

基于潜在语义差异的医学网页聚类

米晓芳, 秦洋, 王立宏, 宋宜斌

(烟台大学计算机学院, 烟台 264005)

摘要: 采用潜在语义索引的全局模型和局部模型表示医学网页时, 模糊聚类结果的类间包含度很大。该文提出一种新的潜在语义差异模型, 将医学网页中的文本抽取出来并分别采用全局模型、局部模型和差异模型进行表示, 利用 FCM 算法进行聚类并计算类间包含度。实验发现, 对给定的 5 类医学网页进行聚类时, 采用差异模型时的类间包含度平均约为全局模型的 85%、局部模型的 80%。

关键词: 潜在语义索引; 差异模型; 文本挖掘; FCM 聚类; 包含度

Medical Webpage Clustering Based on Latent Semantic Difference

MI Xiao-fang, QIN Yang, WANG Li-hong, SONG Yi-bin

(College of Computer, Yantai University, Yantai 264005)

【Abstract】 Fuzzy clustering, two categories of medical Web pages represented by global LSI or local LSI generate two fuzzy sets with a large inclusion degree. A new latent semantic difference model is proposed. The text in medical Webpage is extracted and represented by global LSI, local LSI and difference LSI respectively. FCM algorithm is employed to cluster the feature vectors and inclusion degree between two result fuzzy sets is calculated. Experiments on five given categories of medical Webpages show that, on the average, difference LSI reduces the inclusion degree by a factor of 85% and 80% respectively when compared with global LSI and local LSI.

【Key words】 latent semantic index; difference model; text mining; FCM clustering; inclusion degree

1 概述

随着 Internet 的发展, 互联网出现了海量的、异构的 Web 信息资源, 其中, Web 文本占了主导地位。如何从这些海量的 Web 文本信息中获得有价值的信息, 成为信息处理领域的一个关键问题。人们将数据挖掘技术应用到 Web 的知识发现中, 形成了 Web 挖掘技术。

循证医学网页记录着丰富的循证医学实验信息, 涉及疾病防治、保健、药物等各方面, 典型的网上数据库有: Cochrane, Evidence-Based Medicine, Medline 等, 用户可以免费下载各种报告文章的摘要信息。文献[1]以 Medline 中的摘要信息为研究对象, 抽取出病人群体、对比治疗、疗效评价等信息。本文研究这些医学网页的聚类问题, 主要考查样本网页聚类中涉及的文本表示与降维问题。

潜在语义索引挖掘文本与特征之间潜在的高阶语义结构, 将分解文本特征矩阵, 实现特征空间的降维, 文本和特征被转换到低阶语义空间上进行描述。目前主要的潜在语义模型有全局模型及其改进模型、局部模型及其改进模型^[2]。本文在现有的全局模型和局部模型基础上提出一种潜在语义差异模型。通过对 Medline 数据库下载的医学网页数据进行实验, 结果表明: 利用该模型表示文本向量能有效地改进类间包含度。

2 潜在语义索引差异模型

潜在语义索引(Latent Semantic Indexing, LSI)^[3-4]的基本思想为: 文本中的词与词之间存在某种联系, 即存在某种潜在的语义结构, 可以采用统计的方法来寻找该语义结构, 并且用语义结构来表示词和文本。这样可以消除词之间的相关性, 化简文本向量。

潜在语义分析是基于矩阵的奇异值分解(Single Value

Decomposition, SVD)技术的。给定一组文档, 假定包含 n 篇文档, 其中有 m 个不同的词项, 采用词项作为特征。该文档集可以表示为词项×文档矩阵: $X = [x_{ij}]_{m \times n} \in R^{m \times n}$ 。其中, 矩阵元素 x_{ij} 表示词项 t_i 在文档 d_j 中的权值。矩阵 X 的奇异值分解可以表示为

$$X = UDV^T \quad (1)$$

其中 U 和 V 分别为左奇异矩阵和右奇异矩阵, U 是一个 $m \times m$ 的列标准化正交矩阵; D 是一个 $m \times n$ 的对角矩阵, 其对角线上的元素是按降序排列的非负奇异值; V 是一个 $n \times n$ 的正交矩阵。

