

An Overview of Gene Co-Expression Network Analysis

Cheng He

04/27/2021

Bioinformatics Applications (PLPTH813)

Contents

- Introduction of gene co-expression network (GCN)
- Rationales of gene co-expression network
- Gene co-expression network analysis methods
- Weighted Gene Co-Expression Network Analysis (WGCNA)
- Tutorial: Run WGCNA on your own laptop

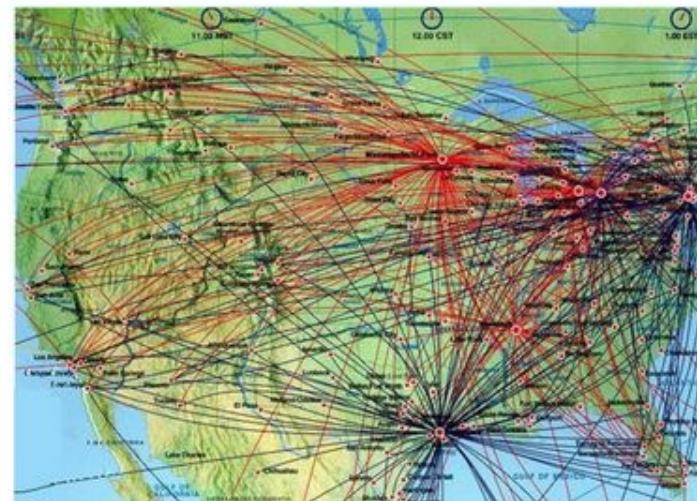
Contents

- Introduction of gene co-expression network (GCN)
- Rationales of a gene co-expression network
- Gene co-expression network analysis methods
- Weighted Gene Co-Expression Network Analysis (WGCNA)
- Tutorial: Run WGCNA on your own laptop

**Does this map tell you
which cities are
important?**



This one does!



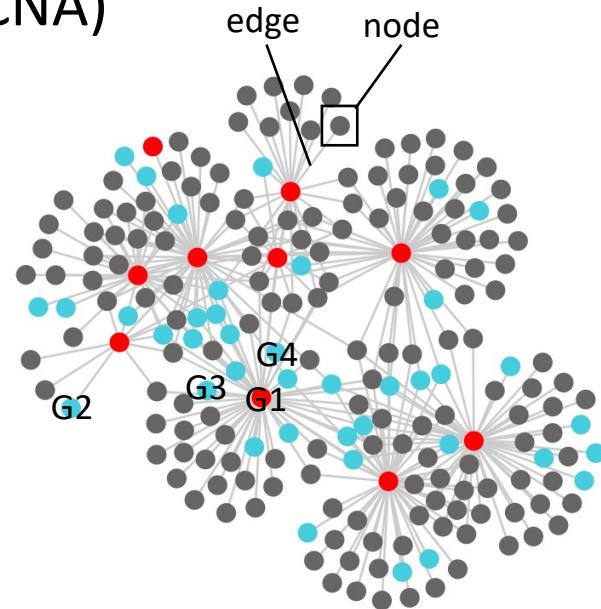
***The nodes with the largest number of links
(connections) are most important!***

Gene co-expression network

- In gene co-expression networks, each gene corresponds to a **node**
- Two genes are connected by an **edge** if their expression values are **highly correlated**
- Definition of "high" correlation can be decided by different methods:
 - (1) Statistical significance (Pearson correlation)
 - (2) Scale free topology criterion (WGCNA)

Gene expression matrix			
Gene_ID	sample_1	sample_2
G1	4	2	xx
G2	0	3	xx
G3	6	3	xx
G4	10	5	xx
⋮	xx	xx	xx
⋮	xx	xx	xx

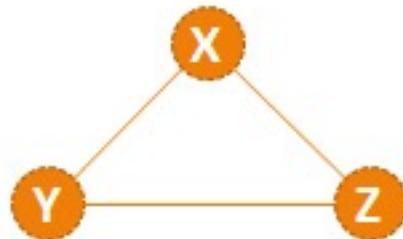
Gene co-expression
network analysis



A gene co-expression network (GCN) is an undirected graph, where each node corresponds to a gene, and a pair of nodes is connected with an edge if there is a significant co-expression relationship between them.

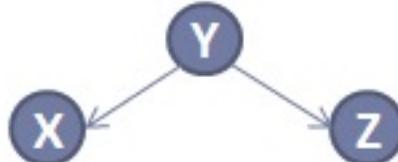
(Stuart, Joshua M; Segal, Eran; Koller, Daphne; Kim, Stuart K (2003). "A gene co-expression network for global discovery of conserved genetic modules". *Science*. **302** (5643): 249–55.)

Gene Co-expression

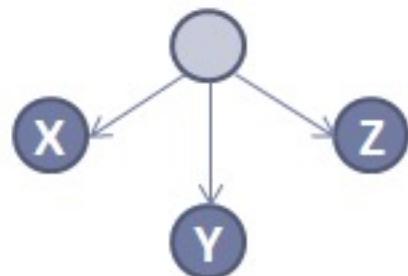


Undirected graph

Gene Regulation



Directed graph



Applications of Gene co-expression network

- Finding clusters or modules of highly relevant genes.
- Summarizing such modules using an intramodular hub gene or the module eigengene.
- Relating modules to one another and to external sample traits with eigengene network methodology.
- Calculating the measures of module membership.
- Gene Network Reverse Engineering.
- Plant Biology - Co-expression analyses have been extensively used to search for novel genes involved in specific plant pathways.

Contents

- Introduction of gene co-expression network (GCN)
- **Rationales of gene co-expression network**
- Gene co-expression network analysis methods
- Weighted Gene Co-Expression Network Analysis (WGCNA)
- Tutorial: Run WGCNA on your own laptop

Gene co-expression similarity measure

- Pearson correlation was used to measure the similarity of 2 gene expressions (s_{ij})

Unsigned network: No gene repression and activation

$$s_{ij}^{unsigned} = |cor(x_i, x_j)|$$

Signed network: Consider gene repression and activation

$$s_{ij}^{signed} = 0.5 + 0.5cor(x_i, x_j)$$

G1 activate G2: $s_{12}^{signed} = 0.5 + 0.5*1 = 1$

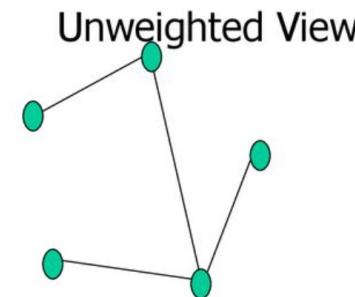
G1 repress G2: $s_{12}^{signed} = 0.5 + 0.5*(-1) = 0$

Network=Adjacency Matrix

- A network can be represented by an adjacency matrix, $A=[a_{ij}]$, that encodes whether/how a pair of nodes is connected.

Unweighted network:

- Entries are 1 or 0 depending on whether or not 2 nodes are adjacent (connected)

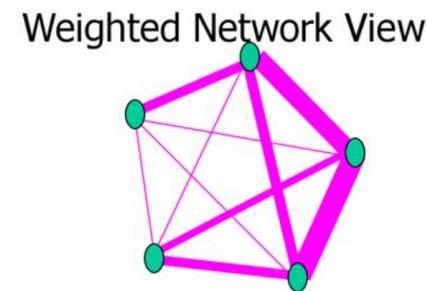


Some genes are connected
All connections are equal



Weighted network:

- The adjacency matrix reports the connection strength between gene pairs



•All genes are connected
•Connection Widths=Connection strengths

Generalized Connectivity

- Gene connectivity = row sum of the adjacency matrix
 - For unweighted networks = number of direct neighbors
 - For weighted networks = sum of connection strengths to other nodes

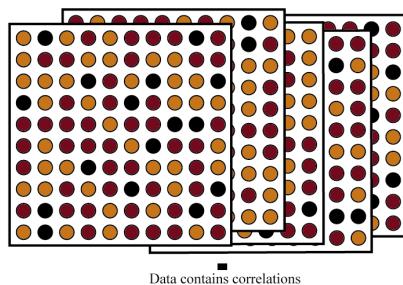
$$k_i = \sum_j a_{ij}$$

k_i : connectivity of gene i
 a_{ij} : adjacency matrix between gene i and gene j

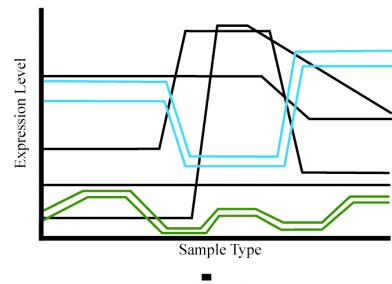
High gene connectivity usually indicates a hub gene

