

文章编号: 1003-0077 (2015) 00-0000-00

基于多标记学习的汽车评论文本多性能识别 *

张晶¹ 李德玉^{1,2} 王素格^{1,2}

(1.山西大学 计算机与信息技术学院, 山西 太原 030006;

2.山西大学 计算智能与中文信息处理教育部重点实验室, 山西 太原 030006)

摘要: 针对汽车产品评论文本中出现的多方面性能, 提出一种基于多标记学习的汽车评论文本多方面性能识别方法。首先, 结合文本挖掘方法, 利用多标记文本特征选择方法选取特征, 将非结构化的文本转化为结构化的多标记数据集。在此基础上, 使用 4 种多标记分类方法, 对待识别的评论文档标注一个或多个方面标记。最后, 以 8 种多标记评价指标评估方面识别的性能。在新浪汽车评论语料上的实验表明, 方面识别的子集准确率达到 95%, 验证了方法的可行性。

关键词: 多标记学习; 文本处理; 汽车评论; 多方面识别

中图分类号: TP391

文献标识码: A

Multiple Performances Identification for Car Review Texts Based on Multi-label Learning

Zhang Jing¹ Li De-yu^{1,2} Wang Su-ge^{1,2}

(1.School of Computer and Information Technology, Shanxi University, Taiyuan, Shanxi 030006, China;

2.Key Laboratory of Computational Intelligence and Chinese Information Processing of Ministry of Education, Shanxi University, Taiyuan, Shanxi 030006, China)

Abstract: Aiming at the characteristics of the multi-aspect performance appeared in the automotive product reviews, this paper proposed a novel method for recognizing the multiple aspects of performance about car comment text based on multi-label learning. Firstly, appropriate words were selected as features by multi-label text feature selection method combined with the text mining technology, and then, the unstructured document corpus are transformed into structured multi-label dataset. After that, we finished marking one or more aspect tags for the unrecognized comment text with four multi-label classification methods. Finally, the recognition accuracy of multiple aspects was assessed by eight multi-label evaluation metrics. On the sina car review corpus, experimental results indicate the subset accuracy reaches up to 95%. Hence, our method was feasible for recognizing the multiple aspects of automobile reviews.

Key words: Multi-label learning; Text processing; Car reviews; Multi-aspect recognition

1 引言

近年来, 评论文本的观点挖掘研究越来越多地受到研究学者的关注^[1]。随着方法研究的深入与实际应用需求的增大, 数据挖掘领域的学者倾向于关注产品的多方面观点挖掘与分析^[2]。例如, 当我们需要购买汽车的时候, 通常会在汽车网站浏览汽车产品的相关评论, 然而我们所看到的一般都是对某种汽车的总体评价及打分。事实上, 人们更想要了解这些评论具体涉及到汽车的哪些方面, 以及各方面的评价信息。此外, 由于评论者对不同方面的认知

收稿日期:

定稿日期:

基金项目: 国家自然科学基金资助项目(No. 61272095, 61175067); 山西省科技攻关项目

(No.20110321027-02); 山西省回国留学人员科研项目(No.2013-014)资助; 山西省科技基础条件平台建设项目(No. 2015091001-0102)

感受存在差异,相同的总体打分往往无法真实体现评论者对各性能方面的重视程度,而企业需要更具体的信息了解消费者的需求。因此,根据评论文本的内容,分析各种品牌不同型号汽车多种性能方面的优劣,对企业和消费者而言显得尤为必要。

对于汽车产品性能的多方面观点挖掘研究,现有的评论文本观点挖掘技术难以解决这类问题。由于这项研究涉及到评论对象的多个方面,对产品的观点挖掘首先需要对产品的方面抽取,而后再结合现有的观点挖掘方法,完成对评论对象的打分及观点分析。以往对于汽车性能的多方面识别,主要是根据汽车本体知识人工构成的用户词典来完成方面的区分。然而,如果将汽车的多种性能(如安全性、操控性等)看作是一篇文档中涉及的多个标记,那么对汽车多方面观点挖掘问题中的方面识别,可以抽象为多标记分类问题,分类结果可以给出性能相关性的排序,为下一步方面的打分奠定基础。目前,多标记学习作为数据挖掘领域一项重要的研究内容,已经得到了广泛的应用。关于多标记文本分类的研究,主要针对多标记文本的特征选取和分类算法两方面。多标记文本分类方法可用于网页文档的主题分类^[3]、微博文本的情感分析^[4]、图书网站的自动标签推荐^[5]等。该项研究不仅能够为评论文本的多方面观点挖掘提供新的思路,而且可以为多标记学习开创新的应用领域。