选择适当 k 值, 将 D 中删除相应的行和列得到 D_k , 删除 U, V 相应的行和列分别得到 U_k, V_k , 得到新的矩阵 $X_k = U_k D_k V_k^T$ 。这个新矩阵是对原矩阵的逼近, 可以用它去近似原始矩阵。LSI 空间是由 U_k 的前 k 维列向量张成的空间, 因此, 矩阵 X 中的文档可以投影到 k 维 LSI 空间, 得到低维表示。LSI 不但可以对文本表示进行有效的降维, 还可以捕捉到文本的潜在语义信息^[2]。

目前主要的潜在语义模型有全局模型、局部模型两种^[2]。假设有两个类, 它们的特征向量分别用 A 和 B 表示(注: 这里的向量是特征词项的集合, 以下相同)。对于全局的潜在语义模型来说, 它的特征向量是取 A 和 B 的并集, 即 $A \cup B$ 组成全局的特征向量, 然后用得到的特征向量对两类的文本集合求词项×文档矩阵。对于局部的潜在语义模型来说, 它的特征向

基金项目: 国家自然科学基金资助项目(60473115); 山东省自然科学基金资助项目(Y2006G22)

作者简介: 米晓芳(1982-), 女, 硕士研究生, 主研方向: Web 信息挖掘; 秦洋, 硕士研究生; 王立宏、宋宜斌, 教授

收稿日期: 2007-11-23 **E-mail:** xiaofangmi@yahoo.com.cn

量是取A或者B的特征向量来表示两类的文本集合的矩阵。全局模型将两类的特征混合起来，综合了两类的特征，但不利于表示类之间的差异；而局部模型用一个类的特征表示两个类，有些重要的特征词项就会无法表示出来，从而影响到两类之间的区分。

本文提出一种新的潜在语义索引差异模型，采用两类的特征向量A和B的对称差来表示新的特征向量C，即

$$C = A \oplus B = (A - B) \cup (B - A) \quad (2)$$

差异模型保留了两类之间的不同点，既保留了A不同于B的特征，同时也保留了B不同于A的特征。可以期望不同类的样本采用该模型表示时能有较大差异，从而使类之间的边界更清晰。该模型处理多个类别时可以采用“1 对其余”的策略。后续实验在对网页进行聚类时发现：采用差异模型时类之间的包含度比采用全局模型和局部模型都有下降，从而证实了差异模型是有效的。

3 医学网页聚类

3.1 医学网页预处理

Web挖掘是在文本挖掘技术上发展起来的，由于网页本身的特征使得Web挖掘比文本挖掘更加难处理，因此本文首先将网页转化为文本文档，再利用文本挖掘技术对其进行处理。实验数据是从Medline医学数据库中按类别下载的英文医学网页，由于是从同一个数据库上下载的，因此网页的结构基本上是相同的。通过分析发现，网页中正文是在标签<D>和</D>之间的，从<D>开始“文章标题”处于第一个和标签之间；“文章作者”处于第2个和标签之间，紧接着“摘要正文”处于

和

标签之间。可以通过C语言编程读出网页中的标题、作者和摘要，并将其保存为与网页同名的文本文档。

在对文档进行特征提取前，需要对文本信息进行预处理。主要完成以下几方面工作：

(1)建立英文停止词表(stop word)，这种词几乎不携带任何信息，如“the”，“a”等连词和虚词对表示网页特征没有意义。

(2)统计每一篇文档的词项和词频，汇总词项得出每一类的特征向量(term₁, term₂, ..., term_n)。以类别为单位，计算每个特征词项term_j出现的文档个数，即文档频率。

