	9/3/3	6/7/2	9/4/2	41/38/11
In Total	56/44/20	48/59/13	24/84/12	77/31/12	56/52/12	69/43/8	330/313/77

- Among all the 720 configurations (15 data sets \times 4 MDC approaches \times 3 metrics \times 2 multi-class classifiers \times 2 versions of KRAM), the KRAM counterpart can achieve superior or at least comparable performance against its corresponding MDC approach in 643 configurations (about 89.3%).
- 335 • BR deals with MDC tasks by independent decomposition, where potential class dependencies are totally neglected in this approach. It is shown in Table 6 that KRAM-BR can achieve superior or at least comparable performance against BR in 167 out of 180 cases. This observation indicates that helpful discriminative information is indeed brought into feature space via the k NN-augmented features generated by KRAM. Specifically, 340 these discriminative information can be regarded as an alternative way of dependency modeling when designing MDC approaches.
- Both ECC and ESC deal with MDC tasks by explicitly modeling class dependencies, which are fulfilled by specifying a chaining structure over class spaces or grouping class spaces into super-classes. It is worth noting that for these two MDC approaches which are designed under the inherent mechanism of dependency modeling, the k NN-augmented features generated by KRAM can also help improve their generalization performance 345

significantly.

- ECP deals with MDC tasks by modeling full-order class dependencies, where all distinct class combinations in training set have been treated as new classes in the learning phase. It is shown that the k NN-augmented features generated by KRAM can generally help improve ECP’s generalization performance, while there are 35 cases where KRAM-ECP performs significantly inferior to ECP (the number of loss cases is a lot more than BR, ECC and ESC relatively). Most of these under-performing cases (29 out of 35) occur on data sets *Flickr*, *Thyroid* and *WaterQuality* (including its two divisions *WQplants* and *WQanimals*), in which the possible number of class combinations is high (e.g., 4^{14} for *WaterQuality*).

Table 7: Win/tie/loss counts of pairwise t -test (at 0.05 significance level) between $\text{KRAM}_d\text{-}\mathcal{A}$ and $\text{KRAM}_c\text{-}\mathcal{A}$ ($\mathcal{A} \in \{\text{BR}, \text{ECC}, \text{ECP}, \text{ESC}\}$) in terms of *hamming score* (HScore), *exact match* (EMatch), and *sub-exact match* (SEMatch).

	multi-class classifier: SVM			multi-class classifier: NB			In Total
	HScore	EMatch	SEMatch	HScore	EMatch	SEMatch	
$\text{KRAM}_d\text{-BR}$ against $\text{KRAM}_c\text{-BR}$	0/13/2	0/12/3	0/14/1	5/4/6	4/6/5	3/7/5	12/56/22
$\text{KRAM}_d\text{-ECC}$ against $\text{KRAM}_c\text{-ECC}$	1/11/3	0/14/1	0/14/1	6/3/6	4/6/5	3/7/5	14/55/21
$\text{KRAM}_d\text{-ECP}$ against $\text{KRAM}_c\text{-ECP}$	0/14/1	0/13/2	0/14/1	6/5/4	5/7/3	4/8/3	15/61/14
$\text{KRAM}_d\text{-ESC}$ against $\text{KRAM}_c\text{-ESC}$	0/14/1	0/14/1	0/13/2	6/4/5	4/8/3	5/5/5	15/58/17
In Total	1/52/7	0/53/7	0/55/5	23/16/21	17/27/16	15/27/18	56/230/74

4.3. Further Analysis

4.3.1. Comparison between two KRAM versions

To show whether the performance between two versions of KRAM, i.e. $\text{KRAM}_d\text{-}\mathcal{A}$ and $\text{KRAM}_c\text{-}\mathcal{A}$, is significantly different, pairwise t -test based on ten-fold cross-validation (at 0.05 significance level) is also conducted, and we use bold face type for the statistically superior one in Tables 2 to 5. Table 7 summarizes the resulting win/tie/loss counts accordingly.

