

学校编码: 10384

分类号_____密级_____

学 号: 23020131153186

UDC_____

厦 门 大 学

硕 士 学 位 论 文

miRNA-disease 关联预测方法研究

The Research of MiRNA-Disease Correlation Prediction

李 金 金

指 导 教 师: 曾 湘 祥 副 教 授

专 业 名 称: 计 算 机 技 术

论 文 提 交 日 期: 2016 年 5 月

论 文 答 辩 日 期: 2016 年 5 月

学 位 授 予 日 期: 2016 年 月

答 辩 委 员 会 主 席: _____

评 阅 人: _____

2016 年 月

miRNA-disease 关联预测方法研究

李 金 金

指 导 教 师 曾 湘 祥 副 教 授

厦 门 大 学

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): 方金金

2016年5月17日

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）： 方金金

2016年5月17日

摘 要

MiRNA 是一类重要的非编码 RNA，它在许多的生命活动中扮演着非常重要的角色。经研究证实，miRNA 与 disease 的产生有着密切的联系。本文主要基于生物异构网络和机器学习的方法对 miRNA-disease 关系进行分析与预测。主要包括以下几个方面的内容：

(1) 本文详细介绍了当前已提出的 miRNA-disease 关联预测的主要方法，并且对比了不同方法之间的差异，为生物信息研究人员后续对 miRNA-disease 关联分析和预测的研究提供参考。现已提出的 miRNA-disease 关联分析和预测方法，大多利用 miRNA 的生物特性构建 miRNA 之间的相似性以及利用文本挖掘的方式研究疾病之间的相似度，从而进一步研究 miRNA-disease 之间的关系。

(2) 引入并改造了成功应用在社交网络中预测社交关系的方法：KATZ 和 CATAPULT 方法，来解决 miRNA-disease 关联分析与预测的问题。KATZ 方法是基于生物异构网络的方法，充分利用异构网络中 miRNA 节点与 disease 节点之间不同步长的不同到达方式的预测方法。CATAPULT 是基于半监督学习的方法，利用偏向的支持向量机预测 miRNA-disease 的关系。然而，对于 miRNA-disease 关联对，正的关联关系是确定的，但是反的关联关系是不确定的。因此，对于问题的这种特性，引入了偏向的支持向量机来做支撑，以准确有效地预测 miRNA-disease 的关联关系。

(3) 提出了基于集成学习的 miRNA-disease 关联关系预测方法。对于 miRNA-disease 关系预测的研究已经涌现了多种方法，但是，较为流行的机器学习方法并没有在 miRNA-disease 关系预测上发挥作用。本文综合多种现有的预测 miRNA-disease 关系的方法作为 miRNA-disease 对的特征向量，再利用集成分类器对 miRNA-disease 对进行分类。

关键词：miRNA-disease 关系；生物异构网络；支持向量机

Abstract

MiRNA is an important class of noncoding RNA and plays a key role in regulating life processes. Researches have shown that miRNA plays an important role in the occurrence and development of disease. In this paper, we analyze and predict the associations between miRNAs and diseases based on the biological heterogeneous network and machine learning method. The main contents are as follows:

(1) In this paper, we introduce in detail the methods that have been proposed to predict the associations between miRNAs and diseases and compare the differences between the methods, which provide the important reference information for the follow-up study of biological information personnel who are interested in the miRNA-disease associations' research. The majority of methods use text mining to calculate the diseases similarity and biological characteristics of miRNA to obtain the miRNA similarity, then, miRNA-disease associations are calculated based on diseases similarity and miRNAs similarity. For example, family and cluster characteristics of miRNA are often used.

(2) We introduce and reform the methods, which succeed in predicting social network association. KATZ based on biological heterogeneous network and synthesizes the different step sizes and different modes of arrival between miRNA and disease in heterogeneous network. CATAPULT based on semi supervised learning method and uses the biased support vector machine to classifier miRNA-disease pairs. For the miRNA-disease pair's classification problems, the positive correlations are determined, but the negative correlations are uncertain. Taking into account the characteristics of the problem, we introduce biased support vector machine to analysis and predict the associations between miRNAs and diseases effectively and accurately.