(3)删除一些低频词，比如文档频率<5的单词，这样就得到了每个类的频繁词特征向量，这是真正使用的特征向量。

3.2 FCM算法和包含度

模糊c-均值(FCM)算法是模式识别的一个重要工具，FCM聚类算法适用于每类样本数相差不大且每类样本点距离相差不大的团状数据。FCM算法需指定数据集的类别数目c，本文只讨论两类之间的聚类问题，故取c=2。模糊聚类问题可表示为以下数学问题^[5]：

$$\min J_m(U, V) = \sum_{i=1}^c \sum_{k=1}^n u_{ik}^m d_{ik}^2 \quad \text{s.t.} \quad \sum_{i=1}^c u_{ik} = 1, \sum_{k=1}^n u_{ik} > 0, u_{ik} \geq 0$$

其中，u_{ik}是第k个样本点对第i个类的隶属度，k=1,2,...,n；U=[u_{ik}]_{c×n}是隶属度矩阵；m通常取2。

FCM算法的实现步骤如下：

(1)初始化。随机给出c个彼此不同的类中心点v_i，i=1,2,...,c。

(2)计算样本点与中心点的距离d_{ik}=||x_k-v_i||，并且求出隶属度u_{ik}。

如果∀k, ∀i, d_{ik} ≠ 0，则

$$u_{ik} = d_{ik}^{1-m} / \left[\sum_{i=1}^c d_{ik}^{1-m} \right]$$

(3)如果存在k, i，使得d_{ik}=0，则令u_{ik}=1，且对j ≠ i u_{jk}=0；计算类中心点的坐标：

$$v_i = \sum_{k=1}^n (u_{ik})^m x_k / \left[\sum_{k=1}^n (u_{ik})^m \right]$$

(4)回到步骤(2)，直到

$$\|V^{(k+1)} - V^k\| < \epsilon$$

其中，V=[v₁, v₂, ..., v_c]是类中心点的集合。

FCM算法能快速、有效地对给定的数据集进行指定类数的聚类，以类中心点的位置和样本点对类的隶属度表示聚类结果，得出c个模糊集合。

通常采用聚类有效性函数对聚类结果进行评价，包含度是一种常用的聚类有效性函数，描述了一个模糊集包含在另一个模糊集的程度。以聚类数c=2为例，设聚类结果是两个模糊集合F₁和F₂，样本x_k对F₁和F₂的隶属度分别为u_{1k}和u_{2k}。定义F₁对F₂的包含度^[5]为

$$\text{sub}(F_1, F_2) = \frac{\sum_{k=1}^n \min(u_{1k}, u_{2k})}{\sum_{k=1}^n u_{1k}} \quad (3)$$

FCM的聚类结果没有空类，所以，上式分母不为0。类似可以定义F₂对F₁的包含度，最后给出聚类结果的有效性函数为^[5]

$$\text{Valid-F}(U, c) = \max(\text{sub}(F_1, F_2), \text{sub}(F_2, F_1)) \quad (4)$$

式(4)以F₁对F₂的包含度和F₂对F₁的包含度中的较大者作为两个模糊集合F₁和F₂互相包含的程度度量。在聚类中，该函数值越小表示两个类别分开得越远，聚类效果越好。

4 实验设计与结果分析

4.1 实验方案设计

实验针对以下几个类进行：二尖瓣闭锁不全(Mitral incompetente)，二尖瓣狭窄(Mitral stenosis)，主动脉闭锁不全(Aortic incompetente)，静脉曲张(Varicose veins)和糖尿病肾病(Diabetic Nephropathy)，其中，前4个类属于心血管类疾病，第5个类是内分泌系统疾病。样本网页来自Medline数据库，以类别名称作关键词，该词出现在网页中即认定和类别相关，再进一步通过阅读网页，剔除同时属于多个类别的网页。各类别的样本情况和预处理后结果如表1所示。