Based on the experimental results mentioned above, it is interesting to obtain the following observations:

- 370
 • Across all the 360 configurations (15 data sets \times 4 MDC approaches \times 3 metrics \times 2 multi-class classifiers), KRAM_d and KRAM_c tie 230 times (about 63.9%), KRAM_d wins 56 times, and KRAM_c wins 74 times. So, two versions of KRAM have comparable performance in general.
- 375
 • When SVM is utilized as the multi-class classifier, performance of KRAM_c is relatively better. Specifically, the performance of KRAM_c is comparable with KRAM_d in most cases and superior in some cases with SVM (with only one loss on HScore with ECC approach).
- 380
 • When NB is utilized as the multi-class classifier, KRAM_d achieves relatively more robust performance than KRAM_c does. Specifically, there are more win and loss cases with NB between KRAM_c and KRAM_d . However, KRAM_c will achieve very poor performance in some configurations compared with KRAM_d (e.g., *Thyroid*) while KRAM_d won't even if it achieves significantly inferior performance to KRAM_c . Possible reason is that NB assumes Gaussian pdf for continuous features, while KRAM_c 's continuous augmented features will unfit this assumption sometimes. Obviously, the more imbalanced class distribution is, the more severe unfit will be,
 385
 then KRAM_c will achieve worse performance, and vice versa.
- 390
 • In summary, KRAM_c should be a better choice when SVM is utilized as the multi-class classifier, while it depends on the characteristics of concrete data set and MDC approach when making a choice between KRAM_d and KRAM_c with NB as the multi-class classifier.

4.3.2. Parameter sensitivity analysis

As shown in Algorithm 1, there is only one parameter to be specified for KRAM, i.e., k , which is fixed to 8 in previous sections. To justify this parameter setting, we also investigate the sensitivity of KRAM w.r.t. the value of k . In
 395
 terms of each evaluation metric, Figure 2 illustrates how the performance of KRAM (with MDC approach BR) changes as the value of k increases from 5

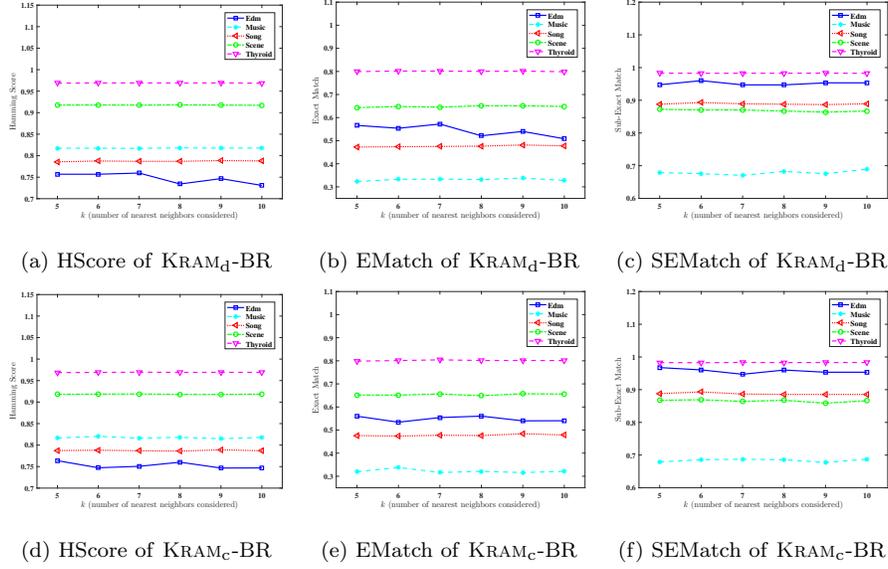


Figure 2: Performance of KRAM-BR changes in terms of *hamming score* (HScore), *exact match* (EMatch), and *sub-exact match* (SEMatch) as k ranges from 5 to 10 on five data sets.

