

# 肿瘤诊断中的应用

健

Department of Computer Science and Engineering, USA. \* 联系人, E-mail: @math.pku.edu.cn

中的一个重要研究领域, 其中最主要的问题... 出了秩和基因选取... 并利用支持向量... 实验表明这种方法... 模型可使得在结肠

选取

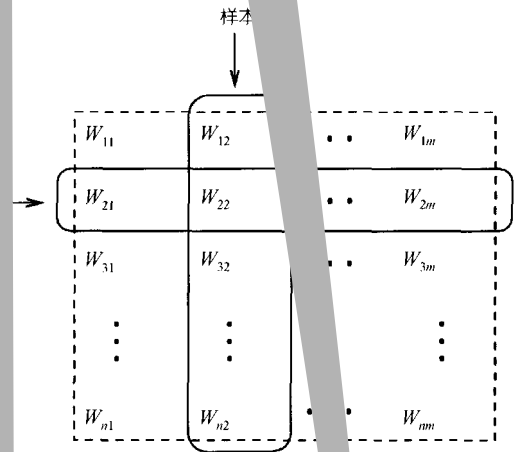


图 1 基因表达谱矩阵

很不理想. 主要表现在肿瘤分类方法的推广不足, 即根据训练样本集所得到的分类规则... 本集上表现出较低的正确率. 即使采用推... 好的 SVM 也是如此. 我们认为其主要症结... 很好的剔除基因表达谱中的噪声. 实际中... 的出现可能仅仅与某些基因的表达水平的... 若笼统地用全部基因表达水平来进行分... 会因数据维数的巨大而难于进行, 而且众多... 将便成噪声大大地干扰分类的结果. 为此... 提出了一些相关基因选取的方法... 应用最广泛...

... 肿瘤诊断... 基因表达... 这些生物... 识别提供... 析出有... 究的主... 因素... 应... 来, 基于... 瘤分类... 其常用的... 特别地, 基... 已成为其中... (1) 首先采... 过程中采... 量选取了... 等对结肠... 表达谱与... 量法进行... [3] 将几种... 分类, 并对... (SVM) 最... 和 G... 等[6] 研究结...

为了避免正态假设, 我们依据非参数统计中的秩和检验理论提出了秩和相关基因选取方法. 然后, 采用 SVM 建立肿瘤诊断模型, 并根据简化后的训练样本数据(相关基因表达谱)进行有监督的学习, 最后在检验样本数据集上进行检验. 通过对两类肿瘤基因表达谱数据的训练和检验, 我们发现这种秩和基因选取方法可以使得 SVM 分类器获得很高的推广能力.

下文中, 我们将在第1节中首先对于相关基因的统计方法进行理论分析, 然后提出了秩和基因选取方法, 并进一步引出了 SVM 作为肿瘤诊断模型. 在第2节中, 我们首先给出了采用秩和方法进行基因选取并应用 SVM 进行肿瘤诊断的一些实验结果, 然后再与 t 统计量方法的结果进行了比较. 最后在第3节给出结论.

## 1 秩和基因选取方法与肿瘤分类模型

### 1.1 相关基因选取的统计分析研究

人们很早便开始了肿瘤相关基因的识别研究, 但基本上是根据生物特性进行的. 随着 DNA 微阵列技术的发展和基于基因表达谱的肿瘤分类方法的研究, 人们提出了一些基于统计分析的肿瘤相关基因选取方法. 这些方法通过引入基因对肿瘤的辨识性度量, 选取出对肿瘤辨识性较大的基因.

特别地, t 统计量及其变形是现今最常用的肿瘤辨识性度量. t 统计量的表达式为  $T = \frac{m_{i,+} - m_{i,-}}{S_w \sqrt{\frac{1}{n_+} + \frac{1}{n_-}}}$ , 其

中  $S_w^2 = \frac{(n_+ - 1)s_{i,+}^2 + (n_- - 1)s_{i,-}^2}{n_+ + n_- - 2}$ ; 式中的  $m_{i,+}$ ,  $m_{i,-}$ ,

$s_{i,+}$  和  $s_{i,-}$  分别为第  $i$  个基因在正、负样本中表达水平的均值和标准差,  $n_+$  和  $n_-$  分别是正负样本的个数. 事实上, t 统计量在统计学的二元 t 检验中可以用于度量两个正态总体的分布差异大小. 因此  $T$  的绝对值越大, 意味着该基因的表达水平的在正负样本中变化越显著, 该基因与样本的肿瘤因素的相关性也就越大. 换句话说该基因对肿瘤的辨识性越强. t 统计量方法<sup>[2]</sup>基于这一统计直观选取出  $T$  的绝对值最大的  $K$  个基因为相关基因.

