



典大小相同的一个数组。

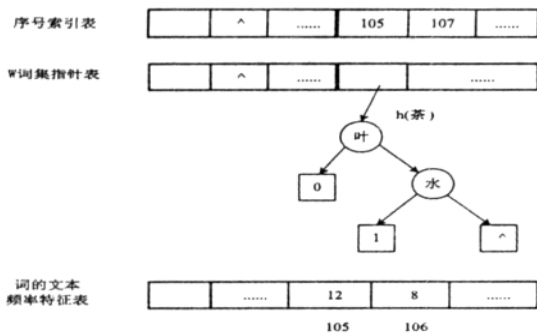


图2 词典的数据结构示意图

(5) 特征项选取:代表该特征项的词的文本频率大于等于1。

(6) 特征项权值计算公式:目前使用较多的是  $tf-idf$  方法,其计算公式为:

$$w_{ik} = \frac{f_{ik} \times \log(N / df_k)}{\sqrt{\sum_{T_k \in D_i} [f_{ik} \times \log(N / df_k)]^2}}$$

其中,  $f_{ik}$  表示特征项  $T_k$  在文本  $D_i$  中出现的次数,即特征项  $T_k$  相对于文本  $D_i$  的文本频率,  $f_{ik}$  越高,意味着特征项  $T_k$  对于文本  $D_i$  越重要。 $df_k$  表示含有特征项  $T_k$  的文本数量,  $df_k$  越高,意味着特征项  $T_k$  在衡量文本相似性方面的作用越低。 $N=||D||$ ,即全部文本的数量,分母归一化因子。 $idf_k = \log(N / df_k)$  称为特征项  $T_k$  逆向文本频率,  $idf_k$  越高,意味着特征项  $T_k$  对于文本的区别作用越大。

(7) 特征项权值表:用于记录向量空间中每个特征项的权值,规模与特征向量空间的大小相同。当某一文本被向量化后,即构成一张赋了值的特征向量权值表。

(8) 特征项文本频率特征表:为一规模与向量空间大小相同的数组,数组元素记录了该特征的文本频率。

#### 4. 实验分析

本实验所用的语料库来自网站“中文自然语言处理开放平台”,由复旦大学李荣陆博士提供,本文选取了其中的2690篇文章,采用台湾大学林智仁(Chih-Jen Lin)博士等开发设计的 LIB-

(上接第78页)

本文方法不仅就有较高的识别率,并且具有较高的识别准确率。图4展示了COLD1\_13\_C08.b1\_A029序列的特征识别结果,图中序列前端识别出来的特征是5TNS,后端是3TNS,两者都得到了精确的识别。识别出来的端非常好的符合了EST序列的期望结构,从而验证了特征识别的准确性。

#### 3. 结论与展望

通过对30多万条序列的实验表明,本文提出的方法能够较好的适应序列的特征识别,具有较高的识别率和准确率。方法基于动态规划算法,能够很好的得到全局最优特征,而不会陷入局部最优解;同时采用扩展识别方法来弥补动态规划算法的速度较慢,并且能够很好的适应不定长度的特征识别。

本文采用动态规划方法进行最短单元识别,能够避免陷入局部最优解,并且采用扩展方法缓解动态规划算法的速度慢缺陷,但如果能够采用其他的智能识别方法[9][10],则能够进一步的加速方法的识别速度,这也是本文方法的进一步研究任务。

#### 参考文献:

- 1.Chun Liang, Yuansheng Liu, Lin Liu, Adam C. Davis, Yingjia Shen, Qingshun Quinn Li. ESTs with cDNA termini: Previously Overlooked Resources for Gene Annotation and Transcriptome Exploration in Chlamydomonas reinhardtii[J]. Genetics, 2008, 179:83-93.
- 2.Chun Liang, Gang Wang, Lin Liu Guoli Ji, Lin Fang, Yuansheng Liu, Kikia Carter, Jason S Webb and Jeffrey FD Dean. ConiferEST: an integrat-