Figure 1

A Array Data



B Correlation Analysis



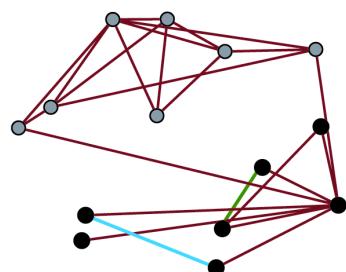
Correlation coefficients for all genes

C Correlation Matrix

	G1	G2	G3	G4	G5	G6	G7	G8	G9	G10	G11	G12	G13	G14
G1	1	0.9	0.9	0.9	0.9	0.8	0.9	0.1	0.9	0.1	0.1	0.8	0.2	0.2
G2	0.9	1	0.9	0.3	0.3	0.7	0.0	0.5	0.3	0.1	0.1	0.2	0.4	0.3
G3	0.9	0.9	1	0.9	0.0	0.2	0.5	0.7	0.6	0.5	0.2	0.6	0.1	0.0
G4	0.9	0.3	0.9	1	0.5	0.3	0.6	0.3	0.0	0.5	0.1	0.2	0.2	0.6
G5	0.9	0.3	0.0	0.5	1	0.1	0.6	0.1	0.3	0.3	0.3	0.5	0.2	0.5
G6	0.8	0.7	0.2	0.3	0.1	1	0.9	0.2	0.1	0.1	0.5	0.3	0.1	0.1
G7	0.9	0.0	0.5	0.6	0.6	0.9	1	0.3	0.1	0.5	0.1	0.3	0.5	0.2
G8	0.1	0.5	0.7	0.3	0.1	0.2	0.3	1	0.9	0.9	0.9	0.8	0.8	0.9
G9	0.9	0.3	0.6	0.0	0.3	0.1	0.1	0.9	1	0.8	0.1	0.3	0.5	0.3
G10	0.1	0.1	0.5	0.5	0.3	0.1	0.5	0.9	0.8	1	0.8	1.0	0.2	0.3
G11	0.1	0.1	0.2	0.1	0.3	0.5	0.1	0.9	0.1	0.8	1	0.5	0.8	0.9
G12	0.8	0.2	0.6	0.2	0.5	0.3	0.3	0.8	0.3	1.0	0.5	1	0.8	0.1
G13	0.2	0.4	0.1	0.2	0.2	0.1	0.5	0.8	0.5	0.2	0.8	0.8	1	0.9
G14	0.2	0.3	0.0	0.6	0.5	0.1	0.2	0.9	0.3	0.3	0.9	0.1	0.9	1

Convert into Adjacency Matrix and Network

D Coexpression Network



Steps for constructing a co-expression network

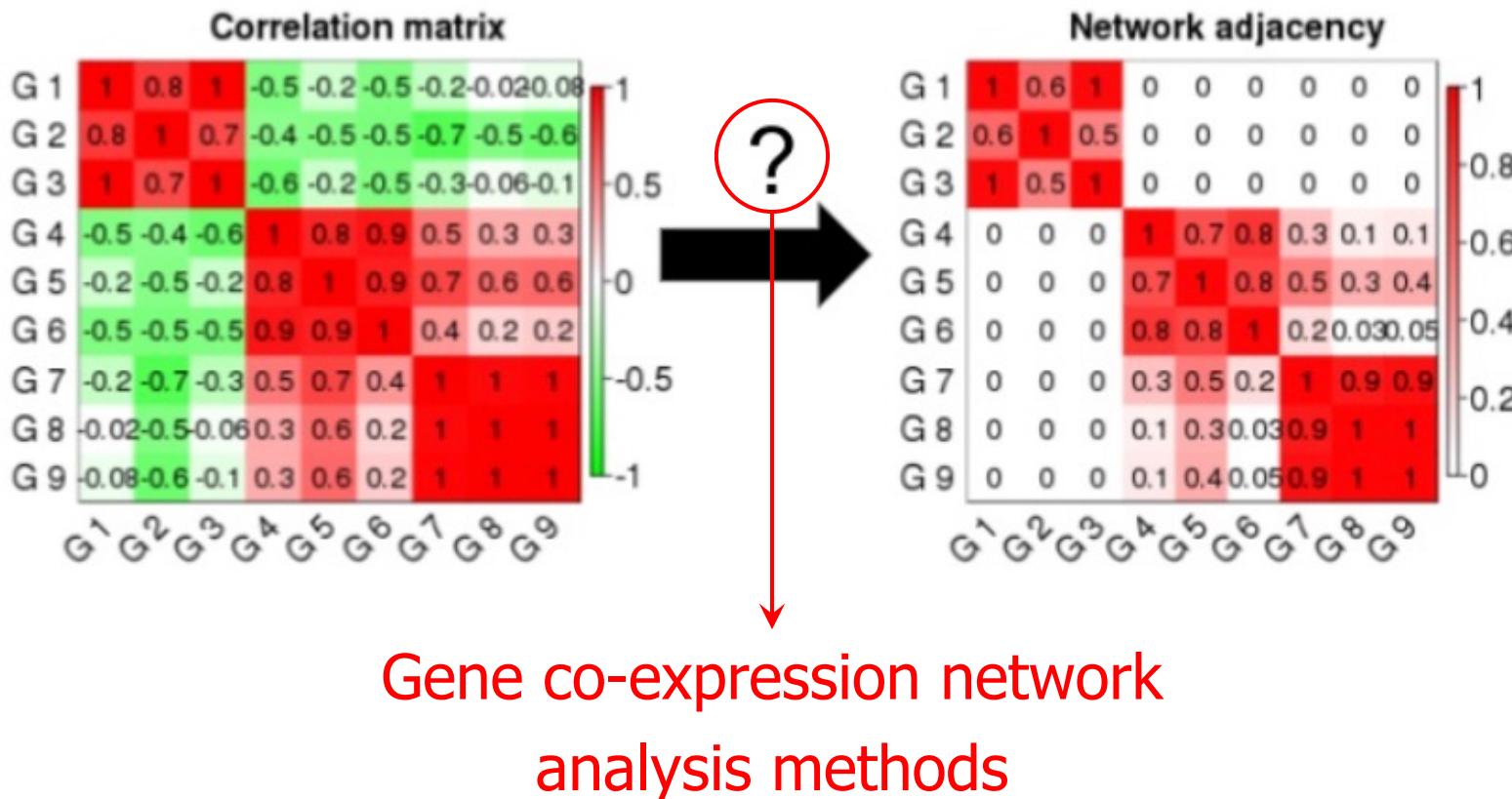
A) Microarray/RNA-Seq gene expression data

B) Measure similarity of gene expression with a Pearson correlation

C) The Pearson correlation matrix is either dichotomized to arrive at an adjacency matrix → unweighted network
Or transformed continuously with the power adjacency function → weighted network

D) The gene connectivity can be calculated from the adjacency matrix

How to turn correlation matrix to network adjacency is the key to construct a good gene co-expression network



Contents

- Introduction of gene co-expression network (GCN)
- Rationales of gene co-expression network
- **Gene co-expression network analysis methods**
- Weighted Gene Co-Expression Network Analysis (WGCNA)
- Tutorial: Run WGCNA on your own laptop

Current Methods

WGCNA based methods

- Weighted Gene Co-expression Network Analysis (Horvath 2008)
- DiffCoEx (Bruno 2010)

Other methods

- THD-Module Extractor (Kakati 2016)

Weighted Gene Co-expression Network Analysis (WGCNA)

Bin Zhang and Steve Horvath

A General Framework for Weighted Gene Co-Expression Network Analysis

De Gruyter | Published online: August 12, 2005

DOI: <https://doi.org/10.2202/1544-6115.1128>

Software | Open Access | Published: 29 December 2008

WGCNA: an R package for weighted correlation network analysis

[Peter Langfelder](#) & [Steve Horvath](#) 

[BMC Bioinformatics](#) **9**, Article number: 559 (2008) | [Cite this article](#)

Now WGCNA has been the most popular and reliable gene co-expression network analysis method

Procedures:

Step 1: Measure network adjacency by the power function

$$a_{ij} = |cor(x_i, x_j)|^\beta \quad \begin{array}{l} a_{ij}: \text{the adjacency matrix between gene i and j} \\ \beta: \text{soft thresholding parameter} \end{array}$$

Step 2: Cluster genes into network modules using TOM method

$$TOM_{ij} = \frac{\sum_u a_{iu} a_{uj} + a_{ij}}{\min(k_i, k_j) + 1 - a_{ij}} \quad \begin{array}{l} k_i: \text{the connectivity of gene i} \\ k_i = \sum_{u \neq i} a_{iu} \end{array}$$

Notes:

1. Often $\beta=6$ or $\beta=12$ works well but in general WGCNA uses the “scale free topology criterion”
2. The topological overlap measure (TOM) combines the adjacency of two genes and the connection strengths these two genes share with other “third party” genes.

Advantages of WGCNA:

- It retains the connectivity of network nodes
- It has powerful analytical efficiency
- Standard data mining methods such as clustering analysis results can be transformed into weight networks (RNA-Seq)

Disadvantages of WGCNA:

- WGCNA doesn't support dataset < 15 samples
- WGCNA takes many parameters with default values to control network construction and distinct module extraction
- WGCNA doesn't assume scale-free topology for differentially expressed genes

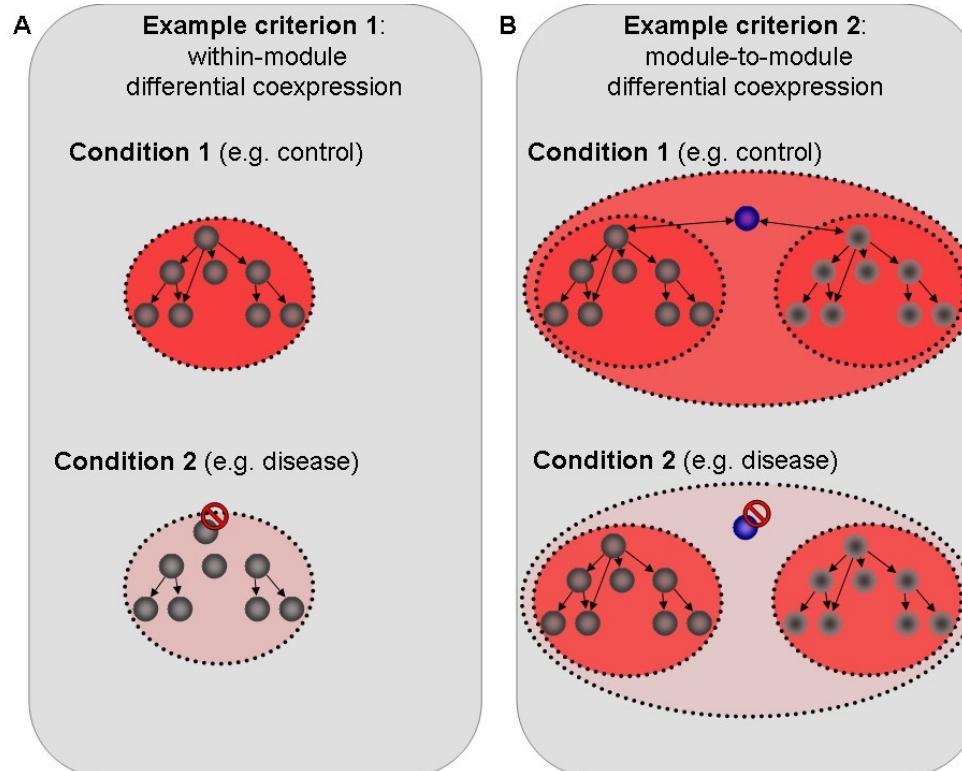
DiffCoEx

Methodology article | [Open Access](#) | Published: 06 October 2010

DiffCoEx: a simple and sensitive method to find differentially coexpressed gene modules