针对汽车评论文本中存在的多方面性能识别问题,本文以文本挖掘技术为基础,结合多标记文本特征选择方法,将非结构化的汽车评论文本表示为结构化的多标记数据。采用四种多标记分类方法,完成对汽车评论文本的多方面性能识别。通过8种多标记评价指标,评估该方法对于汽车评论文本的多性能识别的有效性。

2 相关工作

2.1 评论文本的方面识别与抽取

以往关于评论文本的研究主要在于观点挖掘与倾向性分析。近年来,随着网络文档的多元化发展,评论文本的多方面观点挖掘引起了数据挖掘领域研究者的关注^[2]。而这项工作,首先需要完成对评论文本的多个方面的识别与抽取。目前,对评论文本进行方面识别与抽取的方法主要有基于关联规则、主题模型以及人工定义种子词等方法。

Hu 和 Liu^[6]提出一种基于关联规则的产品评论方面抽取方法。该方法认为评论文本中频繁出现的名词能够很好地体现方面特点,因此,它采用词频统计的方法选取方面表征词,并利用搭配关系找出出现次数较少的特征词,从而完成多个方面的抽取任务。Lu 和 zhai^[7]利用方面的先验知识,提出了两种基于潜在语义模型,分别为非结构化 PLSA 和结构化 PLSA 模型,完成评论文本的方面抽取,将概率主题模型用于对评论短文本的方面摘要的自动生成。Wang 和 Lu^[8]基于 Boot-strapping 方法,提出了一种方面评级回归方法,用于对酒店评论文本的方面识别与抽取。该方法首先人工定义需要抽取的方面,并人工为每个方面定义关键词集合,通过 Boot-strapping 方法对关键词集合的扩充,使用评论文本与关键词集合的匹配识别评论文本的方法。

目前,关于评论文本的方面抽取大多是针对英文语料,且对于汽车领域中汽车性能的方面识别^[9]鲜有相关的研究。因此,本文主要针对汽车这一产品的中文评论文本进行分析,识别其性能的多个方面。

2.2 多标记文本分类方法

在传统的监督学习问题中,单标记分类指一个对象只能被划分到一个类别中,然而,由于存在语义多义性,有的对象需要同时被标记多个类别,这类问题被称为多标记分类。多标记学习技术是数据挖掘领域研究的热点问题^[10]。目前,多标记学习方法主要分为两种:问题转换和算法适应方法。

对多标记文本的分类算法而言,BR(Binary Relevance)^[11]、CC(Classifier Chains)^[12]、

RAKEL(Random k-labelsets)^[13]等多标记分类方法,是问题转换类方法的代表,这类方法旨在将多标记问题转化为一个或多个多标记分类问题,适用于所有类型的多标记分类问题。针对标记间存在的依赖关系,提出了一些基于集成学习和层次分类方法。算法适应方法,是将传统的单标记分类算法进行改进,使其能够直接处理多标记数据。MLkNN^[14]、AdaBoost.MH^[15]、RankSVM^[16]、BPMLL^[17]是这类方法中多标记文本分类的代表。

对多标记文本的特征选择而言,为了减小不相关或冗余特征对分类精度的影响,研究者提出了多种文本特征选择算法。主要是将传统的卡方统计、信息增益、ReliefF等在文本特征选择方面性能较为突出的方法与BR、LP问题转换方法相结合,即可完成多标记文本的特征选择^[18,19]。同时,对现有处理单标记问题的特征选择方法进行改造,也可以对多标记文本实现特征的降维^[20]。

3 汽车多性能识别框架

3.1 问题描述与框架

关于新浪汽车网站上的马自达CX-5品牌车型的用户评论,如图1所示。以此为例可以看出,网友的评论包含汽车的舒适性和经济性等性能方面,针对方面的观点挖掘能够更细致地反应汽车性能特点。



