Learning Hierarchical Bayesian Networks for Genome-Wide Association Studies

Raphaël Mourad¹, Christine Sinoquet² and Philippe Leray¹

KOD team (KnOwledge and Decision), ¹ LINA, UMR CNRS 6241, Ecole Polytechnique de l'Université de Nantes. ² LINA, UMR CNRS 6241, Université de Nantes. FRANCE





Presented by Raphael Mourad PhD student in Bioinformatics raphael.mourad@univ-nantes.fr



Outline

- 1/ Introduction
- 2/ Fondamental concept of association genetics
- 3/ Presentation of genetic data
- 4/ Our approach
- 5/ Results and discussion
- 6/ Conclusion and outlooks



Introduction



• Context:

Complex genetic diseases

= multifactorial genetic diseases caused by a combination of genetic factors (*eg* genes) and environmental factors (*eg* sex, age...).

Examples: diabetes, asthma, hypertension, some cancers...



• Dissect the genetic basis of these diseases:

Genome-wide association studies (GWAS)

 \rightarrow identification of genetic markers associated with common, complex diseases.





Fondamental concept of association genetics



• Linkage disequilibrium (LD):

 \rightarrow dependences generally observed between close SNPs on the chromosome,

 \rightarrow at the basis of GWAS.





Presentation of genetic data



Phenotype

1 binary variable:
1000 non-affected individuals
1000 affected individuals

DNA

> 100k SNP Ternary variables

• Characteristics:

 \rightarrow large number of genetic variables (SNP): combinatorial explosion

 \rightarrow strong dependences among genetic variables



Our approach



• Reduce the data dimension by synthetizing the information of highly dependent SNPs, due to LD.





 Provide a flexible and adapted probabilistic model to reduce dimension for genetic data.





- Advantages of this modelling:
 - \rightarrow hierarchical, thus :
 - various degrees of dimension reduction,
 - various degrees of LD strength,
 - \rightarrow each latent variable can reveal multiple-SNP patterns, potentially relevant to explain the disease,

 \rightarrow contrary to Hierarchical Latent Class model, SNPs are not constrained to be dependent upon one another,

 \rightarrow high-order interactions between SNPs can be taken into account.



• Proposed algorithm to learn both parameters and structure of FHLCMs from data:

CFHLC (Construction of Forests of Hierarchical Latent Class models).

 \rightarrow based on an agglomerative hierarchical procedure to ensure scalability,

 \rightarrow uses clique partitioning methods for an efficient discovery of non-overlapping cliques of dependent SNPs,

 \rightarrow not restricted to binary variables and binary trees, as Hwang *et al.*'s algorithm.



Schema of the algorithm:





Results and discussion



- Protocol testing:
 - \rightarrow C++ implementation,
 - \rightarrow run on a standard pc (3.8 GHz, 3.3 Go RAM),

 \rightarrow tested on simulated unphased genotypic data consisting of 2000 individuals and 1k, 10k or 100k SNPs, generated with the software Hapsimu.



Scalability



Mourad R. et al : Learning Hierarchical Bayesian Networks for GWAS







Conclusion and outlooks



Conclusion:

- CFHLC algorithm have been shown to be efficient on genome-scaled data,
- Can provide a data dimension reduction of 80%.

Perspectives:

- Application on the detection of genetic associations thanks to FHLCM's latent variables,
- Visualization of LD structure through the FHLCM's graph.



Thanks for your attention







Questions

Mourad R. et al : Learning Hierarchical Bayesian Networks for GWAS



Impact of window size on running time



Mourad R. et al : Learning Hierarchical Bayesian Networks for GWAS



Impact of window size on dimension reduction





Bibliography

General on GWASs:

- <u>Balding D. (2006):</u> a tutorial on statistical methods for population association studies.

Specific to probabilistic graphical models:

- <u>Verzilli (2007)</u>: Bayesian graphical models for genome-wide association studies.

- <u>Hwang (2006)</u>: learning hierarchical Bayesian networks for large-scale data analysis.