Statistical learning of epistasis, genetic maps and microbial networks via graphical models

Prof. Ernst C. Wit

Institute of Computational Science Università della Svizzera italiana

Joint work with Pariya Behrouzi, Olaf Schenk and Arnaud Cougoul 16 May 2019



$Correlation \neq causation$

CORRELATION IS NOT CAUSATION!



Both ice cream sales and shark attacks increase when the weather is hot and sunny, but they are not caused by each other (they are caused by good weather, with lots of people at the beach, both eating ice cream and having a swim in the sea)

Events are connected by a common cause: confounding



R.A. Fisher vs. Richard Doll

R.A. Fisher (geneticist and statistician) was a fervent smoker.



"Smoking and lung cancer are confounded"

Smoking Kills

"Control for all possible confounders"

Sir Richard Doll conducted:

- 1950. Lung cancer study in 20 London hospitals.
- 1954–2001 British Doctor Study to eliminate confounders.



Although in a marginal model

lung cancer = $\theta_0 + \theta_1$ smoking + ϵ ,

it may seem like

 $\theta_1 > 0 \implies$ "smoking causes cancer",

but this is the correlation = causation fallacy.

Controlling for confounders

Only if in a complete model

lung cancer = $\theta_0 + \theta_1$ smoking + θ_2 age + θ_3 genetics + ... + ϵ

we still have $\theta_1 > 0$, then we can *hope* for causation.

Eliminating confounders (visually)

Instead of only checking for correlation,

Smoking Lung cancer

we should control for other possible explanations:



No wonder that the British Doctor Study lasted until 2001!

Ps. Philosophically, it is never finished: Think Popper's falsification vs verification.



A **direct influence** network can have a causal interpretation, ... if all potential confounders are included.



Aim of this talk

To estimate this (conditional independence) graph from data



1. Epistatic interactions



Epistasis (= "statistical interaction")

Epistasis

when one gene locus masks or modifies a second gene locus phenotype



Phenotype = squash colour

Arabidopsis Recombinant Inbred Line study

RIL cross made consisting of 367 F_8 generation RIL Arabidopsis plants:



The plants are genotyped at 90 markers.

Svizzera taliana

Epistasis in Arabidopsis



Embryo lethal phenotype





Certain genotype combinations may be functionally incompatible.

Phenotype = survival

Aim of this genotype study

Detection of high-dimensional epistatic selection across Arabidopsis genome.

Approach

Extending graphical model for discrete ordinal data to determine the pattern of conditional independence relationships.



RIL genotype data

For each of i = 1, ..., 367 plants, we obtain genotype at j = 1, ..., 90 loci:

$$L_i^{(i)} =$$
 genotype of marker *j* for plant *i*.

E.g., in heterozygous diploid population:

$$L_{j} = \begin{cases} 0 & \text{homozygous AA from parent 1} \\ 2 & \text{heterozygous AB} \\ 4 & \text{homozygous BB from parent 2} \\ NA & \text{if genotype is missing} \end{cases}$$

PS. Although in RIL genotype is typically homozygous (as in figures), we do have some heterozygous genotypes too.



"Eliminating confounders" idea (1)

Normally, only neighbouring loci are informative about a locus:

$$L_{1} = \theta_{12}L_{2} + \epsilon_{1}$$

$$L_{2} = \theta_{21}L_{1} + \theta_{23}L_{3} + \epsilon_{2}$$

$$L_{3} = \theta_{32}L_{2} + \theta_{34}L_{4} + \epsilon_{3}$$



Genetic laws responsible:

- Genetic linkage
- Independent assortment of chromosomes



Svizzera

IF non-neighbouring locus is predictive, then possibly sign of epistasis:



$$... = ...$$



della Svizzera italiana

Inference idea

The non-zeros in the matrix $\boldsymbol{\Theta}$

$$\Theta = \begin{pmatrix} -- & \theta_{12} & \theta_{13} & \dots \\ \theta_{21} & -- & \theta_{23} & \dots \\ \theta_{31} & \theta_{32} & -- & \dots \\ \dots & \dots & \dots & \dots \end{pmatrix}$$

determine the underlying linkage & epistasis structure.

(A bit) more precisely ... Locus data $L^{(i)} = (L_1^{(i)}, \dots, L_{90}^{(i)})$ can be modelled as $L^{(i)} \sim N(\mu, \Theta^{-1}),$

for $i = 1, \ldots, 367 F_8$ Arabidopsis plants.

