

Statistical Analysis of High-Throughput Metabolomics Data

Differential Signatures Pertaining to Alzheimer's Disease

Carel F.W. Peeters

Dept. of Epidemiology & Biostatistics
VU University medical center, Amsterdam
cf.peeters@vumc.nl

Alzheimercyclus
VUmc Alzheimercentrum
VU University Medical Center, Amsterdam
October 21, 2016

Contributors



Francien de Leeuw
Alzheimer Center, VUMC



Charlotte Teunissen
Neurology Laboratory, Dept. of Clinical Chemistry, VUMC



Wiesje van der Flier
Alzheimer Center, VUMC
Dept. of Epimiology & Biostatistics, VUMC



Wessel N. van Wieringen

Dept. of Epimiology & Biostatistics, VUMC
Dept. of Mathematics, VU University Amsterdam



Anders E. Bilgrau

Novo Nordisk
Dept. of Mathematical Sciences, Aalborg University



Mark A. van de Wiel

Dept. of Epimiology & Biostatistics, VUMC
Dept. of Mathematics, VU University Amsterdam



Thomas Hankemeier
Division for Analytical Biosciences, Leiden University



Herman van Vlijmen
Dept. of Medicinal Chemistry, Leiden University
Molecular Sciences div., Janssen Pharmaceutica



Cornelia van Duijn
Netherlands Institute for Health Sciences
Dept. of Genetic Epidemiology, Erasmus MC

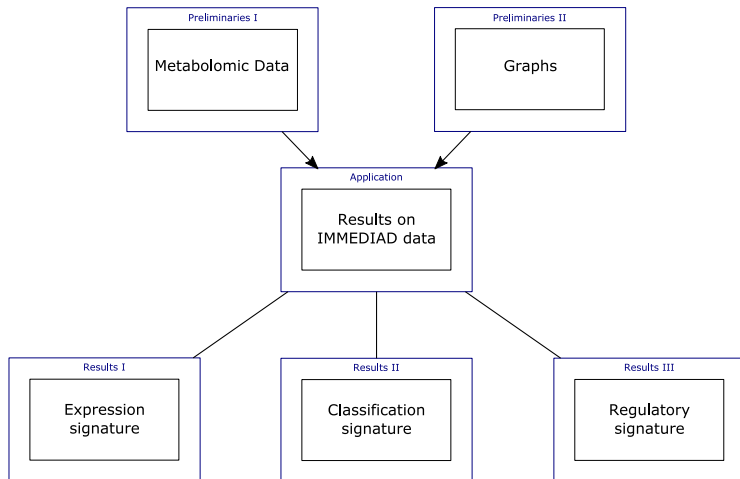
Methodological Developments

European Community Seventh Framework Programme (FP7): grant agreement No. FP7-269553

Data and Analyzes

Janssen Pharmaceutica Stellar funded project (IMMEDIAD): Stellar Neurodegeneration Collaboration Project, Call 2, No. 3

Outline



Omics and omics data

-ome

A totality of some (molecular biological) sort

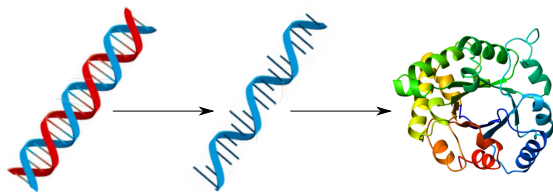
-omics

Collective quantification of some pool of molecular molecules

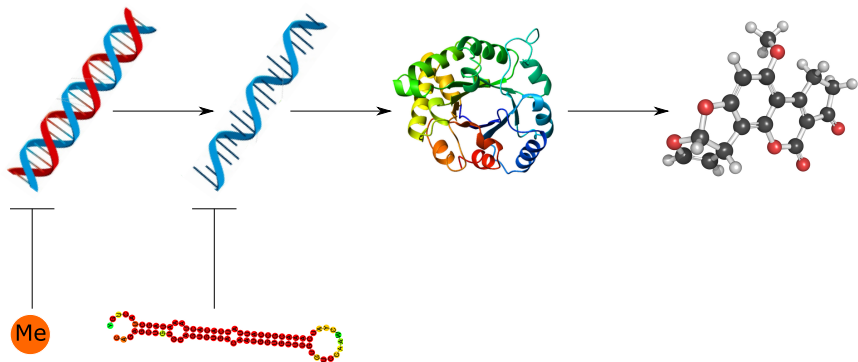
Genomics

The omics of the genome (of some organism)

