

# 基于专家特征的条件互信息多标记特征选择算法

程玉胜<sup>1,2\*</sup>, 宋帆<sup>1</sup>, 王一宾<sup>1,2</sup>, 钱坤<sup>1</sup>

(1. 安庆师范大学 计算机与信息学院, 安徽 安庆 246011; 2. 安徽省高校智能感知与计算重点实验室, 安徽 安庆 246011)

(\* 通信作者电子邮箱 chengyushaq@163.com)

**摘要:** 特征选择对于分类器的分类精度和泛化性能起重要作用。目前的多标记特征选择算法主要利用最大相关性最小冗余性准则在全部特征集中进行特征选择, 没有考虑专家特征, 因此多标记特征选择算法的运行时间较长、复杂度较高。实际上, 在现实生活中专家依据几个或者多个关键特征就能够直接决定整体的预测方向。如果提取关注这些信息, 必将减少特征选择的计算时间, 甚至提升分类器性能。基于此, 提出一种基于专家特征的条件互信息多标记特征选择算法。首先将专家特征与剩余的特征相联合, 再利用条件互信息得出一个与标记集合相关性由强到弱的特征序列, 最后通过划分子空间去除冗余性较大的特征。该算法在7个多标记数据集上进行了实验对比, 结果表明该算法较其他特征选择算法有一定优势, 统计假设检验与稳定性分析进一步证明了所提出算法的有效性和合理性。

**关键词:** 特征选择; 专家特征; 条件互信息; 多标记学习; 局部子空间

中图分类号: TP391 文献标志码: A

## Multi-label feature selection algorithm based on conditional mutual information of expert feature

CHENG Yusheng<sup>1,2\*</sup>, SONG Fan<sup>1</sup>, WANG Yibin<sup>1,2</sup>, QIAN Kun<sup>1</sup>

(1. School of Computer and Information, Anqing Normal University, Anqing Anhui 246011, China;

2. University Key Laboratory of Intelligent Perception and Computing of Anhui Province, Anqing Anhui 246011, China)

**Abstract:** Feature selection plays an important role in the classification accuracy and generalization performance of classifiers. The existing multi-label feature selection algorithms mainly use the maximum relevance and minimum redundancy criterion to perform feature selection in all feature sets without considering expert features, therefore, the multi-label feature selection algorithm has the disadvantages of long running time and high complexity. Actually, in real life, experts can directly determine the overall prediction direction based on a few or several key features. Paying attention to and extracting this information will inevitably reduce the calculation time of feature selection and even improve the performance of classifier. Based on this, a multi-label feature selection algorithm based on conditional mutual information of expert feature was proposed. Firstly, the expert features were combined with the remaining features, and then the conditional mutual information was used to obtain a feature sequence of strong to weak relativity with the label set. Finally, the subspaces were divided to remove the redundant features. The experimental comparison was performed to the proposed algorithm on 7 multi-label datasets. Experimental results show that the proposed algorithm has certain advantages over the other feature selection algorithms, and the statistical hypothesis testing and the stability analysis further illustrate the effectiveness and the rationality of the proposed algorithm.

**Key words:** feature selection; expert feature; conditional mutual information; multi-label learning; local subspace

## 0 引言

多标记学习<sup>[1]</sup>作为机器学习等领域的研究热点之一, 较之传统的单标记学习中的一个对象只能局限于单个标记, 多标记学习框架更具有实际性和广泛性。在真实世界中, 一个对象可能隶属于多个标记<sup>[2]</sup>, 此时, 单个标记就难以表述对象的完整性。所以, 多标记学习的研究对于多义性对象的学习建模更具有实际运用意义, 近年来已成为一个新的研究热点<sup>[3,4]</sup>。

在多标记学习问题中, 由于数据的高维性会引起维数灾难, 导致分类器精度降低<sup>[5]</sup>。而特征选择作为一种普遍的降维手段, 对于分类器的分类精度和泛化性能起重要作用。特征选择的首要目的是在样本数据集中找到一个特征子集, 且使得找到的特征子集蕴含尽可能多的区分类别信息, 同时要考虑子集内部的冗余性尽量小<sup>[6]</sup>。而信息论中的互信息理论, 作为不确定性的一种有效度量方式, 被广泛用于多标记特征选择, 因此许多学者在此方面进行了研究。例如 Zhang 等<sup>[7]</sup>

收稿日期: 2019-08-30; 修回日期: 2019-09-24; 录用日期: 2019-10-09。

基金项目: 安徽省高校重点科研项目(KJ2017A352); 安庆师范大学科研创新团队建设计划项目。

作者简介: 程玉胜(1969—), 男, 安徽安庆人, 教授, 博士, CCF 会员, 主要研究方向: 大数据、粗糙集、特征选择的机器学习; 宋帆(1992—), 男, 安徽铜陵人, 硕士研究生, CCF 会员, 主要研究方向: 多标记学习、神经网络; 王一宾(1970—), 男, 安徽安庆人, 教授, 硕士, 主要研究方向: 多标记学习、机器学习、软件安全; 钱坤(1995—), 男, 安徽滁州人, 硕士研究生, CCF 会员, 主要研究方向: 多标记学习、机器学习、数据统计。

提出的基于最大相关性的多标记维度约简 (Multi-label Dimensionality reduction via Dependence Maximization, MDDM) 算法。Lee 等<sup>[8]</sup>通过多元互信息最大化已选特征与标记集合的相关性,提出了基于多变量互信息的多标记特征选择算法 PMU (Pairwise Multivariate mutual information)。Lin 等<sup>[9]</sup>提出了基于邻域互信息的多标记特征选择。刘景华等<sup>[10]</sup>通过互信息排序已选特征和标记的相关性,提出了基于局部子空间的多标记特征选择算法 (Multi-label Feature Selection algorithm based on Local Subspace, MFSLs)。

上述算法的多标记特征选择算法判定特征是否冗余的标准单一,如信息熵方法仅考虑特征和标记间的相关性,未考虑特征和特征间的关系<sup>[11]</sup>;联合互信息虽然考虑了整体互信息大小,但未考虑单个特征和标记间的相关性。这些算法都没有提取关注专家特征,在整个特征集中进行特征选择,因此时间复杂度很高,如 PMU 算法在大数据集中执行时间很长,基于信息熵的多标签特征选择 (Multi-Label Feature Selection based on Information Entropy, MLFSIE) 算法虽然提高了执行速度,但未考虑特征间的冗余性。

