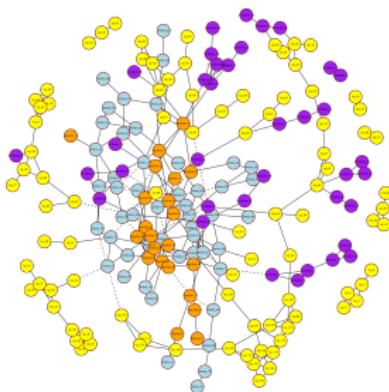


Part I: Single Networks

IBC2022 short course: Network Modeling for High-Dimensional Data

July 10th, 2022, Carel F.W. Peeters

Mathematical & Statistical Methods group — Biometris, Wageningen University & Research



Schedule:

The course is divided a 4 submodules. Each submodule consists of a short lecture and a corresponding hands-on practical.

09:00 – 10:30 Submodule 1: Extracting, visualizing and analyzing single networks

Associated literature: [DOI](#)

R packages used: [rags2ridges](#)

10:30 – 11:00 Break

11:00 – 12:30 Submodule 2: Jointly extracting, visualizing and analyzing multiple networks

Associated literature: [DOI](#)

R packages used: [rags2ridges](#)

12:30 – 13:30 Lunch break

13:30 – 15:00 Submodule 3: Extracting, visualizing and analyzing networks from time-course data

Associated literature: [DOI](#)

R packages used: [ragt2ridges](#)

15:00 – 15:30 Break

15:30 – 17:00 Submodule 4: Miscellanea and extensions

Associated literature:

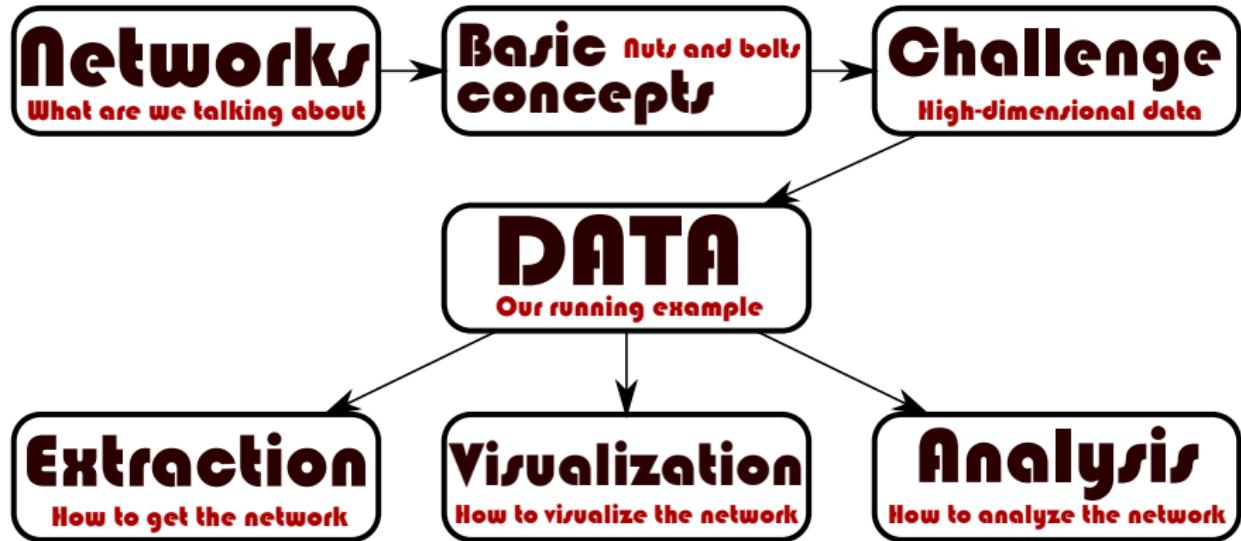
R packages used: [porridge](#)



Carel F.W. Peeters
Wageningen University & Research



Wessel A. van Wieringen
Amsterdam University medical centers
VU University Amsterdam



Part I: Single Networks

IBC2022 short course: Network Modeling for High-Dimensional Data

July 10th, 2022, Carel F.W. Peeters
Mathematical & Statistical Methods group — Biometris, Wageningen University & Research



Preliminaries
Packages
Exercise 1: Get acquainted with the
Exercises
Exercise 2: Find an optimal procedure
Exercise 3: Assess the conditioning
Exercise 4: Extract a network from I.
Visualizations
Exercise 5: Visualize the extracted n.
Analyses
Exercise 6: Find the nodes with the 1.
Exercise 7: Find the nodes with the 1.
Exercise 8: Find and visualize scores
Recap and look ahead
References
Elements

Practical I: Extracting, Visualizing, and Analyzing Single Networks

Carel F.W. Peeters

Mathematical & Statistical Methods group – Biometris
Wageningen University & Research
IBC2022 Short Course on Network Modeling for High-Dimensional Data
carel.peeters.wur.nl

July 10th, 2022

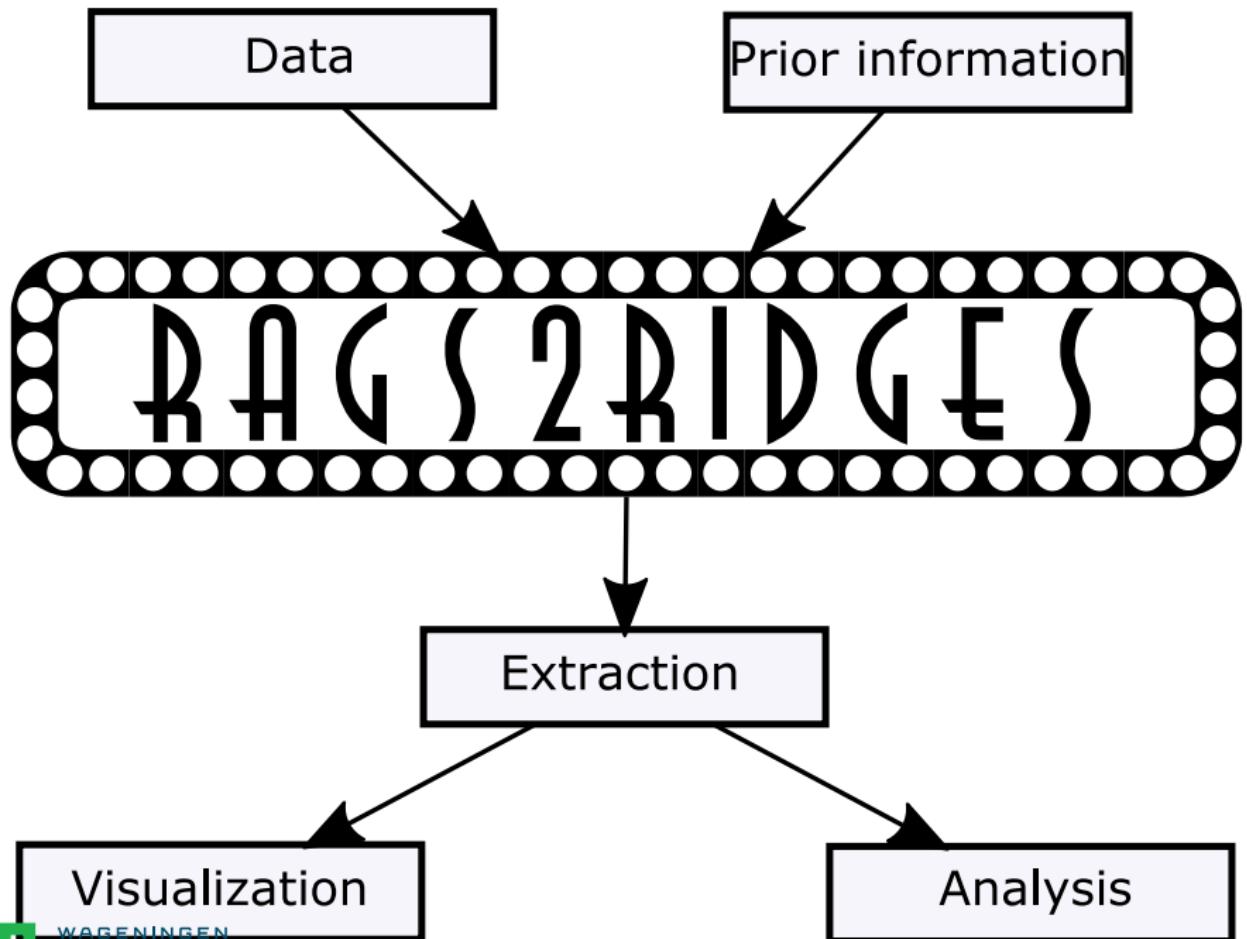
Preliminaries

This document contains a tutorial for Part I of the course on Network Modeling for High-Dimensional Data as given during the 21st International Biometric Conference. It assumes knowledge of R and basic linear algebra. The learning goals for Part I are:

1. To familiarize yourself with basic concepts in network modeling of high-dimensional data;
2. To practice with network modeling of a single high-dimensional data set.

The slides, exercises, and this document form an integrated whole. The slides alternate between the presentation of methodology and corresponding exercises. This document, then, contains possible







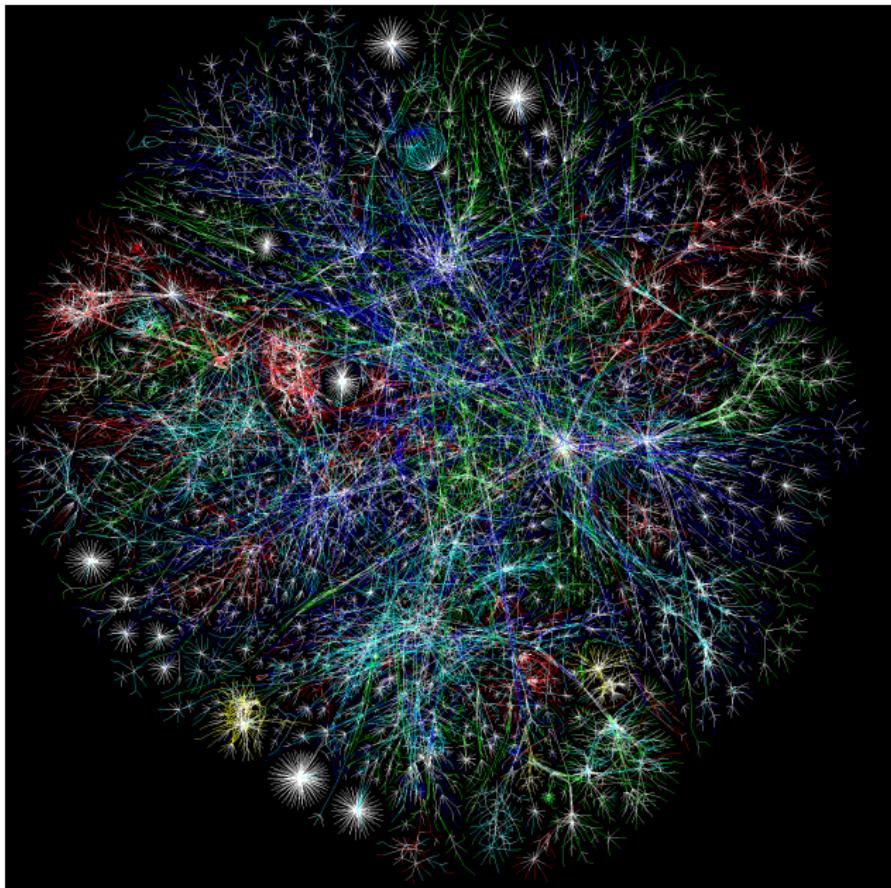
“This has put the ‘IT’ in KITT”
- The Hoff