Bruno M Tesson, Rainer Breitling & Ritseert C Jansen [✉](#)

BMC Bioinformatics 11, Article number: 497 (2010) | [Cite this article](#)



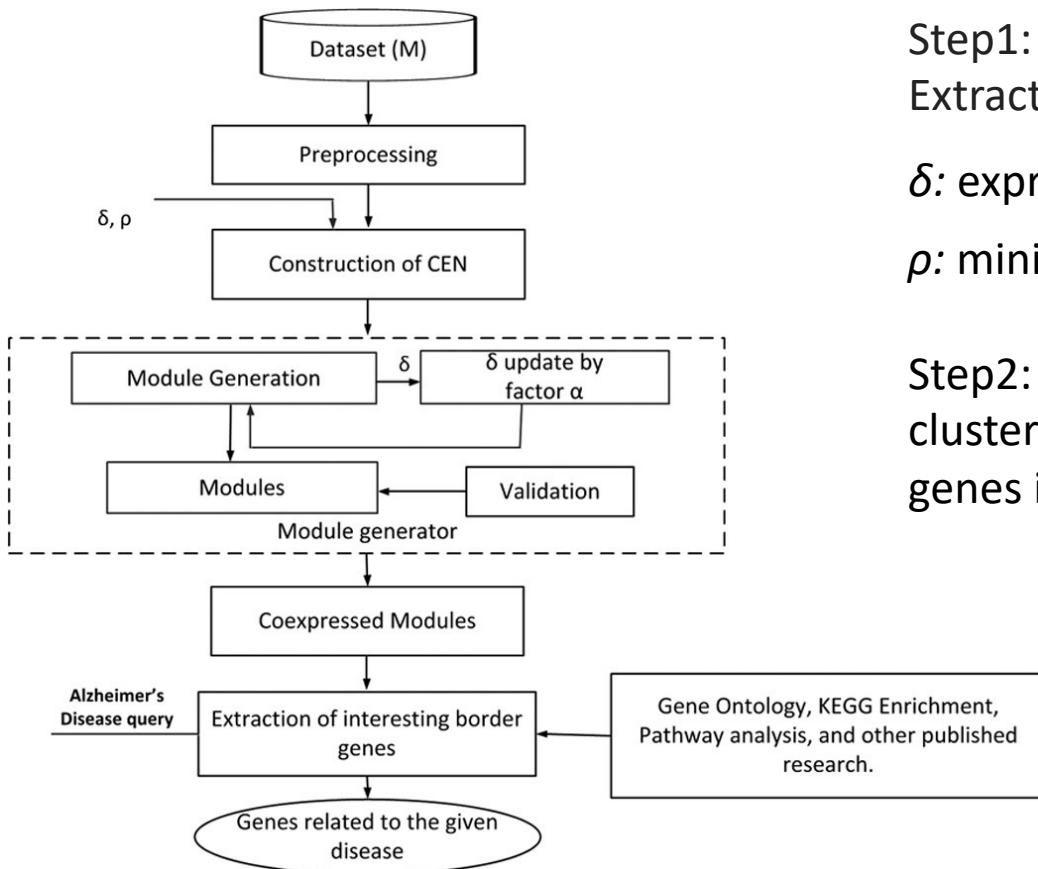
DiffCoEx can identify gene co-expression differences between **multiple conditions**.

THD-Module Extractor

THD-Module Extractor: An Application for CEN Module Extraction and Interesting Gene Identification for Alzheimer's Disease

Tulika Kakati, Hirak Kashyap & Dhruba K. Bhattacharyya

Scientific Reports 6, Article number: 38046 (2016) | Cite this article



Step1: A schematic diagram of THD-Module Extractor framework.

δ : expression similarity threshold

ρ : minimum neighborhood threshold

Step2: A Shifting-and-scaling correlation based clustering algorithm (SSSim) is used to cluster genes into modules

Advantages of THD-Module Extractor:

- THD-Module Extractor considers the important issue of analysis of CEN using both gene expression similarity and **semantic similarity**

Semantic similarity between two genes can be measured on the basis of information content (I_c)

$$I_c = -\ln(P(t))$$

$$P(t) = \frac{\text{number of annotations involving a GO term 't'}}{\text{number of genes}}$$

Comparison of four CEN module extraction techniques

Method	Working principle	Measure used in CEN construction	Parameters	Language	Datasets used
WGCNA	Hierarchical clustering	Dissimilarity measure based on TOM	Soft threshold parameter and tree cutting threshold	R	Microarray and RNA-Seq
DiffCoEx	Differential Co-expression analysis	Dissimilarity measure based on TOM	Soft threshold parameter and tree cutting threshold	R	Microarray
THD-Module Extractor	Density-based clustering	SSSim	Expression similarity threshold (δ), minimum neighborhood threshold (ρ), and δ updating factor α	R and MATLAB	Microarray

Kakati, T., Bhattacharyya, D.K., Barah, P. and Kalita, J.K., 2019. Comparison of methods for differential co-expression analysis for disease biomarker prediction. *Computers in biology and medicine*, 113, p.103380.

Contents

- Introduction of gene co-expression network (GCN)
- Rationales of gene co-expression network
- Gene co-expression network analysis methods
- **Weighted Gene Co-Expression Network Analysis (WGCNA)**
- Tutorial: Run WGCNA on your own laptop

General Framework for WGCNA

(Pearson correlation)

(scale-free topology criterion)

$$a_{ij} = |cor(x_i, x_j)|^\beta$$

$$DistTOM_{ij} = 1 - TOM_{ij}$$

(hierarchical clustering)

Define a Gene Co-expression Similarity

Define a soft thresholding power (β)

β

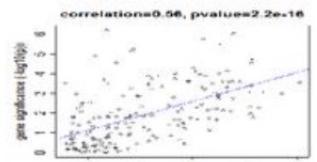
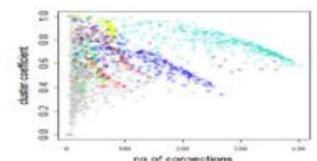
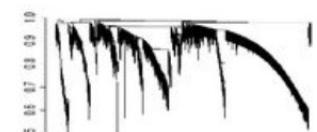
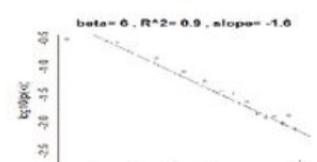
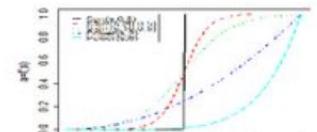
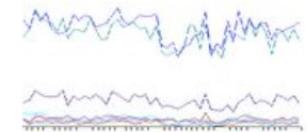
Calculate Adjacency Matrix with β