图1 新浪汽车用户评论示例

从图1可以发现,用户在对汽车评论的过程中,往往会将评论对象的一个或多个方面的观点表述在一句话中,由此,评论文本的多方面识别任务可以描述为,根据评论文档的内容,对每一篇文档标记上相关程度较大的方面标签。因此,本文将汽车评论文本的多方面性能识别构造为多标记文本的分类问题,首先给出如下定义:

由汽车评论文本构成的多标记训练数据集用一个三元组 (D, T, L) 表示,其中, $D = \{D_1, D_2, \dots, D_n\} = \{(d_1, y_1), (d_2, y_2), \dots, (d_n, y_n)\}$,表示由汽车这一实体的 n 篇评论文档构成的多标记数据集, $D_i \in D$,每篇文档 D_i 由特征向量 d_i 和标记向量 y_i 两部分组成, $1 \leq i \leq n$ 。

$T = \{t_1, t_2, \dots, t_p\}$ 表示从 n 篇评论文档中选择的,由 p 项关键词构成的特征集合。具体而言, T 是由能够刻画文档表示方面的词汇(即评论对象的属性等表征词)构成的词汇集合。

$L = \{l_1, l_2, \dots, l_q\}$ 表示由 q 种标签构成的标记集合。具体为,文档集合 D 所涉及的多个性能构成的方面集合,本文使用汽车的舒适性、动力性、操控性、服务性、经济性和安全性六个性能方面^[21]。

特征向量 $d_i = \langle w_{1i}, w_{2i}, \dots, w_{ji}, \dots, w_{pi} \rangle$,其中, w_{ji} 表示关键词 t_j 在文档 D_i 中相应的权值。每篇文档对应于标记集合 L 中的一个或多个性能标签,并由0和1构成一个二值向量 y_i ,如果 D_i 包含类别 $l_j (1 \leq j \leq p)$,则 $y_{ji} = 1$,否则为 $y_{ji} = 0$ 。

该模型的目的是利用给定的特征空间与标记空间的训练数据以及现有的多标记分类方法，通过有监督的学习过程，对未标识方面的多性能评论文档标注一个或多个最接近真实情况的性能标签，并利用多标记分类器给出相关方面的排序。

根据上述定义的问题模型，设计汽车评论文本多性能识别框架如图 2 所示。

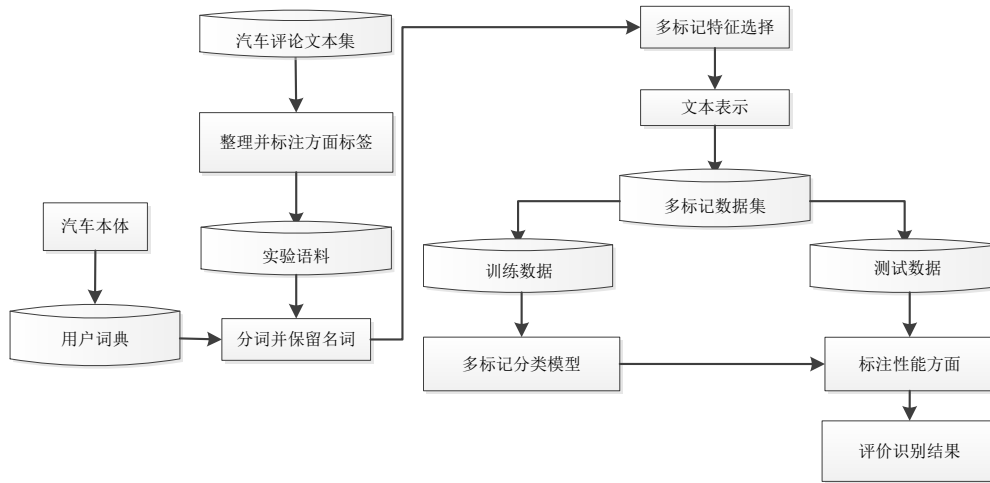


图 2 汽车评论文本多性能识别框架

根据图 2 可得整个框架流程描述为：对从新浪汽车网站上获取的汽车评论文本进行整理，从中随机抽取部分语料，为每篇文档标注出方面标签，删除人工无法标记的文档，构成实验语料；将汽车本体特征加入用户词典，用分词软件对实验数据分词，抽取所有名词性词语作为候选特征集；采用多标记文本特征选择方法，从候选特征集中选取部分词作为分类的特征，用向量空间模型对文档进行表示，构成结构化的多标记数据集；选用多标记分类方法，对训练数据集进行训练，得到每种方法的训练模型；将训练模型用于测试数据集，标注出测试集中汽车评论文本的多个性能方面。