除此之外,包括上述算法在内的大多数多标记特征选择算法均未考虑到优先挑选出专家特征的现实意义,忽略了一个关键的问题:在现实生活中,人们针对分类问题时,通常根据专家经验选取几个或者多个最重要的特征,然后再通过相关评价准则建立特征向任务目标的映射进行多标记分类。例如,在医院中专家医师看病人的病情时,往往根据自己的多年临床经验先确定几个最重要的病症(即对结果不可或缺起到重要作用的特征(专家特征)),然后再在专家特征的基础上进行各种身体检查、血液化验、分析病历,最后分析汇总来确诊。同时也要考虑某些看似不显眼的症状(即与标记空间相关性较弱的特征),因为忽略某些重要性较次要的特征也会产生误诊的可能。

基于此,再结合信息论中的互信息<sup>[12-14]</sup>,本文提出一种基于专家特征的条件互信息多标记特征选择算法 (Multi-label Feature Selection algorithm based on Conditional Mutual Information of Expert Feature, MFSEF)。该算法在最小冗余最大相关性前提下,通过子空间划分,考虑了重要性较次要的特征可能对分类性能产生的影响,受现实生活中实际问题的启发,兼顾考虑了可能决定整体的预测方向的专家特征,从而提升多标记分类性能。首先通过瀑布图联合互信息选出几个最关键的专家特征;再以该专家特征作条件,保持专家特征不变,与余下的特征作并集,构建融合一个新的特征空间,然后计算新特征与标记集合之间的互信息,再进行排序形成新的特征排序集合;借鉴 MFSLs 的思想,最后进行特征选择。实验在 7 个多标记数据集上测试,同其他常用的多标记特征选择算法进行比较,通过 4 个评价指标的结果可以看出,本文算法优于通用的多标记特征选择算法。最后,还通过统计假设检验进一步证实了本文方法的合理性与稳定性。

## 1 理论介绍

### 1.1 多标记学习

由于真实世界的对象具有多义性,多标记学习框架作为一种多义性对象学习建模工具由此产生<sup>[15]</sup>。在该框架下,每个对象由一个示例描述,每个示例具有多个但有限的类别标记,学习的目的是为每个未知示例赋予正确的标记。在数学

语言中,多标记问题可描述为:假定  $X = \{x_1, x_2, \dots, x_n\}^T \in \mathbb{R}^{n \times d}$  表示有  $n$  个样本且每个样本特征维度为  $d$ ,  $Y = \{1, 2, \dots, Q\}$  表示样本对应的标记集合<sup>[16]</sup>。  $T = \{(x_1, Y_1), (x_2, Y_2), \dots, (x_m, Y_m)\} (x_i \in X, Y_i \in Y)$  表示训练集,多标记学习目的是得到映射关系  $f: X \rightarrow \{-1, 1\}^Q$ ,从而对新样本进行标记的预测。

### 1.2 条件互信息

**定义 1** 设集合  $A = \{a_1, a_2, \dots, a_m\}$ , 令  $p(a_i)$  表示元素  $a_i$  的先验概率,则集合  $A$  的信息熵为:

$$H(A) = -\sum_{i=1}^m p(a_i) \log p(a_i) \quad (1)$$

信息熵可以度量集合不确定性的程度,信息熵越大表示集合的不稳定性越大。对于多标记特征选择算法,常通过信息熵来选择特征空间中与标记空间互信息较大的特征。

**定义 2** 设集合  $A = \{a_1, a_2, \dots, a_m\}$ ,  $B = \{b_1, b_2, \dots, b_n\}$ , 则在给定集合  $A$  的条件下集合  $B$  的条件熵为:

$$H(B|A) = -\sum_{i=1}^m \sum_{j=1}^n p(a_i, b_j) \log p(b_j | a_i) \quad (2)$$

条件熵可以度量在集合  $A$  出现的条件下集合  $B$  的不确定程度的大小。

**定义 3** 设集合  $A = \{a_1, a_2, \dots, a_m\}$ , 集合  $B = \{b_1, b_2, \dots, b_n\}$ , 则集合  $A$  与  $B$  的联合熵为:

$$H(A, B) = -\sum_{i=1}^m \sum_{j=1}^n p(a_i, b_j) \log p(a_i, b_j) \quad (3)$$

信息熵、条件熵及联合熵的关系为:

$$H(A, B) = H(A) + H(B|A) = H(B) + H(A|B) \quad (4)$$

**定义 4** 给定集合  $A$  和集合  $B$ , 定义集合  $A$  和  $B$  之间的互信息为:

$$I(A; B) = -\sum_{i=1}^m \sum_{j=1}^n p(a_i, b_j) \log \left[ \frac{p(a_i, b_j)}{p(a_i)p(b_j)} \right] \quad (5)$$

互信息被广泛用于度量随机变量间相关性的大小,即  $I(A; B)$  表示集合  $A$  和集合  $B$  间的相关性大小。 $I(A; B)$  越大,表示两者间的相关性越大。另有  $I(A; B) = I(B; A)$ , 且满足:

$$I(A; B) = H(A) + H(B) - H(A, B) = H(A) - H(A|B) = H(B) - H(B|A) \quad (6)$$

当  $I(A; B) = 0$  时,集合  $A$  和集合  $B$  无相关性,集合  $A$  和集合  $B$  之间未提供任何信息。

**定义 5** 设集合  $A = \{a_1, a_2, \dots, a_m\}$ ,  $B = \{b_1, b_2, \dots, b_n\}$ ,  $C = \{c_1, c_2, \dots, c_l\}$ , 则在集合  $C$  条件下集合  $A$  和  $B$  间的条件互信息<sup>[17]</sup>为:

$$I(A; B|C) = H(A|B, C) - H(A|B) \quad (7)$$

联合互信息可由式(6)和式(7)得出:

$$I(A, C; B) = I(A; B|C) + I(B; C) \quad (8)$$

联合互信息是考虑  $A, C$  整体同  $B$  之间的关系,由上式可知条件互信息和互信息之和为联合互信息,根据式(5)得出联合互信息为:

$$I(A, C; B) = -\sum_{i=1}^m \sum_{j=1}^n \sum_{k=1}^l p(a_i, b_j, c_k) \log \left[ \frac{p(a_i, b_j, c_k)}{p(a_i, c_k)p(b_j)} \right] \quad (9)$$