表1 实验类别基本情况

类别名称	样本网页个数	频繁特征词数量
二尖瓣闭锁不全	288	942
二尖瓣狭窄	293	968
主动脉闭锁不全	242	841
静脉曲张	231	905
糖尿病肾病	314	1798

实验研究任意两个类别在全局模型、局部模型和差异模型表示下的边界情况，将两个类别的样本混在一起，利用FCM算法进行c=2的聚类，利用包含度来评价聚类结果的有效性。

实验步骤如下：

(1)将网页转换成文本文件，按3.1节进行预处理，得出每个类的频繁词向量。

(2)任选两个类别，将样本混合起来，采用全局模型表示两个类中的每个文档，得到词项×文档矩阵X，矩阵元素x_{ij}表示词项i在文档d_j中的出现次数。

(3)对矩阵X进行SVD分解，设置k=20，取U的前20维U_k，然后将文档集合映射到低维空间表示为X'=U_k^TX。

(4)对 X' 进行FCM聚类,按式(4)计算相应的包含度值^[5]。

(5)将步骤(2)中的全局模型换成局部模型和差异模型重复该实验。

步骤(3)中的SVD分解是利用Matlab中的SVD()来完成的,其余代码使用C语言完成。

4.2 实验结果及分析

表2是对全局模型、局部模型和差异模型的对比描述。局部模型是用一个类别的特征向量来表示这两个类别混合到一起后文本集的特征向量(A向量是指用A的特征向量来表示文本集合A-B)。全局模型的特征向量是A、B特征向量的并集,而差异模型的特征向量是A、B特征向量的对称差。通过表2中3种模型的包含度比较可以看出:用差异模型来表示文档时,类之间的包含度更低,两个类别的边界更清晰。通过统计发现:采用差异模型时的类间包含度平均约为全局模型时的85%,约为局部模型时的80%(按A类向量计算)。

表2 不同模型表示下的聚类包含度比较

类别/包含度		差异模型	全局模型	局部模型	
A类	B类			A向量	B向量
主动脉闭锁不全	糖尿病肾病	0.575 407	0.661 739	0.758 599	0.755 267
主动脉闭锁不全	二尖瓣闭锁不全	0.980 027	0.996 279	0.990 824	0.989 191
主动脉闭锁不全	二尖瓣狭窄	0.634 628	0.895 548	0.894 961	0.890 452
主动脉闭锁不全	静脉曲张	0.647 109	0.743 331	0.785 942	0.839 902
糖尿病肾病	二尖瓣闭锁不全	0.591 232	0.659 152	0.749 092	0.739062
糖尿病肾病	二尖瓣狭窄	0.625 148	0.704 036	0.789 944	0.749 229
糖尿病肾病	静脉曲张	0.630 191	0.758 795	0.802 587	0.908 257
二尖瓣闭锁不全	二尖瓣狭窄	0.610 342	0.861 538	0.908 661	0.855 062
二尖瓣闭锁不全	静脉曲张	0.619 515	0.742 879	0.784 234	0.844 447
二尖瓣狭窄	静脉曲张	0.778 981	0.820 774	0.847 693	0.883 508

包含度越大说明两个类的包含程度越高,边界就越模糊。从二尖瓣闭锁不全和主动脉闭锁不全两类的包含度看出,这两个类是很相近的。实际上通过阅读网页发现,这两个类的网页讨论的内容基本上都涉及肺动脉回流、瓣膜性回流、肥

原型心肌、心室间隔缺损、左主动脉、急性风湿等内容,具有很大的相似性。糖尿病肾病与其他类别的包含度都很低,事实上这个类与其他类别的疾病分别属于两个不同种类的疾病,糖尿病肾病是属于内分泌系统方面的疾病,而其余4类是属于心血管方面的疾病,它们本身的相似性就很低,实验也验证了这一点。