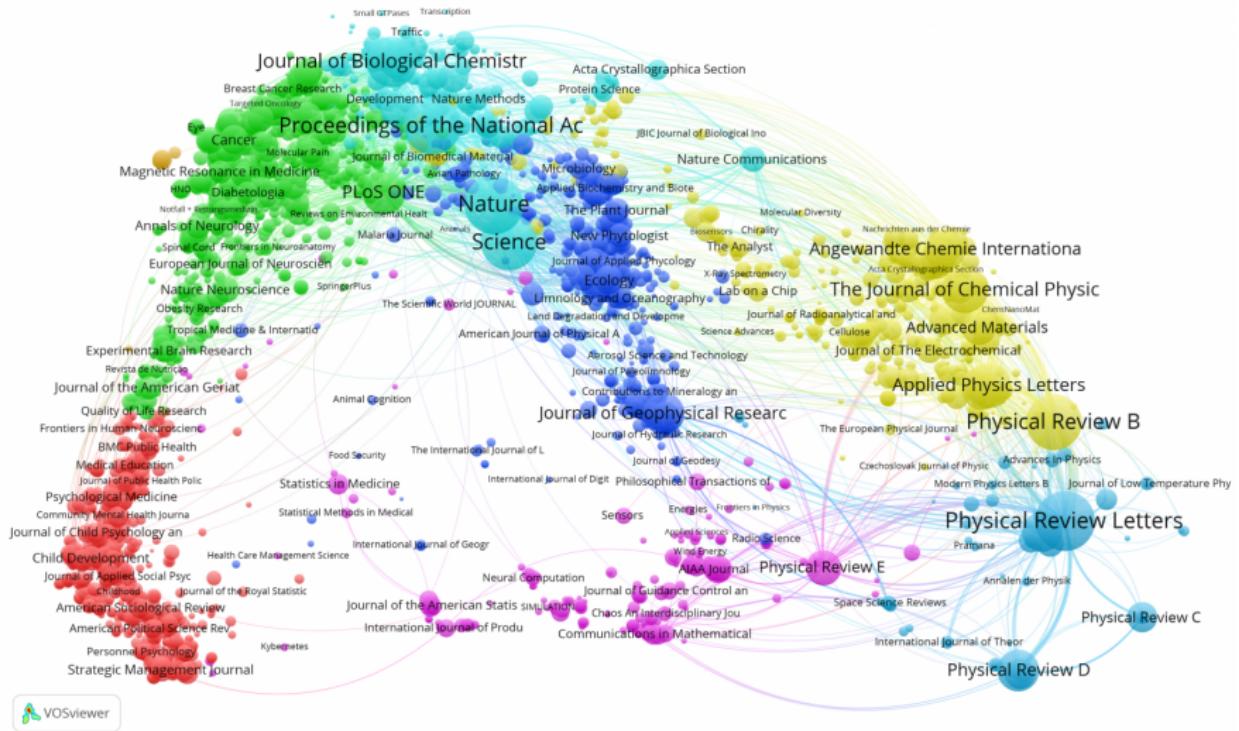
networks

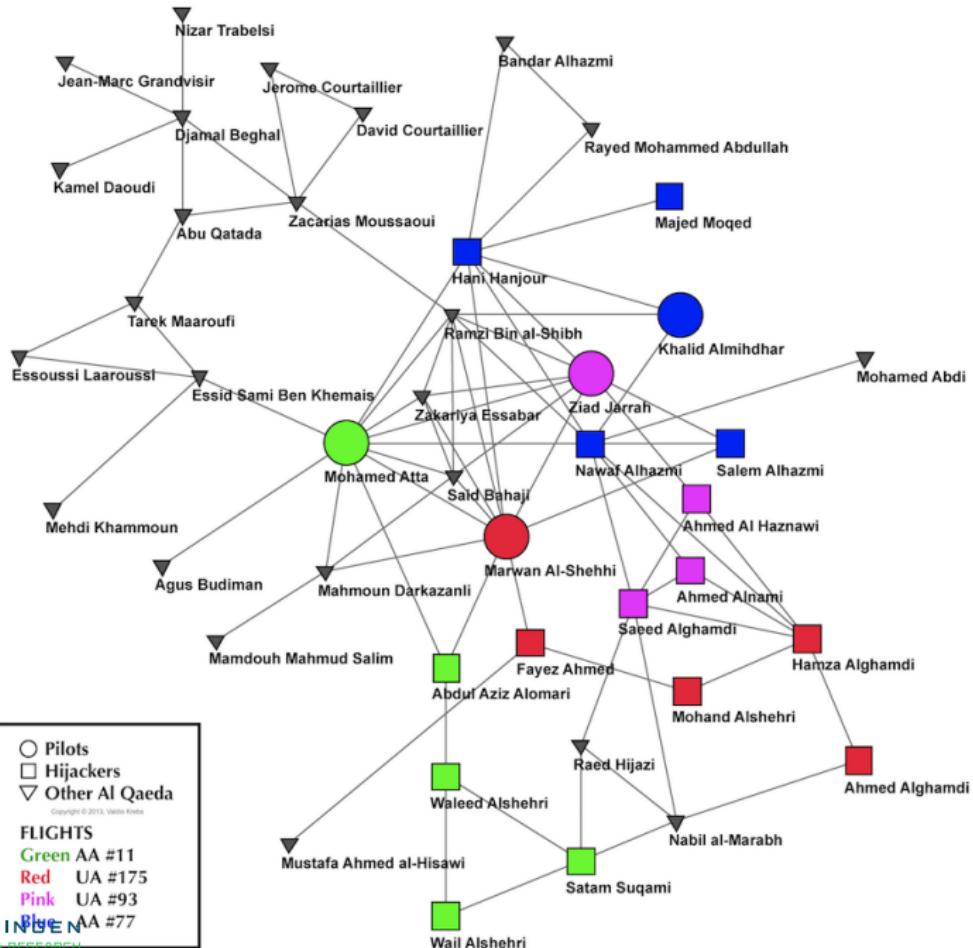
What are we talking about

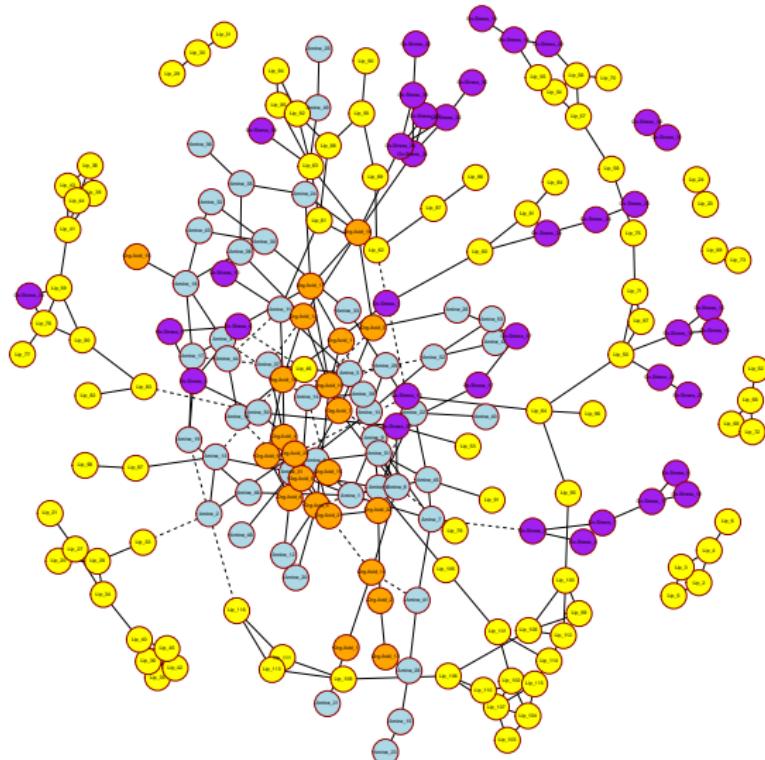












Basic concepts

nuts and bolts

Network

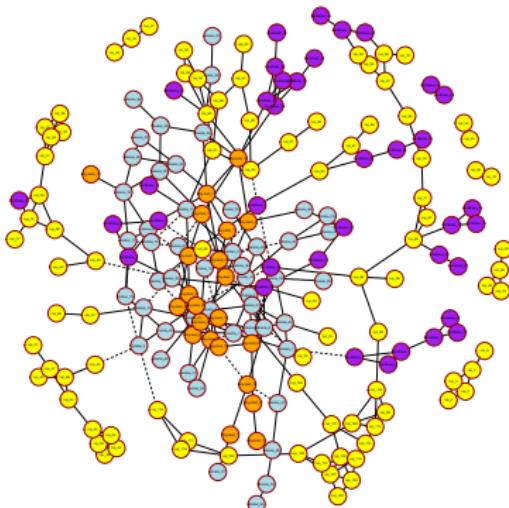
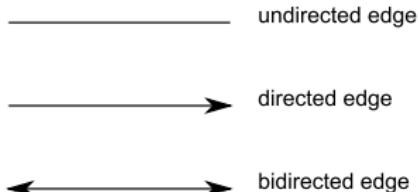
Collection of nodes that may have pairwise relationships. Represented by *graph*

Vertices

○ Node or vertex represents feature

Edges

Edge or connection represents some functional pairwise relation



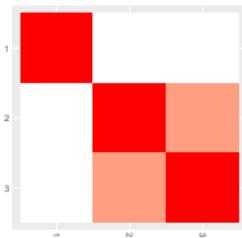
Example

Three variables: Y_1 , Y_2 , and Y_3

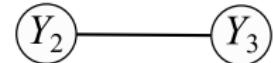
$$\text{cor}(Y_1, Y_2) = 0$$

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

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

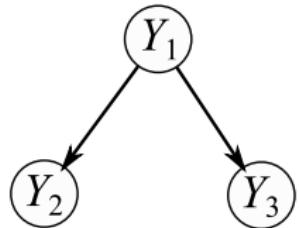


(Y_1)

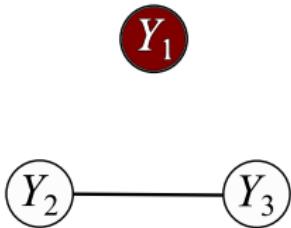


Marginal dependence

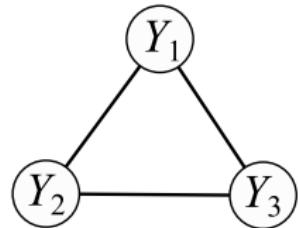
Undirected edge represents marginal dependence



True mechanism



Not observing Y_1 :
Spurious association



Observing Y_1 :
Saturated graph

Partial correlation

Measures degree of association between two random variables when controlling for remaining variables