The omic cascade



The omic cascade



Metabolite quantification

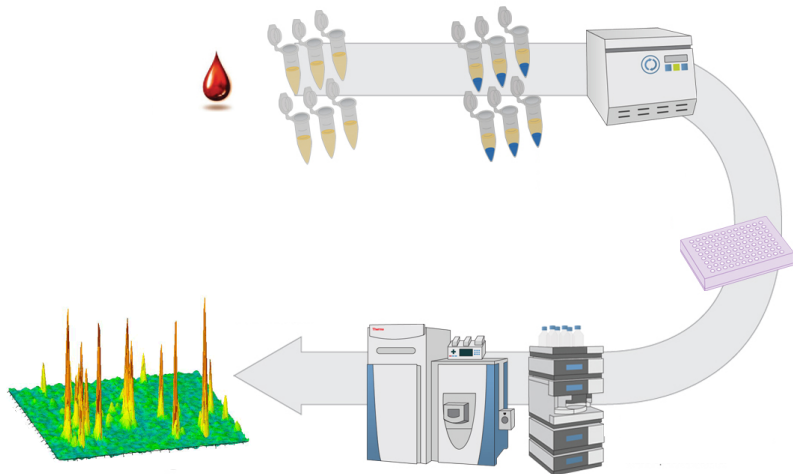


Illustration adapted from: <http://planetorbitrap.com/untargeted-metabolomics#.Vzw6yfmlRaQ> &

<http://metabolomicsplatform.com/metabolomics-overview/>

Challenge: Dimensionality metabolomic data

		Variables				
		1	2	3	p
Observations	1	Y_{11}	Y_{12}	Y_{13}	Y_{1p}
	2	Y_{21}	Y_{22}	Y_{23}	Y_{2p}
	3	Y_{31}	Y_{32}	Y_{33}	Y_{3p}
	4	Y_{41}	Y_{42}	Y_{43}	Y_{4p}
	5	Y_{51}	Y_{52}	Y_{53}	Y_{5p}

n	Y_{n1}	Y_{n2}	Y_{n3}	Y_{np}	

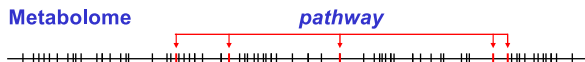
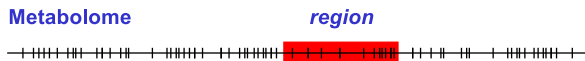
Regular data: $n > p$

		Variables (features)						
		1	2	3	4	5	p
Observations	1	Y_{11}	Y_{12}	Y_{13}	Y_{14}	Y_{15}	Y_{1p}
	2	Y_{21}	Y_{22}	Y_{23}	Y_{24}	Y_{25}	Y_{2p}
	3	Y_{31}	Y_{32}	Y_{33}	Y_{34}	Y_{35}	Y_{3p}

	n	Y_{n1}	Y_{n2}	Y_{n3}	Y_{n4}	Y_{n5}	Y_{np}

Metabolomic data: $p > n$ or $p \gg n$

Unit of analysis



Graphs

Representation

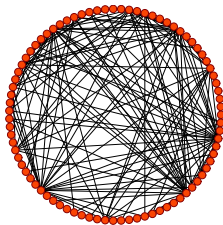
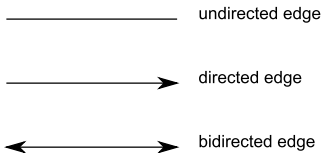
Pathways are represented by a *graph* (or *network*)

Vertices

○ *Node* or *vertex* represents molecular feature

Edges

Edge or *arrow* represents some functional relation



Nested models

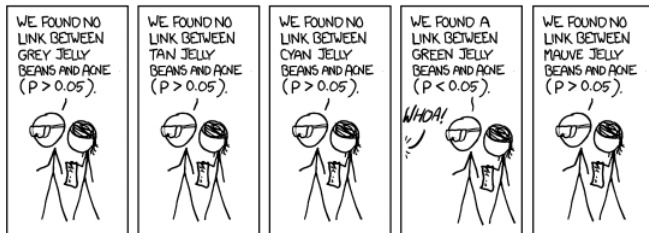


Nested models

$$\text{Metabolite expression} = \beta_0 + \beta_1 \text{SBP} + \beta_2 \text{ApoE} + \epsilon$$

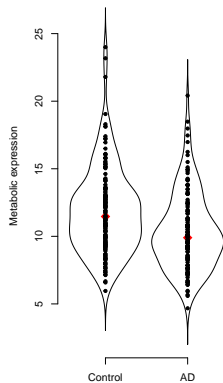
$$\text{Metabolite expression} = \beta_0 + \beta_1 \text{SBP} + \beta_2 \text{ApoE} + \beta_3 \text{AD} + \epsilon$$

Multiple testing

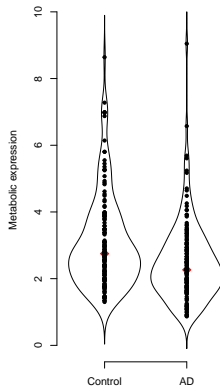


Results

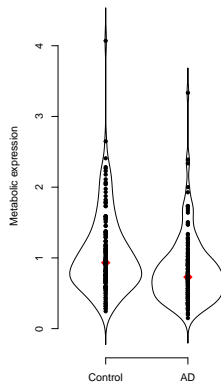
Metabolite: Am.L.Tyrosine



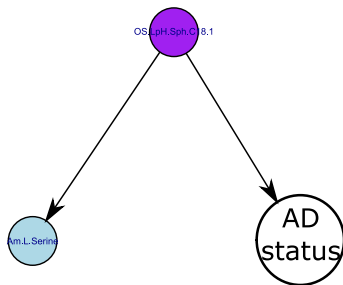
Metabolite: Lip.TG.54.6.



