

学校编码：10384

分类号__密级__

学号：23120090153690

UDC__

廈門大學

博士学位论文

基因芯片表达数据分析及基因调控网络 模型研究

**Data Analysis of Expression with Gene Microarray and
Investigation for Gene Regulatory Networks**

郭顺

指导教师姓名：王守觉院士
郭东辉教授

专业名称：电路与系统

论文提交日期：2017年 月

论文答辩时间：2017年 月

学位授予日期：2017年 月

答辩委员会主席：_____

评阅人：_____

2017年

厦门大学博硕士学位论文摘要库

厦门大学学位论文原创性声明

本人呈交的学位论文是本人在导师指导下,独立完成的研究成果。本人在论文写作中参考其他个人或集体已经发表的研究成果,均在文中以适当方式明确标明,并符合法律规范和《厦门大学研究生学术活动规范(试行)》。

另外,该学位论文为()课题(组)的研究成果,获得()课题(组)经费或实验室的资助,在()实验室完成。(请在以上括号内填写课题或课题组负责人或实验室名称,未有此项声明内容的,可以不作特别声明。)

声明人(签名):

年 月 日

厦门大学博硕士学位论文摘要库

厦门大学学位论文著作权使用声明

本人同意厦门大学根据《中华人民共和国学位条例暂行实施办法》等规定保留和使用此学位论文,并向主管部门或其指定机构送交学位论文(包括纸质版和电子版),允许学位论文进入厦门大学图书馆及其数据库被查阅、借阅。本人同意厦门大学将学位论文加入全国博士、硕士学位论文共建单位数据库进行检索,将学位论文的标题和摘要汇编出版,采用影印、缩印或者其它方式合理复制学位论文。

本学位论文属于:

- () 1. 经厦门大学保密委员会审查核定的保密学位论文, 于
年 月 日解密, 解密后适用上述授权。
- () 2. 不保密, 适用上述授权。

(请在以上相应括号内打“√”或填上相应内容。保密学位论文应是已经厦门大学保密委员会审定过的学位论文, 未经厦门大学保密委员会审定的学位论文均为公开学位论文。此声明栏不填写的, 默认为公开学位论文, 均适用上述授权。)

声明人(签名):

年 月 日

厦门大学博硕士学位论文摘要库

摘要

基因芯片表达数据可以揭示在各种不同条件下基因的活动情况及基因之间相互作用关系等，因此，对该数据的分析有着非常广泛的应用前景。近年来，随着基因芯片表达数据快速增长，传统的分析方法已经不能适应现代生物迅速发展的要求，迫切需要更加高效的方法来分析处理这些数据。

本文首先简述了基因芯片表达数据分析的关键技术及其研究进展，指出了现有方法存在处理冗余和噪声数据能力不足、结果缺乏生物学解释以及调控网络识别率有限等问题。针对这些问题，我们研究了基因芯片表达数据分析三种新方法，包括：基因芯片表达数据的特征选择方法、维度约简方法以及重构基因调控网络方法，为建立分类模型辅助诊断以及识别疾病相关基因等提供了解决方案，为生物进程分析、发现潜在药物靶点以及基因功能预测等提供了有效路径。本文的主要工作和创新点如下：

1. 提出了基于核函数为类分布中心的特征选择方法，即以类的分离性准则来定义求解的目标函数，给出具有较好的处理噪声和冗余数据的数据分析方法，提高特征选择的分类识别率，并且证明其在初始值不为零的条件下可获得全局最优解的结论；
2. 提出了基于 l_1 正则化的两阶段的局部维度约简方法，该方法不仅可以提高基因的分别识别率，且所获得的满足生物学解释的结果，其优点是所使用特征选择方法可去除冗余和不相关特征，并可实现对所筛选的特征进行具有生物学解释的特征提取；
3. 提出基于偏最小二乘(PLS)的基因调控网络重构方法，该方法将基因调控网络的模型重构问题分解成多个特征选择的子问题，每个子问题可通过基于 PLS 的特征选择方法来求解，并结合误差统计来进一步优化重构的调控网络，提高调控网络重构的准确率。

最后，基于本论文所提的技术与方法，我们设计一个基因芯片数据分析系统软件，该系统主要包含特征选择、维度约简以及重构调控网络这三种功能模块，可以直观的反映出本文工作成果，体现了本课题研究的实用性。

关键词：基因芯片；数据分析；特征选择；特征提取；调控网络

Abstract

The gene microarray data is able to reveal the gene activity in various conditions and the interactions between different genes, therefore, it has a very broad application prospects for analyzing gene microarray data. Recently, the amount of gene microarray data is growing faster than the rate at which it can be analyzed and more effective techniques and methods are needed to analysis these data.