联合互信息  $I(A, C; B)$  越大,则表示  $A, C$  整体同  $B$  间的

相关性越强。另外关于条件互信息还可变形表示如下:

$$I(A; B|C) = H(A, C) + H(B, C) - H(A, B, C) - H(C) \quad (10)$$

## 2 MFSEF

在通过互信息考虑特征与标记之间的相关性来进行多标记特征选择中,先给定  $f$  表示描述样本的特征,  $l$  表示样本的类别标记, 则  $I(f; l)$  虽然仅可以在单标记中描述在样本中特征和类别标记之间的相关性程度,而在多标记中,一个样本是由多个特征向量表示且隶属于多个类别标记,故给出以下定义。

**定义 6** 给定特征  $f$  和标记空间  $L = \{l_1, l_2, \dots, l_n\}$ ,  $\forall l_i \in L (i = 1, 2, \dots, n)$ ,  $I(f; l_i)$  为特征  $f$  和标记  $l_i$  的互信息,那么特征  $f$  和标记空间集  $L$  的互信息可定义为:

$$FMI(f; L) = \sum_{i=1}^n I(f; l_i) \quad (11)$$

**定义 7** 给定一个特征子集为  $S = \{f_1, f_2, \dots, f_m\}$ , 特征  $f_i$  与特征子集空间的互信息定义为:

$$FMI(f_i; S) = \sum_{j=1, j \neq i}^m I(f_i, f_j) \quad (12)$$

特征与标记集合之间的互信息大小描述了两个集合间的相关性程度,特征和标记集合的互信息越大,表明该特征越重要;反之,表明该特征重要性越弱,当特征和标记集合的互信息为零时,表明该特征和每个类别标记相互独立,此时特征和标记集合的互信息也取得最小值。

给定训练样本  $U = \{x_1, x_2, \dots, x_n\}$  和其构成样本的特征集合  $F = \{f_1, f_2, \dots, f_j\}$ , 以及标记空间集合  $L = \{l_1, l_2, \dots, l_i\}$ 。

由于专家特征在现实生活中是通过人的经验主观性选定,而本文实验所采用数据集为常用的多标记数据集,若单纯地人为指定专家特征,可能会影响实验的可靠性和有效性,故可事先通过数据画画出对应的特征值瀑布图,然后根据瀑布图联合互信息理论挑选出几个特征作为专家特征。图 1 和图 2 为常用多标记数据集所画的部分代表性瀑布图。

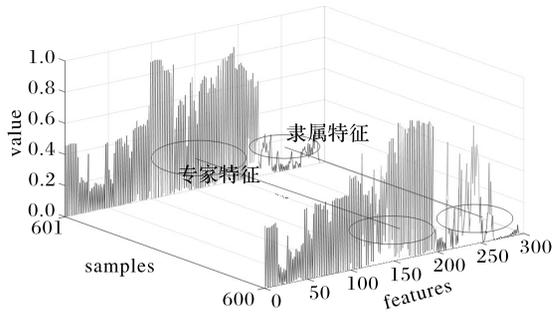


图 1 第 600 和 601 个样本中所有特征对应的特征值构成的瀑布图

Fig. 1 Waterfall plot of eigenvalues corresponding to all features in the 600th and 601th samples

其中图 1 展示了在第 600 个样本到第 601 个样本中,每个特征所对应的特征值大小的对比,由图可看出:两个样本对应的前 100 个特征的特征值基本较小且小于 0.5,两样本特征值基本相同,表明该段特征对标记的影响甚微;而在第 200 到 300 之间的特征所对应的特征值数值有大有小,分布极其不均衡,波动较大,说明该段特征对样本类别的划分起到决定性作用,即称隶属特征;在第 100 和 200 之间的特征所对应的特征值基本数值较大,且分布均衡,无较大波动,表明该段特征

对标记的影响至关重要,即称为专家特征。

在图 2 中展示了在所有样本中,每个特征所对应的特征值大小的对比。针对 100~200 的专家特征,明显看出所有样本的特征值基本较大接近 1,这表明这部分特征对于标记的划分不可或缺,但这部分特征值基本相同,表明全部样本基本上均具有该特征,所以如果事先把这些特征挑选部分出来,作为条件与剩余特征相联合,然后再针对剩余特征进行后续特征选择操作,无疑可以避免重复计算,提高特征选择速度,而且更符合现实生活中面对分类问题时,常优先选出专家特征的操作习惯,因此更具有实际应用价值。

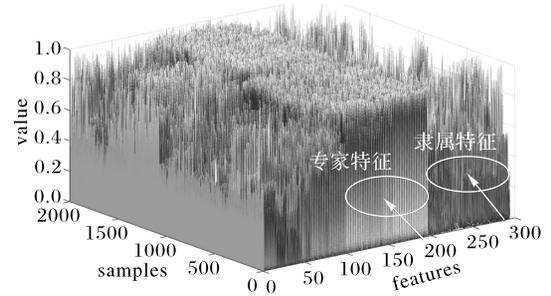


图 2 全部样本所有特征对应的特征值构成的整体瀑布图

Fig. 2 Overall waterfall plot composed of eigenvalues corresponding to all features of all samples

先通过瀑布图观察出专家特征的大致分布,再通过互信息考虑特征空间和标记空间的相关性大小,由互信息大小降序挑选靠前的特征,最后综合考虑两要素选出若干个特征作为专家特征。本文实验取前四个特征作专家特征,记作  $E = \{e_1, e_2, e_3, e_4\}$ 。

传统的基于互信息的多标记特征选择算法仅考虑特征空间  $F$  与标记空间  $L$  的相关性:  $F \rightarrow L$ , 本文以专家特征  $E$  为条件,保持专家特征不变,将专家特征和每个原始特征作并集构建新的特征空间,再考虑其与标记空间的相关性:  $E \cup F \rightarrow L$ 。