Metabolite: Lip.TG.56.8.



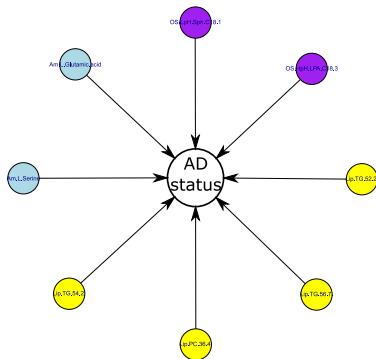
Not the complete picture



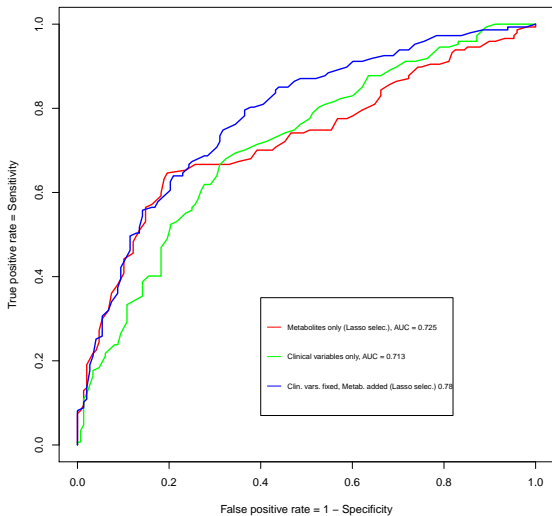
The prediction model

Metabolome

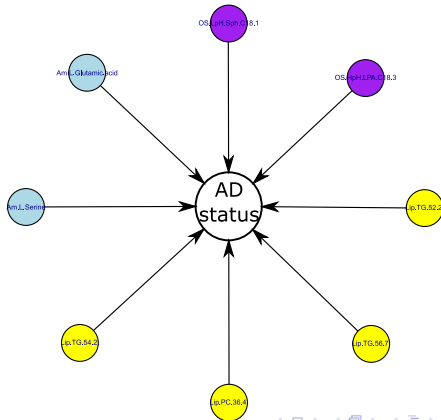
region



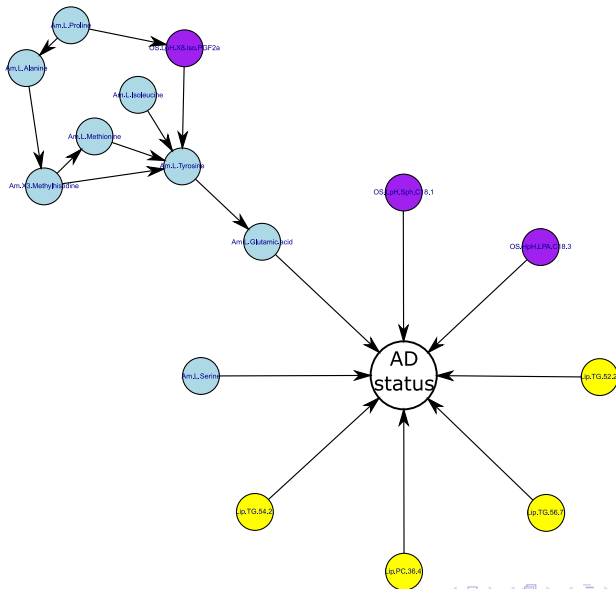
Results



Not the complete picture



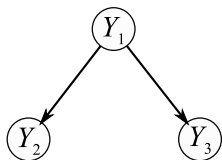
Not the complete picture



Graphical modeling

Metabolome

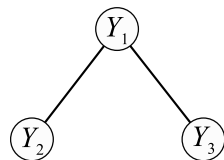
pathway



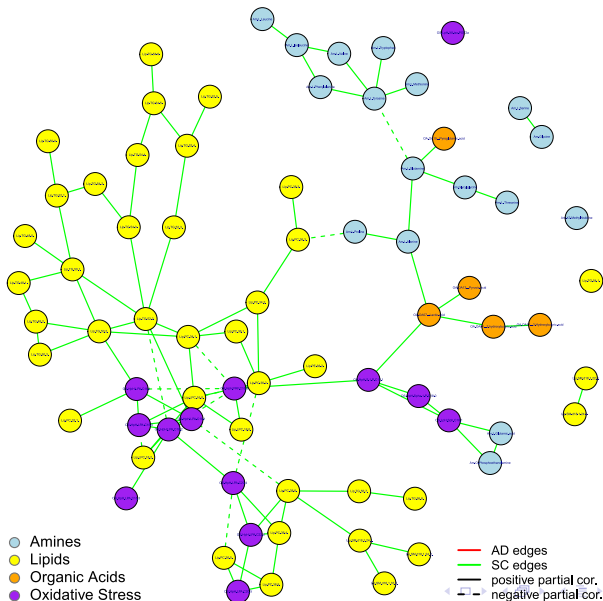
$$\text{cor}(Y_1, Y_2 | Y_3) \neq 0$$

$$\text{cor}(Y_1, Y_3 | Y_2) \neq 0$$

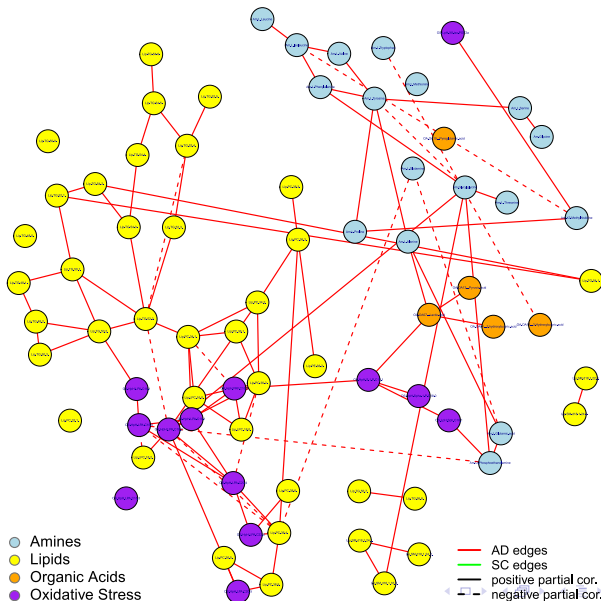
$$\text{cor}(Y_2, Y_3 | Y_1) = 0$$



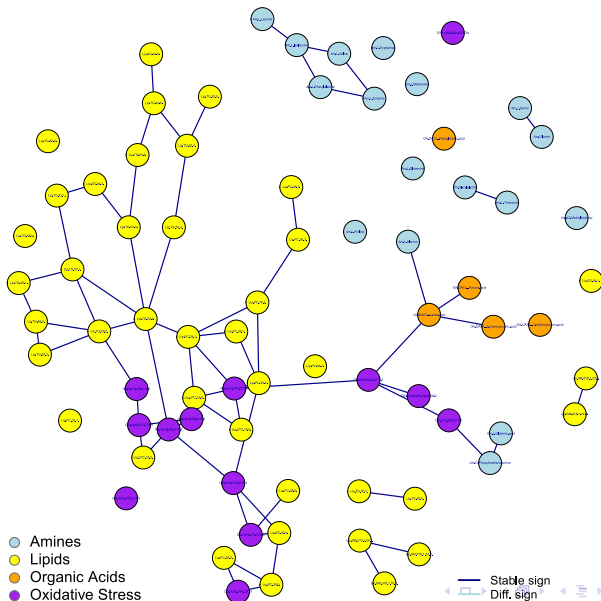
Control connections



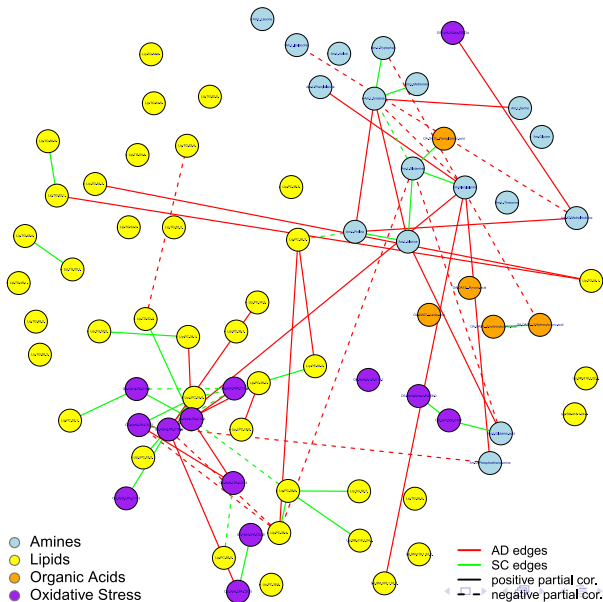
AD connections



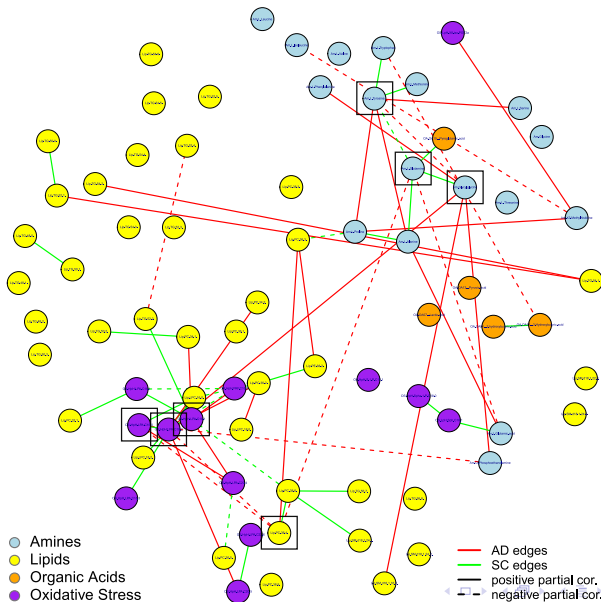
Shared connections



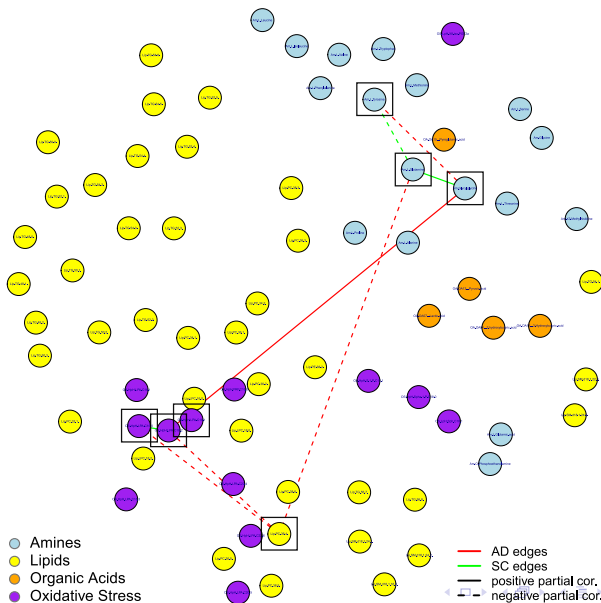
Differential connections



