

DOCTORAL THESIS

Statistical methods for integrative analysis of genomic data

Ming, Jingsi

Date of Award:
2018

[Link to publication](#)

General rights

Copyright and intellectual property rights for the publications made accessible in HKBU Scholars are retained by the authors and/or other copyright owners. In addition to the restrictions prescribed by the Copyright Ordinance of Hong Kong, all users and readers must also observe the following terms of use:

- Users may download and print one copy of any publication from HKBU Scholars for the purpose of private study or research
- Users cannot further distribute the material or use it for any profit-making activity or commercial gain
- To share publications in HKBU Scholars with others, users are welcome to freely distribute the permanent URL assigned to the publication

Abstract

Thousands of risk variants underlying complex phenotypes (quantitative traits and diseases) have been identified in genome-wide association studies (GWAS). However, there are still several challenges towards deepening our understanding of the genetic architectures of complex phenotypes. First, the majority of GWAS hits are in non-coding region and their biological interpretation is still unclear. Second, most complex traits are suggested to be highly polygenic, i.e., they are affected by a vast number of risk variants with individually small or moderate effects, whereas a large proportion of risk variants with small effects remain unknown. Third, accumulating evidence from GWAS suggests the pervasiveness of pleiotropy, a phenomenon that some genetic variants can be associated with multiple traits, but there is a lack of unified framework which is scalable to reveal relationship among a large number of traits and prioritize genetic variants simultaneously with functional annotations integrated. In this thesis, we propose two statistical methods to address these challenges using integrative analysis of summary statistics from GWASs and functional annotations.

In the first part, we propose a latent sparse mixed model (LSMM) to integrate functional annotations with GWAS data. Not only does it increase the statistical power of identifying risk variants, but also offers more biological insights by detecting relevant functional annotations. To allow LSMM scalable to millions of variants and hundreds of functional annotations, we developed an efficient variational expectation-maximization (EM) algorithm for model parameter estimation and statistical inference. We first conducted comprehensive simulation studies to evaluate the performance of LSMM. Then we applied it to analyze 30 GWASs of complex phenotypes integrated with nine genic category annotations and 127 cell-type specific functional annotations from the Roadmap project. The results demonstrate that our method possesses more statistical power than conventional methods, and can help researchers achieve deeper understanding of genetic architecture of these complex phenotypes.

In the second part, we propose a latent probit model (LPM) which combines summary statistics from multiple GWASs and functional annotations, to characterize relationship and increase statistical power to identify risk variants. LPM can also perform hypothesis testing for pleiotropy and annotations enrichment. To enable the scalability of LPM as the number of GWASs increases, we developed an efficient parameter-expanded EM (PX-EM) algorithm which can execute parallelly. We first validated the performance of LPM through comprehensive simulations, then applied it to analyze 44 GWASs with nine genic category annotations. The results demonstrate the benefits of LPM and can offer new insights of disease etiology.

Table of Contents

Declaration	i
Abstract	ii
Acknowledgements	iii
Table of Contents	iv
List of Figures	v
List of Tables	vi
Chapter 1 Introduction	1
1.1 Complex Traits are Highly Polygenic	1
1.2 Integrative Analysis of GWAS with Functional Annotations	2
1.3 Integrative Analysis of Multiple GWASs	3
1.4 Major Contributions	5
1.5 Outline of the Thesis	5
Chapter 2 LSMM: A statistical approach to integrating functional annotations with genome-wide association studies	7
2.1 Introduction	7
2.2 Model	8
2.3 Algorithm	10
2.3.1 The Variational EM Algorithm	10
2.3.2 Details of the Four-stage Algorithm	17
2.4 Inference	29

2.4.1	Identification of Risk SNPs	29
2.4.2	Detection of Relevant Cell-type Specific Functional Annotations	29
2.5	Simulation	30
2.5.1	Ideal Case	30
2.5.2	More simulations	45
2.5.3	Comparison with other methods	56
2.6	Real Data Analysis	60
Chapter 3	LPM: a latent probit model to characterize relationship among complex traits using summary statistics from multiple GWASs and functional annotations	69
3.1	Introduction	69
3.2	Model	70
3.3	Algorithm	72
3.3.1	The PX-EM Algorithm	72
3.3.2	Details of the Three-stage Algorithm	77
3.3.3	Method to get the parameter estimation in LPM	86
3.4	Inference	87
3.4.1	Identification of risk SNPs	87
3.4.2	Relationship test among traits	90
3.4.3	Hypothesis testing of annotation enrichment	90
3.5	Theorem based on composite likelihood approach	92
3.6	Simulation	93
3.6.1	Simulation of eight traits	93
3.6.2	Simulation of eight traits without annotation	102
3.6.3	Simulations of two traits	107
3.6.4	Computational time	113
3.6.5	More simulations	114
3.7	Real Data Analysis	115
Chapter 4	Conclusion	131
	Curriculum Vitae	140