

# 改进的 $\chi^2$ 统计文本特征选择方法

肖 婷,唐 雁

XIAO Ting,TANG Yan

西南大学 计算机与信息科学学院,重庆 400715

School of Computer & Information Science, Southwest University, Chongqing 400715, China

E-mail:xiaott2006@163.com

**XIAO Ting,TANG Yan.** Improved  $\chi^2$  statistics method for text feature selection. *Computer Engineering and Applications*, 2009, 45(14):136–137.

**Abstract:** Feature selection is a hot topic in current search field, especially in the field of text categorization. In this paper,  $\chi^2$  statistical method has two defects. One is reducing the weight of the low-frequency words. The other is increasing the weight of the characteristics in the designated class. The characteristics little appear in designated class but other classes. Through simulation and comparison experiment, the result is better than before.

**Key words:** text categorization; feature selection;  $\chi^2$  statistics

**摘要:** 特征选择是当今研究领域的一个热点,尤其是文本分类领域中的热点。针对 $\chi^2$ 统计方法的两个缺陷:降低了低频词的权重和提高了很少在指定类中出现但普遍存在于其他类的特征在该类中的权重,对 $\chi^2$ 统计方法进行改进,并通过做模拟和对比实验,对比改进前后的方法对文本分类的影响。在模拟和对比实验中,改进后方法的分类效果要好于传统的方法。

**关键词:** 文本分类;特征选择; $\chi^2$ 统计

**DOI:** 10.3777/j.issn.1002-8331.2009.14.041   **文章编号:** 1002-8331(2009)14-0136-02   **文献标识码:** A   **中图分类号:** TP39

## 1 引言

文本挖掘是近几年来数据挖掘领域的一个研究热点<sup>[1]</sup>,文本分类作为文本挖掘的重要手段和工具之一,其目标是对未知类别的文档进行自动处理、判断它们所属的类别。特征选择作为文本分类的前提,其重要性和意义就可想而知。一个合理、有效的特征选择方法可以在数据预处理阶段去掉数据中的冗余,降低特征空间的维数,提高分类的效率。

文本分类中,特征选择的基本思想通常是构造一个评价函数,对特征集的每个特征进行分别评估,后对所有的特征按照其评估分的大小进行排序,选取预定数目的最佳特征作为结果的特征子集<sup>[2]</sup>。

常用的特征选择方法有:文档频率(DF)、信息增益(IG)、互信息(MI)、 $\chi^2$ 统计方法(CHI)、期望交叉熵(Expected Cross Entropy)、文本证据权(Weight of Evidence text)等。但是在众多的特征选择方法中,Yang 指出 IG(信息增益)和 CHI 的效果最好<sup>[3]</sup>,而 IG 计算量相对其它几种方法较大,因此本文主要针对目前特征选择方法中效果最好 $\chi^2$ 统计方法进行研究和改进。

## 2 $\chi^2$ 统计方法相关研究

### 2.1 $\chi^2$ 统计方法

$\chi^2$ 统计方法度量词条  $t$  和文档类别  $c$  之间的相关程度,假

设词条  $t$  和类别  $c$  之间符合具有一阶自由度的 $\chi^2$  分布。词条对于某类的 $\chi^2$  统计值越高,它与该类之间的相关性就越大,携带的类别信息也就越多。令  $N$  表示训练语料库中的文档总数,  $c$  为某一特定类别,  $t$  表示特定的词条,  $A$  表示属于  $c$  类且包含  $t$  的文档频数,  $B$  表示不属于  $c$  类但包含  $t$  的文档频数,  $C$  表示属于类别  $c$  但不包含  $t$  的文档频数,  $D$  表示不属于  $c$  也不包含  $t$  的文档频数。则  $t$  对于  $c$  的 CHI 值由下式计算:

$$\chi^2(t,c) = \frac{N(AD-BC)^2}{(A+C)(B+D)(A+B)(C+D)} \quad (1)$$

当特征  $t$  与类别  $c$  的相关性越强,  $\chi^2(t,c)$  的值就越大,代表特征  $t$  包含与特定类别  $c$  之间的独立性就越小,对分类的贡献就越大。有关的信息就越多。当特征  $t$  与类别  $c$  相互独立时,  $\chi^2(t,c)=0$ , 代表特征  $t$  不包含任何与特定类别  $c$  有关的信息。