NOTE:

$$\mathsf{Cor}(\mathit{L}_1, \mathit{L}_4| \; \mathsf{rest}) = rac{- heta_{14}}{\sqrt{ heta_{11} heta_{44}}}.$$

della Svizzera

Gaussian Copula

Assume latent variable $Z \sim N(0, \Theta^{-1})$ in the following way:

Relationship between latent and observed variables



NOTE: Observations i = 1, ..., 367, Markers j = 1, ..., 90. Ernst C. Wit epistasis, genetic maps

•
$$L_j = F_j^{-1}(\Phi(Z_j))$$

epistasis, genetic maps and microbial networks

15 <u>/ 41</u>

Penalized graphical Gaussian models

To achieve "sparse" graph structure $\widehat{\Theta},$ introduce an ℓ_1 penalty:

$$\widehat{\Theta} = ext{argmax } \log | \mathbf{\Theta} | - \mathsf{Tr}(S\mathbf{\Theta})$$
 $ext{subject to} \sum_{i
eq j} | \mathcal{K}_{ij} | \leq lpha$

This problem, a.k.a. Graphical Lasso, has been considered by

- Meinshausen and Buehlman (2006): implemented in huge,
- Banerjee (2007) and
- Friedman et al. (2008): implemented in glasso.

Or a Bayesian version:

• Mohammadi and Wit (2015): R-package BDgraph.



L_1 penalty = sparsity

Maximize
$$\ell(\beta) = -(y - X\beta)^t (y - X\beta)$$
, subject to $||\beta|||_1 \le 3$.



Ernst C. Wit

17 / 41

Sparse inference of GCGM

Likelihood:

$$\ell_{Y}(\Theta) \approx \frac{n}{2} \log |\Theta| - \frac{1}{2} \sum_{i=1}^{n} \int_{c_{p-1}^{(i)}}^{c_{p}^{(i)}} \dots \int_{c_{1}^{(i)}}^{c_{1}^{(i)}} z^{(i)T} \Theta z^{(i)} dz_{1} \dots dz_{p}$$

Penalized EM algorithm

E-step: Compute $Q_{\lambda}(\Theta|\Theta^*) = E[\ell_{Y,Z}^{p}(\Theta)|Y,\Theta^{(m)}]$

$$Q_{\lambda}(\Theta|\Theta^{(m)}) = -\frac{np}{2}\log 2\pi + \frac{n}{2}\{\log|\Theta| - tr\{\left[\frac{1}{n}\sum_{i=1}^{n}E(Z^{(i)}Z^{(i)T}|Y^{(i)}\Theta^{(m)})\right]\Theta\} - \lambda||\Theta||_{1}$$

M-step: $\widehat{\Theta} = \arg_{\Theta} \max Q(\Theta | \Theta^*)$ subject to $||\Theta||_1 \leq \lambda$.



E-step: Approximating conditional expectations fast

$$E(z_k^{(i)} z_l^{(i)^T} \mid y^{(i)}, \widehat{}) \approx \begin{cases} E(z_k^{(i)} \mid y^{(i)}, \widehat{}) E(z_l^{(i)} \mid y^{(i)}, \widehat{}) & \text{if } 1 \le k \ne l \le p \\ E(z_k^{(i)^2} \mid y^{(i)}, \widehat{}) & \text{if } k = l \end{cases}$$

$$E(z_{j}^{(i)} \mid y^{(i)}; \widehat{,}, \widehat{C}) \approx \widehat{}_{j,-j} \widehat{}_{-j,-j}^{\circ-1} E(z_{-j}^{(i)^{T}} \mid y^{(i)}; \widehat{,}, \widehat{C}) + \frac{\phi(\widehat{\delta}_{j,y_{j}^{(i)}-1}^{(i)} - \phi(\widetilde{\delta}_{j,y_{j}^{(i)}}^{(i)})}{\Phi(\widetilde{\delta}_{j,y_{j}^{(i)}}^{(i)}) - \Phi(\widetilde{\delta}_{j,y_{j}^{(i)}-1}^{(i)})} \widetilde{\sigma}_{j}^{(i)}$$