Conditioned correlation

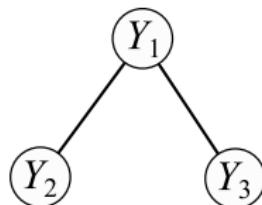
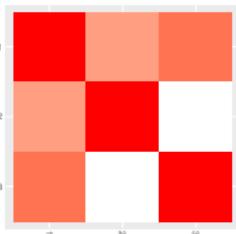
$$\text{cor}(Y_1, Y_2|Y_3)$$

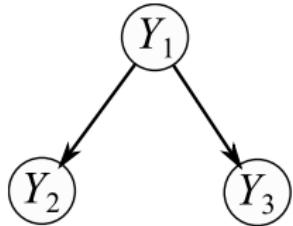
$$\text{cor}(Y_1, Y_3|Y_2)$$

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

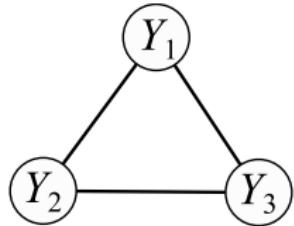
If, e.g., $\text{cor}(Y_2, Y_3|Y_1) = 0$, we say Y_2 and Y_3 are independent given Y_1

$$\begin{aligned}\text{cor}(Y_1, Y_2|Y_3) &\neq 0 \\ \text{cor}(Y_1, Y_3|Y_2) &\neq 0 \\ \text{cor}(Y_2, Y_3|Y_1) &= 0\end{aligned}$$

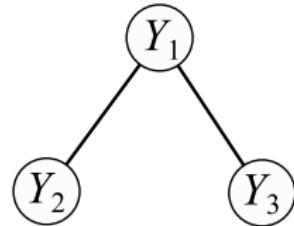




True mechanism



Correlation graph



Conditional
independence graph

Graphical modeling

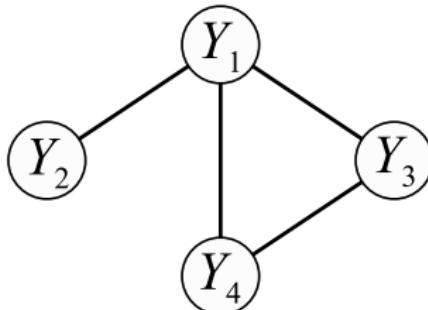
Class of models using graphs to express conditional (in)dependence relations between random variables

Gaussian setting

- Vertices: Correspond to random variables with normal distribution
- Edges: Correspond to the dependence structure
- Say $\mathbf{y} \sim \mathcal{N}_p(\mathbf{0}, \boldsymbol{\Sigma})$, and define $\boldsymbol{\Sigma}^{-1} \equiv \boldsymbol{\Omega}$. Then, for $a, b \in$ vertex set V , $a \neq b$

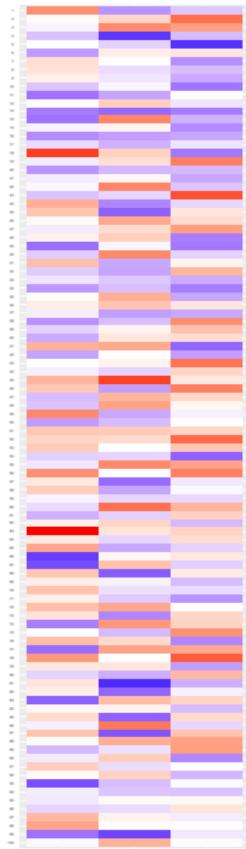
$$-\frac{\omega_{ab}}{\sqrt{\omega_{aa}\omega_{bb}}} = 0 \iff \omega_{ab} = 0 \iff a \perp\!\!\!\perp b | V \setminus \{a, b\} \iff a \not\sim b$$

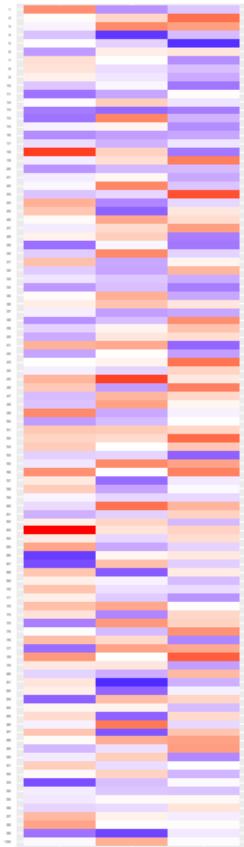
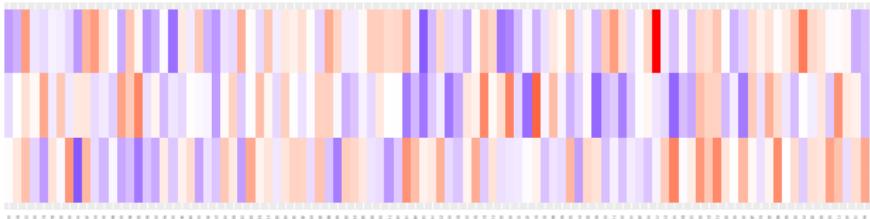
$$\begin{bmatrix} \omega_{11} & \omega_{12} & \omega_{13} & \omega_{14} \\ \omega_{21} & \omega_{22} & 0 & 0 \\ \omega_{31} & 0 & \omega_{33} & \omega_{34} \\ \omega_{41} & 0 & \omega_{43} & \omega_{44} \end{bmatrix}$$

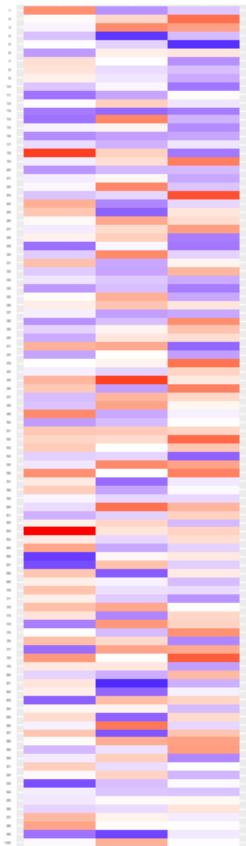
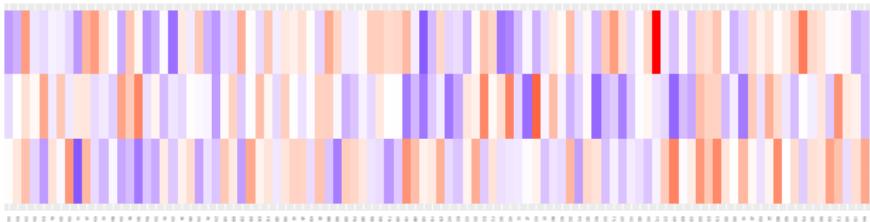


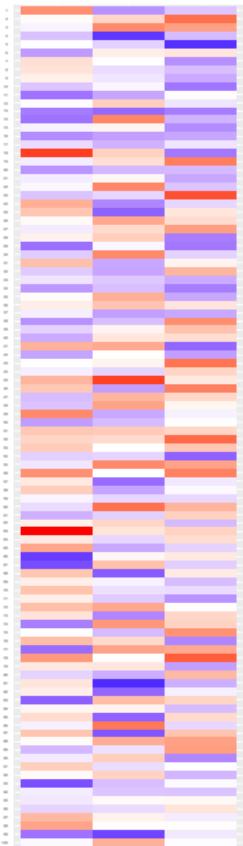
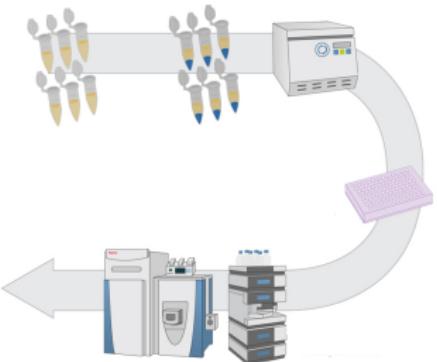
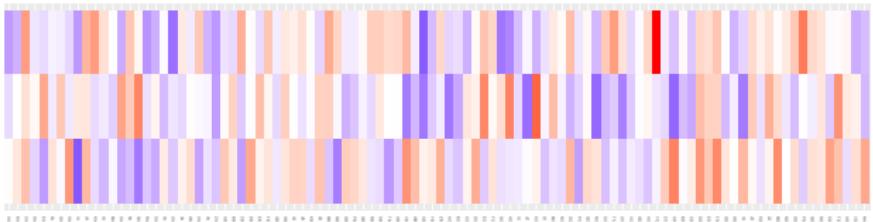
Challenge

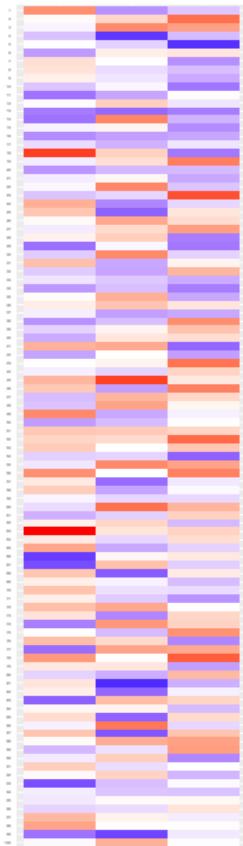
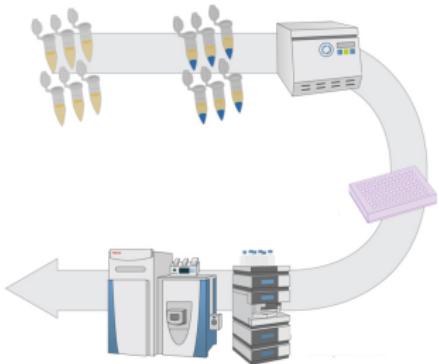
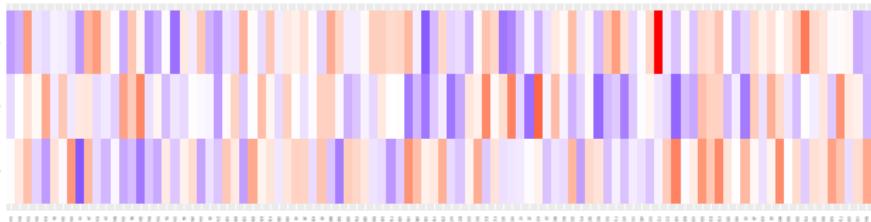
High-dimensional data

$n > p$ 

$n > p$  $p > n$ 

$n > p$  $p > n$ 

$n > p$  $p > n$ 

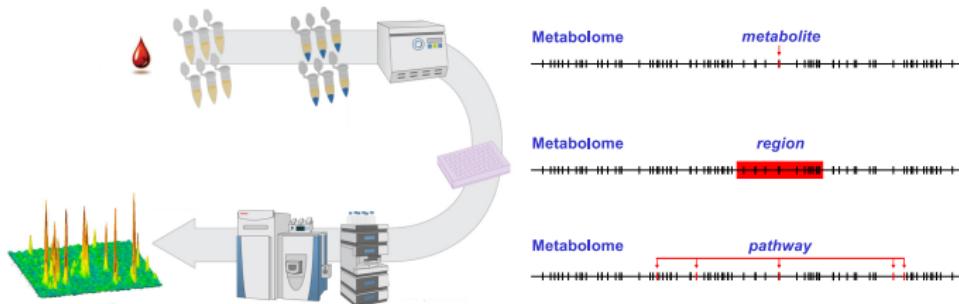
$n > p$  $p > n$ 

Inverse correlation when $n > p$

```
> ## Obtain some data
> p = 3
> n = 20
> X = matrix(rnorm(n*p), nrow = n, ncol = p)
> R <- cor(X)
>
> ## Obtain inverse and check
> Ri <- solve(R)
> R %*% Ri
      [,1]           [,2]           [,3]
[1,] 1.000000e+00  2.515349e-17  0.000000e+00
[2,] 3.469447e-18  1.000000e+00 -2.775558e-17
[3,] 6.938894e-18 -2.775558e-17  1.000000e+00
```