对于多个类别的问题,分别计算  $t$  对于每个类别的 $\chi^2$  值,再选择出整个语料库中 $\chi^2$  值最大的作为整个语料库的 $\chi^2$  值。可以事先假定一个阈值,从原始的特征空间中删除那些低于特定阈值的词条,保留大于该阈值的特征词条作为文档表示的特征词。

### 2.2 存在的问题

Yang 的研究表明, $\chi^2$ 统计方法虽然是目前对于中文文本分类效果最好的一个特征选择方法之一<sup>[3]</sup>。大多数中文分类系

统都采用这种方法<sup>[4-5]</sup>。但是,根据分析,它依然存在着缺陷,主要表现在两个方面。

### (1)降低了低频词的权重。

$\chi^2$ 统计方法不能准确地保留往往是某类特有的,具有很强代表性的低频词。这里的低频词是指文档频数,但某些低频词往往是某类特有的特征词,这些特征词具有很强的代表性,虽然只出现在指定类的少量文章中,但在这少量文章中频繁的出现,因此对分类的贡献很大。但是通过公式 $\chi^2(t, c)$ 计算出来的 $\chi^2$ 统计值相对来说比较小,而根据TFIDF理论<sup>[6]</sup>,权重低的特征项将被过滤掉。于是,由于低频词的权重比较低, $\chi^2$ 统计就删除了这本该保留的具有很强代表性的特征项。

(2)提高了很少在指定类中出现但普遍存在于其他类的特征在该类中的权重。

在整个训练集语料库中一些出现频率较高的词,而这些词语在指定类中出现得很少甚至几乎是不出现的( $A$ 值很低),显然这些词是不能很好地代表该指定类,应该被过滤掉。但是 $BC >> AD$ 导致用公式 $\chi^2(t, c)$ 计算出来的 $\chi^2$ 统计值相对来说比较大,而根据TFIDF理论<sup>[6]</sup>,权重高的特征项将被保留,作为该指定类的特征项。

## 3 $\chi^2$ 统计方法的改进

特征选择方法的好坏直接影响文本分类的效果。前面分析的 $\chi^2$ 统计方法有两个方面的缺陷:第一,降低了低频词的权重;第二,提高了很少在指定类中出现但普遍存在于其他类的特征在该类中的权重。这里认为降低了低频词的权重这一点,对分类效果的影响相对较为明显。所以主要针对降低了低频词的权重这一问题,引入了文档内频度的概念,并且兼顾第二个问题引入类内正确度的概念,对 $\chi^2$ 统计方法进行改进。

### (1) 文档内频度

针对问题1,把特征项在具体的文档中出现的频度考虑进 $\chi^2$ 统计的计算公式。设训练集中类别为 $C_i$ 的文本有 $d_{i1}, d_{i2}, \dots, d_{ik}, \dots, d_{ij}$ ,特征 $t$ 在文本 $d_{ik}$ ( $1 \leq k \leq j$ )中出现的频度为 $tf_{ik}$ ,则特征 $t$ 在类别 $C_i$ 中出现的频度记为 $\alpha$ ,即文档内频度表示为: $\alpha = \sum_{k=1}^j tf_{ik}$

### (2) 类内正确度

针对问题2,由于是 $A$ 值太小,导致 $BC >> AD$ 造成的,所以引入一个调节函数 $\beta = \frac{A}{A+B}$ ( $0 \leq \beta \leq 1$ )即词条出现在指定类中的比例。当 $\beta$ 靠近0时,即在整个语料库出现频繁,而在指定类出现较少。当 $\beta$ 靠近1时,即在指定类出现较多。由于 $\beta$ 的范围是[0,1],这样的 $\beta$ 值对计算出的 $\chi^2$ 值影响有限,因此,对 $\beta$ 函数变换到 $\beta'$ ,把[0,1]映射到[M<sup>-1</sup>, M]中去,M是正整数。即当 $A > B$ 时(即词条能很好地代表指定类), $\beta' > 1$ ,参数 $\beta$ 放大计算出的 $\chi^2$ 值;当 $A < B$ 时(词条指定类出现很少), $\beta' < 1$ ,参数 $\beta$ 减小计算出的 $\chi^2$ 值。因此 $\beta'$ 记作类内正确度,表示为:

$$\beta' = M^{\beta-1} \quad (2)$$

鉴于以上分析,考虑将 $\alpha$ 和 $\beta'$ 加入原来的 $\chi^2$ 统计方法的计算公式,使得通过 $\chi^2$ 公式计算出来的词条 $t$ 和类别 $c$ 之间的相关程度更为准确,这样提取出来的特征项更能代表指定类。将改进后的 $\chi^2$ 统计方法的计算公式表示如下:

$$\chi^2(t, c) = \frac{N(AD - BC)}{(A+C)(B+D)(A+B)(C+D)} \times \alpha \times \beta' \quad (3)$$