In this dissertation, we introduce the key techniques and research progress of analysis of gene microarray data, and point out that many existed methods have some weakness including low capability of dealing with redundant and noisy data, unexplained mining results and limited accuracy of inferred regulatory gene networks. In light of these problems, we investigate three new gene microarray data mining methods, including feature selection and dimension reduction methods for gene microarray data, and the methods for the reconstruction of gene regulatory networks, which provide solutions for establishing classification model and identifying disease genes, and provide an effective way for biological process analysis, identifying the potential drug targets and predicting the gene function. The main work and innovation of this dissertation are as follows:

1. Proposing a centroid-based feature selection method. Theoretical analysis indicates that the global optimal solution of the proposed method can be reached with a non-zero initial point. In the proposed method, a kernel-based approach is used to estimate the class centroid to define the class separability criterion, based on which the objective function is formulated;
2. Proposing a two-stage l_1 -regularized local dimension reduction method, which firstly uses a feature selection method to remove the redundant and irrelevant features, and then implements feature extraction on the selected features. The proposed method not only improves the classification accuracy, but also can generate biologically interpretable results;
3. Proposing a PLS-based method for the reconstruction of gene regulatory networks (GRNs). The proposed method decomposes the GRNs inference problem into multiple feature selection subproblems, each of which is solved by a PLS-based

feature selection method. Then, a statistical technique is used to refine the predictions of the inference network.

Finally, Based on the technologies and methods proposed in this dissertation, a Gene Microarray Data Analysis System is designed and implemented. The system contains three main functional modules, including Feature Selection, Dimension Reduction and Gene Regulatory Network Inference, which reflect the results of our work and the practicality of this research.

Key Words: Gene Microarray Data; Data Analysis; Feature Selection; Feature Extraction; Gene Regulatory Networks

目 录

第一章 绪论	1
1.1 引言	1
1.2 基因数据分析相关技术	4
1.3 关键技术及其研究进展	6
1.2.1 基因选择技术	6
1.2.2 特征提取技术	7
1.2.3 重构调控网络	8
1.4 主要内容与章节安排	10
第二章 基因芯片数据分析方法及技术原理	13
2.1 引言	13
2.2 特征选择	13
2.2.1 Filter 类特征选择	14
2.2.2 Wrapper 类特征选择	17
2.2.3 Embedded 类特征选择	18
2.2.4 l_1 正则化的特征选择	20
2.2.5 多类特征选择框架	23
2.3 特征提取	24
2.3.1 主成分分析	25
2.3.2 线性判别分析	26
2.3.3 偏最小二乘	27
2.4 重构基因调控网络	28
2.4.1 布尔网络模型	28
2.4.2 贝叶斯网络模型	29
2.4.3 微分方程模型	31
2.4.4 集成方法模型	32
2.5 本章小结	34

第三章	分布式数据中心定位特征选择方法	35
3.1	相关背景	35
3.2	研究方法	37
3.2.1	相关定义	37
3.2.2	目标函数	39
3.2.3	求解算法	41
3.2.4	多类问题	43
3.2.5	计算复杂度	44
3.3	实验结果与分析	44
3.3.1	实验数据	45
3.3.2	实验设置	46
3.3.3	比较方法	46
3.3.4	分类性能分析	47
3.4	本章小结	57
第四章	基因芯片数据维度约简范数优化方法	58
4.1	相关背景	58
4.1.1	维度约简方法	58
4.1.2	特征提取方法	59
4.2	研究方法	62
4.2.1	l_1 正则化特征选择	62
4.2.2	局部维度约简方法	65
4.2.3	计算复杂度	66
4.3	实验结果与分析	67
4.3.1	实验数据	67
4.3.2	实验设置	68
4.3.3	比较方法	69
4.3.4	分类性能分析	70
4.3.5	生物学验证	79

4.4	本章小结.....	84
第五章	基因调控网络误差优化重构方法.....	85
5.1	相关背景	85
5.2	研究方法	87
5.2.1	特征选择方法.....	87
5.2.2	重构网络优化.....	88
5.2.3	参数设置.....	89
5.2.4	计算复杂度.....	90
5.3	实验结果与分析.....	91
5.3.1	实验数据.....	91
5.3.2	评价方法.....	92
5.3.3	比较方法.....	93
5.3.4	识别性能分析.....	93
5.4	本章小结.....	100
第六章	基因芯片数据分析系统设计	101
6.1	引言	101
6.2	系统设计方案	101
6.3	系统功能展示	103
6.4	本章小结	109
第七章	总结与展望.....	110
7.1	总结	110
7.2	展望	111
	参考文献.....	112
	在学期间取得的科研成果.....	120
	致 谢.....	121