SVM工具,选择分类问题(C-SVC;选择惩罚参数为C的Support Vector Classification),核函数类型设置为RBF核。为了获取最优的训练参数,本文使用LibSVM提供的交叉验证工具对训练集进行了10重交叉验证,得到的结果为,C=5,g=0.05。在此情况下,对SVM进行模拟实验获得的平均实验数据如表1所示:

类别	计算机	体育	农业	政治	文学	经济
查准确率	0.832	0.825	0.763	0.698	0.738	0.786
查全率	0.855	0.815	0.687	0.754	0.672	0.815
F1值	0.842	0.835	0.782	0.712	0.679	0.821

表1 在SVM下的分类效果

根据以上实验结果进行分析:

(1)SVM有着良好的分类性能、泛化能力,同时,SVM进行分类决策的时间相对较快。这些优点使得SVM技术应用在信息过滤系统中,有着一定的优势。

(2)由于没有一个标准化的向量空间词典,虽然分类结果在一定程度上有所提高,但还不是完全精确。

#### 5. 小结

本章重点讨论了普遍认为性能较好的SVM分类方法中的词典结构问题,SVM在特征独立性的基础上进行了模拟实验,实验证明,基于SVM的分类算法精确度高,分类速度快,但由于没有一个标准化的向量空间词典在文本分类处理时也有错分的词。

#### 参考文献:

- 1.李荣陆.文本分类及其相关技术研究[D].上海:复旦大学,2005.
- 2.林鸿飞,李业丽,姚天顺.中文文本过滤的信息分流机制[J].计算机研究与发展,2000,37(4):470-476.
- 3.林鸿飞.中文文本过滤的逻辑模型.东北大学博士论文,2000,6-7
- 4.林鸿飞,战学钢,姚天顺.基于概念扩充的中文文本过滤模型[J].计算机科学,2000,27(2):88-90.
- 5.邓乃扬,田英杰.数据挖掘中的新方法——支持向量机.北京:科学出版社,2004.
- 6.Nello Cristianini, John Shawe-Taylor.支持向量机导论.李国正等译.北京:电子工业出版社,2004,3
- 7.黄董菁,吴立德.独立于语种的文本分类方法[D].中文信息学报,2000,14(6):1-7.

- ed bioinformatics system for data reprocessing and minig of conifer expressed sequence tags(ESTs)[J]. BMC Genomics 2007, 8:134.83-93
- 3.BLAST <http://blast.ncbi.nlm.nih.gov/Blast.cgi>[Z].
- 4.Andress Wilm,Desmond G. Higgins and Cedric Notredame R-Coffee: a method for multiple alignment of non-coding [J] RNA Nucleic Acids Research, 2008, Vol. 36, No. 9
- 5.James Robert White, Michael Roberts, James A. Yorke and Mihai Pop. Figaro: a novel statistical method for vector sequence removal [J]. BMC Bioinformatics 2008 Vol. 24 no. 4.
- 6.Arthur L. Delcher, Adam Phillippy, Jane Carlton and Steven L. Salzberg Fast algorithms for large-scale genome alignment and comparison[J]. Nucleic Acids Research 2002 Vol. 30 No. 11.
- 7.Liang C, Sun F, Wang H, Qu J, Freeman RM Jr, Pratt LH, Cordonnier-Pratt MM: MAGIC-SPP: a database-driven DNA sequence processing package with associated management tools[J]. BMC Bioinformatics 2006, 7: 115.
- 8.Liang C, Wang G, Liu L, Ji G, Liu Y, Chen J, Webb JS, Reese G, Dean JF. WebTraceMiner: a web service for processing and mining EST sequence trace files[J]. Nucleic Acids Res (2007) 35:W137-W142
- 9.M Dorigo, V Maniezzo and A Colomi. The Ant System: Optimization by a colony of cooperation agnets [J]. IEEE Transactions on Systems, 1996, 26(1):1-13.
- 10.Viharos Z J, etal. Traning and Application of Artificial Neural Networks with Incomplete Data[J]. In:LNAI 2358,2002.649-659.