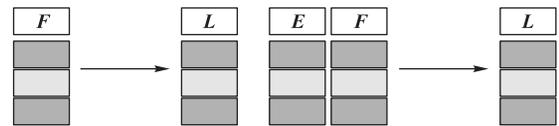


图 3 专家特征联合的图解描述

Fig. 3 Graphical description of combining expert features

由定义 6 可知特征和标记集合的互信息越大,表明该特征越重要;反之,表明该特征重要性越弱。故用互信息大小进行降序排列得到一组新的特征序列:

$$S = \{f'_1, f'_2, \dots, f'_d\}$$

对于多标记特征选择,由于每个标记隶属于不同特征空间,因此对于特征和标记集合的相关性通过互信息大小进行计算时,相关性强的特征之间可能有比较大的冗余性,而相关性弱的特征也不一定对判别标记类别不起作用,也有可能某个相关性弱的特征往往对最后分类结果起决定性作用。考虑到此情景,可以通过建立局部子空间模型来解决此问题<sup>[10]</sup>。文献[18]中说明,当数据集的特征维度较小的时候,子空间个数可以划分 2、3、4,考虑到较多保留相关性强的特征,同时兼顾对某些类别标记贡献较大的但是特征与标记相关性较弱的特征,故每个特征子空间的采样比例可以设置为由大到小,例如对于 2 个子空间,采样比例可为:  $\{0.6, 0.4\}$ 、 $\{0.7, 0.3\}$ 、

$\{0.9, 0.2\}$ 。又由文献[10]中实验证明 3 个子空间的预测效果最好,故可以将已经通过专家特征进行互信息大小降序排列后的特征序列划分为三个子空间,再通过  $\{0.6, 0.3, 0.1\}$  的采样比例进行进一步的特征选择。

局部子空间的详细过程如下:

给定有  $d$  维特征空间,三个子空间,故每个子空间特征维数为  $\lceil d/3 \rceil$ ,由定义 6 可计算出在子空间中每个特征和剩余特征的互信息大小:此时互信息越大,表明其特征的冗余性越高;反之,特征间互信息越小,冗余性越低。因此将通过定义 6 新得到的三个子空间中的特征进行升序,三个子空间关于特征间的冗余性排列分别为:  $S'_1 = \{f'_{11}, f'_{12}, \dots, f'_{1[\lceil d/3 \rceil]}\}$ ,  $S'_2 = \{f'_{21}, f'_{22}, \dots, f'_{2[\lceil d/3 \rceil]}\}$ ,  $S'_3 = \{f'_{31}, f'_{32}, \dots, f'_{3[\lceil d/3 \rceil]}\}$ 。

再通过采样比例分别在三个子空间中选择冗余性比较小的特征,由于采样比例为  $\{0.6, 0.3, 0.1\}$ ,故三个子空间通过比例选择出的特征个数分别为:  $N_1 = \lceil 0.6 * \lceil d/3 \rceil \rceil$ ,  $N_2 = \lceil 0.3 * \lceil d/3 \rceil \rceil$ ,  $N_3 = \lceil 0.1 * \lceil d/3 \rceil \rceil$ 。

通过采样比例分别挑选出  $S'_1$ 、 $S'_2$  和  $S'_3$  的前  $N_1$ 、 $N_2$  和  $N_3$  个特征,记作:  $S''_1 = \{f''_{11}, f''_{12}, \dots, f''_{1N_1}\}$ ,  $S''_2 = \{f''_{21}, f''_{22}, \dots, f''_{2N_2}\}$ ,  $S''_3 = \{f''_{31}, f''_{32}, \dots, f''_{3N_3}\}$ 。最后将  $S''_1$ 、 $S''_2$  和  $S''_3$  合并,就成了新的特征选择后的最终特征子集。

模拟现实世界中分类问题,引入专家特征作为条件,同原始特征作并集,通过信息论中的互信息理论背景,计算特征与标记空间的相关性大小,再结合局部子空间模型的划分,最后实现特征选择。这样既考虑了多标记特征选择在现实社会中的合理性和实用性,也避免了传统特征选择只是根据相关准则选择较强的特征导致的特征间的冗余性,最后实验结果也显示了较好的分类性能。

算法 1 MFSEF。

输入 多标记数据集  $D$ , 专家特征  $E$ ;

输出 特征子集  $Sub$ 。

- 1)  $FMI = \emptyset$ ;
- 2)  $E = \{e_1, e_2, e_3, e_4\}$ ;
- 3) for each  $f_i \in L$
- 4)  $S = \emptyset$ ;  $Sub = \emptyset$ ;
- 5) for each  $l_j \in L$
- 6) 通过定义 6 计算  $FMI(f_i \cup E; l_j)$
- 7) end
- 8) 根据  $\sum_{j=1}^n FMI(f_i \cup E; l_j)$  计算每个特征  $f_i$  在专家特征  $E$  的条件下对所有标记互信息大小;
- 9) end
- 10) 通过第 8) 步计算出来的特征空间和标记空间互信息大小,对特征进行一个降序,从而得到新特征序列  $S$ ;
- 11) 将特征集合  $S$  均分成 3 个子空间  $S_1$ 、 $S_2$  和  $S_3$ ;
- 12) 对三个子空间分别通过定义 7 计算特征和特征的互信息大小,然后进行升序排列,再通过采样比例  $\{0.6, 0.3, 0.1\}$  在三个子空间分别挑选出新的特征子集;
- 13) 将四个专家特征和新得到的三个特征子集合并,按顺序依次放入  $Sub$ 。

### 3 实验及其结果分析

#### 3.1 实验数据集

为验证本文算法的有效性,选取了 Entertainment、Recreation、Artificial、Reference、Health、Business、Computer 共 7

个数据集,详细信息见表 1。

表 1 多标记数据集

Tab. 1 Multi-label datasets

数据集	样本数	特征数	类别数	训练数	测试数
Health	5 000	612	32	2 000	3 000
Recreation	5 000	606	22	2 000	3 000
Artificial	5 000	462	26	2 000	3 000
Reference	5 000	793	33	2 000	3 000
Entertainment	5 000	640	21	2 000	3 000
Business	5 000	438	30	2 000	3 000
Compute	5 000	681	33	2 000	3 000

#### 3.2 实验环境及评价指标

本实验代码在 Matlab 2016a 中运行;硬件环境为 Intel Core i5-2525M 2.50 GHz CPU, 8 GB 内存;操作系统为 Windows 10。实验选取多标记常用的 4 种性能评价指标<sup>[19-20]</sup>,即平均精度(Average Precision, AP)、海明损失(Hamming Loss, HL)、排序损失(Ranking Loss, RL)和 1-错误率(One Error, OE)来综合评价多标记学习算法性能,且分别简称为:  $AP \uparrow$ 、 $HL \downarrow$ 、 $RL \downarrow$  和  $OE \downarrow$ ,其中:  $\uparrow$  代表指标数值越高越好,  $\downarrow$  代表指标数值越低越好。设:多标记分类器  $h(\cdot)$ , 预测函数  $f(\cdot, \cdot)$ , 排序函数  $rank_f$ , 多标记数据集  $D = \{(x_i, Y_i) | 1 \leq i \leq n\}$ 。上述 4 种评价指标 AP、HL、RL 和 OE 形式化定义如下:

1) Average Precision: 评估在特定标记  $y \in Y_i$  排列的正确标记的平均分数。

$$AP_D(f) = \frac{1}{n} \sum_{i=1}^n \frac{1}{|Y_i|} \sum_{y \in Y_i} \frac{|\{rank_f(x_i, y') \leq rank_f(x_i, y), y' \in Y_i\}|}{rank_f(x_i, y)} \quad (13)$$

2) Hamming Loss: 用于度量样本在单个标记的真实标记和预测标记的错误匹配情况。

$$HL_D(h) = \frac{1}{n} \sum_{i=1}^n \frac{1}{|Y_i|} |h(x_i) \neq Y_i| \quad (14)$$

3) One Error: 评估对象最高排位标记并未正确标记的次數情况。

$$OE_D(f) = \frac{1}{n} \sum_{i=1}^n \left[ \arg \max_{y \in Y} f(x_i, y) \notin Y_i \right] \quad (15)$$

4) Ranking Loss: 用来考察样本的不相关标记的排序低于相关标记的排序的情况。

$$RL_D(f) = \frac{1}{n} \sum_{i=1}^n \frac{1}{|Y_i| |\bar{Y}_i|} \left| \left\{ (y_1, y_2) \mid f(x_i, y_1) \leq f(x_i, y_2), (y_1, y_2) \in Y_i \times \bar{Y}_i \right\} \right| \quad (16)$$

#### 3.3 算法选择与相关参数设置

为验证新提出的特征选择算法性能,实验将 MFSEF 算法与 4 个经典多标记特征选择算法进行对比,分别是 MDDM<sub>spc</sub>、MDDM<sub>proj</sub>、PMU 及 MFSLs。表 2 到表 5 中的第 2 列 Original 表示在原始特征空间下仅通过基本的经典多标记分类器 ML- $k$ NN 的分类性能;MDDM 是基于最大相关性的多标记维度约简算法,而 MDDM 又可划分为 MDDM<sub>spc</sub> 和 MDDM<sub>proj</sub>,PMU 是通过多元互信息最大化已选特征与标记集合的相关性,提出基于多变量互信息的多标记特征选择算法。其中 MDDM<sub>spc</sub> 和 MDDM<sub>proj</sub> 算法需先进行原始数据归一化,再进行特征选择,MFSLs 和 PMU 是针对离散型数据进行处理,鉴于此,为了实验的严谨和合理性,以 MFSLs 和 PMU 离散

化方法为基准,对本文实验数据先两折离散化。由文献[18]可知,当选择多标记数据集的特征维度不高时,将子空间划分为3个,特征采样比例设置为{0.6,0.3,0.1},专家特征个数 $k_1$ 设为4。另本实验最后分类器采用ML- $k$ NN,故相关参数选择默认值,近邻个数 $k$ 取10,平滑参数 $s$ 取1。

### 3.4 实验结果

表2到表5给出了本文算法和其他4种算法在7个多标记数据集上实验结果,最好的结果加粗表示;同时,每种方法在所有数据集上的平均排序结果列在最后一行;数据右下角括号数字表示6种算法分别在评价指标下的排序序号。

表2 各算法在7个数据集上的平均精度测试结果

Tab. 2 AP(↑) results of each algorithm on 7 datasets

数据集	Original	MDDMspc	MDDMproj	PMU	MFSLS	MFSEF
Health	0.6812 <sub>(3)</sub>	0.6794 <sub>(4)</sub>	0.6516 <sub>(6)</sub>	0.6709 <sub>(5)</sub>	0.7237 <sub>(2)</sub>	<b>0.7280</b> <sub>(1)</sub>
Recreation	0.4547 <sub>(4)</sub>	0.4497 <sub>(5)</sub>	0.4628 <sub>(3)</sub>	0.4441 <sub>(6)</sub>	0.5102 <sub>(2)</sub>	<b>0.5225</b> <sub>(1)</sub>
Artificial	0.5094 <sub>(3)</sub>	0.4974 <sub>(4)</sub>	0.4848 <sub>(6)</sub>	0.4909 <sub>(5)</sub>	0.5363 <sub>(2)</sub>	<b>0.5364</b> <sub>(1)</sub>
Reference	0.6194 <sub>(3)</sub>	0.6014 <sub>(5)</sub>	0.5992 <sub>(6)</sub>	0.6145 <sub>(4)</sub>	0.6304 <sub>(2)</sub>	<b>0.6347</b> <sub>(1)</sub>
Entertainment	0.6023 <sub>(3)</sub>	0.5513 <sub>(6)</sub>	0.5588 <sub>(5)</sub>	0.5671 <sub>(4)</sub>	0.6032 <sub>(2)</sub>	<b>0.6040</b> <sub>(1)</sub>
Business	<b>0.8798</b> <sub>(1)</sub>	0.8707 <sub>(5)</sub>	0.8731 <sub>(4)</sub>	0.8628 <sub>(6)</sub>	0.8762 <sub>(3)</sub>	0.8765 <sub>(2)</sub>
Computer	0.6335 <sub>(3)</sub>	0.6319 <sub>(4)</sub>	0.6250 <sub>(6)</sub>	0.6252 <sub>(5)</sub>	0.6405 <sub>(2)</sub>	<b>0.6424</b> <sub>(1)</sub>
均值	0.6258	0.6117	0.6079	0.6108	0.6458	0.6492