3.2 文本预处理

中文文本的预处理过程与英文相比较为复杂。针对本文需要处理的汽车评论文本，首先需要将中文文档进行分词处理。为避免破坏部分具有代表性的词语，文本首先将汽车本体特征词^[21]添加到用户词典中，然后使用中科院 ICTCLAS 分词软件对全部文档进行分词，并标注出词性，将结果以文本形式存储。

在评论文本的方面抽取中，通常名词性词语能更有效地表示方面特征，因此，本文对分词结果去除非名词性词语，即保留的词语中包含普通名词以及动名词等。如标注文档“最低/a 配/v 增加/v 点/qt 必要/a 配置/vn 吧/y ! /wt”中，“配置/vn”一词显然能表示文本方面特征，需要将其作为候选特征。

3.3 特征选择

在利用多标记分类器对文档标注方面标记之前，需要确定特征集合 T 以及每篇评论文档对应的特征向量 d_i 。中文分词后的词汇量比较大，而部分词对分类作用不大甚至会影响分类的性能。因此，在抽取的名词性词语的基础上，为了得到更好的分类性能，需要进一步精简候选特征。由于本文主要是在多标记文本分类方法的基础上完成评论文本的方面识别任务，因此，我们依据多标记文本的特征选择方法，选取与方面标记具有较强相关性的词，作为分类学习过程中的特征集合 T 中的关键词，从而达到降维的目的。

单标记文本分类的特征选择主要采用信息增益、卡方统计、文档频率等方法。考虑到多标记数据与单标记数据在标签方面的区别，根据 2.2 节关于多标记文本特征选择方法的介绍，

本文采用文献[19][22]中提及的 Chi-BR 多标记特征选择方法完成对文本特征的降维。该方法是在 χ^2 统计的基础上,结合 BR 多标记转换方法,将一个汽车评论多标记数据集转换为多个单标记数据集,每个单标记数据集的标签分为正反两种,即将文档在某一性能方面上分成正例和反例。Chi-BR 方法的 χ^2 统计度量了一个词 $word$ 和某一方面 $l_j(1 \leq j \leq q)$ 之间的相关性,见公式 (1)。

$$\chi^2(word, l_j) = \frac{n[P(word, l_j)P(\overline{word}, \overline{l_j}) - P(word, \overline{l_j})P(\overline{word}, l_j)]^2}{P(word)P(\overline{word})P(l_j)P(\overline{l_j})} \quad (1)$$

其中, n 表示文档总数, $P(word, l_j)$ 表示词 $word$ 在文档 D_i 中且 $l_{ij} = 1$ 时组合出现的次数,

同理, \overline{word} 表示词 $word$ 不在文档 D_i 中, 用 $\overline{l_j}$ 表示 $l_{ij} = 0$ 。

文档集 D 中, 词 $word$ 的 Chi-BR 方法中 χ^2 值的度量采用平均 χ^2 值和最大 χ^2 值两种聚合策略, 见公式(2)-(3)。

$$\chi^2(word) = \frac{1}{q} \sum_{i=1}^q \chi^2(word, l_i) \quad (2)$$

$$\chi^2(word) = \max_{1 \leq i \leq q} \chi^2(word, l_i) \quad (3)$$

利用 Chi-BR 方法对每个特征词的 χ^2 值排序, 可从高到低选取部分词作为特征项, 以词频作为特征的权值, 采用常用的向量空间模型对文本进行表示, 即可得到每篇评论文档对应的特征向量 d_i 。