Differential connections



Differential connections



Concluding

Messages

- 3 metabolic signatures
- Each signature comes with different (complementary) information
- Tyrosine influential in all signatures

Now what?

- Follow-up studies
- Validation of findings
- Translational efforts

Manual/Download

- Peeters, C.F.W., Bilgrau, A.E., & van Wieringen, W.N. (2016). “rags2ridges: Ridge Estimation of Precision Matrices from High-Dimensional Data”. R package, version 2.1.1 URL: <https://cran.r-project.org/package=rags2ridges>.

Theory/Methodology

- Peeters*, C.F.W., Bilgrau*, A.E., Eriksen, P.S., Boegsted, M., & van Wieringen, W.N. (2015). “Targeted Fused Ridge Estimation of Inverse Covariance Matrices from Multiple High-Dimensional Data Classes”. [arXiv:1509.07982v1](https://arxiv.org/abs/1509.07982v1) [stat.ME].
- Peeters, C.F.W., van Wieringen, W.N., & van de Wiel, M.A. (in preparation). “Directed Cyclic Mixed Graph Modeling for High-Dimensional Genomic Data Integration”.
- van Wieringen, W.N. & Peeters, C.F.W. (2016). “Ridge Estimation of Inverse Covariance Matrices from High-Dimensional Data”. *Computational Statistics & Data Analysis*, 103: 284-303. [arXiv:1403.0904v3](https://arxiv.org/abs/1403.0904v3) [stat.ME].

Software

- Peeters, C.F.W., van de Wiel, M.A., & van Wieringen, W.N. (2016) “The Spectral Condition Number Plot for Regularization Parameter Determination”. [arXiv:1608.04123v1](https://arxiv.org/abs/1608.04123v1) [stat.CO].
- van Wieringen, W.N. & Peeters, C.F.W. (2015). “Application of a New Ridge Estimator of the Inverse Covariance Matrix to the Reconstruction of Gene-Gene Interaction Networks”. In: di Serio, C., Lio, P., Nonis, A., and Tagliaferri, R. (Eds.) ‘Computational Intelligence Methods for Bioinformatics and Biostatistics’. *Lecture Notes in Computer Science*, vol. 8623. Springer, pp. 170–179.