Inverse correlation when $p > n$

```
> ## Obtain some (high-dimensional) data
> p = 3
> n = 2
> X = matrix(rnorm(n*p), nrow = n, ncol = p)
> R <- cor(X)
>
> ## Try to obtain inverse
> Ri <- solve(R)
Error in solve.default(R) :
  Lapack routine dgesv: system is exactly singular
```

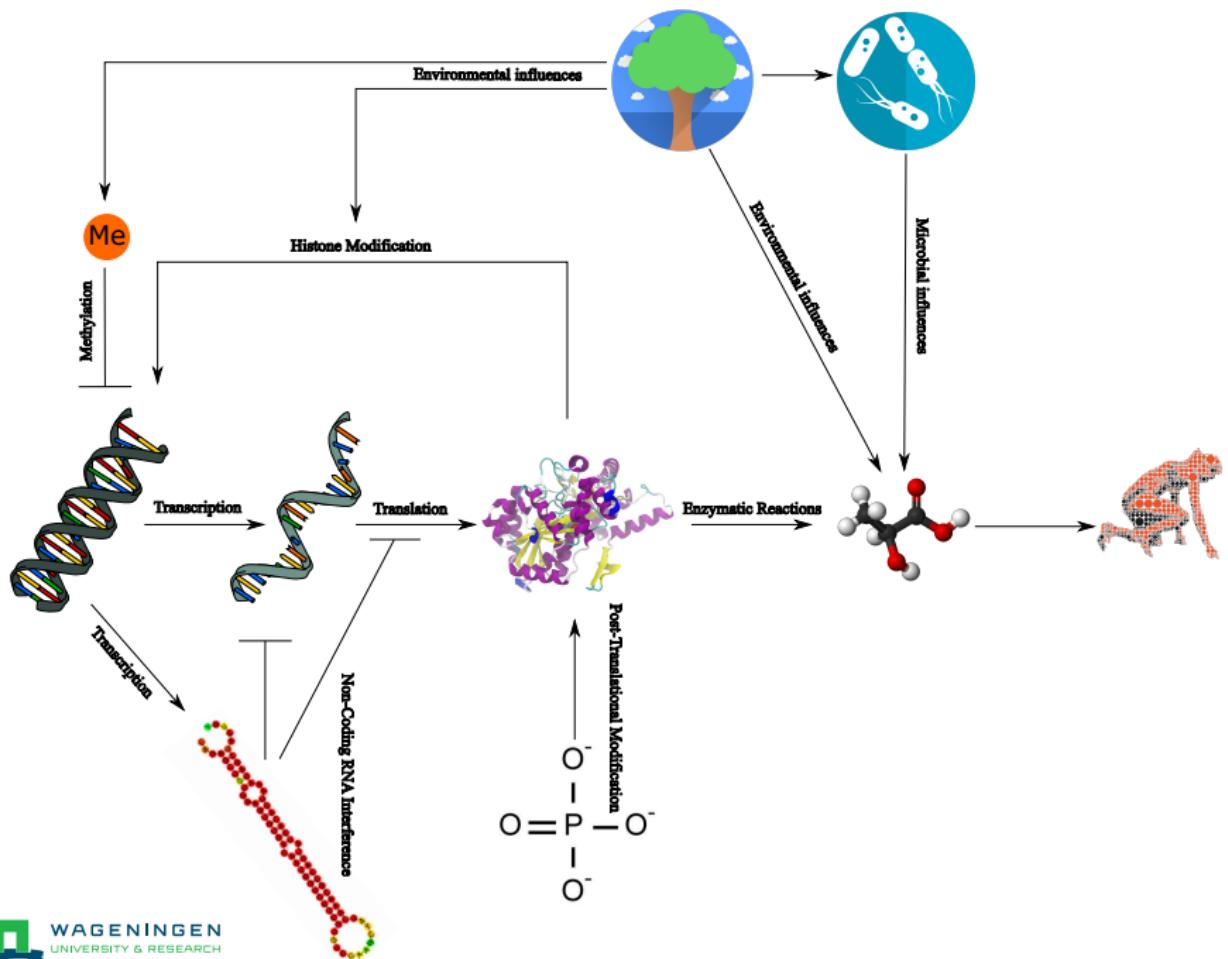


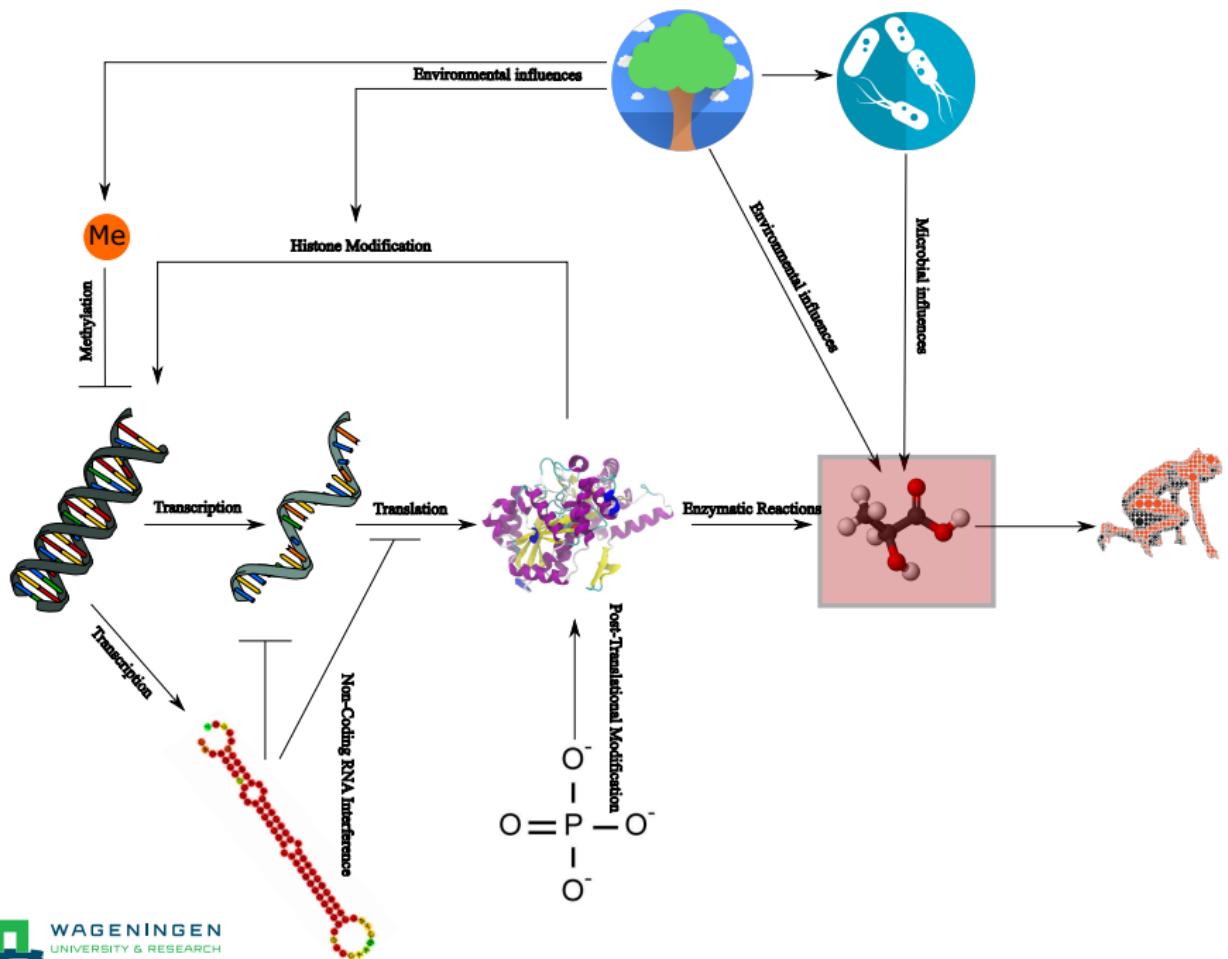
Problem

- Let $\mathbf{S} \equiv \hat{\Sigma}$ denote the sample covariance matrix on n realizations \mathbf{y}_i
- When $p \approx n$ or $p > n$, \mathbf{S} is ill-behaved or singular
- The empirical precision $\mathbf{S}^{-1} \equiv \hat{\Omega}$ is then undefined