Define a Measure of Node Dissimilarity

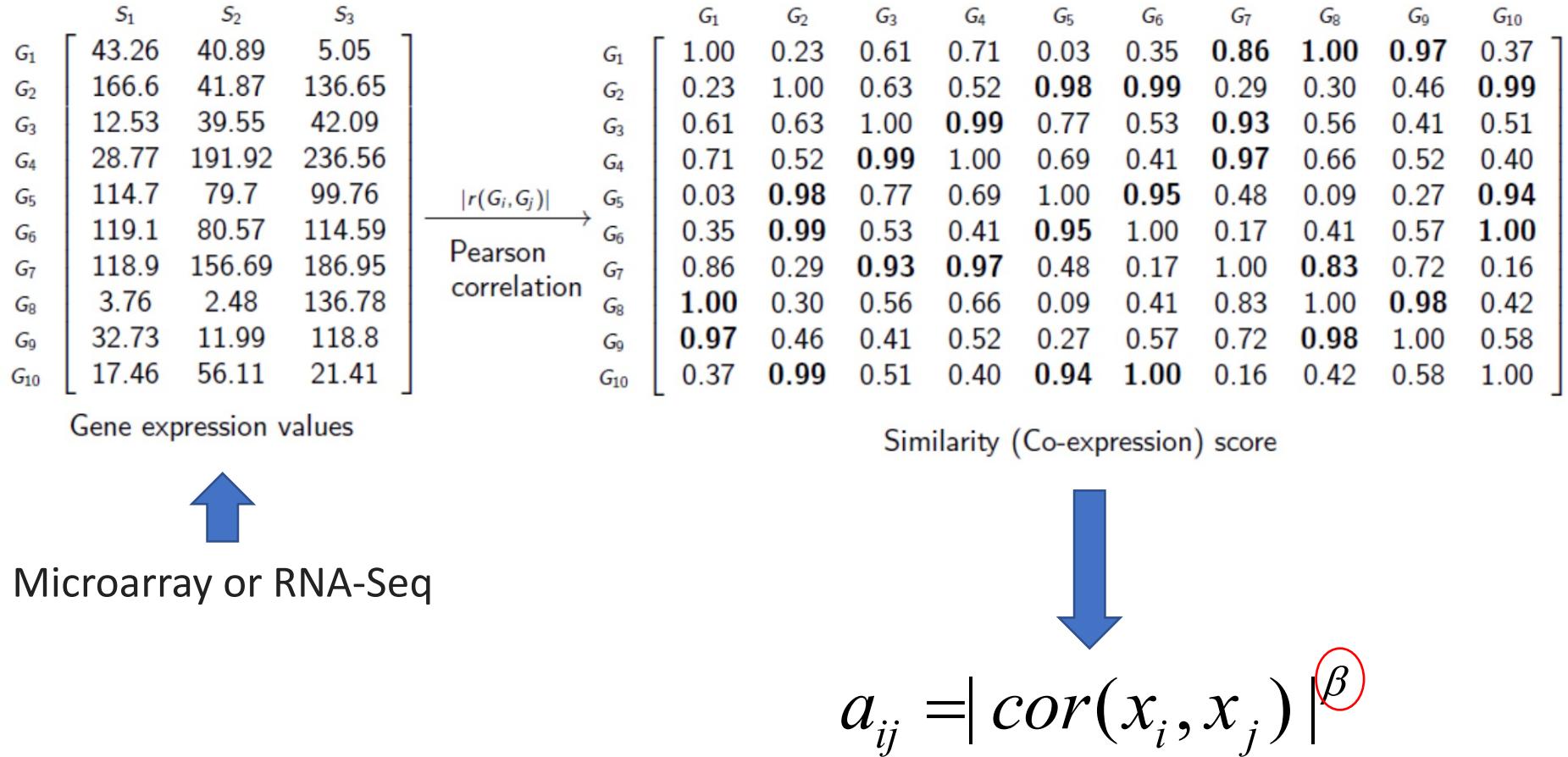
Identify Network Modules (Clustering)

Visualization of Co-expression network

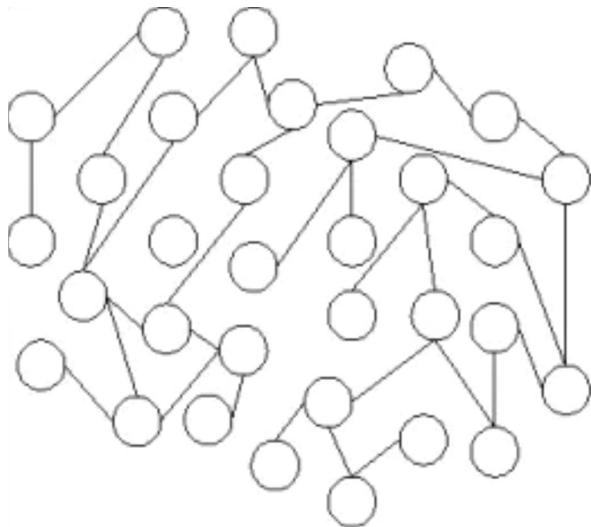
Relate the Network Concepts to External Gene or Sample Information



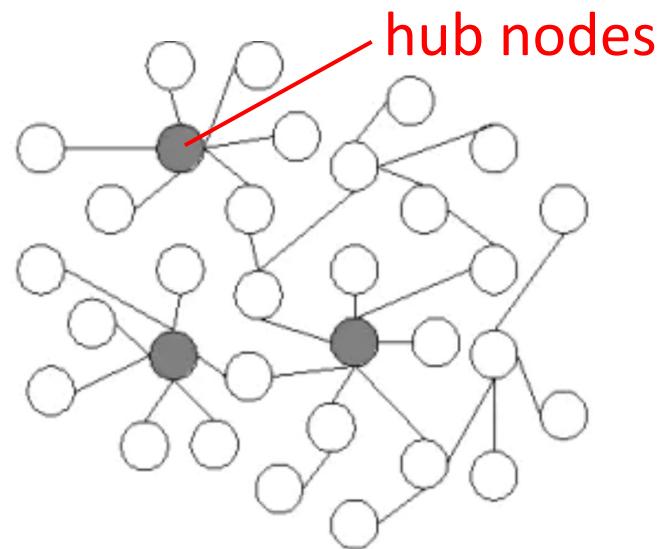
Define a Gene Co-expression Similarity



Define a soft thresholding power (β)



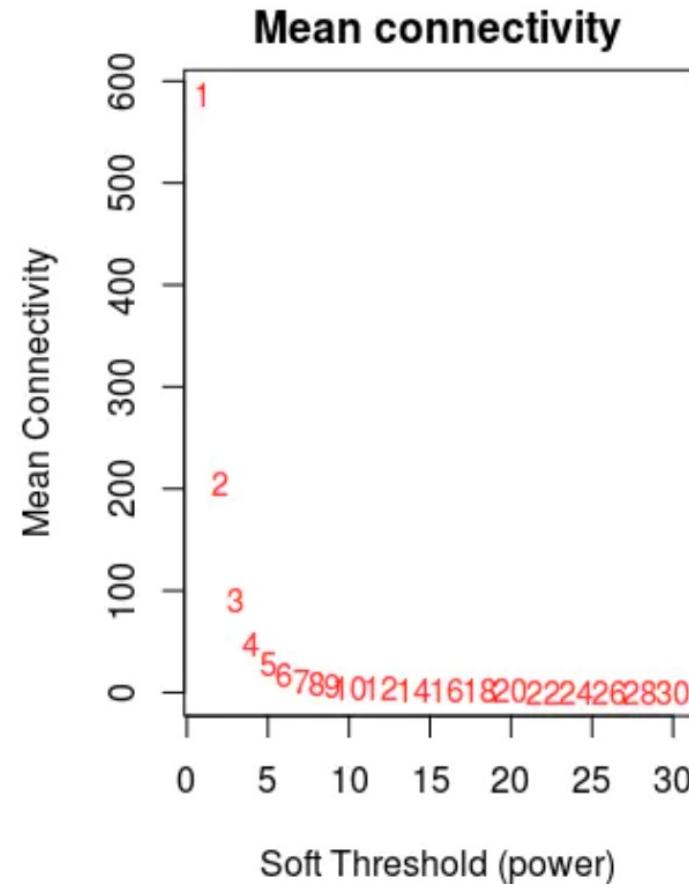
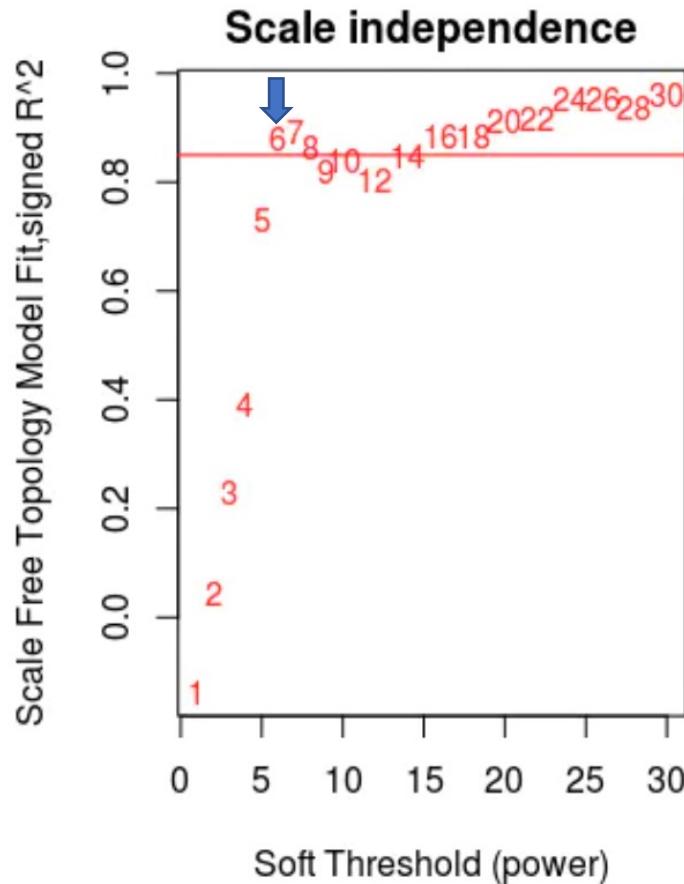
(a) Random network



(b) Scale-free network

- A scale-free network is a network whose degree distribution follows a power law.
- Unlike random network, some nodes in scale-free network have many more connections than others, which is called “hubs”
- Real-world networks are often claimed to be scale free

scale-free topology criterion



We need to find a best β that leads to a network that satisfies scale free topology most (R^2). Consequently, the **lowest β value** that leads to a R^2 of 0.8 (or the highest possible model fit) was used.