## 4 实验与结论

特征选择是文本分类的关键技术,特征提取的好与坏直接关系到文本分类的效果,将特征选择方法作为重点研究对象,对 $\chi^2$ 统计方法提出了改进措施,用文本分类的结果来验证改进的有效性。

### 4.1 朴素贝叶斯分类算法

在文本分类算法的选择中,这里选择最为常见的分类贝叶斯算法。在朴素贝叶斯分类(简称NB)基于贝叶斯定理,可以用来预测类成员关系的可能性,给出文本属于某特定类别的概率。分类时根据预测结果该样本分到概率最高的类别中去即可。在这个方法中,有一个“独立性假设”<sup>[7]</sup>:在给定的文本类语境下,文本中每个特征词是相互独立的。

设 $d$ 为一任意文本,它属于文档类 $C = \{c_1, c_2, \dots, c_j\}$ 中的某一类 $c_k$ 。根据NB分类方法有:

$$p(c_k|d) = \frac{p(c_k)p(d|c_k)}{p(d)} \quad (4)$$

对文本 $d$ 进行分类,就是按照上述公式计算所有文本类在给定 $d$ 情况下的概率,概率值最大的那个类就是文本 $d$ 所属的类。其中的 $p(c_k) = \frac{N_k}{N}$ , $p(d|c_k) = \frac{1+N_{ik}}{M + \sum_{k=1}^j N_{kj}}$ , $N_{ik}$ 表示特征 $i$ 在训练类别 $c_k$ 在文本出现的次数, $N_k$ 表示类别 $k$ 包含的训练文本数, $i$ 表示特征项数, $d_i$ 表示特征 $i$ 。

### 4.2 实验及结果

从来源于搜狐新闻网站的大量经过编辑手工整理与分类的新闻语料中抽取部分短文共计1500多篇,采用交叉验证的方法对5个类的文档进行实验。本实验公式(2)中的 $M$ 值取4。

一般用准确率和召回率来衡量信息挖掘系统的结果。因此,最终的实验结果采用通用的召回率(recall)和准确率(precision)两个指标作为评测标准:

$$\text{准确率(precision)} = \frac{\text{分类正确的文本数}}{\text{实际分类的文本数}}$$

$$\text{召回率(recall)} = \frac{\text{分类正确的文本数}}{\text{应有的文本数}}$$

下面3个表是特征向量空间维数为500、1000、1200维的实验数据。

通过对比实验,从表1~表3可以看出,改进后的召回率和准确率普遍高于改进前的。这就证明,本文对于 $\chi^2$ 统计方法的改进是有较好的效果的。

表1 特征向量空间500维

TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area	Class
0.767/0.737	0.190/0.136	0.335/0.403	0.767/0.737	0.466/0.521	0.874/0.894	computer
0.733/0.76	0.028/0.023	0.764/0.806	0.733/0.76	0.748/0.782	0.940/0.959	health
0.830/0.880	0.009/0.009	0.922/0.926	0.830/0.880	0.874/0.903	0.977/0.989	sports
0.847/0.900	0.018/0.008	0.858/0.931	0.847/0.900	0.852/0.915	0.975/0.992	train
0.563/0.687	0.014/0.025	0.833/0.777	0.563/0.687	0.672/0.729	0.899/0.939	travel

表2 特征向量空间1000维

TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area	Class
0.707/0.750	0.100/0.081	0.468/0.537	0.707/0.750	0.563/0.626	0.902/0.925	computer
0.747/0.830	0.026/0.020	0.783/0.836	0.747/0.830	0.765/0.833	0.950/0.962	health
0.863/0.920	0.012/0.012	0.902/0.998	0.863/0.920	0.882/0.914	0.985/0.994	sports
0.880/0.907	0.016/0.006	0.874/0.948	0.880/0.907	0.877/0.927	0.988/0.992	train
0.723/0.757	0.024/0.015	0.792/0.860	0.723/0.757	0.756/0.805	0.935/0.955	travel

(下转140页)