5 结束语

随着Internet和Web挖掘技术的快速发展,数据挖掘技术不再局限于结构化数据的挖掘。如何从Web文本中发现隐含的、有意义的信息,已成为数据挖掘的重要目标。在本文实验中,每一类的特征向量都是由文档频率>4的特征词组成的。实验中将两个类的数据混合在一起进行聚类,并计算类间的包含度,从而比较不同的潜在语义模型得到的向量对聚类的影响。实验结果表明,用潜在语义差异模型得到的特征向量使类别之间的边界更清晰,因此,可以认为用潜在语义差异模型来表示文本的特征效果更好。下一步将此模型用于Web文本分类,并与全局模型和局部模型的改进模型进行比较。

参考文献

- [1] Kazuo H, Matsumoto Y. Information Extraction from Medline Abstracts of Clinical Trials[J]. Technical Report of IEICE AI, 2004, 42(12): 45-49.
- [2] 孙建涛. Web挖掘中的降维和分类方法研究[D]. 北京: 清华大学, 2005.
- [3] Berry M W, Dumais S T, O'brien G W. Using Linear Algebra for Intelligent Information Retrieval[J]. SIAM Review, 1995, 37(4): 573-595.
- [4] Deerwester S, Dumais S T, Furnas G W, et al. Indexing by Latent Semantic Analysis[J]. Journal of the American Society for Information Science, 1990, 41(6): 391-407.
- [5] 范九伦, 吴成茂. 用于聚类有效性判定的包含度公式[J]. 模糊系统与数学, 2002, 16(1): 80-86.

(上接第63页)

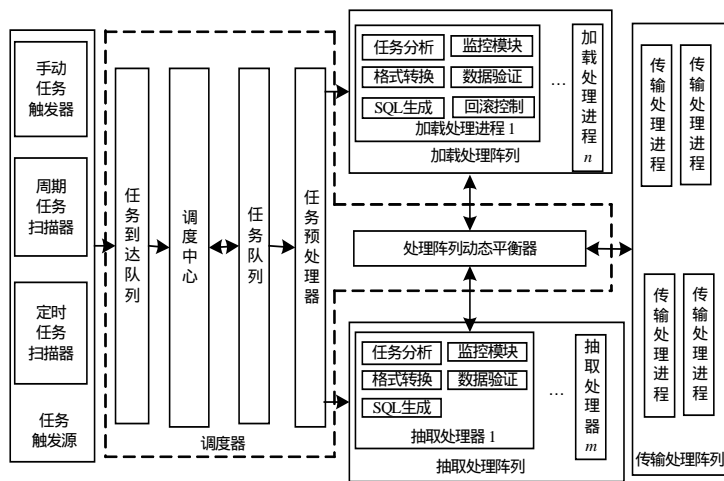


图5 交换引擎结构

关系数据库间的信息交换,以后可以考虑:(1)其他格式化文档与数据库间的交换;(2)在双方数据库都符合信息结构体系的情况下的自动数据交换。

参考文献

- [1] 中共中央组织部. 全国组织干部人事管理信息系统信息结构体系[Z]. 2006.
- [2] Fuxman A, Phokion G. Peer Data Exchange[C]//Proc. of TODS'06. [S. l.]: ACM Press, 2006: 1454-1498.
- [3] Lucian P. Translation Web Data[C]//Proc. of VLDB'02. Hong Kong, China: [s. n.], 2002: 598-609.
- [4] Arenas M, Libkin L. XML Data Exchange: Consistency and Query Answering[C]//Proceedings of the 24th ACM Sigmod-sigact-sigart Symposium on Principles of Database Systems. New York, USA: ACM Press, 2005.
- [5] Audsley N C. Deadline Monotonic Scheduling[D]. Charlottesville, Canada: Department of Computer Science, University of Virginia, 1990.

5 结束语

以修订的信息结构体系^[1]为标准,该平台方便了组织、干部系统异构数据库间的数据交换。目前平台主要考虑的是