Calculate Adjacency Matrix with β

Correlation matrix

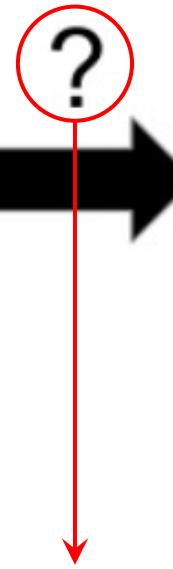
G 1	1	0.8	1	-0.5	-0.2	-0.5	-0.2	-0.02	0.08
G 2	0.8	1	0.7	-0.4	-0.5	-0.5	-0.7	-0.5	-0.6
G 3	1	0.7	1	-0.6	-0.2	-0.5	-0.3	-0.06	-0.1
G 4	-0.5	-0.4	-0.6	1	0.8	0.9	0.5	0.3	0.3
G 5	-0.2	-0.5	-0.2	0.8	1	0.9	0.7	0.6	0.6
G 6	-0.5	-0.5	-0.5	0.9	0.9	1	0.4	0.2	0.2
G 7	-0.2	-0.7	-0.3	0.5	0.7	0.4	1	1	1
G 8	-0.02	-0.5	0.06	0.3	0.6	0.2	1	1	1
G 9	-0.08	-0.6	-0.1	0.3	0.6	0.2	1	1	1

G¹ G² G³ G⁴ G⁵ G⁶ G⁷ G⁸ G⁹

Network adjacency

G 1	1	0.6	1	0	0	0	0	0	0
G 2	0.6	1	0.5	0	0	0	0	0	0
G 3	1	0.5	1	0	0	0	0	0	0
G 4	0	0	0	1	0.7	0.8	0.3	0.1	0.1
G 5	0	0	0	0.7	1	0.8	0.5	0.3	0.4
G 6	0	0	0	0.8	0.8	1	0.2	0.03	0.05
G 7	0	0	0	0.3	0.5	0.2	1	0.9	0.9
G 8	0	0	0	0.1	0.3	0.03	0.9	1	1
G 9	0	0	0	0.1	0.4	0.05	0.9	1	1

G¹ G² G³ G⁴ G⁵ G⁶ G⁷ G⁸ G⁹



$$a_{ij} = |cor(x_i, x_j)|^\beta$$

a_{ij} : adjacency matrix
between gene i and gene j



k_i : connectivity of gene i

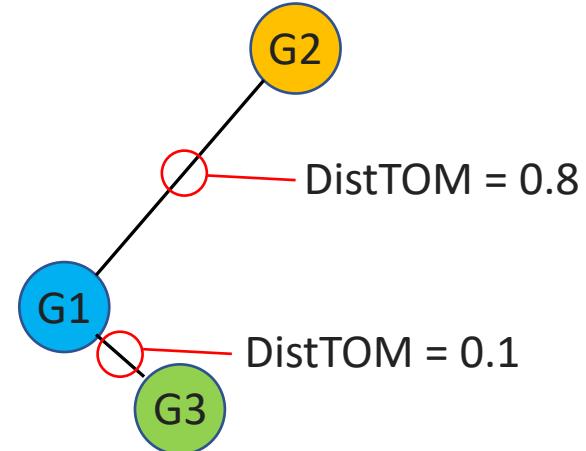
$$k_i = \sum_j a_{ij}$$

Define a Measure of Mode Dissimilarity

Topological Overlap leads to a network distance measure
(Ravasz et al., 2002)

$$TOM_{ij} = \frac{\sum_u a_{iu} a_{uj} + a_{ij}}{\min(k_i, k_j) + 1 - a_{ij}}$$

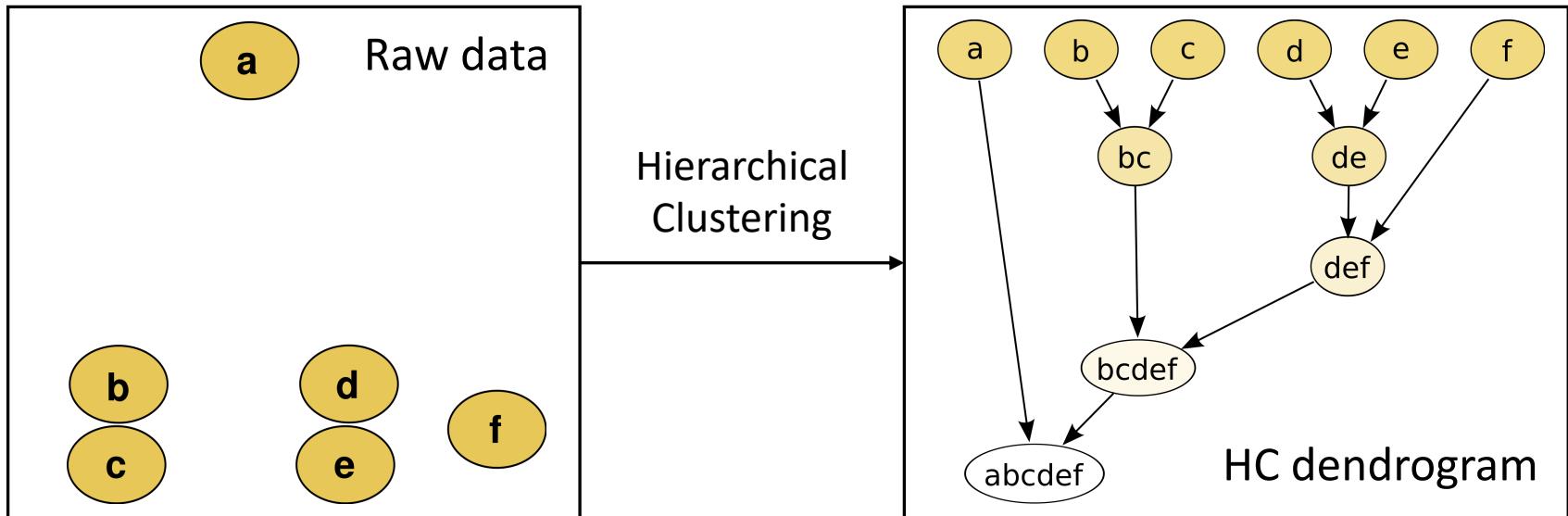
$$DistTOM_{ij} = 1 - TOM_{ij}$$



- The TOM describes how well connected the genes are in respect of how many neighbors they share
- The TOM matrix was transformed to a dissimilarity TOM (distTOM, 1-TOM), where a high connectivity produces a low number and no connectivity gives us a value of 1, or close.

Identify Network Modules

Hierarchical Clustering (HC) is a method of cluster analysis which seeks to build a hierarchy of clusters (Rokach 2005. 321-352.)



(https://en.wikipedia.org/wiki/Hierarchical_clustering)

Euclidean distance:

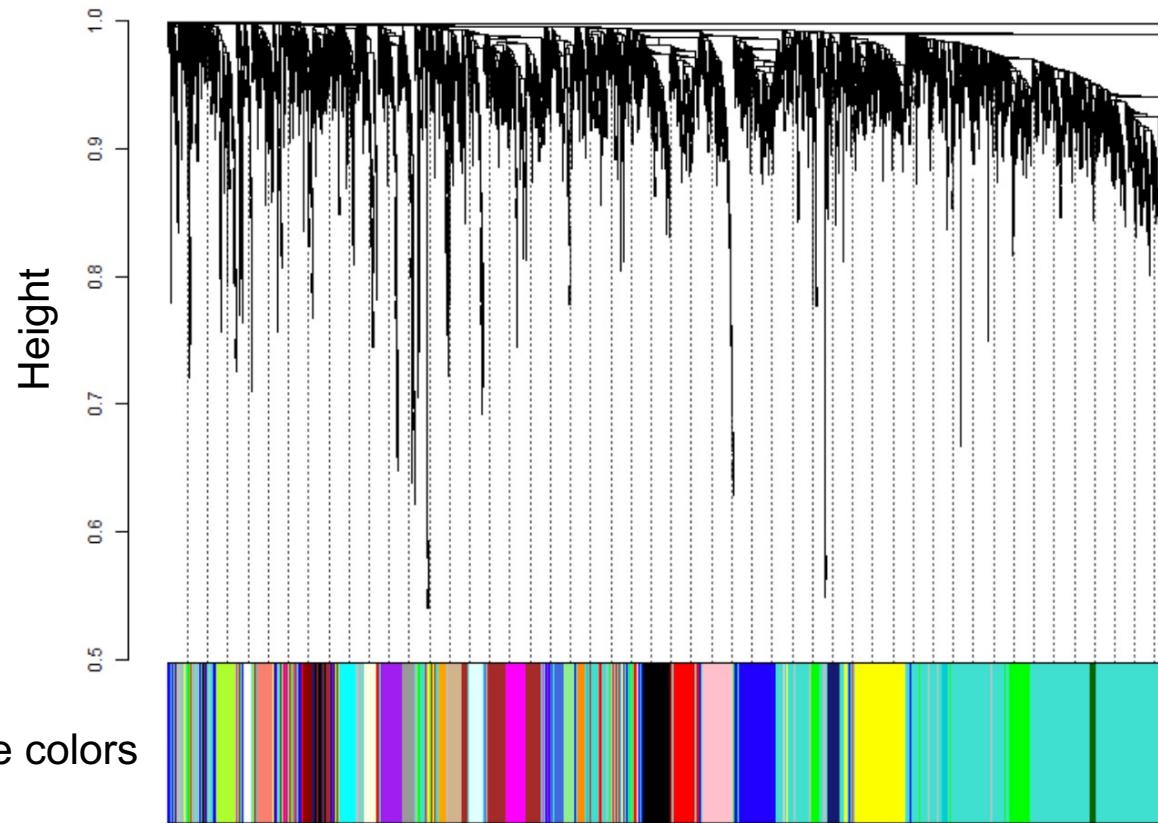
$$\|a - b\|_2 = \sqrt{\sum_i (a_i - b_i)^2}$$

$a_i : DistTOM$ of gene a

$b_i : DistTOM$ of gene b

Clustering of the genes in the distTOM matrix produces clusters of genes sharing many common neighbors into the same module.