由于 t 统计量的表达式较复杂, 人们也提出了一些简化的表达式. 例如 Golub 等<sup>[1]</sup>采用的辨识性度量

为  $W_i = \frac{m_{i,+} - m_{i,-}}{s_{i,+} + s_{i,-}}$ , 并要求所选取基因的  $W_i$  取正值和

负值的个数相同. 另外, Furey 等<sup>[5]</sup>采用了  $W_i$  的绝对值, Pavlidis 等<sup>[7]</sup>采用了  $W_i$  二次值作为辨识性度量. t 统计量及其变形的表达式本质上是一致的, 其分子

都是正负样本基因表达水平的均值的差, 分母都是正负样本基因表达水平的方差的函数, 用于正规化辨识性度量的表达式.

此外, Ding<sup>[8]</sup>提出了 t 统计量的一般化形式 F 统计量作为辨识性度量, 可以处理肿瘤类数多于 2 的相关基因选取问题. 假定某基因的表达水平为  $g = (g_1, g_2, \dots, g_n)$ , 肿瘤类别数为  $K$ , 则 F 统计量可

表达为  $F = \left[ \sum_k n_k (\bar{g}_k - \bar{g})^2 / (K - 1) \right] / s^2$ , 其中  $\bar{g}$  和

$\bar{g}_k$  分别是该基因在全体样本和第  $k$  类样本中的平均表达水平,  $n_k$  和  $s_k$  表示第  $k$  类的样本数和方差,

$s^2 = \left[ \sum_k (n_k - 1) s_k^2 \right] / (n - K)$ . 当  $K = 2$  时,  $F = t^2$ ,

$t = \sqrt{\frac{n_1 n_2}{n_1 + n_2} \frac{\bar{g}_1 - \bar{g}_2}{s}}$ , F 统计量退化为 t 统计量.

如前所述, t 统计量方法及其变形都是以 t 检验作为其统计依据的. 然而 t 检验是一种参数检验方法, 以样本服从正态总体的假设为前提. 那么在基因表达谱不服从正态分布的情况下, 使用 t 统计量方法选取相关基因能否得到好的分类结果呢? 从以前的研究成果和我们的实验(见实验部分)可以看出: 一般来说, 即使样本不完全满足正态条件, 利用 t 统计量方法进行相关基因选取依然是有效的. 即利用 t 统计量方法进行基因选取的分类器, 比起不做基因选取的分类器在正确率上能有较大提高. 因此, 即使正态假设不完全满足, t 统计量仍能在一定程度上反映出基因表达水平对肿瘤因素的辨识性.

然而从理论上讲, t 统计量在正态假设不满足的条件下对基因的辨识性度量则是不精确的, 缺乏牢靠统计依据的. 这时候采用 t 统计量方法会产生两个问题: 第一, 利用 t 统计量对基因的排序与真实按照基因对肿瘤辨识性大小的排序可能出现不一致, 即基因错位. 例如, 假设有两个基因 A 和 B, A 的表达水平服从正态分布, B 的表达水平服从某一区间内的均匀分布, A 的 t 统计量的值大于 B 的 t 统计量的值. 因此, A 被排在了 B 的前面. 然而问题的关键在于, 真实反映 B 基因对肿瘤因素辨识性大小的 p 值(p-value, 或称临界值)<sup>[9]</sup>应该按照均匀分布而非正态分布的条

表1 在0.05显

	结肠正常样本	结肠肿瘤样本
数	2000	2000
被否定的基因数	730	148

数据集的正态性检验结果

本	乳腺肿瘤样本	本	急性淋巴白血病	急性髓性白血病
	5776		7129	7129
	74		4542	2558

算. 因此按照 A, B 两基因真实的 p 值排序, 有可能排在 A 基因的前面. 这就是 t 统计量在非正态条件下造成了基因的错位. 第二, 当基因表达谱从正态分布时, 计算出的 t 统计量也不再服从 t 分布, 因此无法通过查统计表或利用统计软件计算基因真实的 p 值. 得不到 p 值就不能利用显著性水平阈值确定合适的相关基因数目, 这将给研究中确定合适的相关基因数目带来不便. 因此可以说, t 统计量方法只有在正态性条件满足时才能表现得最好. 然而, 正态性假设在许多基因表达数据中的应用范围.

为了验证实际问题中基因表达谱是否服从正态分布, 我们在结肠数据<sup>[1]</sup>中利用峰度和偏度<sup>[9]</sup>来检验样本是否服从正态分布, 所犯第一类错误的概率小. 结果表明, 约半数基因的表达谱服从正态分布, 约半数基因的表达谱不服从正态分布. 因此显然不能从正态假设出发来选取相关基因. 对于肿瘤分类问题的基因表达谱不服从正态分布的情况应该是普遍的. 在正态假设不满足时, 需要一种有统计依据的方法替代 t 统计

率高的基因. 在样本量之下, 由秩和方法得到相关基因列表对肿瘤因素的辨识性. 秩和方法在思想上是对每个基因的表达水平在肿瘤和正常样本的分布间是否存在显著差异. 是, 表明肿瘤因素对该基因的表达水平有显著影响. 因此该基因与肿瘤关系紧密, 将其选入相关基因列表. 秩和方法直接利用基因表达水平的值(秩)先对表达水平排序得到“秩”(秩的位置), 再对秩和统计量进行分析. 基因表达数据中含有大量的噪声和异常值, 对利用 t 统计量和对利用秩和统计量的值造成很大影响. 因此秩和方法更适用于基因表达谱数据.