$$\begin{split} E(z_{j}^{(i)^{2}} \mid y^{(i)}; \widehat{\widehat{\ }}, \widehat{C}) &\approx \widehat{\widehat{\ }}_{j,-j} \widehat{\widehat{\ }}_{-j,-j}^{-1} E(z_{-j}^{(i)^{T}} z_{-j}^{(i)} \mid y^{(i)}; \widehat{\widehat{\ }}, \widehat{C}) \widehat{\widehat{\ }}_{-j,-j}^{-1} \widehat{\widehat{\ }}_{j,-j}^{T} + \widetilde{\sigma}_{j}^{(i)^{2}} \\ &+ \ldots + \frac{\delta_{j,y_{j}^{(i)}-1}^{(i)} \phi(\widetilde{\delta}_{j,y_{j}^{(i)}-1}^{(i)}) - \widetilde{\delta}_{j,y_{j}^{(i)}}^{(i)} \phi(\widetilde{\delta}_{j,y_{j}^{(i)}}^{(i)})}{\Phi(\widetilde{\delta}_{j,y_{j}^{(i)}}^{(i)}) - \Phi(\widetilde{\delta}_{j,y_{j}^{(i)}-1}^{(i)})} \widetilde{\sigma}_{j}^{(i)^{2}} \end{split}$$

Take-away message: all explicit, so

M-step: Large scale maximization via SQUIC

The orginal QUIC algorithm considers an approximation

$$Q(\Theta + \Delta | \Theta^*) \approx -tr((S - W)\Delta) - \frac{1}{2}tr(W\Delta W\Delta) - \lambda ||\Theta + \Delta ||_1,$$

for one-dimensional maximization steps $\Delta = c(e_i e'_j + e_j e'_i)$. Bottlenecks of QUIC algorithm:

- dense empirical covariance matrix S is reference for each S_{ij}.
- $\Theta + \Delta$ has to be checked for positive-definiteness.
- $W = \theta^{-1}$ is required.

We now work with Olaf Schenk to extend apply SQUIC method:

- method is fast
- on laptop can deal with 100,000 variables (on HPC up to 10 million)



$$eBIC_{\gamma}(\lambda) = -2Q(\widehat{\Theta}_{\lambda}|\widehat{\Theta}_{\lambda}^{(m)}) + 2H(\widehat{\Theta}_{\lambda}|\Theta_{\lambda}^{(m)}) + df(\widehat{\Theta}_{\lambda})$$

where

$$\begin{array}{lll} H(\widehat{\Theta}_{\lambda}|\widehat{\Theta}_{\lambda}^{(m)}) &=& E[logL_{Z|z\in\mathcal{D}}(\widehat{\Theta}_{\lambda})|z\in\mathcal{D};\widehat{\Theta}_{\lambda}^{(m)}]\\ df(\widehat{\Theta}_{\lambda}) &=& (\log n+4\gamma\log p)d\\ d &=& \sum_{1\leq k< l\leq p} l(\widehat{\Theta}_{\lambda}\neq 0) \ \text{and} \ \gamma\in[0,1]. \end{array}$$

Alternatively, a CV based estimator is given by

$$df(\widehat{\Theta}_{\lambda}) = 2 \frac{\sum_{i=1}^{n} \mathsf{c}[(\widehat{K}_{\lambda}^{-1} - S_{i}) \circ I_{\lambda}]^{\top} (\widehat{K}_{\lambda} \otimes \widehat{K}_{\lambda}) \mathsf{c}[(S - S_{i}) \circ I_{\lambda}]}{(n-1)},$$

• $S_k = x_k x_k^T$, • $I_{\lambda} = 1 * (\widehat{K}_{\lambda} != 0)$, an indicator matrix.

della Svizzera italiana

Arabidopsis thaliana experiment

A RIL cross between two A.thaliana lines

- Columbia (Col-0) and Cape Verde Island (Cvi-0)
- p = 90 SNP markers, n = 367 individuals
- Heterozygous population: $Y_j^{(i)} \in \{0, 1, 2\}$
- Contains missing genotypes
- Its genome has 5 chromosomes





Epistatic selection in A.thaliana RIL



Epistatic selection in A.thaliana RIL



Genetic inbreeding experiment in Maize

p = 1106 SNP markers ; n = 193



Existence of such trans-chromosomal edges reveals "aberrant" marker-marker associations that are due to epistatic selection.



2. Genetic map construction



From this picture, it is clear that...

epistasis is a little boat in the ocean of genetic linkage.



... but sampling design is important!

Consider a two parents with genotypes of two loci very close together:

 $AB \times Ab$ and $AB \times aB$,

So all off-spring will have genotype

 $A. \times .B$

In fact, each of the following will have probability 1/4:

 $Ab \times aB$, $AB \times aB$, $Ab \times AB$, $AB \times AB$.