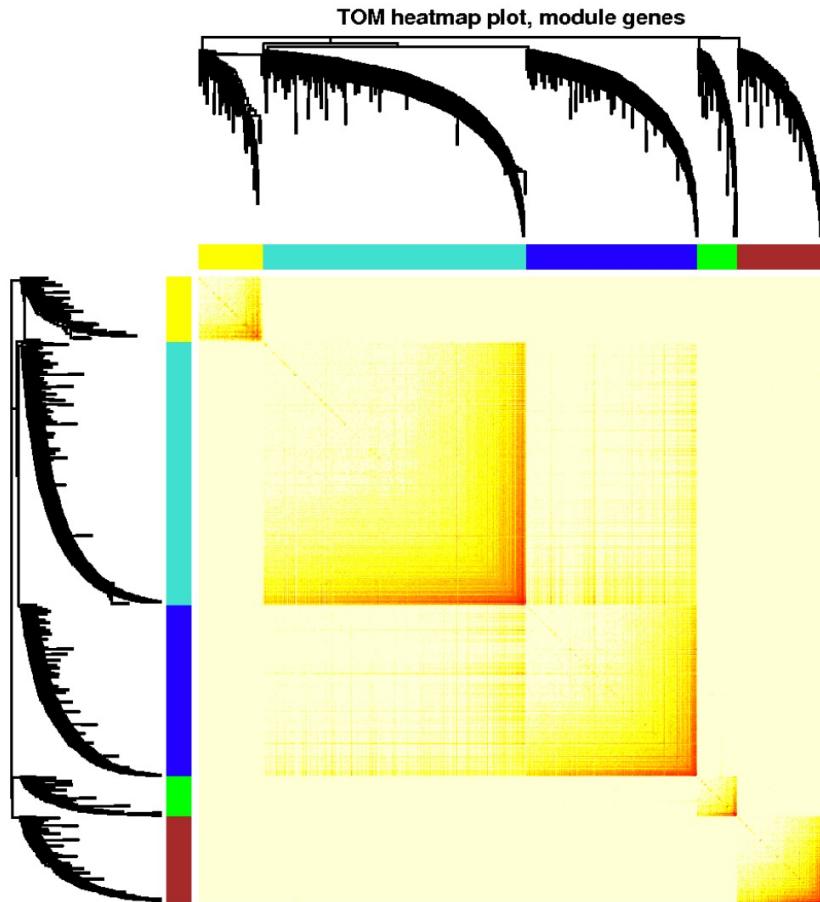
Hierarchical Cluster Dendrogram



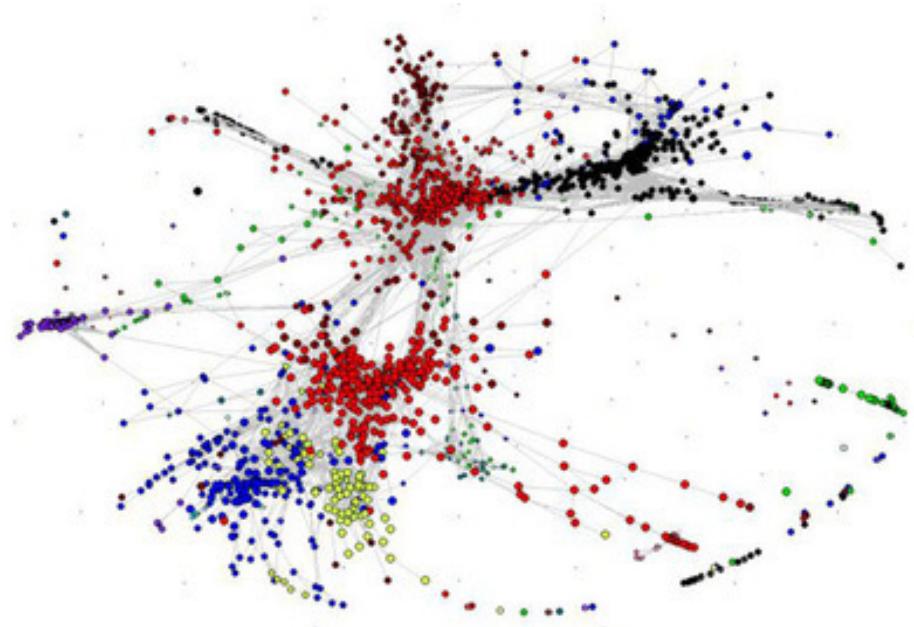
The 1st principal component for each co-expression module was defined as module **eigengene**, which represents the whole gene expression patterns for each module

Visualization of Co-expression network

Heatmap plot of the TOM



Network Visualization



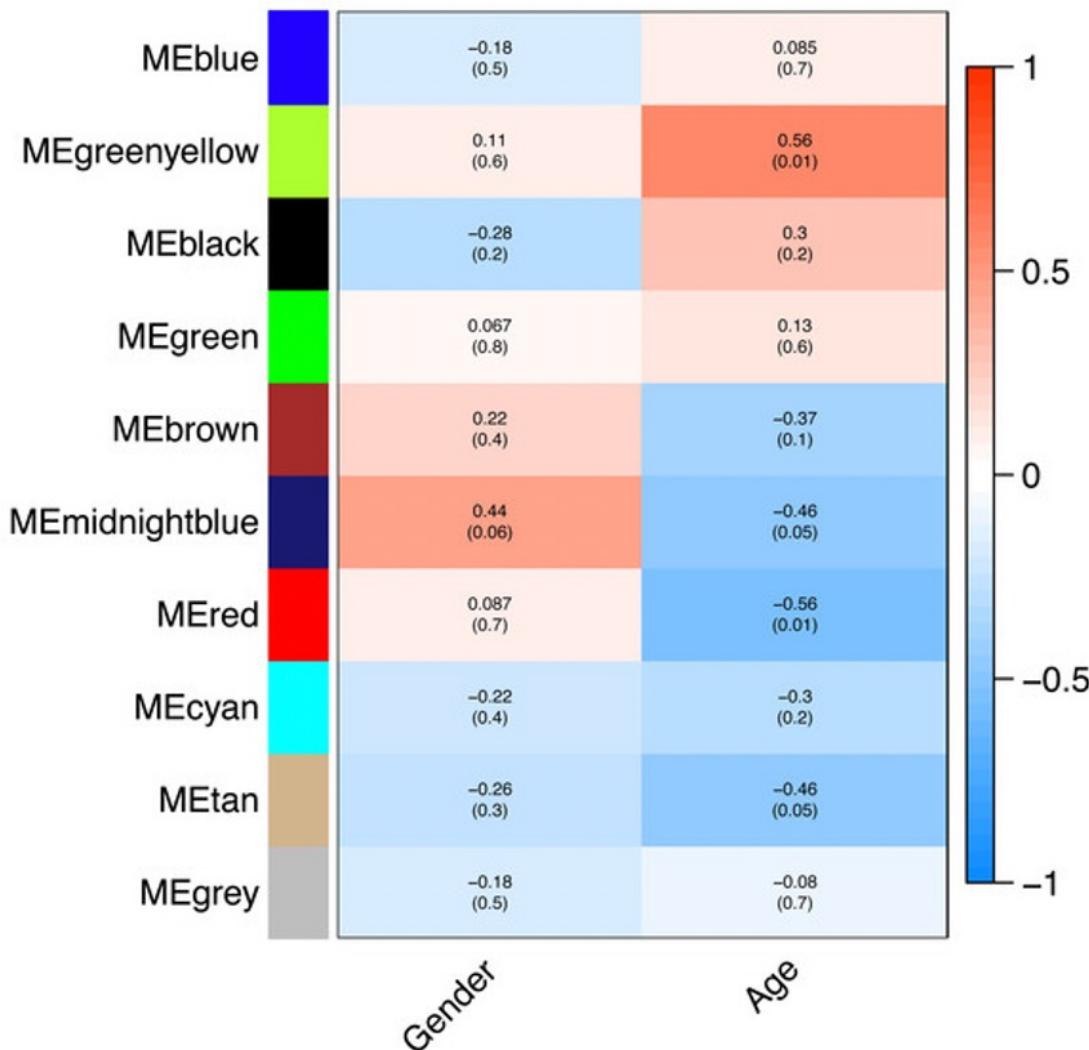
Cytoscape (<https://cytoscape.org>)

Gephi (<https://gephi.org/>)

(<https://Horvath.genetics.ucla.edu/html/CoexpressionNetwork/Rpackages/WGCNA/Tutorials/>)

Relate the Network Concepts to External Information

Module–trait relationships



- Module-trait relationships analysis
- GO terms analysis of genes in specific modules
- Functional analysis of hub genes
-

Contents

- Introduction of gene co-expression network (GCN)
- Steps to construct a gene co-expression network
- Gene co-expression network analysis methods
- Weighted Gene Co-Expression Network Analysis (WGCNA)
- Tutorial: Run WGCNA on your own laptop

Software requirement

R project or R Studio (<https://www.rstudio.com/>)

Gephi (<https://gephi.org/>)

Tutorial input data

Gene expressions of 56 breast cancer cell lines
(GSE48213_gene_expression_fpkm.txt)

Traits of 56 breast cancer cell lines (GSE48213_trait.txt)

Tutorial codes

WGCNA_GSE48213_coexp.R

(github)