3.4 分类方法

多标记分类方法是一种有监督的机器学习方法。在方面识别的任务中, 它主要根据已有的经验知识, 采用特定的分类学习系统中的函数关系, 依据给定的刻画文档表示方面的特征向量 d_i , 预测与该文档相关的方面标记向量 y_i , 从而识别相应的方面标记集合 L 。

由于目前还未看到有关多标记分类方法处理评论文本多方面识别的研究。因此, 本文选用四种典型的多标记分类方法完成方面识别任务, 并比较各方法的优劣。四种分类方法分别为 BR^[11]、CC^[12]、MLkNN^[14]以及 rankSVM^[16]。其中, BR 和 CC 方法是基于问题转化方法; 而 MLkNN 和 rankSVM 两种方法是基于算法适应方法, 分别根据传统的 k 近邻和支持向量机(SVM)进行改进。同时, 在标记相关性方面, BR 和 MLkNN 方法对多标记数据分类时, 不考虑标记之间的相关性, 即假定各标记之间是相互独立的。BR 方法将多标记数据转化为等同于标记数量的多个单标记数据集, 适用于标记数量少的情况, MLkNN 方法是基于最大化后验概率的原则确定标记集合; 而 rankSVM 方法考虑两两标记之间的相关性, 是在“间隔最大化”策略的基础上, 最小化 RankingLoss 损失函数, 对于 CC 方法, 它是考虑全局的相关性, 是对 BR 方法的一种改进, 将 BR 方法训练的多个二分类器连成一条链, 每一个二分类器的训练结果作为下一个分类器训练过程中的样本属性参与训练。

4 实验结果及分析

4.1 实验数据

本文对从新浪汽车网站上搜集的评论文本进行整理, 随机抽取部分评论文本, 选出 2000 条有一个或多个性能方面的评论, 标注方面标签, 构成实验语料。标记集合 L 中共包含 6 个方面。通常, 对于多标记数据, 需要统计其标记密度 LD(Label Density)与标记集合的势 LC(Label Cardinality)。这两个指标表示如下:

$$LD = \frac{1}{n} \sum_{i=1}^n \frac{|Y_i|}{|L|} \quad (4)$$

$$LC = \frac{1}{n} \sum_{i=1}^n |Y_i| \quad (5)$$

其中, Y_i 表示文档 D_i 的方面标签集合, 即 L 的子集。

对多标记实验数据进行统计, 给出了该语料中方面数量与文档数量的对应关系及 LD 统计量和该语料中各方面对应的文档数量及 LC 统计量, 见表 1 所示。

表 1 实验语料中方面及其统计量

方面数量与文档数量关系以及 LD 统计量			各方面对应的文档数量及 LC 统计量		
方面数量	文档数量	LD	方面	文档数量	LC
1	1142	0.27	舒适性	694	1.51
2	587		动力性	558	
3	207		操控性	478	
4	52		服务性	681	
5	11		经济性	616	
6	1		安全性	179	

表 1 中 LD 和 LC 两个对于标记数量的统计信息, 表明该汽车评论文本数据集达到了多标记数据集的标准, 因此, 本文的实验可以在该数据集上完成。

4.2 评价指标

本文采用 8 种多标记评价指标^[9]评估多性能识别的准确率。这 8 种指标可以分为 4 种基于样本的评价和 4 种基于排序的评价两类。由于篇幅限制, 文本仅给出各种指标的简单介绍。

基于样本的评价指标:

Hamming Loss(HL): 计算每个样本被错误标记的平均次数, 将与样本相关的标签未被标记和与样本不相关的标签被标记两种情况计算在内。

Subset Accuracy(SA): 计算预测的标记集合与真实的标记集合完全匹配的次数占样本总数的比例。

F1 值: 表示每个样本的 precision 与 recall 的调和平均数。

Accuracy(A) 根据预测标记集合与真实标记集合之间的 Jaccard 系数定义的。

基于排序的评价指标:

Average Precision(AP): 计算预测出的排序位于某个真实相关标记 L_i 之前的标记确实属于该样本的比率。

Coverage(C): 计算平均意义下, 找到一个样本的所有真实标记, 需要遍历的预测标记有序列中的标记个数。

Ranking Loss(RL) 计算没有被正确排序的标记对的平均比率。

One Error(OE) 计算在预测标记有序列中, 排名第一的标记与实际对象不相关的次数。

以上 8 种指标中, HL、C、RL 和 OE 的值越小, 表示分类性能越好; SA、F1、A 和 AP 的值越大, 表示分类性能越好。

4.3 实验结果与分析

4.3.1 实验设置

在文本分类中, 支持向量机(SVM)常被用作分类器, 表现了较好的性能。因此, 对于两

种问题转换方法 BR 和 CC，本文采用 SVM 作为基础分类器，该分类器是在 LibSVM¹上实现的，并使用线性核函数，完成数据转化后的二类分类；对于 MLkNN 方法，参数 k 设置为 10，平滑参数选用默认值²；RankSVM 方法，采用线性核函数，其他均采用默认参数²。

本次实验中，我们对 2000 个实验数据随机采样，采取五次交叉的方式完成实验，实验结果以“平均值±标准差”的形式表示。

4.3.2 χ^2 值的不同聚合策略对分类的影响

为分析不同的聚合策略对文本中的评论文本多性能识别的影响，首先将候选特征分别按 CHI-BR-mean 和 CHI-BR-max 特征选择方法得到的 χ^2 值以降序排列，以此选取前 20% 的特征，作为本次实验的特征集合。以 BR 分类方法为例，分类结果如图 3 所示。

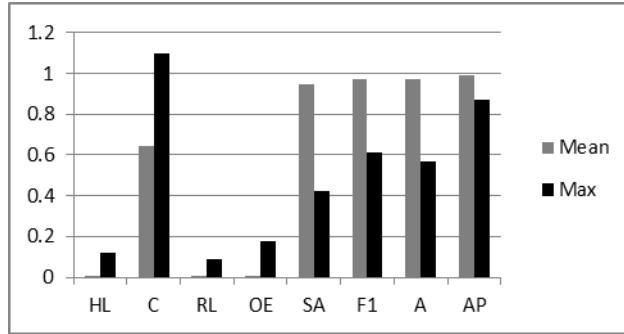


图 3 BR 分类方法中， χ^2 值的不同聚合策略选取特征对分类的影响

图 3 中，前四项 HL、C、RL、OE 的值越小越好，后四项 SA、F1、A、AP 的值越大越好。可以看出，在 CHI-BR 特征选择方法中，采用平均聚合策略的实验结果在各项指标上都明显优于最大聚合策略。因此，在之后的方法对比试验中，采用平均聚合策略对特征进行降维。

4.3.3 不同分类方法的分类结果比较

对于相同的数据，不同的分类方法往往会影响分类的性能。本文选用平均聚合策略对候选特征进行降维，选取 χ^2 值降序排序后的前 20% 的特征作为特征项，采用向量空间模型对其进行表示，使用 4 种方法进行分类用于分类模型的训练与测试。表 2 给出了 4 种分类器在 8 种多标记评价指标上的分类结果，最优结果用加粗方式表示。

表 2 4 种分类器在 8 种多标记评价指标上的分类结果

	BR	CC	MLkNN	RankSVM
HL	0.0085±0.0019	0.0084±0.0020	0.1436±0.0109	0.0603±0.0128
SA	0.9500±0.0121	0.9505±0.0124	0.4515±0.0299	0.7155±0.0635
F1	0.9749±0.0033	0.9750±0.0034	0.6643±0.0317	0.8772±0.0378
A	0.9699±0.0052	0.9701±0.0053	0.6090±0.0311	0.8396±0.0441
AP	0.9929±0.0016	0.9929±0.0016	0.8517±0.0159	0.9716±0.0150
C	0.6445±0.0683	0.6450±0.0680	1.3440±0.0475	0.7455±0.0405
RL	0.0066±0.0012	0.0066±0.0011	0.1306±0.0113	0.0227±0.0079
OE	0.0075±0.0031	0.0075±0.0031	0.1836±0.0284	0.0345±0.0287

由表 2 可以看出：

(1) 四种方法中，BR 和 CC 两种方法对于评论文本的多性能识别效果最好，两者基本相等。CC 方法是 BR 的改进，考虑了方面标签与评论文档之间的依赖关系，相似的结果表明该数据集在标记上与样本之间的依赖关系较弱。这两种方法都是在 LibSVM 的基础上实现的。该基础分类器最初是针对二值分类器设计的，因此在多标记数据转换后的正反例分类