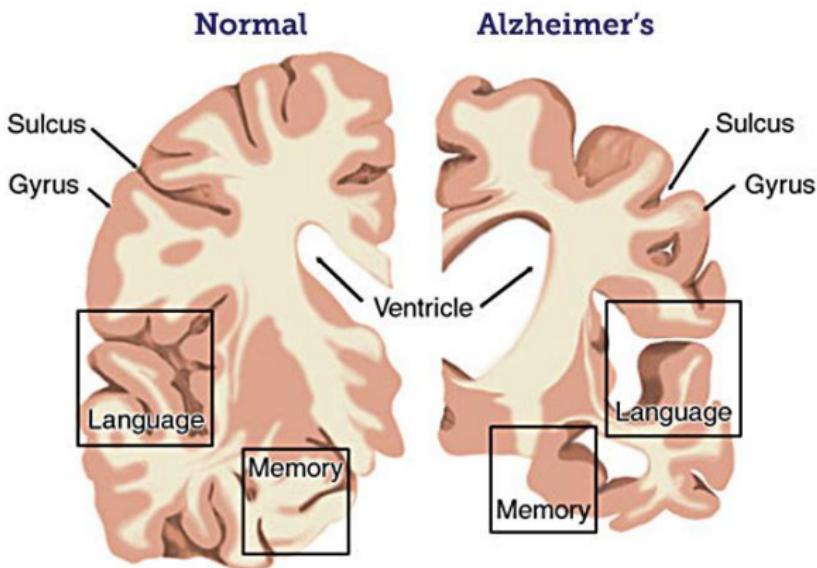
DATA

Our running example





Brain Cross-Sections



© 2000 by BrightFocus Foundation



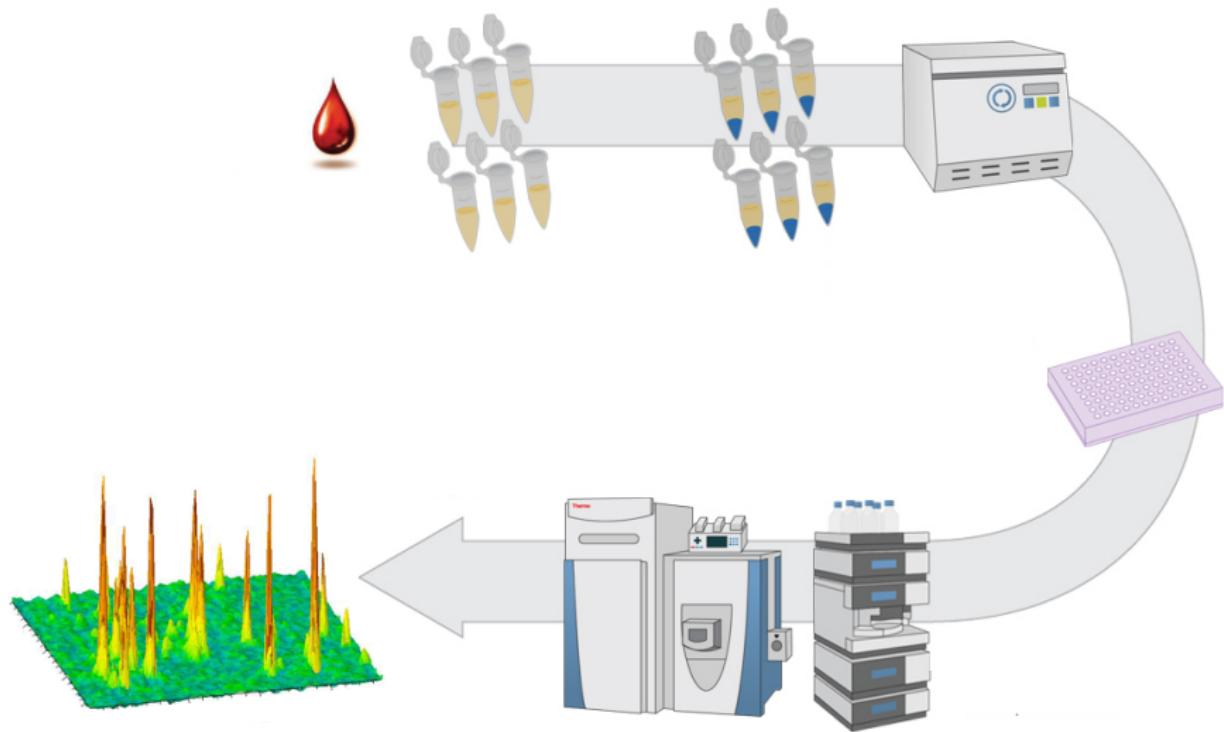
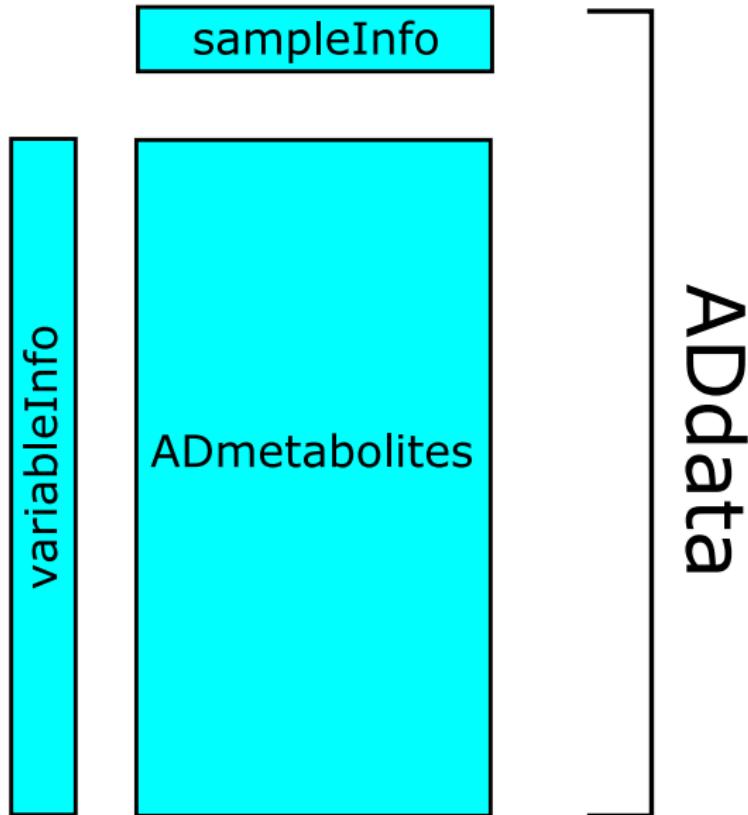


Illustration adapted from: <http://planetorbitrap.com/untargeted-metabolomics#.Vzw6yfmLRaQ> & <http://metabolomicsplatform.com/metabolomics-overview/>



Exercise 1: Get acquainted with the data

Get acquainted with the data

```
## Set working directory to convenience
setwd("")

## Needed package
library("rags2ridges")

## Data packaged as 3 data objects in ADdata object
data("ADdata", package = "rags2ridges")

## To probe objects one could use:
objects()
head(ADmetabolites)
head(sampleInfo)
colnames(variableInfo); table(variableInfo)
```

	Variables (features)						
	1	2	3	4	5	p
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}

- Amines
- Lipids
- Organic Acids
- Oxidative Stress

$p = 230$: 53 Amines, 116 Lipids, 22 Organic acids, 39 Oxidative stress compounds
 $n = 127$: 40 AD class 1, 87 AD class 2

Extraction

How to get the network

Maximize

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

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

Analytic penalized ML estimator

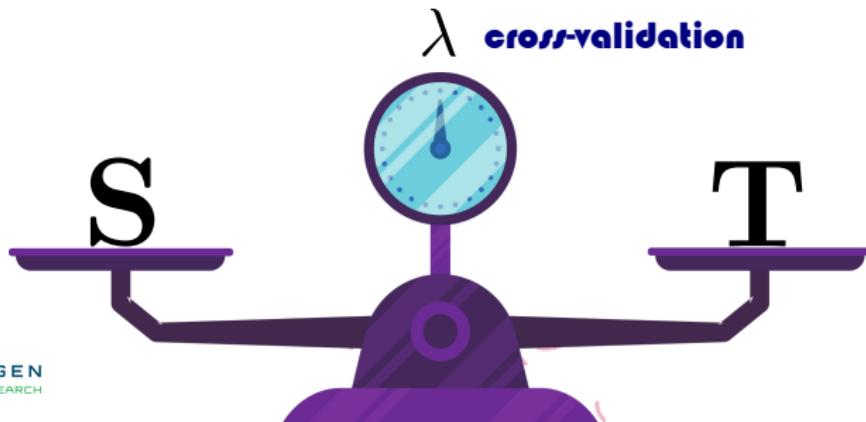
$$\hat{\Omega}(\lambda) = \left\{ \left[\lambda \mathbf{I}_p + \frac{1}{4} (\mathbf{S} - \lambda \mathbf{T})^2 \right]^{1/2} + \frac{1}{2} (\mathbf{S} - \lambda \mathbf{T}) \right\}^{-1}$$

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. $\lim_{\lambda \rightarrow \infty} \hat{\Omega}(\lambda) = \mathbf{T}$.

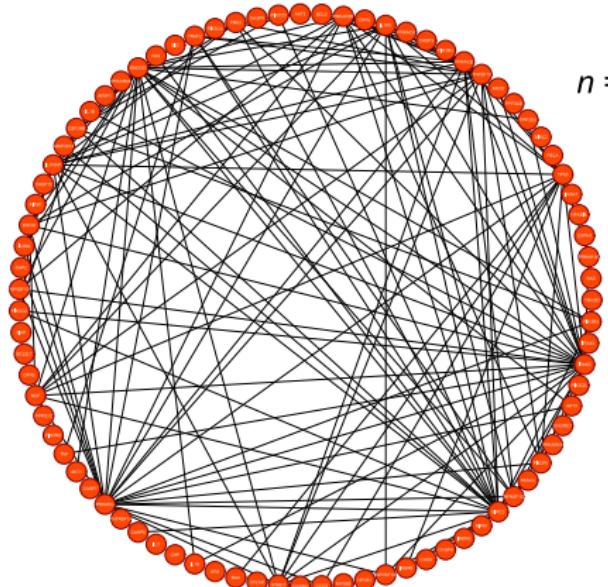
Consistency

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

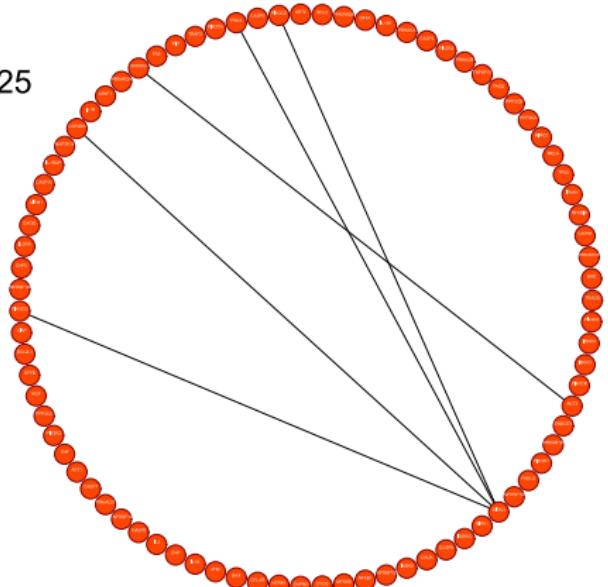


Why not maximize the ℓ_1 -Penalized log-likelihood?

$$\underbrace{\ln |\Omega| - \text{tr}(\hat{\Sigma}\Omega)}_{\text{log-likelihood}} - \underbrace{\lambda \|\Omega\|_1}_{\ell_1\text{-penalty}}$$



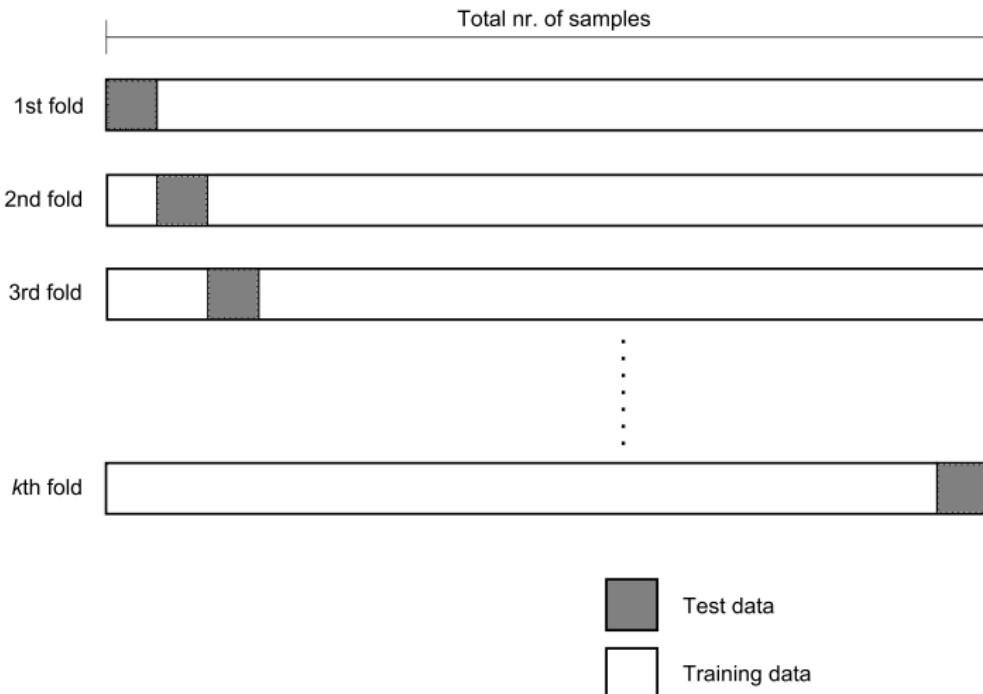
$n = 25$



graphical lasso

K-fold cross-validation (CV)