表3 各算法在7个数据集上的海明损失测试结果

Tab. 3 HL(↓) results of each algorithm on 7 datasets

数据集	Original	MDDMspc	MDDMproj	PMU	MFSLS	MFSEF
Health	0.0458 <sub>(5)</sub>	0.0438 <sub>(4)</sub>	0.0462 <sub>(6)</sub>	0.0435 <sub>(3)</sub>	0.0410 <sub>(2)</sub>	<b>0.0386</b> <sub>(1)</sub>
Recreation	0.0618 <sub>(3)</sub>	0.0633 <sub>(5)</sub>	0.0619 <sub>(4)</sub>	0.0637 <sub>(6)</sub>	0.0598 <sub>(2)</sub>	<b>0.0588</b> <sub>(1)</sub>
Artificial	0.0612 <sub>(4)</sub>	0.0616 <sub>(6)</sub>	0.0609 <sub>(3)</sub>	0.0615 <sub>(5)</sub>	<b>0.0587</b> <sub>(1)</sub>	0.0594 <sub>(2)</sub>
Reference	0.0314 <sub>(4)</sub>	0.0324 <sub>(6)</sub>	0.0311 <sub>(3)</sub>	0.0307 <sub>(2)</sub>	0.0315 <sub>(5)</sub>	<b>0.0288</b> <sub>(1)</sub>
Entertainment	0.0612 <sub>(4)</sub>	0.0624 <sub>(6)</sub>	0.0620 <sub>(5)</sub>	0.0607 <sub>(3)</sub>	0.0594 <sub>(2)</sub>	<b>0.0591</b> <sub>(1)</sub>
Business	<b>0.0269</b> <sub>(1)</sub>	0.0280 <sub>(4,5)</sub>	0.0280 <sub>(4,5)</sub>	0.0285 <sub>(6)</sub>	0.0274 <sub>(3)</sub>	0.0272 <sub>(2)</sub>
Computer	0.0412 <sub>(6)</sub>	0.0408 <sub>(4,5)</sub>	0.0408 <sub>(4,5)</sub>	0.0405 <sub>(3)</sub>	0.0401 <sub>(2)</sub>	<b>0.0400</b> <sub>(1)</sub>
均值	0.0471	0.0475	0.0473	0.0470	0.0454	0.0446

表4 各算法在7个数据集上的排序损失测试结果

Tab. 4 RL(↓) results of each algorithm on 7 datasets

数据集	Original	MDDMspc	MDDMproj	PMU	MFSLS	MFSEF
Health	0.0605 <sub>(3)</sub>	0.0652 <sub>(5)</sub>	0.0699 <sub>(6)</sub>	0.0635 <sub>(4)</sub>	<b>0.0563</b> <sub>(1)</sub>	0.0565 <sub>(2)</sub>
Recreation	0.1914 <sub>(5)</sub>	0.1892 <sub>(3)</sub>	0.1924 <sub>(6)</sub>	0.1895 <sub>(4)</sub>	0.1770 <sub>(2)</sub>	<b>0.1730</b> <sub>(1)</sub>
Artificial	0.1520 <sub>(3)</sub>	0.1539 <sub>(5)</sub>	0.1576 <sub>(6)</sub>	0.1530 <sub>(4)</sub>	0.1468 <sub>(2)</sub>	<b>0.1457</b> <sub>(1)</sub>
Reference	0.0919 <sub>(4)</sub>	0.0925 <sub>(5)</sub>	0.0933 <sub>(6)</sub>	0.0870 <sub>(2)</sub>	0.0883 <sub>(3)</sub>	<b>0.0855</b> <sub>(1)</sub>
Entertainment	0.1154 <sub>(3)</sub>	0.1264 <sub>(6)</sub>	0.1258 <sub>(5)</sub>	0.1226 <sub>(4)</sub>	0.1150 <sub>(2)</sub>	<b>0.1131</b> <sub>(1)</sub>
Business	<b>0.0374</b> <sub>(1)</sub>	0.0433 <sub>(5)</sub>	0.0416 <sub>(4)</sub>	0.0459 <sub>(6)</sub>	0.0402 <sub>(2)</sub>	0.0407 <sub>(3)</sub>
Computer	0.0922 <sub>(4)</sub>	0.0919 <sub>(3)</sub>	0.0945 <sub>(5)</sub>	0.0952 <sub>(6)</sub>	0.0896 <sub>(2)</sub>	<b>0.0894</b> <sub>(1)</sub>
均值	0.1058	0.1089	0.1107	0.1081	0.1019	0.1006

表5 各算法在7个数据集上的1-错误率测试结果

Tab. 5 OE(↓) results of each algorithm on 7 datasets

数据集	Original	MDDMspc	MDDMproj	PMU	MFSLS	MFSEF
Health	0.4207 <sub>(4)</sub>	0.4053 <sub>(3)</sub>	0.4463 <sub>(6)</sub>	0.4313 <sub>(5)</sub>	0.3477 <sub>(2)</sub>	<b>0.3400</b> <sub>(1)</sub>
Recreation	0.7063 <sub>(4)</sub>	0.7143 <sub>(5)</sub>	0.6883 <sub>(3)</sub>	0.7190 <sub>(6)</sub>	0.6163 <sub>(2)</sub>	<b>0.6083</b> <sub>(1)</sub>
Artificial	0.6327 <sub>(3)</sub>	0.6470 <sub>(4)</sub>	0.6670 <sub>(6)</sub>	0.6570 <sub>(5)</sub>	0.5840 <sub>(2)</sub>	<b>0.5837</b> <sub>(1)</sub>
Reference	0.4730 <sub>(3)</sub>	0.4973 <sub>(6)</sub>	0.4960 <sub>(5)</sub>	0.4910 <sub>(4)</sub>	0.4617 <sub>(2)</sub>	<b>0.4540</b> <sub>(1)</sub>
Entertainment	0.5297 <sub>(2)</sub>	0.6020 <sub>(6)</sub>	0.5997 <sub>(5)</sub>	0.5847 <sub>(4)</sub>	<b>0.5257</b> <sub>(1)</sub>	0.5330 <sub>(3)</sub>
Business	<b>0.1213</b> <sub>(1)</sub>	0.1283 <sub>(5)</sub>	0.1263 <sub>(4)</sub>	0.1360 <sub>(6)</sub>	0.1227 <sub>(2,5)</sub>	0.1227 <sub>(2,5)</sub>
Computer	0.4363 <sub>(3)</sub>	0.4453 <sub>(4)</sub>	0.4537 <sub>(6)</sub>	0.4497 <sub>(5)</sub>	0.4323 <sub>(2)</sub>	<b>0.4307</b> <sub>(1)</sub>
均值	0.4743	0.4914	0.4968	0.4955	0.4415	0.4389