¹ <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>

² http://cse.seu.edu.cn/people/zhangml/Resources.htm#codes_mll

中，效果较为突出。

(2) RankSVM 方法的结果相比于 BR 和 CC 方法较低。这种方法考虑的是两两标记之间的相关性，该结果表明对于本文所用评论文本的多方面性能而言，其方面间的相关性较弱。

(3) MLkNN 方法的分类效果最差，主要原因是其基础分类器是 kNN，且该方法中使用欧氏距离作为度量相似度，不适合特征稀疏的文本数据的分类。

4.3.4 每种标记上的方面识别结果

文献[21]建立了面向汽车领域观点挖掘的本体库，本文将文档的属性特征对应于本体库中的特征词，从而检索该词所对应的性能方面的方法称为基于领域本体的特征匹配 (Ontology-Match)。本次实验将该方法与 4.3.3 节中四种多标记识别方法进行对比，比较评论文本中识别单个方面的平均正确率，结果见图 4。

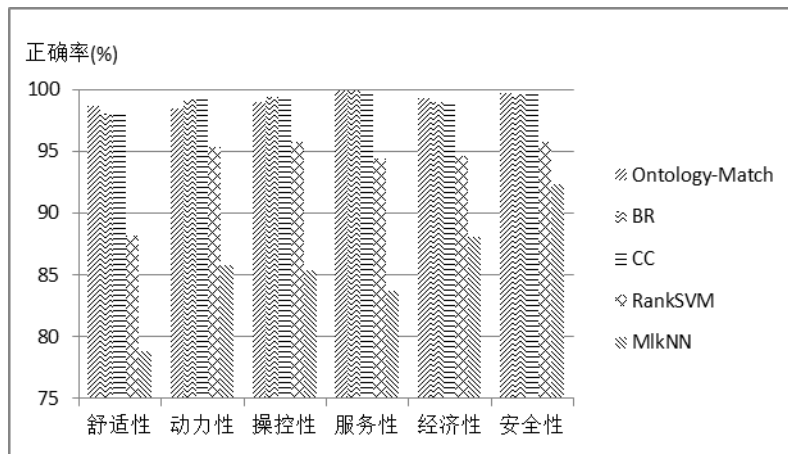


图 4 各分类方法中，各性能方面识别的平均正确率

由图 4 可知，在 BR 和 CC 分类方法下，各性能方面的识别正确率达到了很高的标准，服务性方面识别正确率高达 99.85%，其他方面的识别正确率均在 95% 以上。与基于领域本体的特征匹配结果相比，在动力性与操控性方面，多标记分类方法的识别准确率明显较高，在服务性方面两者相等，在其他方面与多标记分类方法的识别结果基本持平。由此可以表明，以 BR 和 CC 方法分类的汽车评论文本多性能识别，可为下一步基于方面的观点挖掘提供准确性的保证。

5 结论与展望

文本针对汽车评论文本中存在的多方面性能特点，以文本挖掘方法为基础，将多标记学习方法用于对汽车评论文本的多性能识别。在 8 种多标记评价指标上，给出了四种类型的多标记分类方法的实验结果，其中，基于 LibSVM 的 BR 和 CC 方法取得了高达 95% 的子集准确率。

此次工作仅完成了对汽车评论文本的多性能识别，下一步，我们将考虑利用倾向性分析，完成对多方面的观点打分，针对多方面的评论文本观点挖掘问题进行研究。

参考文献

- [1] Dave K, Lawrence S, Pennock D. Minging the peanut gallery: opinion extraction and semantic classification of product reviews[C] // Proceedings of the 12th international conference on World Wide Web (WWW). New York: ACM, 2003, 519-529.
- [2] Qiu G, Liu B, Bu J, et al. Expanding domain sentiment lexicon through double propagation[C] // Proceedings of the 21st International Joint Conference on Artificial Intelligence (IJCAI). 2009, 1199-1204.