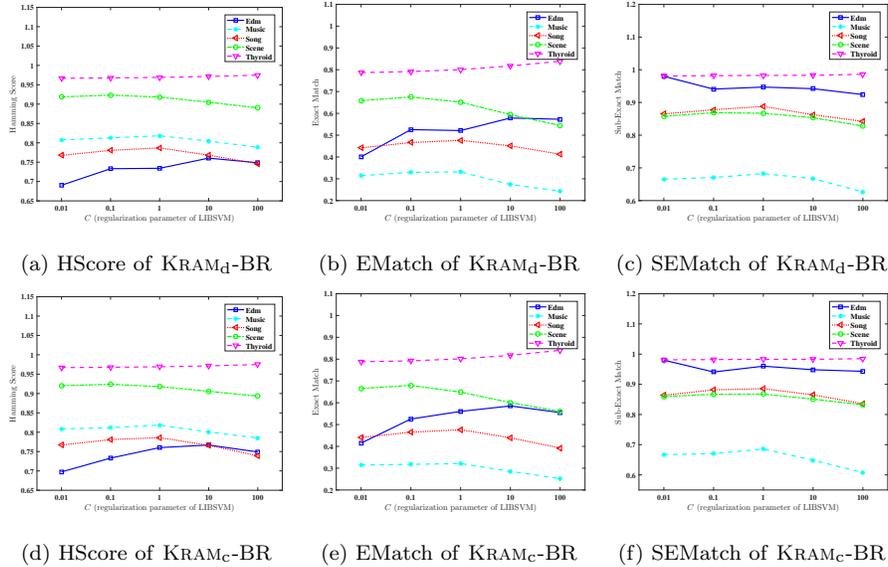


Figure 3: Performance of KRAM-BR changes in terms of *hamming score* (HScore), *exact match* (EMatch), and *sub-exact match* (SEMatch) as the LIBSVM regularization parameter C ranges in $\{0.01, 0.1, 1, 10, 100\}$ on five data sets.

to 10. It is shown that both two KRAM versions can achieve relatively stable performance when the value of k varies. Therefore, in this paper, the value of k is moderately set to 8 which is also the recommended parameter setting of KRAM. Moreover, parameter insensitivity serves as a desirable property which keeps KRAM away from sophisticated parameter-tuning issue for practical use.

In this paper, both SVM and NB are utilized as the multi-class classifier to implement each MDC approach. Here, we also investigate the sensitivity of KRAM w.r.t. the regularization parameter C of LIBSVM [35] (The common NB implementation doesn't involve specific parameters to be tuned). In terms of each evaluation metric, Figure 3 illustrates how the performance of KRAM (with MDC approach BR) changes as the regularization parameter C ranges in $\{0.01, 0.1, 1, 10, 100\}$. It is shown that both of the two KRAM versions can achieve relatively better performance when $C = 1$, which is also the default parameter setting of LIBSVM.

4.3.3. Alternative experimental analysis

The four selected compared approaches (i.e. BR, ECC, ECP, ESC) all necessitate a multi-class classifier for implementation, in this subsection, we also investigate the gMML approach [13] as compared approach which does not necessitate a base multi-class classifier. Specifically, gMML alternately learns linear regression models for each class label as well as a Mahalanobis distance metric to solve MDC problem effectively, where the Mahalanobis distance metric can make the distance between linear regression outputs of one example and its ground-truth class label vector closer. The detailed experimental results are reported in Table 8. It is shown that the k NN-augmented features generated by KRAM can also help improve the generalization performance of gMML and the performance of the KRAM_c version is relatively better.

5. Conclusion

Most existing works for multi-dimensional classification focus on modeling class dependencies in output space, while this paper extends our preliminary

Table 8: Predictive performance of gMML and its two version KRAM counterparts (mean \pm std. deviation) in terms of all three evaluation metrics. In addition, \bullet / \circ indicates whether KRAM_d-gMML or KRAM_c-gMML (denoted as KRAM_d or KRAM_c for short) is statistically superior/inferior to gMML on each data set, and results in bold face of KRAM_d or KRAM_c indicate that current KRAM version is statistically superior to the other version on each data set (pairwise t -test at 0.05 significance level).