(3) We proposed the prediction method based on ensemble learning. A number of approaches have emerged for miRNA-disease association problem, but more popular machine learning method does not play a role in the issue. In this method, we

integrated multiple methods predicting the miRNA-disease association to act as the features. This not only synthesizes the advantages of various proposed methods, but also introduces the machine learning methods.

Key words: miRNA-disease association; Biological heterogeneous network; SVM

廈門大學博碩

目 录

摘 要.....	I
Abstract.....	II
第一章 绪 论.....	1
1.1 课题研究背景及意义.....	1
1.1.1 课题研究背景.....	1
1.1.2 课题研究意义.....	2
1.2 miRNA 相关知识介绍.....	3
1.2.1 miRNA 的产生和作用机制.....	3
1.2.2 miRNA-disease 的关联关系.....	4
1.3 国内外研究现状.....	5
1.3.1 生物异构网络的研究现状.....	5
1.3.2 miRNA 研究现状.....	7
1.3.3 miRNA-disease 关系预测中存在的问题.....	9
1.4 本文的主要研究内容.....	10
第二章 miRNA 与 disease 关联预测算法介绍.....	13
2.1 引言.....	13
2.1.1 相关生物网络.....	13
2.1.2 致病 miRNA 排序问题的定义.....	14
2.2 现有致病 miRNA 的预测算法.....	15
2.2.1 基于本地网络的预测算法.....	16
2.2.2 基于全局网络的预测算法.....	21
2.2.3 其他的预测算法.....	28
2.3 本章小结.....	29
第三章 基于异构网络和半监督学习预测致病 miRNA.....	31
3.1 引言.....	31
3.2 预测方法介绍.....	31

3.2.1 基于异构网络拓扑结构的关联预测算法介绍.....	31
3.2.2 基于半监督学习的关联预测算法介绍.....	34
3.3 实验结果与分析.....	36
3.3.1 数据集.....	36
3.3.2 预测评估方法.....	37
3.3.3 性能比较.....	41
3.4 本章小结.....	42
第四章 基于集成学习的 miRNA-disease 关系预测方法.....	45
4.1 引言.....	45
4.1.1 机器学习.....	45
4.1.2 集成学习.....	45
4.2 数据集与方法.....	47
4.2.1 数据集介绍.....	47
4.2.2 实验方法与结果.....	49
4.3 实验结果与分析.....	52
4.3.1 预测评估方法.....	52
4.3.2 集成学习方法的预测性能与对比.....	53
4.4 本章小结.....	57
第五章 总结与展望.....	59
5.1 本文总结.....	59
5.2 未来展望.....	60
参 考 文 献.....	61
攻读学位期间发表的学术论文.....	67
致 谢.....	69

CONTENTS

Abstract(CN)	I
Abstract(EN)	II
Chapter 1 Introduction	1
1.1 Background	1
1.1.1 Background.....	1
1.1.2 Significance of subject.....	2
1.2 miRNA	3
1.2.1 Generation and mechanism of miRNA.....	3
1.2.2 miRNA-disease association.....	4
1.3 Reasarch Status	5
1.3.1 Biological heterogeneous network.....	5
1.3.2 miRNA Reasarch Status.....	7
1.3.3 miRNA-disease association problems.....	9
1.4 Main Research Contents	10
Chapter 2 Introduction MiRNA-disease Association Prediction	
Algorithm	13
2.1 Introduction	13
2.1.1 Related Biological network.....	13
2.1.2 Definition ranking miRNA for the disease.....	14
2.2 Existing prediction algorithms for miRNA-disease association	15
2.2.1 Prediction algorithms based on local network.....	16
2.2.2 Prediction algorithms based on global network.....	21
2.2.3 Other prediction algorithms.....	28
2.3 Conclusion	29