Table: List of clinical confounders/variables

Anthropometric:

Age

Sex

ApoE ϵ 4 allele status (at least one allele ϵ 4 yes,no)

Systolic blood pressure

Diastolic blood pressure

Height

Weight

Smoking (yes, no, quit)

Alcohol (yes, no)

Comorbidities (binary):

Hypertension

Diabetes Mellitus

Hypercholesterolemia

Medication use (binary):

Cholesterol lowering medications

Antidepressant medications

Antiplatelet medications

Antidiabetic medication

Antiepileptic medication

Antiparkinson medication

Antipsychotics

Table: Metabolites selected for prediction model

Am.X1.Methylhistidine
 Am.X3.Methylhistidine
 Am.Citrulline
 Am.Cysteine
 Am.gamma.L.glutamyl.L.alanine
 Am.Histamine
 Am.L.carnosine
 Am.L.Tyrosine
 Am.Methyldopa
 Am.O.acetyl.L.serine
 Am.Putrescine
 OA.OA14Methylmalonic.acid
 OA.OA173.Hydroxybutyric.acid
 OA.OA263.Hydroxyisobutyric.acid
 OA.OA273.hydroxyisovaleric.acid
 OA.OA28Glyceric.acid
 Lip.TG.51.3.
 Lip.TG.56.8.
 Lip.TG.58.10.
 Lip.LPC.18.1.
 Lip.LPC.20.3.
 Lip.LPC.20.4.
 Lip.PC.O.34.3.
 Lip.PC.O.36.5.
 Lip.PC.O.38.6.
 Lip.SM.d18.1.18.2.
 Lip.SM.d18.1.20.1.
 Lip.SM.d18.1.23.0.
 OS.LpH.X8.12.iPF2a.IV
 OS.LpH.X..5.iPF2a.VI
 OS.LpH.NO2.aLA
 OS.LpH.NO2.OA
 OS.LpH.PGD2
 OS.LpH.Spha.C18.0
 OS.cLpH.PGA2
 OS.cLpH.X2.3.dinor.8.iso.PGF2a
 OS.HpH.LPA.C14.0

Explaining the inverse

The scalar inverse

- Let a denote a number (excluding 0)
- The inverse is then the number b such that $a \times b = 1$
- Clearly, $b = \frac{1}{a}$

Matrix

A matrix is a generalization of a number, an array of numbers

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1p} \\ a_{21} & a_{22} & \cdots & a_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ a_{p1} & a_{p2} & \cdots & a_{pp} \end{bmatrix}$$

Explaining the inverse

The Matrix Inverse

Consider the matrix \mathbf{A} . Its inverse $\mathbf{B} = \mathbf{A}^{-1}$ is defined such that

$$\mathbf{AB} = \mathbf{I},$$

where

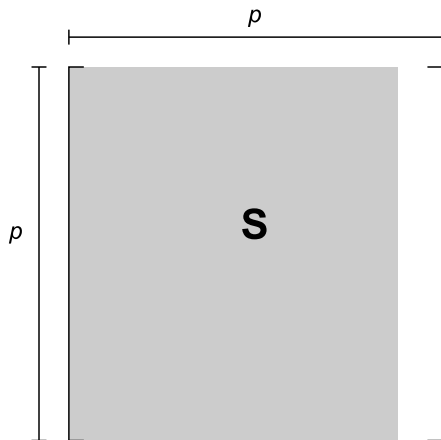
$$\mathbf{I} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}$$

Solution

$$\mathbf{A}^{-1} = \begin{bmatrix} \mathbf{A}_{11}^{-1} + \mathbf{A}_{11}^{-1} \mathbf{A}_{12} \mathbf{Q}^{-1} \mathbf{A}_{21} \mathbf{A}_{11}^{-1} & -\mathbf{A}_{11}^{-1} \mathbf{A}_{12} \mathbf{Q}^{-1} \\ -\mathbf{Q}^{-1} \mathbf{A}_{21} \mathbf{A}_{11}^{-1} & \mathbf{Q}^{-1} \end{bmatrix},$$

with $\mathbf{Q} = \mathbf{A}_{22} - \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12}$ denoting the Schur complement.

Singularity



Ridge estimation

Maximize

$$\underbrace{\ln |\mathbf{\Omega}| - \text{tr}(\mathbf{S}\mathbf{\Omega})}_{\text{log-likelihood}} - \underbrace{\frac{\lambda}{2} \|\mathbf{\Omega} - \mathbf{T}\|_2^2}_{\ell_2\text{-penalty}}$$

- \mathbf{T} denotes a p.d. symmetric target matrix
- $\lambda \in (0, \infty)$ denotes a penalty parameter

Analytic penalized ML estimator

$$\hat{\mathbf{\Omega}}(\lambda) = \left\{ \left[\lambda \mathbf{I}_p + \frac{1}{4}(\mathbf{S} - \lambda \mathbf{T})^2 \right]^{1/2} + \frac{1}{2}(\mathbf{S} - \lambda \mathbf{T}) \right\}^{-1}$$