Data Set	Hamming Score			Exact Match			Sub-Exact Match		
	gMML	KRAM _d	KRAM _c	gMML	KRAM _d	KRAM _c	gMML	KRAM _d	KRAM _c
Edm	.714 \pm .083	.770 \pm .090 \bullet	.766 \pm .082 \bullet	.487 \pm .145	.586 \pm .188 \bullet	.579 \pm .166 \bullet	.941 \pm .065	.954 \pm .055	.954 \pm .055
Flare1	.923 \pm .033	.923 \pm .035	.924 \pm .034	.821 \pm .073	.818 \pm .075	.821 \pm .073	.951 \pm .036	.954 \pm .039	.954 \pm .039
Song	.788 \pm .027	.787 \pm .024	.786 \pm .025	.484 \pm .059	.481 \pm .054	.481 \pm .054	.883 \pm .041	.883 \pm .040	.882 \pm .041
WQpla.	.655 \pm .015	.663 \pm .018	.662 \pm .016	.092 \pm .035	.098 \pm .039	.096 \pm .036	.286 \pm .053	.295 \pm .043	.290 \pm .040
WQani.	.630 \pm .015	.642 \pm .014 \bullet	.642 \pm .013 \bullet	.062 \pm .023	.062 \pm .010	.062 \pm .013	.227 \pm .033	.240 \pm .033	.248\pm.033
WQ	.643 \pm .013	.649 \pm .012	.649 \pm .012	.006 \pm .008	.008 \pm .006	.008 \pm .006	.049 \pm .024	.054 \pm .018	.057 \pm .020
Voice	.842 \pm .009	.945 \pm .010 \bullet	.946 \pm .010 \bullet	.699 \pm .017	.892 \pm .020 \bullet	.893 \pm .018 \bullet	.985 \pm .011	.998 \pm .003 \bullet	.998 \pm .003 \bullet
Thyroid	.960 \pm .003	.967 \pm .003 \bullet	.967 \pm .003 \bullet	.741 \pm .015	.790 \pm .018 \bullet	.790 \pm .010 \bullet	.982 \pm .005	.981 \pm .004	.981 \pm .004
Flickr	.779 \pm .004	.782 \pm .005 \bullet	.782 \pm .006 \bullet	.287 \pm .009	.296 \pm .011 \bullet	.295 \pm .030 \bullet	.690 \pm .016	.696 \pm .015 \bullet	.696 \pm .016 \bullet
Music	.801 \pm .018	.815 \pm .023 \bullet	.817 \pm .027 \bullet	.254 \pm .057	.320 \pm .095 \bullet	.333 \pm .027 \bullet	.652 \pm .040	.681 \pm .033 \bullet	.677 \pm .049
Enron	.832 \pm .009	.833 \pm .011	.833 \pm .011	.197 \pm .027	.202 \pm .038	.206\pm.024\bullet	.477 \pm .042	.476 \pm .042	.478 \pm .041
Image	.811 \pm .010	.839 \pm .010 \bullet	.840 \pm .010 \bullet	.289 \pm .025	.448 \pm .028 \bullet	.454 \pm .019 \bullet	.787 \pm .027	.781 \pm .021	.779 \pm .024
Scene	.893 \pm .009	.917 \pm .008 \bullet	.919\pm.007\bullet	.457 \pm .046	.646 \pm .025 \bullet	.657\pm.011\bullet	.908 \pm .017	.858 \pm .023 \circ	.861 \pm .018 \circ
Yeast	.800 \pm .005	.811 \pm .006 \bullet	.811 \pm .006 \bullet	.134 \pm .018	.210 \pm .020 \bullet	.211 \pm .035 \bullet	.266 \pm .026	.311 \pm .016 \bullet	.309 \pm .018 \bullet
Mediamill	.811 \pm .001	.847 \pm .001 \bullet	.850\pm.001\bullet	.111 \pm .006	.250 \pm .005 \bullet	.258\pm.005\bullet	.342 \pm .007	.497 \pm .006 \bullet	.505\pm.007\bullet

work [12] which solves MDC problems by manipulating the input space. The major contribution of our research is to propose a feature augmentation strategy for multi-dimensional classification that augments feature space with augmented features generated via combining with class space, which suggests an alternative solution to induce MDC models.

To justify the proposed strategy, a simple yet effective approach named KRAM is designed based on k NN techniques, and comprehensive comparative studies have been conducted to validate its effectiveness accordingly. Experimental results show that: (a) Both versions of KRAM, i.e., discrete version KRAM_d and continuous version KRAM_c, can improve predictive performance of existing MDC approaches; (b) KRAM_d has more stable performance than KRAM_c when the base multi-class classifier is sensitive to the feature type (discrete or con-

tinuous), e.g., Naïve Bayes classifier. For other base multi-class classifier less sensitive to feature type, e.g., SVM, KRAM_c is likely to be a better choice; (c) In light of the effectiveness of KRAM , feature augmentation technique can be further studied as an alternative strategy for modeling class dependencies.