实验结果显示:MFSEF在Health、Recreation、Artificial、Reference、Entertainment、Business、Computer等7个多标记数据集上实验结果的平均排序都占优,其中,对于AP指标,平均数值越大,算法性能越优,对于其他三个评价指标,平均数值越小,算法性能越优。从表2到表5可发现:

① 在AP指标中,MFSEF仅在Business数据集中AP不是最优,排名第二。对比4种算法和原始特征空间,MFSEF在其他数据集中AP值最大,即分类性能达到最优。

② 在HL指标中,MFSEF在Business和Artificial数据集中HL值排名第二,对比4种算法和原始特征空间,MFSEF在其

他数据集中 HL 值最小,即分类性能达到最优。

③ 在 RL 指标中, MFSEF 在 Health 和 Business 数据集中分别排第二和第三,在其他 5 个数据集对比 4 种算法和原始特征空间, MFSEF 的 RL 值最小,即分类性能达到最优。

④ 在 OE 指标中, MFSEF 在 Entertainment 数据集上排第三,在 Business 数据集和 MFSLS 对比算法并列排第二,在其他 5 个数据集对比 4 种算法和原始特征空间, MFSEF 的 OE 值最小,即分类性能达到最优。

上述实验分析充分表明,通过本文算法特征选择后的特征子集在后续分类性能上,对比其他 4 种算法和原始特征空间在 7 个多标记数据集中多数占优,验证了本文算法的有效性和鲁棒性。

图 4 是对比其他 4 种算法,随着选择后的特征数目的逐渐变多,其分类性能的变化情况。针对每一种算法,都有 28 种对比结果。介于篇幅所限,本文只选取了 Artificial 数据集的曲线图进行分析,分别展示了在 AP、HL、RL 和 OE 四种评价指标时,5 种算法在特征数逐渐变大时分类性能的变化情况。对比原始特征空间和 PMU、MDDM<sub>spsc</sub>、MDDM<sub>proj</sub>、MFSLS 这 4 种算法的分类性能,在 Artificial 数据集上, MFSEF 基本上占优。基本上在前 80 个特征范围类,本文算法在四个评价指标上均明显优于其他对比算法,并且往往能通过较少的特征数更快地达到更好的分类性能。另外,在其他未展示的数据集上,本文算法的分类性能曲线变化也多数占优,这充分地表明 MFSEF 的有效性和合理性。

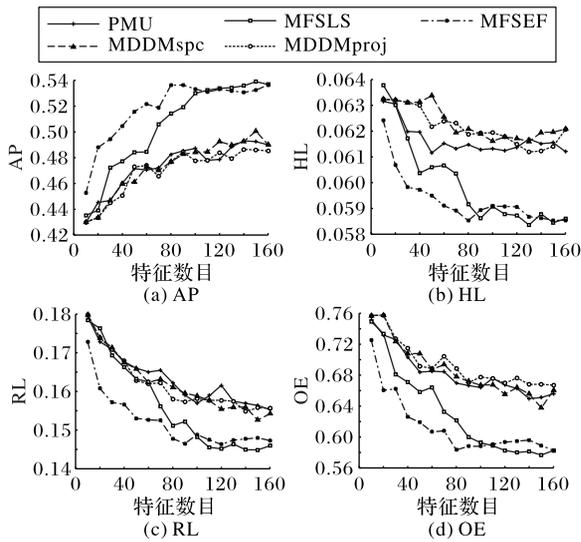


图 4 Artificial 数据集的各个评价指标的性能变化

Fig. 4 Changes in performance of various evaluation indicators on Artificial dataset

### 3.5 统计假设检验

在上述实验分析中,本文算法达到了显而易见的效果,下面将运用统计学中统计假设检验<sup>[21]</sup>进一步说明本文算法的有效性和合理性。

统计假设检验:在上述 7 个数据集上采用显著性水平为 5% 的 Nemenyi 检验<sup>[22-23]</sup>来对比 MFSEF 算法与其他对比算法。如果两个算法在所有数据集上的平均排序的差值小于或者等于临界差值(Critical Difference, CD),那么这两个算法之间没有显著性差异,反之存在显著性差异。如图 5 所示,在最上行为临界值 CD=2.850 0 时,若两个算法之间没有显著性差异则

用实线连接。在图 5 中,随着坐标轴上的数值增大其算法性能依次降低。

图 5 展示了各算法在 AP、RL、HL、OE 四个指标上的 CD 图,从中可以看出,本文算法在 4 个指标上性能均处于首位。具体在平均精度指标上,如图 5(a)、(b)、(d)所示,在 Average Precision、Ranking Loss 和 One Error 三个指标上, MLSEF 与 MDDM<sub>spsc</sub>、PMU、MDDM<sub>proj</sub> 均有显著差异,且优于这三种算法;如图 5(c)在 Hamming Loss 指标上,本文算法与 MDDM<sub>proj</sub> 和 MDDM<sub>spsc</sub> 有显著差异,且优于这两种算法。与其他算法相比,在统计上,本文算法有 45% 的概率与其他算法无显著性差异。

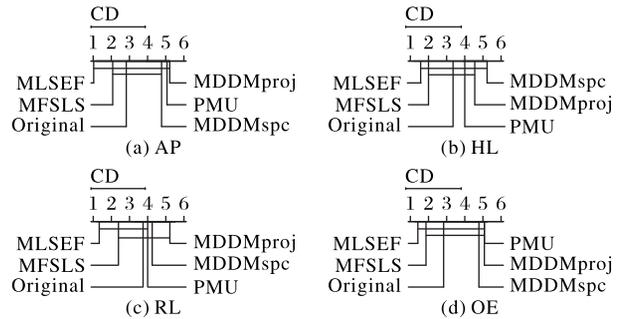


图 5 算法综合性能比较

Fig. 5 Performance comparison of algorithms.