Properties

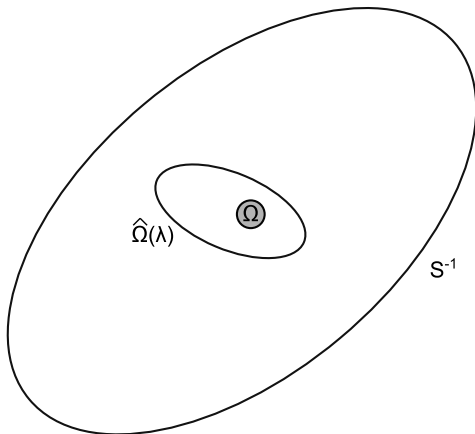
Behavior

- i. $\hat{\Omega}(\lambda) \succ 0$, for all $\lambda \in (0, \infty)$;
- ii. $\lim_{\lambda \rightarrow 0^+} \hat{\Omega}(\lambda) = \mathbf{S}^{-1}$ if $p < n$;
- iii. $\lim_{\lambda \rightarrow \infty} \hat{\Omega}(\lambda) = \mathbf{T}$.

Consistency

- i. $\lim_{n \rightarrow \infty} \mathbb{E} \left[\hat{\Omega}_n(\lambda_n) \right] \longrightarrow \lim_{n \rightarrow \infty} \mathbb{E} \left(\mathbf{S}_n^{-1} \right) = \Omega$;
- ii. $\lim_{n \rightarrow \infty} \mathbb{E} \left(\left\| \hat{\Omega}_n(\lambda_n) - \Omega \right\|_F^2 \right) = 0$.

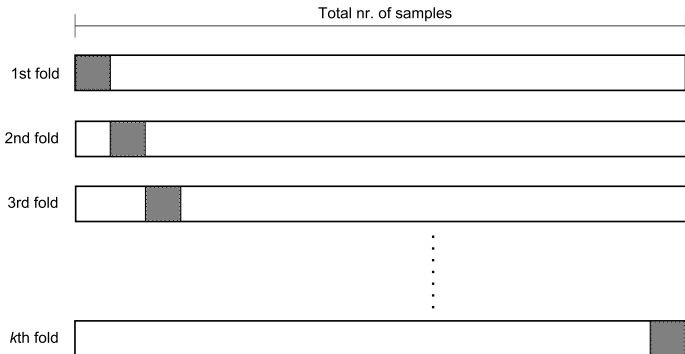
Visual explanation



Choosing the penalty value

K -fold cross-validation (CV)

Single iteration of K -fold CV



Test data



Training data

Choosing the penalty value

K -fold CV score

$$\varphi^K(\lambda) = \sum_{k=1}^K n_k \left\{ -\ln |\hat{\Omega}(\lambda)_{-k}| + \text{tr}[\hat{\Omega}(\lambda)_{-k} \mathbf{S}_k] \right\},$$

n_k is the size of subset k , for $k = 1, \dots, K$ disjoint subsets;

\mathbf{S}_k denotes the sample covariance matrix on k th test set;

$\hat{\Omega}(\lambda)_{-k}$ denotes the estimated regularized precision matrix on k th training set

Highest predictive accuracy

Choose $n_k = 1$, such that $K = n$ (known as leave-one-out CV - LOOCV)

Problem

K -fold CV is computationally demanding for large p and/or large K

Solution

Computationally efficient approximate LOOCV score

Support determination

Scaling

$\hat{\mathbf{P}}(\lambda)$: Regularized precision estimate scaled to partial correlation form

Assume

Nonredundant off-diagonal partial correlation coefficients (indexed by $j < j'$) follow a mixture distribution:

$$f \left\{ [\hat{\mathbf{P}}(\lambda^*)]_{jj'} \right\} = \eta_0 f_0 \left\{ [\hat{\mathbf{P}}(\lambda^*)]_{jj'}; \kappa \right\} + (1 - \eta_0) f_{\mathcal{E}} \left\{ [\hat{\mathbf{P}}(\lambda^*)]_{jj'} \right\}$$

- $\eta_0 \in [0, 1]$ is the mixture weight
- $f_0 \{ \cdot \}$ denotes the distribution of a null-edge
- $f_{\mathcal{E}} \{ \cdot \}$ denotes the distribution of a present edge

Determine

$$P \left(Y_j \neq Y_{j'} \mid [\hat{\mathbf{P}}(\lambda^*)]_{jj'} \right)$$

Situation

Data

- G classes of $(n_g \times p)$ -dimensional data
- Classes defined by data sets and/or (subtypes of) diseases

Assumption

Precision matrices of constituent classes chiefly share the same structure but potentially differ in a number of locations of interest

Desire

Integrative or meta-analytic Gaussian graphical modeling

Targeted fused ridge estimation: General Formulation

Maximize

$$\underbrace{\mathcal{L}(\{\Omega_g\}; \{S_g\})}_{\text{log-likelihood}} - \sum_g \underbrace{\frac{\lambda_{gg}}{2} \|\Omega_g - T_g\|_F^2}_{\text{ridge-penalty}} - \sum_{g_1, g_2} \underbrace{\frac{\lambda_{g_1 g_2}}{4} \|(\Omega_{g_1} - T_{g_1}) - (\Omega_{g_2} - T_{g_2})\|_F^2}_{\text{fusion-penalty}}$$