As an initial attempt towards solving MDC problem with feature augmentation, there are several potential ways that the current KRAM instantiation can be improved: (a) In addition to the simple counting statistics derived from k nearest neighbors, more advanced information could be utilized for feature augmentation by trying to exploit available domain knowledge [36]; (b) Other than the meta-strategy for feature augmentation, customized feature augmentation techniques can be investigated for given MDC approaches; (c) Similar to the label-specific features techniques for multi-label classification [37, 38], it is worthwhile to investigate the feasibility of generating specific augmented features w.r.t. each class space.

References

- [1] J. Ortigosa-Hernández, J. D. Rodríguez, L. Alzate, M. Lucania, I. Inza, J. A. Lozano, Approaching sentiment analysis by using semi-supervised learning of multi-dimensional classifiers, *Neurocomputing* 92 (2012) 98–115.
- [2] C. Tu, Z. Liu, H. Luan, M. Sun, PRISM: Profession identification in social media, *ACM Transactions on Intelligent Systems and Technology* 8 (6) (2017) Article 81.
- [3] J. D. Rodríguez, A. Pérez, D. Arteta, D. Tejedor, J. A. Lozano, Using multidimensional Bayesian network classifiers to assist the treatment of multiple sclerosis, *IEEE Transactions on Systems, Man, and Cybernetics-Part C: Applications and Reviews* 42 (6) (2012) 1705–1715.
- [4] H. Borchani, C. Bielza, C. Toro, P. Larrañaga, Predicting human immun-

- 465 odeficiency virus inhibitors using multi-dimensional Bayesian network clas-
sifiers, *Artificial Intelligence in Medicine* 57 (3) (2013) 219–229.
- [5] F. Serafino, G. Pio, M. Ceci, D. Malerba, Hierarchical multidimensional
classification of web documents with multiwebclass, in: *Lecture Notes in
Computer Science* 9356, Springer, Berlin, 2015, pp. 236–250.
- 470 [6] J. Arias, J. A. Gamez, T. D. Nielsen, J. M. Puerta, A scalable pairwise
class interaction framework for multidimensional classification, *Internation-
al Journal of Approximate Reasoning* 68 (2016) 194–210.
- [7] J. H. Zaragoza, L. E. Sucar, E. F. Morales, C. Bielza, P. Larrañaga,
Bayesian chain classifiers for multidimensional classification, in: *Proceed-
ings of the 22nd International Joint Conference on Artificial Intelligence,*
475 *Barcelona, Spain, 2011*, pp. 2192–2197.
- [8] J. Read, L. Martino, D. Luengo, Efficient monte carlo methods for multi-
dimensional learning with classifier chains, *Pattern Recognition* 47 (3)
(2014) 1535–1546.
- 480 [9] C. Bielza, G. Li, P. Larrañaga, Multi-dimensional classification with
Bayesian networks, *International Journal of Approximate Reasoning* 52 (6)
(2011) 705–727.
- [10] J. H. Bolt, L. C. van der Gaag, Balanced sensitivity functions for tun-
ing multi-dimensional Bayesian network classifiers, *International Journal
of Approximate Reasoning* 80 (2017) 361–376.
485
- [11] J. Read, C. Bielza, P. Larrañaga, Multi-dimensional classification with
super-classes, *IEEE Transactions on Knowledge and Data Engineering*
26 (7) (2014) 1720–1733.
- 490 [12] B.-B. Jia, M.-L. Zhang, Multi-dimensional classification via kNN feature
augmentation, in: *Proceedings of the 33rd AAAI Conference on Artificial
Intelligence, Honolulu, HI, 2019*, pp. 3975–3982.