秩和检验分为 Wilcoxon 秩和检验和 Kruskal-Wallis 秩和检验. 前者用于两类别(或两样本与正常)的相关基因选取问题, 后者适用于多类别的基因选取. 本文着重于两肿瘤类别的基因选取问题, 以下简称 Wilcoxon 秩和检验用于基因选取的基本步骤:

- (1) 建立假设  
一般建立如下假设: 肿瘤和正常样本(或两

非参数统计中的秩和检验应用于相关基因的选取, 建立了基因表达谱的秩和方法. 非参数统计方法的优点是具有样本分布的无关性, 即不需要假设样本分布的类型. 其中秩和检验是一种常用的、检验效率很高的非参数检验方法. 非参数统计的理论证明, 在样本总体不满足正态条件的情况下, 秩和检验的 Pitman 渐进效率远高于 t 检验<sup>[11]</sup>. 对于相关基因选取问题而言, 检验效

小到... 秩表... 到 n... 示.

统计量... 该类...  $H_0$ ... 的期望值是

1) 结肠数据, 40 个肿瘤样本, 22 个正常样本, 2000 个基因. 来自 <http://www.genepattern.org/data/colorectal.html>  
白血病数据, 47 个急性淋巴白血病样本(ALL), 25 个急性髓性白血病样本(AML). 来自 [http://www.genepattern.org/data/ALL\\_AML.html](http://www.genepattern.org/data/ALL_AML.html)

2) 乳腺肿瘤数据, 14 个正常样本, 14 个肿瘤样本, 5776 个基因. 来自 <http://genepattern.org/data/breast.html>



文

结肠  
白血病

函数  
Matlab  
light  
hims  
下,  
选  
然  
的  
/M  
验式  
4

在有限元数值分析中



数据集	基因选取
结肠数据	秩和方 t 统计量
数据集	基因选取
白血病数据	秩和方 t 统计量

比较

	34 个基因	8 个基因
	96.2%	88.8%
	93.6%	88.8%
	844 个基因	398 个基因
	100%	98.6%
	95.4%	94.4%

一优势还是十分显著的。结肠和白血病这两个数据集的正态分布，因此实验结果在数据集上优于 t 统计量。

### 3 总结

本文对于基于基因选择的研究，而肿瘤诊断问题。本文针对传统的秩和方不足，提出了秩和方在正态条件的缺陷，另一方面，秩和方来确定出与肿瘤基因数的确定更为科学。来确定相关基因。机(SVM)在选择诊断模型。实验 SVM 的肿瘤数据集和白血病数据集。实验也证明了方法的。这些诊断模型是

致谢 本文  
助项目。

1999, 96: 6745~6750

P S, Grundy W N, Lin D, et al. Knowledge-based analysis of microarray gene expression data by using support vector machines. Proc Nat'l Acad Sci, 2000, 97(1): 262~267

Fridyand J, Speed T P. Comparison of discrimination methods for the classification of tumor using gene expression data. American Statistical Association, 2002, 97(457): 77~87

Christianini N, Duffy N, et al. Support vector machine classification and validation of cancer tissue samples using gene expression data. Bioinformatics, 2000, 16(10): 1053~1062

Stanton J, Barnhill S, et al. Gene selection for cancer diagnosis using support vector machine. Machine Learning, 2000, 38: 389~422

Stanton J, Cai J, et al. Gene functional Analysis from microarray data. Proc Fifth Int. Conf. on Computational Intelligence in Bioinformatics. New York: ACM Press, 2001. 249~255

Gene expression analysis of gene expression profiles: class discovery and class prediction. In: Proc RECOMB, 2002. 127~136

Methods of Statistical Analysis. (2<sup>nd</sup> edition). New York: Wiley & Sons, 1956

Van der Waerden S G. Statistical Inference Based on Ranks, New York: Wiley & Sons, Inc, 1984

Van der Waerden S G. Asymptotic efficiency of non-parametric tests. University Press, 1995

Learning Theory. New York: Wiley, 1998

Large-scale SVM learning practical. In: Proceedings in Kernel Methods-Support Vector Machines. New York: Wiley, 1999

收稿, 2004-03-16 收第 1 次修改稿, 2004-04-29 收第 2 次修改稿)

- 1 Golub
- 2 /