Conclusion

Even though loci are really close together, the off-spring genotype contains NO information of this fact.

3 different population types

Consider

- genotype data Y
- on 10 markers
- of a polyploid species
- from 3 different population types



(a) homozygous, (b) inbred, (c) outcrossing (outbred) populations

Mapping algorithm



| | ۰. |
|---|----|
| n | |
| | |
| | |

| | | M1 | M2 | M3 | M4 | M5 | M6 | M7 | M8 | M9 | M10 | M11 | M500 |
|----------|---------|----|----|----|----|----|----|----|----|----|-----|-----|----------|
| | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 1 | 2 |
| | 2 | 0 | 0 | 0 | 0 | 0 | - | 2 | 2 | 2 | 1 | - | 2 |
| → | 3 | 2 | 2 | 2 | 2 | 1 | 1 | 0 | 0 | - | 0 | 0 | - |
| <i>´</i> | 4 | 0 | -* | 0 | 1 | 1 | 1 | 1 | - | 0 | 0 | 0 | 0 |
| | 1 | | | | | | | | | | | | |
| | 200 | 2 | 2 | 2 | 1 | 1 | 1 | 1 | 2 | 2 | 2 | 0 | 0 |
| | A A _ 0 | A | | | | | | | | | | | |

* Missing genotype







It beats the competition





Mapping diploid Barley genome

Estimated number of linkage groups (LGs) for OWB data set

| | Estimated # | Size of the LGs |
|---------|-------------|--------------------------------------|
| netgwas | 7 | 140, 199, 211, 187, 236, 182, 173 |
| MSTMap | 1 | 1328 |

Comparison of ordering accuracy between netgwas and MSTMap. In this Table assumed MSTMAP has estimated correctly the number of LGs in the OWB data set.

| Linkage Group | Sensitivity Score | | |
|---------------|-------------------|--------|--|
| (LG) | netgwas | MSTMap | |
| 1 | 0.86 | 0.96 | |
| 2 | 0.78 | 0. 52 | |
| 3 | 0.78 | 0.92 | |
| 4 | 0.74 | 0.49 | |
| 5 | 0.71 | 0.38 | |
| 6 | 0.61 | 0.50 | |
| 7 | 0.70 | 0.61 | |
| Average | 0.74 | 0.63 | |





3. Microbial interaction networks



symbiotic relationship

Microorganisms play a central role in many biological processes.



Ernst C. Wit

34 / 41

Metagenomic data from 16S rRNA sequencing



The study collected microbiomes of healthy individual (Methe et al. 2012)



We focus on 306 most prevalent bacteria (OTUs) in 360 stool samples.



Sequence read data

For each stool sample i = 1, ..., 360, we measure for OTU j = 1, ..., 306:

$$O_j^{(i)} =$$
 read count of OTU *j* for stool sample *i*.

It is well-known that the read count $O_i^{(i)}$:

- depends on sequencing depth of sample *i*,
- depends particular peculiarities of OTU j
- has a lot of zero counts
- is overdispersed

We control for all these nuisance effects,

$$O_j^{(i)} \sim ext{Zero-inflated negative binomial}$$

where
 $\log E_{\beta} O_j^{(i)} = ext{seq. depth}_i + \beta_0 + \beta_j + \dots$

Microbial network as Copula Graphical Model

But most importantly, the OTUs depend on each other:



Copula graphical model for interactions:

- **1** Precision matrix Θ is associated with interaction graph;
- 2 Latent Gaussians are generated via $Z^{(i)} \sim N(0, \Theta^{-1})$.
- Observed count $O_j^{(i)}$ is generated by transforming $Z_j^{(i)}$ via ZINB.



Svizzera

MAGMA: inferring microbial interactions



Compared to other methods, our method MAGMA

- Finds fewer spurious links
- Is able to insert real biology in the read count distributions
- Accounts for variability in sequencing depth between samples
- resulting in a cleaner interpretation of the results



Wrapping up

Conclusions

- Confounding is the real enemy of causality.
- Networks account for confounding
 - Detection of epistatic selection
 - Construct genetic map in in any polyploid species.
 - Reconstruct microbial interaction networks from sequence read counts.

Software

netgwas: Our multi-core R package is available on CRAN **rMAGMA**: Our microbial network R package is available on github.