- [13] Z. Ma, S. Chen, Multi-dimensional classification via a metric approach, *Neurocomputing* 275 (2018) 1121–1131.
- [14] 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.
- [15] E. Gibaja, S. Ventura, A tutorial on multilabel learning, *ACM Computing Surveys* 47 (3) (2015) Article 52.
- [16] M.-L. Zhang, Y.-K. Li, X.-Y. Liu, X. Geng, Binary relevance for multi-label learning: An overview, *Frontiers of Computer Science* 12 (2) (2018) 191–202.
- [17] D. Xu, Y. Shi, I. W. Tsang, Y. Ong, C. Gong, X. Shen, Survey on multi-output learning, *IEEE Transactions on Neural Networks and Learning Systems* (2020) in press.
- [18] B.-B. Jia, M.-L. Zhang, Maximum margin multi-dimensional classification, in: *Proceedings of the 34th AAAI Conference on Artificial Intelligence*, New York, NY, 2020, in press.
- [19] B.-B. Jia, M.-L. Zhang, Multi-dimensional classification via stacked dependency exploitation, *Science China Information Sciences* (2020) in press.
- [20] C. Liu, P. Zhao, S.-J. Huang, Y. Jiang, Z.-H. Zhou, Dual set multi-label learning, in: *Proceedings of the 32th AAAI Conference on Artificial Intelligence*, New Orleans, LA, 2018, pp. 3635–3642.
- [21] Z. Ma, S. Chen, A convex formulation for multiple ordinal output classification, *Pattern Recognition* 86 (2019) 73–84.
- [22] J. Read, B. Pfahringer, G. Holmes, E. Frank, Classifier chains for multi-label classification, *Machine Learning* 85 (3) (2011) 333–359.

- [23] I. Batal, C. Hong, M. Hauskrecht, An efficient probabilistic framework for multi-dimensional classification, in: Proceedings of the 22nd ACM International Conference on Information & Knowledge Management, San Francisco, CA, 2013, pp. 2417–2422.
- [24] M. Zhu, S. Liu, J. Jiang, A hybrid method for learning multi-dimensional Bayesian network classifiers based on an optimization model, Applied Intelligence 44 (1) (2016) 123–148.
- [25] M. Benjumbeda, C. Bielza, P. Larrañaga, Tractability of most probable explanations in multidimensional Bayesian network classifiers, International Journal of Approximate Reasoning 93 (2018) 74–87.
- [26] A. Karalič, I. Bratko, First order regression, Machine Learning 26 (1997) 147–176.
- [27] D. Dheeru, E. Karra Taniskidou, UCI machine learning repository, <http://archive.ics.uci.edu/ml> (2017).
- [28] S. Džeroski, D. Demšar, J. Grbović, Predicting chemical parameters of river water quality from bioindicator data, Applied Intelligence 13 (1) (2000) 7–17.
- [29] D. Kocev, C. Vens, J. Struyf, S. Džeroski, Ensembles of multi-objective decision trees, in: Lecture Notes in Computer Science 4701, Springer, Berlin, 2007, pp. 624–631.
- [30] M. J. Huiskes, M. S. Lew, The MIR flickr retrieval evaluation, in: Proceedings of the 1st ACM International Conference on Multimedia Information Retrieval, Vancouver, Canada, 2008, pp. 39–43.
- [31] M.-L. Zhang, Z.-H. Zhou, ML-KNN: A lazy learning approach to multi-label learning, Pattern Recognition 40 (7) (2007) 2038–2048.
- [32] M. R. Boutell, J. Luo, X. Shen, C. M. Brown, Learning multi-label scene classification, Pattern Recognition 37 (9) (2004) 1757–1771.

- [33] A. Elisseeff, J. Weston, A kernel method for multi-labelled classification,
545 in: *Advances in Neural Information Processing Systems*, MIT Press, Cambridge, MA, 2001, pp. 681–687.
- [34] C. G. Snoek, M. Worring, J. C. Van Gemert, J.-M. Geusebroek, A. W. Smeulders, The challenge problem for automated detection of 101 semantic concepts in multimedia, in: *Proceedings of the 14th ACM International*
550 *Conference on Multimedia*, Santa Barbara, CA, 2006, pp. 421–430.
- [35] C.-C. Chang, C.-J. Lin, LIBSVM: A library for support vector machines, *ACM Transactions on Intelligent Systems and Technology* 2 (3) (2011) Article 27, software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- 555 [36] Z.-H. Zhou, Abductive learning: Towards bridging machine learning and logical reasoning, *Science China Information Sciences* 62 (7) (2019) Article 076101.
- [37] M.-L. Zhang, L. Wu, LIFT: Multi-label learning with label-specific features, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 37 (1)
560 (2015) 107–120.
- [38] X.-Y. Jia, S.-S. Zhu, W.-W. Li, Joint label-specific features and correlation information for multi-label learning, *Journal of Computer Science and Technology* 35 (2) (2020) 247–258.