1. Install & load WGCNA package

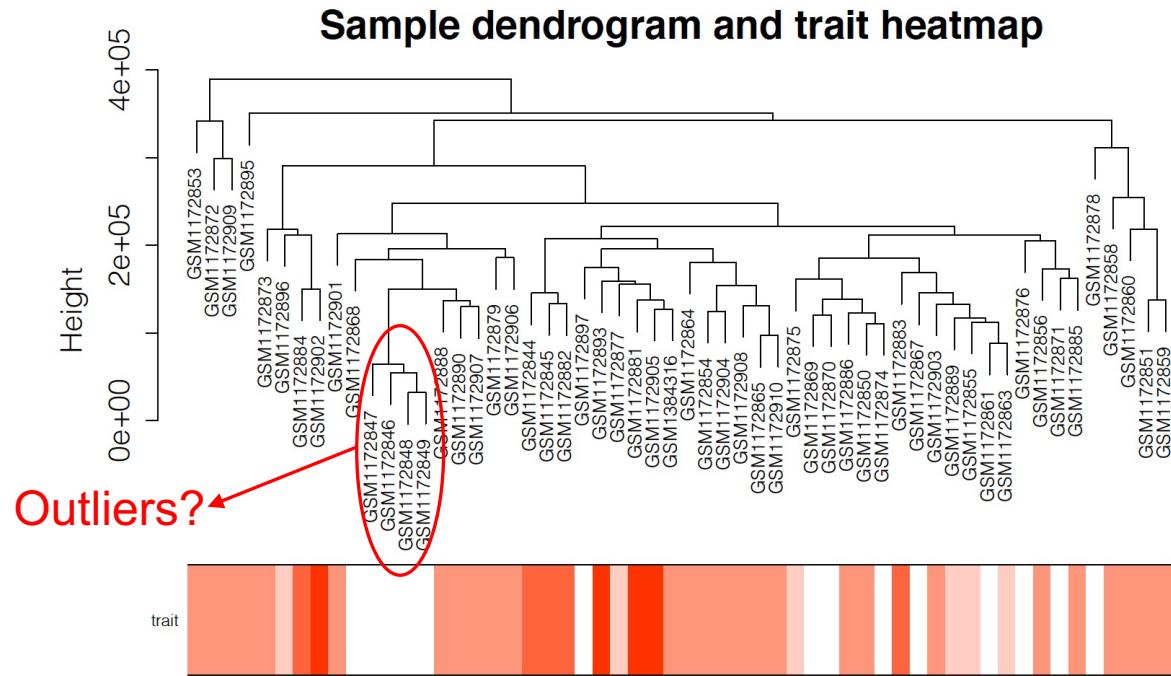
Code Line 1 - Line 27

2. Load tutorial input data

Code Line 29 – Line 44

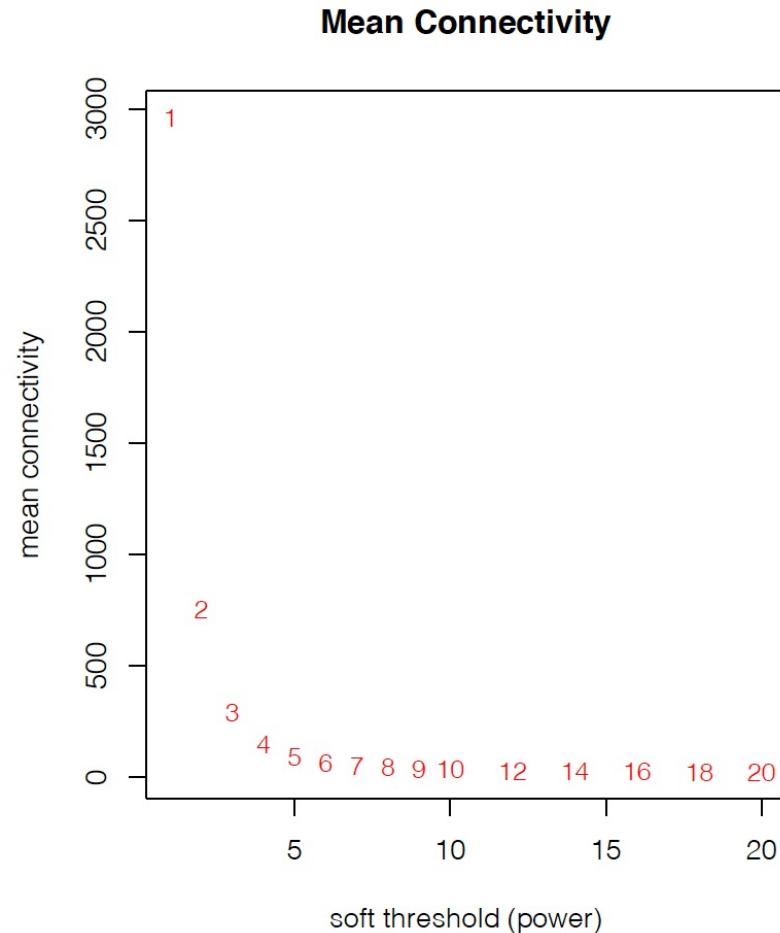
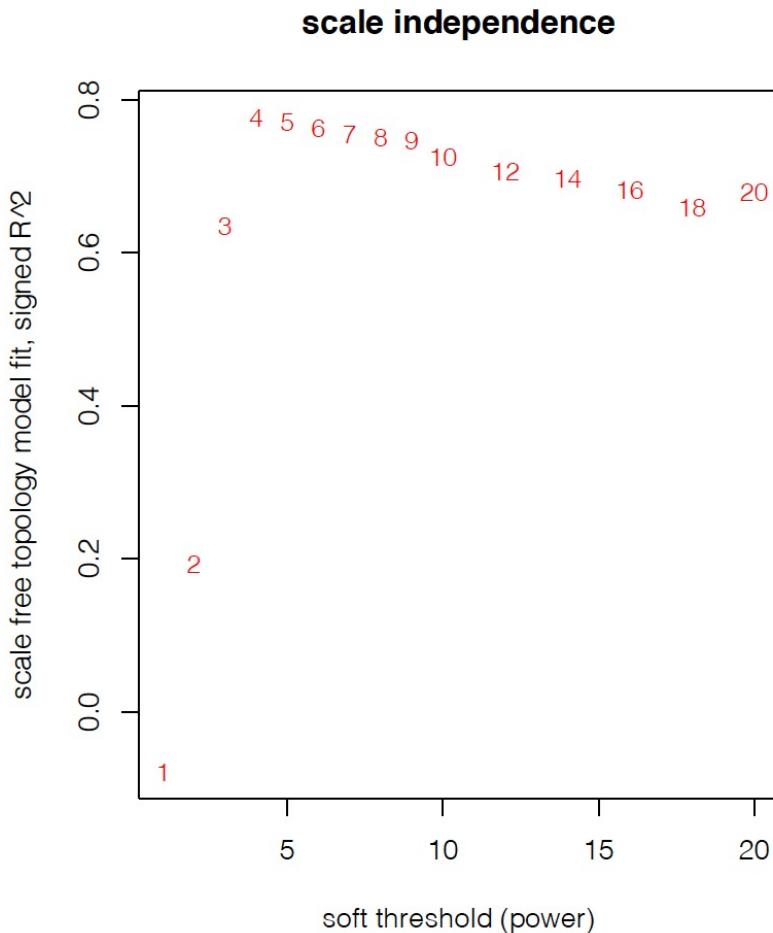
3. Sample clustering to identify outliers

Code Line 44 – Line 54



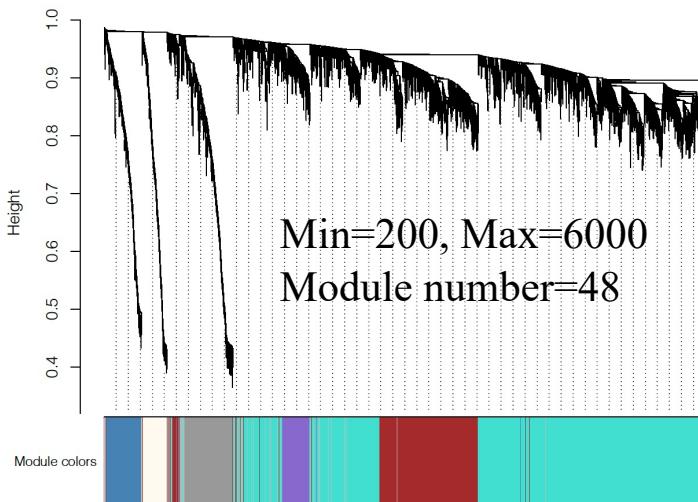
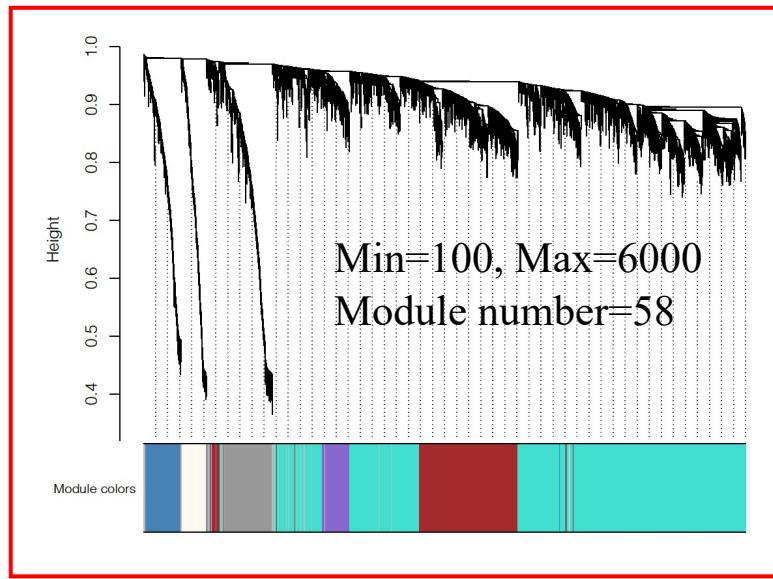
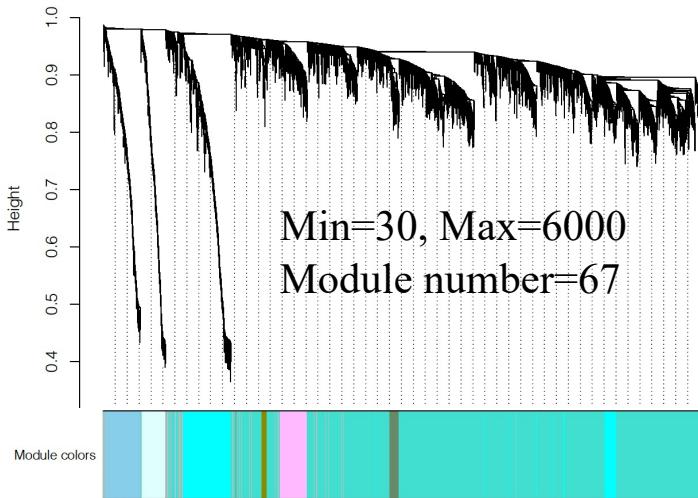
4. Define the soft thresholding power (β)