通过以上针对图 5 的分析,可知本文算法综合性能最为优秀,在统计上也优于其他对比算法。基于以上的实验结果和统计分析再次充分表明本文算法的优越性。

### 4 结语

与通过相关准则挑选特征和标记相关性强的多标记特征选择算法相比,本文不仅考虑了重要性较次要的特征可能对分类性能产生的影响,还考虑了可能决定整体预测方向的关键特征。模拟现实生活中的实际情况,通过经验优先挑选出部分专家特征与剩余特征相联合,利用条件互信息和联合互信息理论得出一个与标记集合相关性由强到弱的特征序列,再通过划分子空间去除冗余性较大的特征,最后保留专家特征和挑选出的新的特征作为最后的特征子集。在已公开的多个基准多标记数据集上的实验结果表明,该算法在实验中较其他对比的多标记特征选择算法有一定优势和较好的稳定性,且更具有实际应用价值。

本文算法在专家特征的选取上还可以进一步探讨,目前只是通过瀑布图联合互信息理论模拟选出几个专家特征,所以最后结果可能由于个人选取专家特征的不同,实验结果和预期效果存在一定的误差,但是针对具体问题分析数据,合理选择专家特征,还是可以有效减少误差。

### 参考文献 (References)

- [1] GIBAJA E, VENTURA S. A tutorial on multilabel learning [J]. ACM Computing Surveys, 2015, 47(3): 1-38.
- [2] 何志芬, 杨明, 刘会东. 多标记分类和标记相关性的联合学习 [J]. 软件学报, 2014, 25(9): 1967-1981. (HE Z F, YANG M, LIU H D. Joint learning of multi-label classification and label correlations [J]. Journal of Software, 2014, 25(9): 1967-1981.)
- [3] WANG Z, CHEN T, LI G, et al. Multi-label image recognition by recurrently discovering attentional regions [C]// Proceedings of the 2017 IEEE International Conference on Computer Vision. Piscataway:

- IEEE, 2017: 464-472.
- [4] OZONAT K, YOUNG D. Towards a universal marketplace over the web: statistical multi-label classification of service provider forms with simulated annealing [C]// Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. New York: ACM, 2009: 1295-1304.
- [5] SHU X, LAI D, XU H, et al. Learning shared subspace for multi-label dimensionality reduction via dependence maximization [J]. *Neurocomputing*, 2015, 168: 356-364.
- [6] PEREIRA R B, PLASTINO A, ZADROZNY B, et al. Categorizing feature selection methods for multi-label classification[J]. *Artificial Intelligence Review*, 2018, 49(1): 57-78.
- [7] ZHANG Y, ZHOU Z. Multilabel dimensionality reduction via dependence maximization [J]. *ACM Transactions on Knowledge Discovery from Data*, 2010, 4(3): No. 14.
- [8] LEE J, KIM D W. Feature selection for multi-label classification using multivariate mutual information [J]. *Pattern Recognition Letters*, 2013, 34(3): 349-357.
- [9] LIN Y, HU Q, LIU J, et al. Multi-label feature selection based on neighborhood mutual information [J]. *Applied Soft Computing*, 2016, 38: 244-256.
- [10] 刘景华,林梦雷,王晨曦,等. 基于局部子空间的多标记特征选择算法[J]. *模式识别与人工智能*, 2016, 29(3): 240-251. (LIU J H, LIN M L, WANG C X, et al. Multi-label feature selection algorithm based on local subspace [J]. *Pattern Recognition and Artificial Intelligence*, 2016, 29(3): 240-251.)
- [11] 王晨曦,林耀进,唐莉,等. 基于信息粒化的多标记特征选择算法[J]. *模式识别与人工智能*, 2018, 31(2): 123-131. (WANG C X, LIN Y J, TANG L, et al. Multi-label feature selection based on information granulation [J]. *Pattern Recognition and Artificial Intelligence*, 2018, 31(2): 123-131.)
- [12] LEE J, LIM H, KIM D W. Approximating mutual information for multi-label feature selection [J]. *Electronics Letters*, 2012, 48(15): 929-930.
- [13] YU S, HUANG T. Exponential weighted entropy and exponential weighted mutual information [J]. *Neurocomputing*, 2017, 249: 86-94.
- [14] LI F, MIAO D, PEDRYCZ W. Granular multi-label feature selection based on mutual information[J]. *Pattern Recognition*, 2017, 67: 410-423.
- [15] ZHANG M, ZHOU Z. ML- $k$ NN: a lazy learning approach to multi-label learning[J]. *Pattern recognition*, 2007, 40(7): 2038-2048.
- [16] ZHANG M, ZHOU Z. A review on multi-label learning algorithms [J]. *IEEE Transactions on Knowledge and Data Engineering*, 2014, 26(8): 1819-1837.
- [17] FLEURET F. Fast binary feature selection with conditional mutual information [J]. *Journal of Machine Learning Research*, 2004, 5: 1531-1555.
- [18] 杨明,王飞. 一种基于局部随机子空间的分类集成算法[J]. *模式识别与人工智能*, 2012, 25(4): 595-603. (YANG M, WANG F. A classifier ensemble algorithm based on local random subspace [J]. *Pattern Recognition and Artificial Intelligence*, 2012, 25(4): 595-603.)
- [19] TSOUMAKAS G, VLAHAVAS I. Random  $k$ -labelsets: an ensemble method for multilabel classification [C]// Proceedings of the 18th European Conference on Machine Learning, LNCS 4701. Berlin: Springer, 2007: 406-417.
- [20] ZHANG M, PEÑA J M, ROBLES V. Feature selection for multi-label naive Bayes classification [J]. *Information Sciences*, 2009, 179(19): 3218-3229.
- [21] DEMŠAR J. Statistical comparisons of classifiers over multiple data sets [J]. *Journal of Machine Learning Research*, 2006, 7: 1-30.
- [22] ZHANG M, WU L. LIFT: multi-label learning with label-specific features [J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2015, 37(1): 107-120.
- [23] 程玉胜,钱坤,王一宾,等. 融合萤火虫方法的多标签懒惰学习算法[J]. *计算机应用*, 2019, 39(5): 1305-1311. (CHENG Y S, QIAN K, WANG Y B, et al. Multi-label lazy learning algorithm based on firefly method [J]. *Journal of Computer Applications*, 2019, 39(5): 1305-1311.)

This work is partially supported by the Key University Natural Science Research Project of Anhui Province (KJ2017A352), the Program for Innovative Research Team in Anqing Normal University.

**CHENG Yusheng**, born in 1969, Ph. D, professor. His research interests include big data, rough set, machine learning for feature selection.

**SONG Fan**, born in 1992, M. S. candidate. His research interests include multi-label learning, neural network.

**WNAG Yibin**, born in 1970, M. S., professor. His research interests include multi-label learning, machine learning, software security.

**QIAN Kun**, born in 1995, M. S. candidate. His research interests include multi-label learning, machine learning, data statistics.