Single iteration of K-fold CV



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

Choose

$$\lambda^* = \arg \min_{\lambda \in \mathbb{R}^+} \varphi^K(\lambda)$$

Efficiency

- Implementation uses C++ at its core
- Makes use of a root-finding (Brent) algorithm
- Utilizes rotational equivariance property when possible

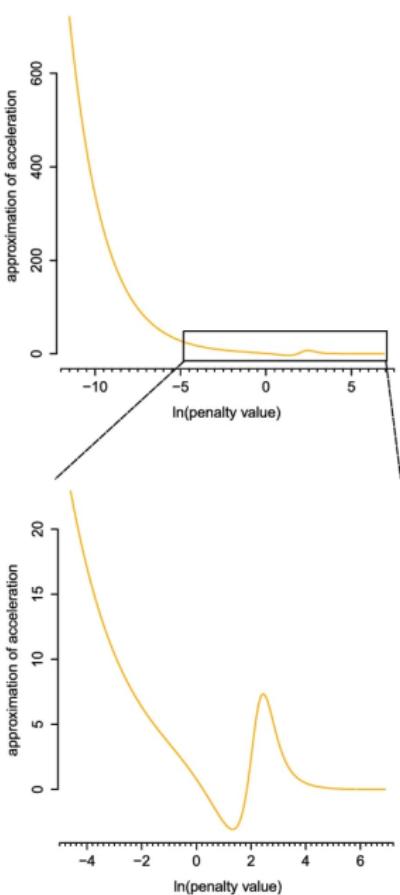
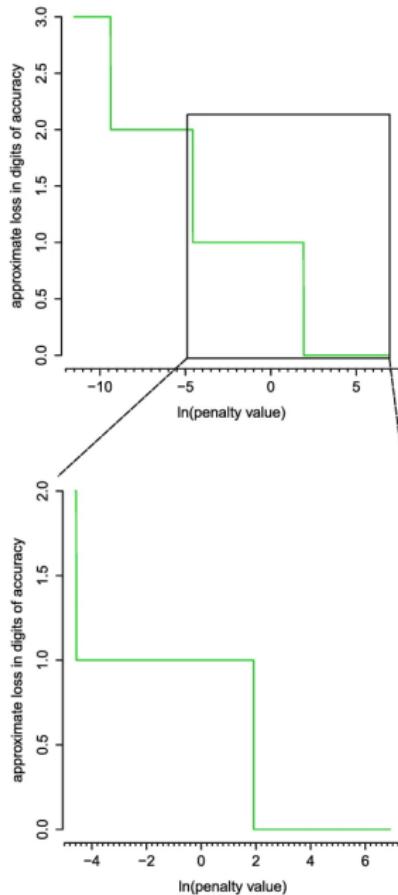
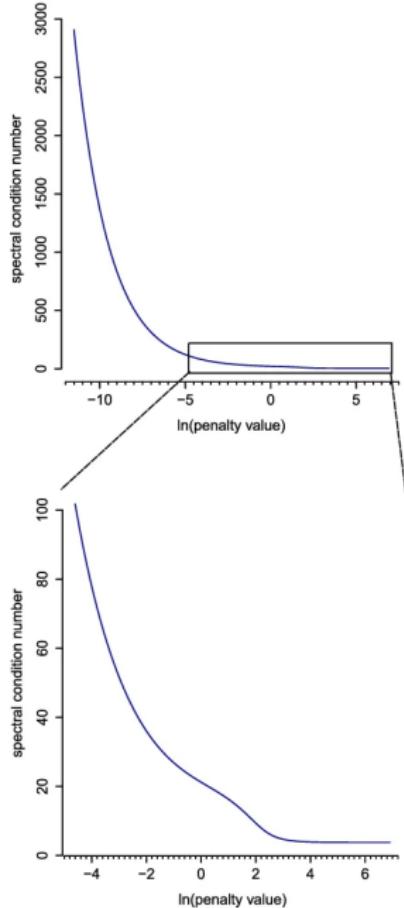
Exercise 2: Find an optimal precision matrix for the Class 2 AD data

CV function

```
optPenalty.kCVauto(  
  Y,          ## data  
  fold,       ## number of folds  
  lambdaMin, ## minimum penalty value  
  lambdaMax, ## maximum penalty value  
  target      ## target matrix T: use default.target()  
)
```

Returns list object

- \$optLambda: Optimal penalty parameter
- \$optPrec: Precision estimate under optimal penalty parameter



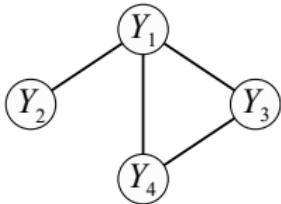
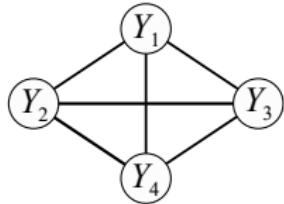
Exercise 3: Assess the conditioning of the optimal precision matrix

Function for condition number plot

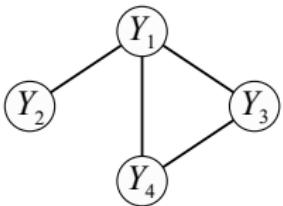
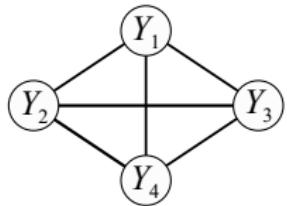
```
CNplot(S,          ## covariance matrix data
        lambdaMin, ## minimum penalty value
        lambdaMax, ## maximum penalty value
        step,       ## coarseness grid
        target,     ## target matrix T
        Iaids,      ## logical for interpretational aids
        vertical,   ## logical for inclusion vertical bar
        value,      ## value if vertical = TRUE
        verbose     ## logical controlling on-screen printing
    )
```

Returns list object

- \$optLambda: Optimal penalty parameter
- \$optPrec: Precision estimate under optimal penalty parameter



$$\begin{bmatrix} \omega_{11} & \omega_{12} & \omega_{13} & \omega_{14} \\ \omega_{21} & \omega_{22} & \omega_{23} & \omega_{24} \\ \omega_{31} & \omega_{32} & \omega_{33} & \omega_{34} \\ \omega_{41} & \omega_{42} & \omega_{43} & \omega_{44} \end{bmatrix} \rightarrow \begin{bmatrix} \omega_{11} & \omega_{12} & \omega_{13} & \omega_{14} \\ \omega_{21} & \omega_{22} & 0 & 0 \\ \omega_{31} & 0 & \omega_{33} & \omega_{34} \\ \omega_{41} & 0 & \omega_{43} & \omega_{44} \end{bmatrix}$$



$$\begin{bmatrix} \omega_{11} & \omega_{12} & \omega_{13} & \omega_{14} \\ \omega_{21} & \omega_{22} & \omega_{23} & \omega_{24} \\ \omega_{31} & \omega_{32} & \omega_{33} & \omega_{34} \\ \omega_{41} & \omega_{42} & \omega_{43} & \omega_{44} \end{bmatrix} \rightarrow \begin{bmatrix} \omega_{11} & \omega_{12} & \omega_{13} & \omega_{14} \\ \omega_{21} & \omega_{22} & 0 & 0 \\ \omega_{31} & 0 & \omega_{33} & \omega_{34} \\ \omega_{41} & 0 & \omega_{43} & \omega_{44} \end{bmatrix}$$



Scaling

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

Assume

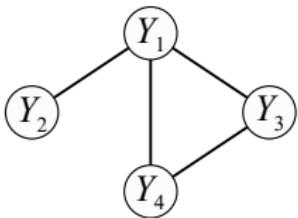
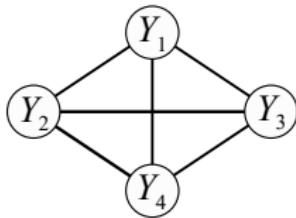
Nonredundant coefficients (indexed by $j < j'$) follow a mixture distribution:

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

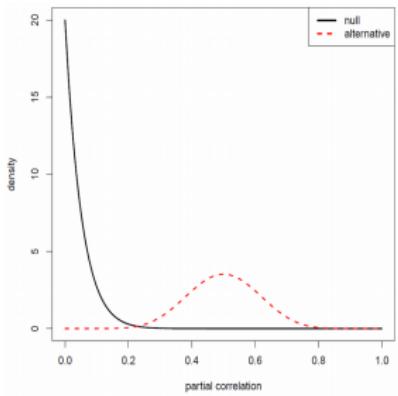
- $\eta_0 \in [0, 1]$ is the mixture weight
- $f_0\{\cdot\}$ denotes the distribution of a null-edge
- $f_{\mathcal{E}}\{\cdot\}$ denotes the distribution of a present edge
- κ denotes degrees of freedom

Determine

$$P(Y_j \neq Y_{j'} | [\hat{\mathbf{P}}(\lambda^*)]_{jj'}) = \frac{\hat{\eta}_0 f_0 \left\{ [\hat{\mathbf{P}}(\lambda^*)]_{jj'}; \hat{\kappa} \right\}}{\hat{\eta}_0 f_0 \left\{ [\hat{\mathbf{P}}(\lambda^*)]_{jj'}; \hat{\kappa} \right\} + (1 - \hat{\eta}_0) \hat{f}_{\mathcal{E}} \left\{ [\hat{\mathbf{P}}(\lambda^*)]_{jj'} \right\}}$$



$$\begin{bmatrix} \omega_{11} & \omega_{12} & \omega_{13} & \omega_{14} \\ \omega_{21} & \omega_{22} & \omega_{23} & \omega_{24} \\ \omega_{31} & \omega_{32} & \omega_{33} & \omega_{34} \\ \omega_{41} & \omega_{42} & \omega_{43} & \omega_{44} \end{bmatrix} \rightarrow \begin{bmatrix} \omega_{11} & \omega_{12} & \omega_{13} & \omega_{14} \\ \omega_{21} & \omega_{22} & 0 & 0 \\ \omega_{31} & 0 & \omega_{33} & \omega_{34} \\ \omega_{41} & 0 & \omega_{43} & \omega_{44} \end{bmatrix}$$