Table of Contents

Chapter1 Introduction.....	1
1.1 Research Background.....	1
1.2 Gene Microarray Data Analysis Techniques.....	4
1.3 Key Techniques and Research Progress.....	6
1.3.1 Feature Selection Techniques	6
1.3.2 Feature Extraction Techniques	7
1.3.3 Gene Regulatory Networks Inference	8
1.4 Mainly content and Structure of This Thesis.....	10
Chapter 2 Gene Microarray Data Analysis Methods and Technical Principle.....	13
2.1 Introduction.....	13
2.2 Feature Selection.....	13
2.2.1 Filter Methods	14
2.2.2 Wrapper Methods	17
2.2.3 Embedded Methods	18
2.2.4 l_1 -Regularized Methods.....	20
2.2.5 Framework for Multi-category Problems	23
2.3 Feature Extraction.....	24
2.3.1 Principal Component Analysis	25
2.3.2 Linear Discriminant Analysis	26
2.3.3 Partial Least Squares	27
2.4 Reconstruction of Gene Regulatory Networks.....	28
2.4.1 Boolean Networks	28
2.4.2 Bayesian Networks	29
2.4.3 Differential Equations Model	31
2.4.4 Methods using Ensembles of Feature Selection Algorithms	32
2.5 Summary.....	34

Chapter 3 Distributed Data Centroid Location Feature Selection

Method	35
3.1 Relevent Background	35
3.2 Research Method	37
3.2.1 Definitions.....	37
3.2.2 Objective Function.....	39
3.2.3 An Algorithm For the Objective Function.....	41
3.2.4 Feature Selection for Multiclass Problems.....	43
3.2.5 Computational Complexity.....	44
3.3 Experiments and Analysis	44
3.3.1 Experimental Datasets.....	45
3.3.2 Experimental Setting.....	46
3.3.3 Compared Methods and Parameter Selection.....	46
3.3.4 Experimental Results and Analysis.....	47
3.4 Summary	57

Chapter 4 Gene Microarray Data Dimension Reduction using Norm

Optimaiton based Method	58
4.1 Relevent Background	58
4.1.1 Dimension Reduction.....	58
4.1.2 PLS-based Feature Extraction.....	60
4.2 Research Method	62
4.2.1 l_1 -Regularized Feature Selection	62
4.2.2 Local Dimension Reduction.....	65
4.2.3 Computational Complexity.....	66
4.3 Experiments and Analysis	67
4.3.1 Experimental Datasets.....	67
4.3.2 Experimental Setting.....	68
4.3.3 Compared Methods and Parameter Selection.....	69

4.3.4	Experimental Results and Analysis.....	70
4.3.5	Biological Interpretation of Results.....	79
4.5	Summary.....	84
Chapter 5	Gene Regulatory Networks Inference using Error Optimization based Method.....	85
5.1	Relevant Background.....	85
5.2	Research Method.....	87
5.2.1	Feature Selection Method.....	87
5.2.2	Refining the Inferred Regulatory Network.....	88
5.2.3	Parameter Selection.....	83
5.2.4	Computational Complexity.....	84
5.3	Experiments and Analysis.....	86
5.3.1	Experimental Datasets.....	86
5.3.2	Evaluation.....	93
5.3.3	Compared Methods and Parameter Selection.....	94
5.3.4	Experimental Results and Analysis.....	94
5.4	Summary.....	100
Chapter 6	Gene Microarray Data Analysis System.....	101
6.1	Introduction.....	101
6.2	Design Scheme.....	101
6.3	Function Demonstration.....	102
6.4	Summary.....	109
Chapter 7	Conclusions and Future Work.....	110
7.1	Conclusions.....	110
7.2	Future Work.....	111
	References.....	112
	Publication List During Post Graduation.....	120

Acknowledgements 121

厦门大学博硕士学位论文摘要库

Degree papers are in the "[Xiamen University Electronic Theses and Dissertations Database](#)". Full texts are available in the following ways:

1. If your library is a CALIS member libraries, please log on <http://etd.calis.edu.cn/> and submit requests online, or consult the interlibrary loan department in your library.
2. For users of non-CALIS member libraries, please mail to etd@xmu.edu.cn for delivery details.

厦门大学博硕士论文摘要库