Chapter 3 MiRNA-disease Association Prediction Algorithm Based on heterogeneous network and semi-supervised learning.....	31
3.1 Introduction.....	31
3.2 Introduction prediction algorithm.....	31
3.2.1 Prediction Algorithm Based on Heterogeneous Network.....	31
3.2.2 Prediction Algorithm Based on Semi-supervised Learning.....	34
3.3 Comparison and Discussion.....	36
3.3.1 Dataset.....	36
3.3.2 Evaluation.....	37
3.3.3 Comparison.....	41
3.4 Conclusion.....	42
Chapter 4 Based on ensemble learning algorithm to predict miRNA-disease association.....	45
4.1 Introduction.....	45
4.1.1 Machine Learning.....	45
4.1.2 Ensemble Learning.....	45
4.2 Dataset and Methods.....	47
4.2.1 Dataset.....	47
4.2.2 Methods and Results.....	49
4.3 Comparison and Discussion.....	52
4.3.1 Evaluation.....	52
4.3.2 Comparison.....	53
4.4 Conclusion.....	57
Chapter 5 Conclusion and Future Work.....	59
5.1 Conclusions.....	59
5.2 Future Work.....	60
References.....	61

Publications.....	67
Acknowledgement.....	69

廈門大學博碩

第一章 绪论

1.1 课题研究背景及意义

1.1.1 课题研究背景

1985年，美国科学家提出了一个大规模跨国跨学科的科学探索工程-人类基因组计划。人类基因组计划是生命科学中的一个重大课题，它的目的在于获得人类基因组完整的核苷酸序列，用于鉴别人类基因组中的所有基因。人类将借此重大的科学计划破译人类的生命，解读人类自身的奥秘。人类基因组计划于1990年10月1日正式启动，预计人类基因组计划15年内将耗资30亿美元解读人类基因。2000年6月，人类基因组计划参与国的协作组宣布：人类生命蓝图已基本完成。公告意味着人类从自己基因的角度对人类生老病死的客观规律、疾病诊断和治疗有了更深入的了解和掌握，对人类生命的起源有了更准确的认识。同时，标识着人类在摸索生命奥秘的道路上树立了新的里程碑。

人类基因组计划的完成，意味着全新的“后基因组时代”即将开启。迄今为止，人类已经拥有了大量的生物数据，加上大数据时代的带来，生物数据仍然在迅猛地增长。从海量的生物数据中挖掘出生命的信息是人类基因组计划的真正目的。后基因组时代即将遇到的问题主要包括：基因组的表达与调控、蛋白质产物的功能、基因组的多样性、模式生物基因组的研究等。后基因组时代的研究将让人类更加深入地了解遗传语言功能的逻辑架构；基因的功能与结构之间的关系；人类生长、发育、衰老和死亡的原理；脑功能的表现与神经活动之间的关系；疾病的发生与发展；信息传递和作用机理；基因后机理以及关于人类生命科学的各种问题等。人类基因组研究的目的是从根底开始探索人类生命的起源；人类物种间以及不同个体之间的差别；导致疾病发生与发展的原因；长寿与衰老的原因等。美国国立卫生院和能源部实施了两个重大科研计划是：基因组到生命(Genomes to Life, GTL)^[1]与 Roadmap。GTL 科研计划的基础是准确地刻画出生命系统的所有“分子机器”，认清“分子机器”在生命体中是如何协调工作的。其目标是：“分子机器”的辨别，其中，分子机器是蛋白质的复合物产物，其扮演着执行生命系统的功能；充分地认识与理解调控“分子机器”的原理；发展计算机技术以

Degree papers are in the “[Xiamen University Electronic Theses and Dissertations Database](#)”.

Fulltexts 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.

廈門大學博碩