Exercise 4: Extract network from the optimal precision matrix

Function for support determination

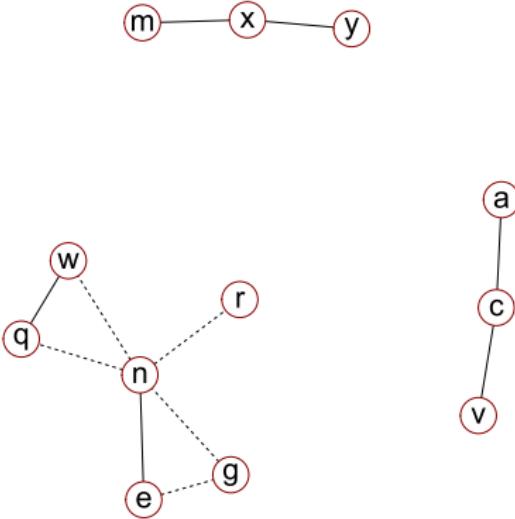
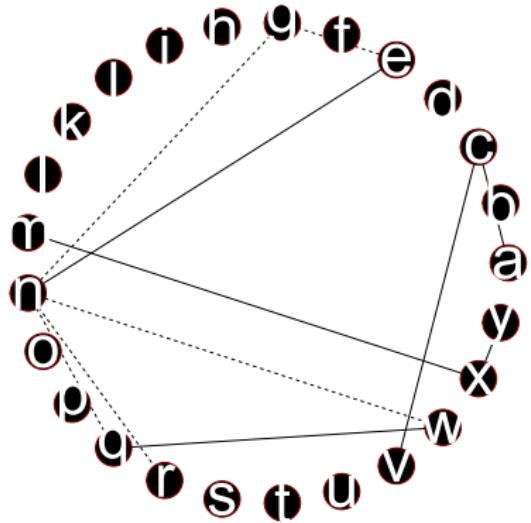
```
sparsify(P,          ## estimated precision matrix
         threshold,    ## type of thresholding: we use "localFDR"
         FDRcut,       ## cut-off for 1 - 1FDR
         verbose,      ## logical controlling on-screen printing
         )
```

Returns list object

- `$sparsePrecision`: Sparsified precision matrix
- `$sparseParCor`: Sparsified partial correlation matrix

Visualization

How to visualize the network



Exercise 5: Visualize the extracted network

Function for visualization

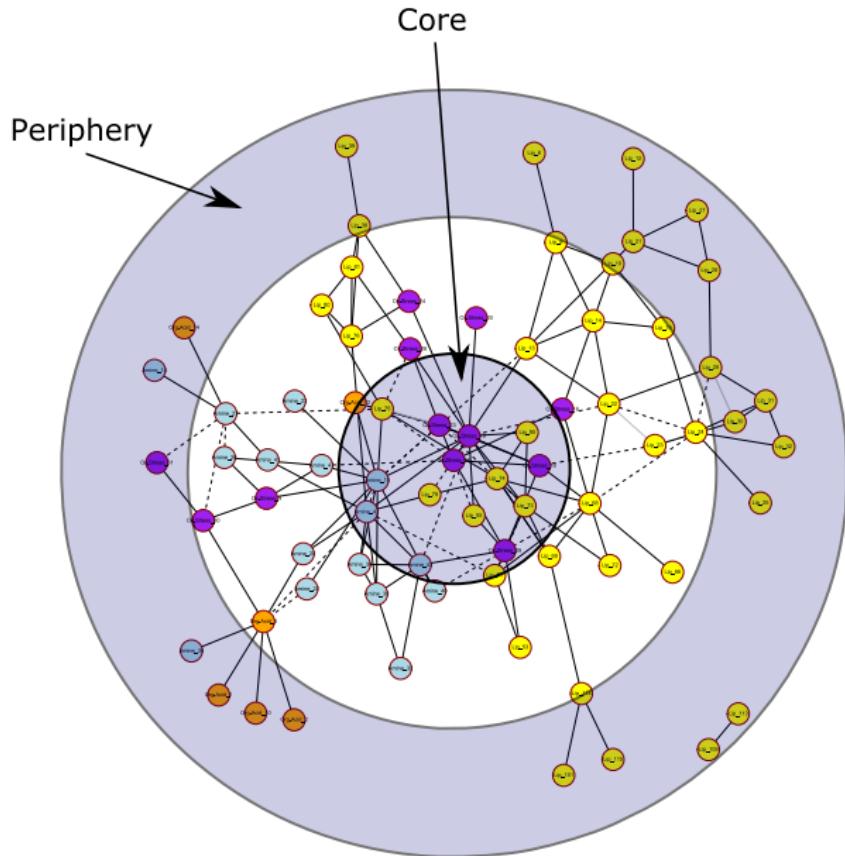
```
Ugraph(M,      ## matrix to be visualized
       type,    ## "plain", "fancy", "weighted"
       lay,     ## specifies layout
       coords,  ## gives control over node coordinates
       Vcolor,  ## vertex color
       Vsize,   ## vertex size
       Vcex,    ## size of labels within vertices
       ...
       )
```

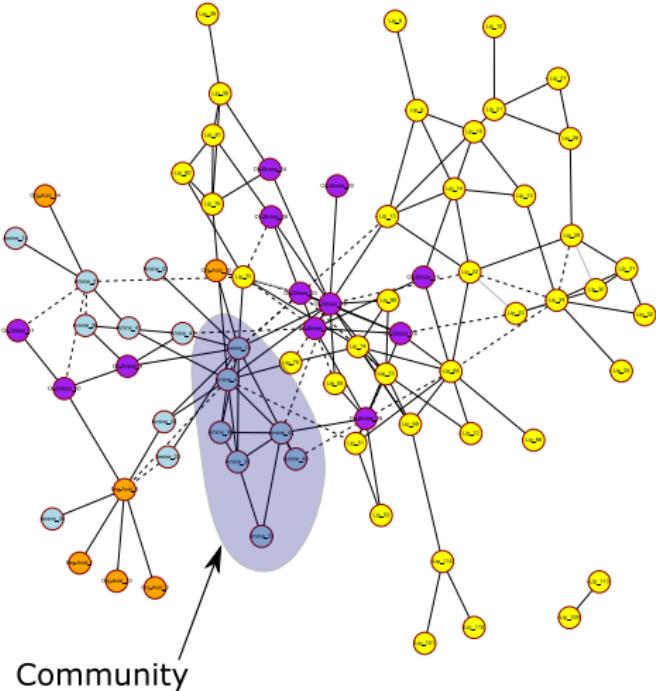
Returns matrix object

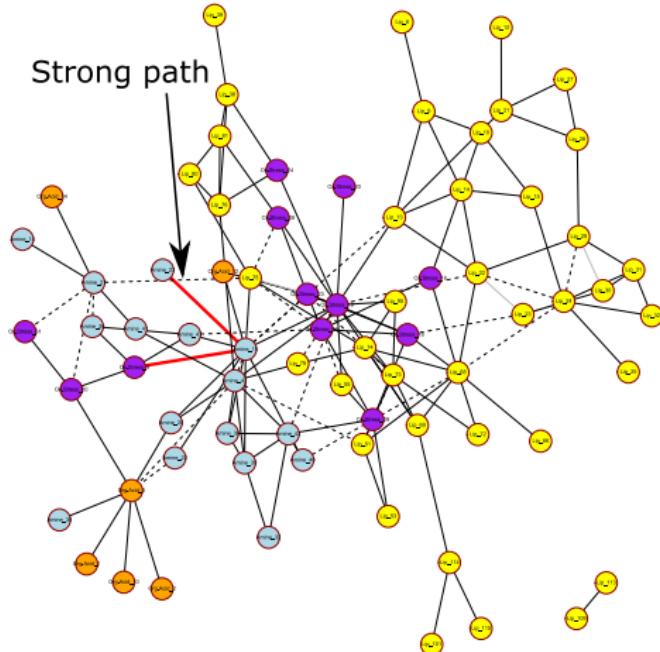
Containing the coordinates of the vertices in the given graph/network