- [3] Huang S, Peng W, Li J, et al. Sentiment and topic analysis on social media: A multi-task multi-label classification approach[C] // Proceedings of the 5th Annual ACM Web Science Conference (WebSci). New York: ACM, 2013, 172-181.
- [4] Shuhua Monica L, Jiun-Hung C. A multi-label classification based approach for sentiment classification[J]. Expert Systems with Applications, 2015, 42(3):1083-1093.
- [5] Katakis I, Tsoumakas G, Vlahavas I. Multi-label text classification for automated tag suggestion[C] // Proceedings of the ECML/PKDD 2008 Discovery Challenge. Belgium, 2008.
- [6] Hu M, Liu B. Mining and summarizing customer reviews[C] // Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining(KDD). New York: ACM, 2004, 168-177.
- [7] Lu Y, Zhai C X, Sundaresan N. Rated aspect summarization of short comments[C]. // Proceedings of the 18th International Conference on World Wide Web (WWW). New York: ACM, 2009: 131-140.
- [8] Wang H, Lu Y, Zhai C. Latent aspect rating analysis on review text data: a rating regression approach[C] // Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD). York: ACM, 2010, 783-792.
- [9] 王素格, 尹学倩, 李茹, 等. 基于非完备信息系统的评价对象情感聚类[J]. 中文信息学报, 26(4): 98-102.
- [10] Zhang M L, Zhou Z H. A review on multi-label learning algorithms[J]. IEEE Transactions on Knowledge and Data Engineering, 2014, 26(8): 1819-1837.
- [11] Shen X, Boutell M, Luo J, et al. Mutil-abel machine learning and its application to semantic scene classification[C] // Proceedings of the 2004 International Symposium on Electronic Imaging. San Jose, California, USA, 2004, 18-22.
- [12] Read J, Pfahringer B, Holmes G et al. Classifier chains for multi-label classification[J]. Machine Learning, 2011, 85(3): 333-359.
- [13] Tsoumakas G, Vlahavas I. Random k-labelsets: An ensemble method for multi-label classification[C] // Proceedings of the 18th European Conference on Machine Learning (ECML). Berlin, Heidelberg, 2007, 406-417.
- [14] Zhang M L, Zhou Z H. ML-*k*NN: A lazy learning approach to multi-label learning[J]. Pattern Recognition, 2007, 40(7): 2038-2048.
- [15] Schapire R, Singer Y. BoosTexter: A boosting-based system for text categorization[J]. Machine Learning, 2000, 39 (2):135-168.
- [16] Elisseeff A, Weston J. A kernel method for multi-labelled classification[C]. Advances in Neural Information Processing Systems 14. Cambridge, MA: MIT Press, 2002: 681-687.
- [17] Zhang M L, Zhou Z H. Multi-label neural net-works with applications to functional genomics and text categorization[J]. IEEE Transactions on Knowledge and Data Engineering, 2006, 18: 1338-1351.
- [18] Newton Spolaôr, Everton Alvares Cherman, et al. A comparison of multi-label feature selection methods using the problem transformation approach[J]. Electronic Notes in Theoretical Computer Science, 2013, 292(5):135-151.
- [19] Newton Spolaôr, Grigorios Tsoumakas. Evaluating feature selection methods for multi-label text classification[C] // Proceedings of the first workshop on bio-medical semantic indexing and question answering. Valencia, Spain, 2013.
- [20] Zhang M L, JM Peña, V Robles. Feature selection for multi-label Naive Bayes classification[J]. Information Sciences. 2009, 179(19): 3218-3229.
- [21] 冯淑芳, 王素格. 面向观点挖掘的汽车本体知识库的构建[J]. 计算机应用与软件, 2011, 28(5): 45-47.
- [22] Plaban Kumar Bhowmick, Anupam Basu, et al. Reader perspective emotion analysis in text through ensemble based multi-label classification framework[J]. Computer and Information Science. 2009, 2(4): 64-74.

作者简介: 张晶 (1990-), 女, 硕士生, 主要研究方向为数据挖掘; 李德玉(1965-), 男, 教授, 博士生导师, 主要研究方向为人工智能等, E-mail: lidy@sxu.edu.cn(通信作者); 王素格(1964-), 女, 教授, 博士生导师, 主要研究方向为自然语言处理、智能检索等。