- T_g indicate class-specific target matrices
- $\lambda_{gg} \in (0, \infty)$ denote class-specific ridge penalty parameters
- $\lambda_{g_1 g_2} \in [0, \infty)$ denote pair-specific fusion penalty parameters, $\lambda_{g_1 g_2} = \lambda_{g_2 g_1}$

Penalty matrix

All penalties can be collected into a non-negative symmetric matrix $\Lambda = [\lambda_{g_1 g_2}]$

Targeted fused ridge estimation

Maximizing argument for class g_0

$$\hat{\Omega}_{g_0}(\Lambda, \{\Omega_g\}_{g \neq g_0}) = \left\{ \left[\bar{\lambda}_{g_0} \mathbf{I}_p + \frac{1}{4} (\bar{\mathbf{S}}_{g_0} - \bar{\lambda}_{g_0} \mathbf{T}_{g_0})^2 \right]^{1/2} + \frac{1}{2} (\bar{\mathbf{S}}_{g_0} - \bar{\lambda}_{g_0} \mathbf{T}_{g_0}) \right\}^{-1},$$

where

$$\bar{\mathbf{S}}_{g_0} = \mathbf{S}_{g_0} - \sum_{g \neq g_0} \frac{\lambda_{gg_0}}{n_{g_0}} (\Omega_g - \mathbf{T}_g), \quad \text{and} \quad \bar{\lambda}_{g_0} = \frac{\sum_g \lambda_{gg_0}}{n_{g_0}}$$

Properties

Behavior

- i. $\hat{\Omega}_g \succ \mathbf{0}$ for all $\lambda_{gg} \in (0, \infty)$;
- ii. $\lim_{\lambda_{gg} \rightarrow 0^+} \hat{\Omega}_g = \mathbf{S}_g^{-1}$ if $\sum_{g' \neq g} \lambda_{gg'} = 0$ and $p \leq n_g$;
- iii. $\lim_{\lambda_{gg} \rightarrow \infty} \hat{\Omega}_g = \mathbf{T}_g$ if $\lambda_{gg'} < \infty$ for all $g' \neq g$;
- iv. $\lim_{\lambda_{g_1 g_2} \rightarrow \infty} (\hat{\Omega}_{g_1} - \mathbf{T}_{g_1}) = \lim_{\lambda_{g_1 g_2} \rightarrow \infty} (\hat{\Omega}_{g_2} - \mathbf{T}_{g_2})$ if $\lambda_{g'_1 g'_2} < \infty$ for all $\{g'_1, g'_2\} \neq \{g_1, g_2\}$.

Block coordinate ascent

- 1: **Input:**
- 2: *Sufficient data:* $(\mathbf{S}_1, n_1), \dots, (\mathbf{S}_G, n_G)$
- 3: *Penalty matrix:* $\mathbf{\Lambda}$
- 4: *Convergence criterion:* $\varepsilon > 0$
- 5: **Output:**
- 6: *Estimates:* $\hat{\mathbf{\Omega}}_1, \dots, \hat{\mathbf{\Omega}}_G$
- 7: **procedure** RIDGEP.FUSED($\mathbf{S}_1, \dots, \mathbf{S}_G, n_1, \dots, n_G, \mathbf{\Lambda}, \varepsilon$)
- 8: *Initialize:* $\hat{\mathbf{\Omega}}_g^{(0)}$ for all g .
- 9: **for** $c = 1, 2, 3, \dots$ **do**
- 10: **for** $g = 1, 2, \dots, G$ **do**
- 11: Update $\hat{\mathbf{\Omega}}_g^{(c)} := \hat{\mathbf{\Omega}}_g(\mathbf{\Lambda}, \hat{\mathbf{\Omega}}_1^{(c)}, \dots, \hat{\mathbf{\Omega}}_{g-1}^{(c)}, \hat{\mathbf{\Omega}}_{g+1}^{(c-1)}, \dots, \hat{\mathbf{\Omega}}_G^{(c-1)})$
- 12: **end for**
- 13: **if** $\max_g \left\{ \frac{\|\hat{\mathbf{\Omega}}_g^{(c)} - \hat{\mathbf{\Omega}}_g^{(c-1)}\|_F^2}{\|\hat{\mathbf{\Omega}}_g^{(c)}\|_F^2} \right\} < \varepsilon$ **then**
- 14: **return** $(\hat{\mathbf{\Omega}}_1^{(c)}, \dots, \hat{\mathbf{\Omega}}_G^{(c)})$
- 15: **end if**
- 16: **end for**
- 17: **end procedure**