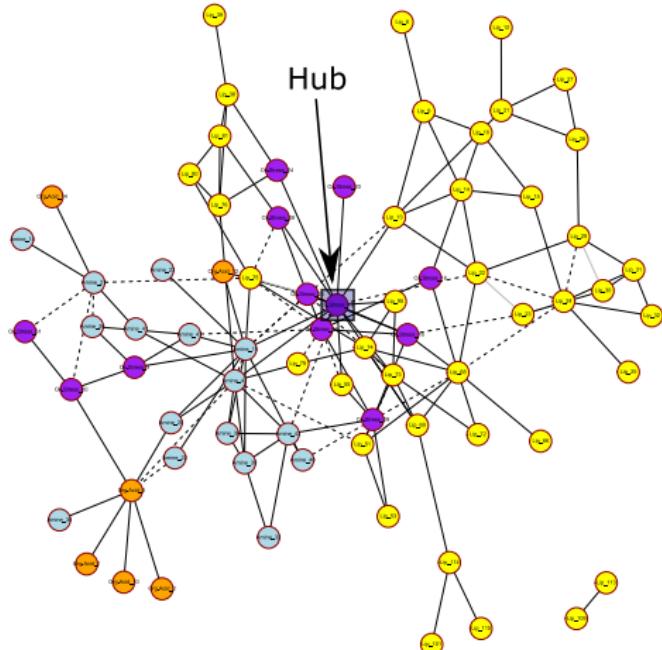
Analysis

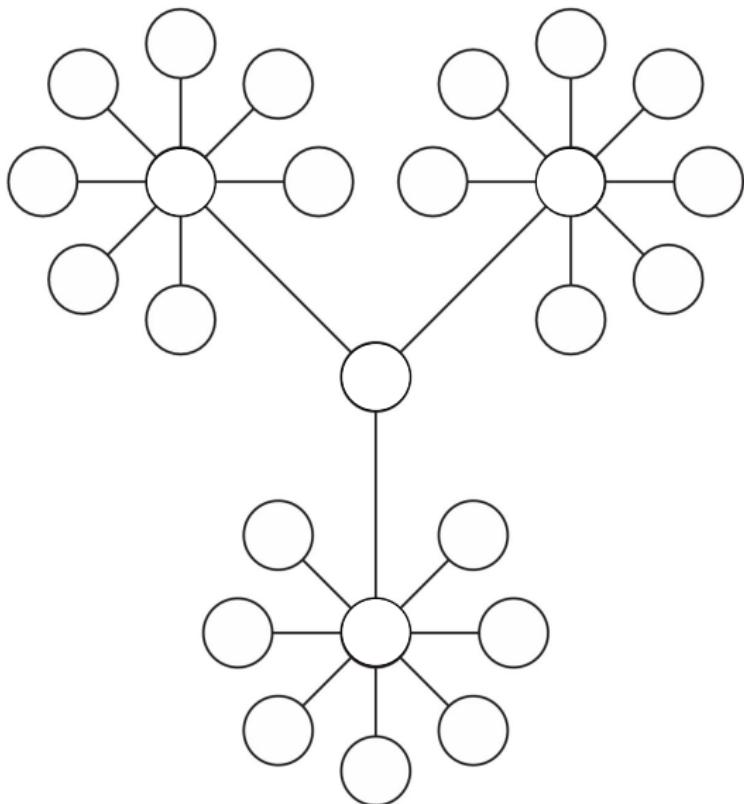
How to analyze the network

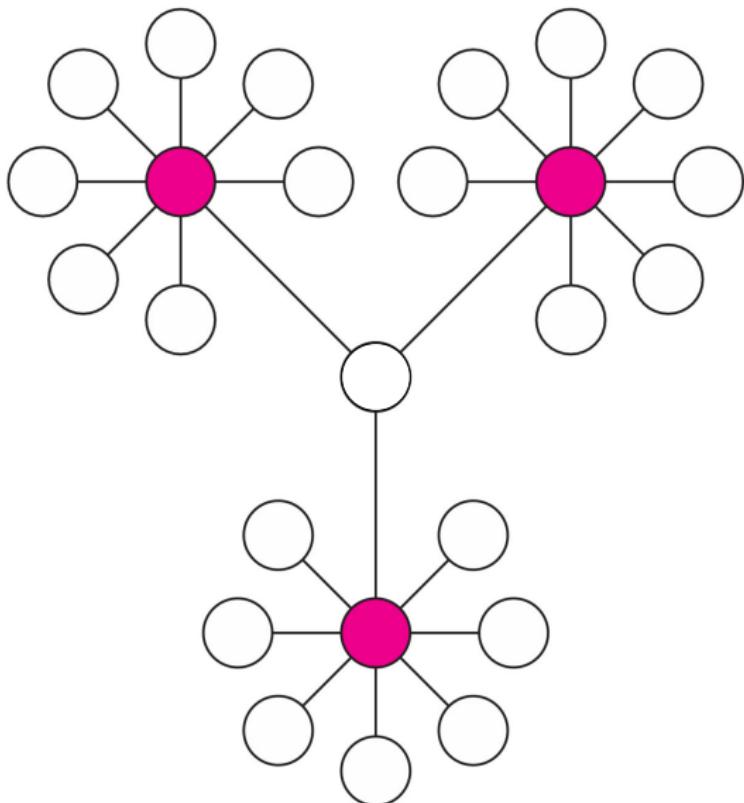


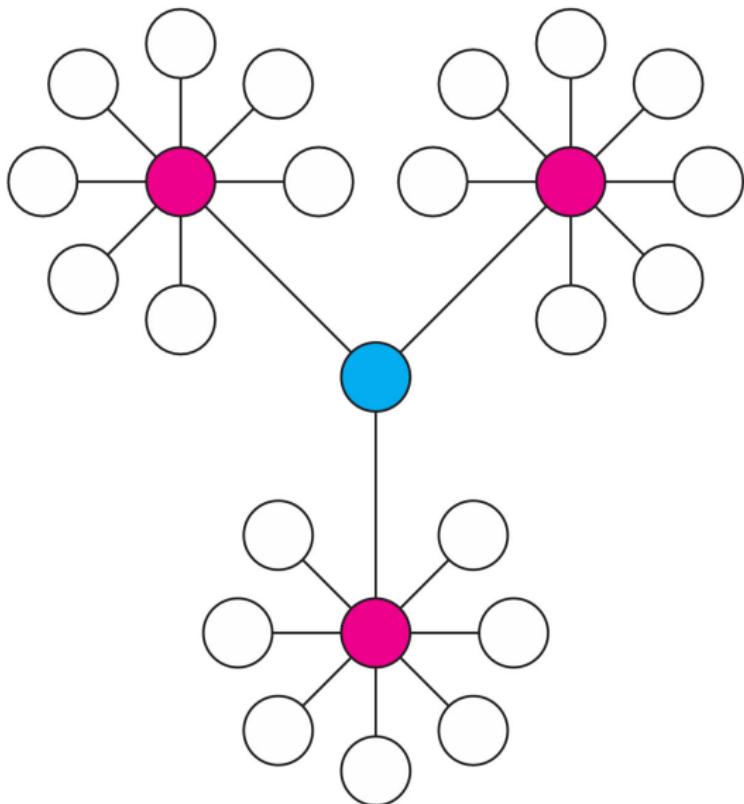












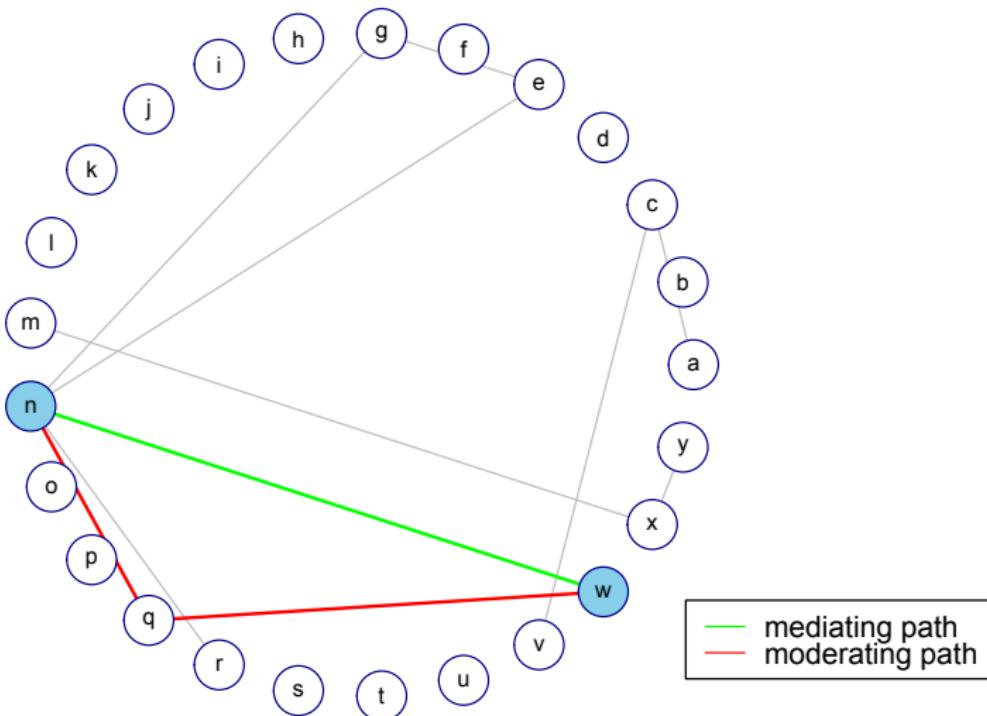
Exercise 6: Find the nodes with the top degree and betweenness centralities

Function for simple network statistics

```
GGMnetworkStats(sparseP ## sparse matrix  
    )
```

Returns list object

- \$degree: Degree centrality
- \$betweenness: Betweenness centrality
- ...



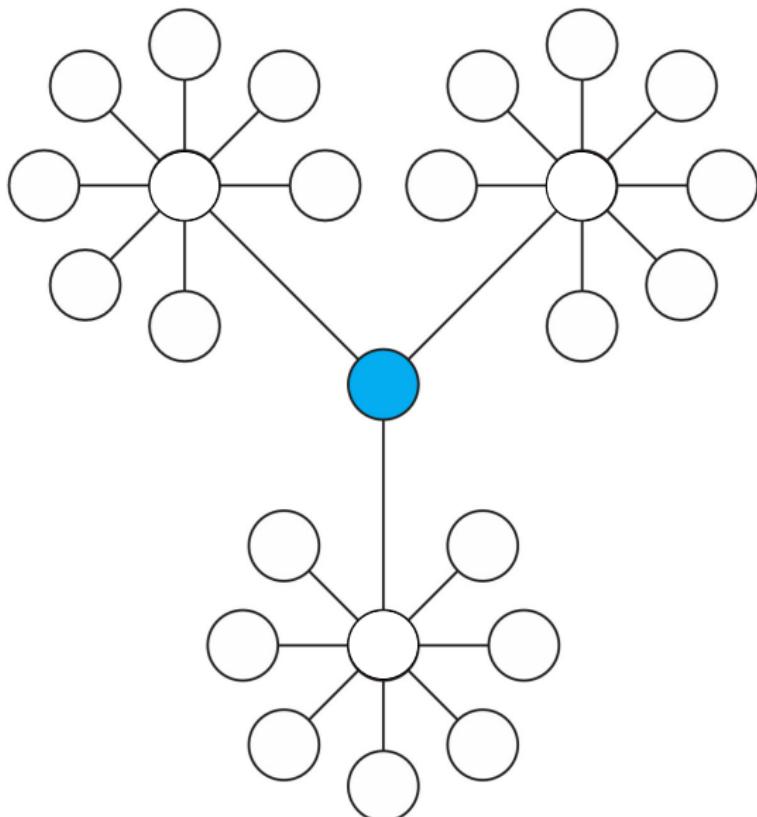
Exercise 7: Find the 2 strongest paths between Amines 1 and 2

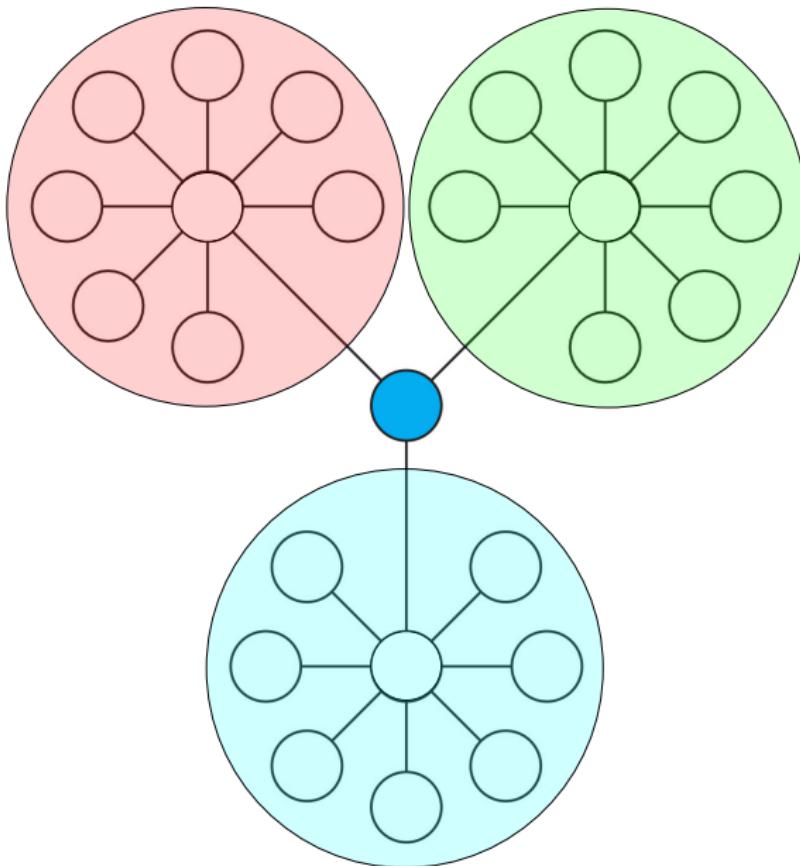
Function for path statistics

```
GGMpathStats(P0,      ## sparse precision matrix
              node1,    ## endpoint 1
              node2,    ## endpoint 2
              graph,    ## logical, should graph be produced
              ...       ## arguments to Ugraph
            )
```

Returns list object

- \$pathStats: Matrix specifying paths
- ...





Exercise 8: Find and visualize communities

Function for community detection

```
Communities(P,      ## sparse matrix
            graph, ## logical, graph is produced when TRUE
            ...
            )      ## arguments to Ugraph
```

Returns list object

- **\$membership:** Community membership for each feature
- **\$modularityscore:** Modularity score

Part II: Jointly extracting, visualizing, and analyzing multiple networks