Code Line 56 – Line 76



4. Construction gene co-expression modules

Code Line 78 – Line 89



The parameter **minModuleSize** and **maxBlockSize** can be used to limit the size of each module, which will result in different module numbers finally

5. Save module results

Code Line 91 – Line 101

Coexp_modules_for_each_gene_P4M100.txt ————— Module information for each gene

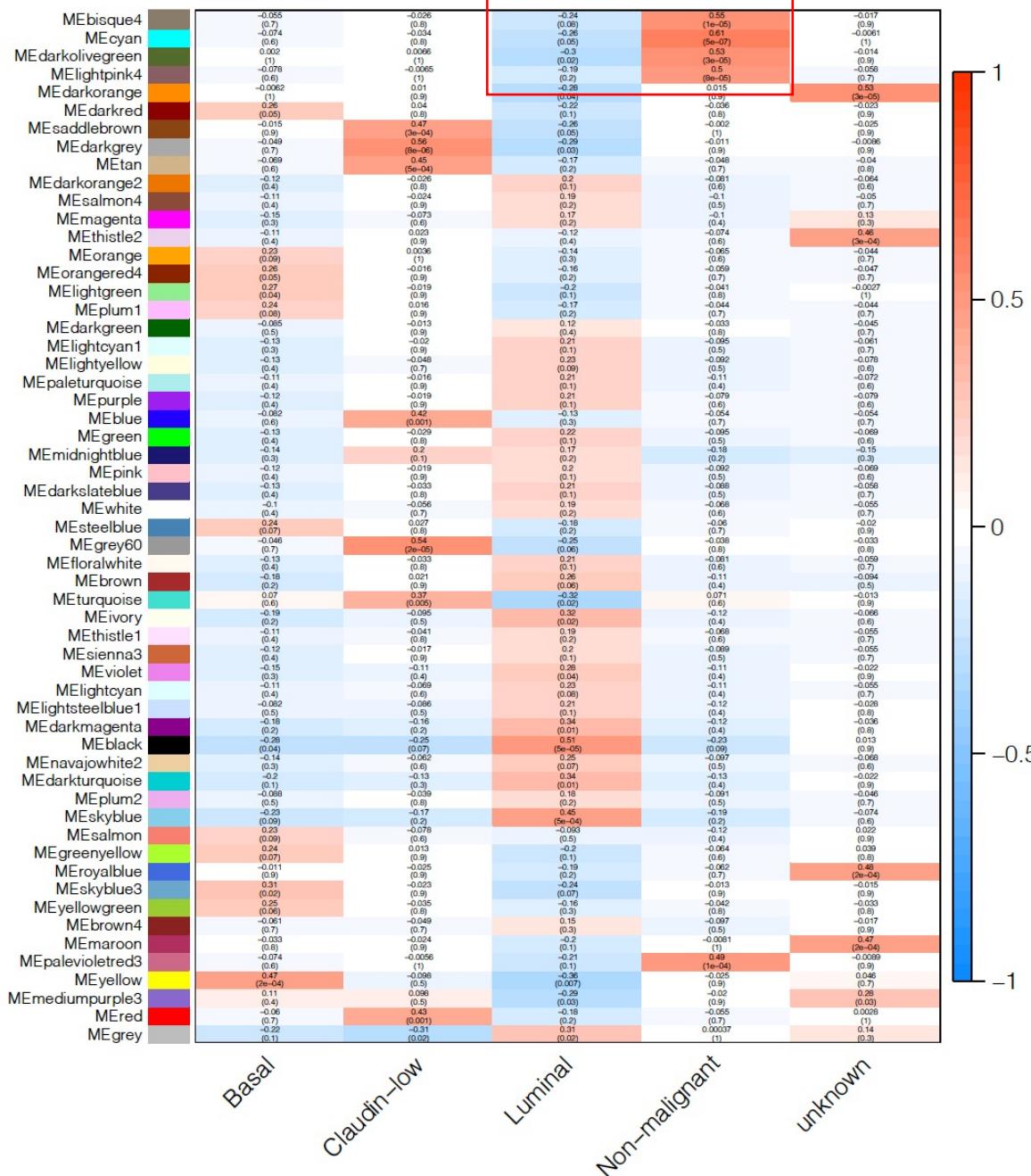
Summary_of_coexp_modules_P4M100.txt ————— Summary of gene numbers in each module

Eigengene_for_each_module_P4M100.txt ————— Eigengene expressions for each module (ME)

6. Quantifying module–trait associations

Code Line 103 – Line 136

Module-trait associations



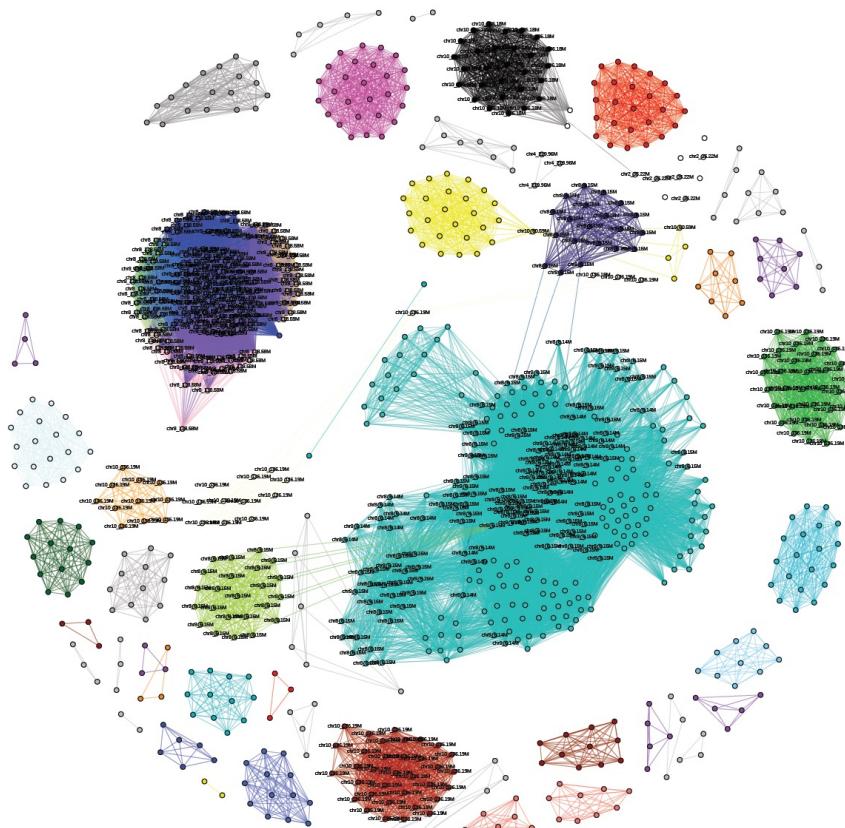
7. Exporting Files for Cytoscape Visualization

Code Line 138 – Line 148

Cytoscape (<https://cytoscape.org>)

8. Network visualization with Gephi

Folder: `gephi_visualization`



A short methodological summary of the publications.

- WGCNA methods
 - Horvath S (2011) Weighted Network Analysis. Applications in Genomics and Systems Biology. Springer Book. ISBN: 978-1-4419-8818-8
 - Zhang B, Horvath S (2005) "A General Framework for Weighted Gene Co-Expression Network Analysis", Statistical Applications in Genetics and Molecular Biology: Vol. 4: No. 1, Article 17
 - Langfelder P, Horvath S (2008) WGCNA: an R package for Weighted Correlation Network Analysis. BMC Bioinformatics. 2008 Dec 29;9(1):559. PMID: 19114008 PMCID: PMC2631488
 - Langfelder P et al (2011) Is my network module preserved and reproducible? PLoS Comp Biol. 7(1): e1001057. PMID: 21283776
- Math and WGCNA:
 - Horvath S, Dong J (2008) Geometric Interpretation of Gene Co-Expression Network Analysis. PLoS Computational Biology. 4(8): e1000117. PMID: 18704157
- Empirical evaluation of WGCNA
 - Langfelder P, et al (2013) When Is Hub Gene Selection Better than Standard Meta-Analysis? PLoS ONE 8(4): e61505.
 - Song L, Langfelder P, Horvath S. (2012) Comparison of co-expression measures: mutual information, correlation, and model based indices. BMC Bioinformatics;13(1):328. PMID: 23217028
- What is the topological overlap measure? Empirical studies of the robustness of the topological overlap measure:
 - Yip A, Horvath S (2007) Gene network interconnectedness and the generalized topological overlap measure. BMC Bioinformatics 8:22
- Dynamic branch cutting:
 - Langfelder P, Zhang B, Horvath S (2008) Defining clusters from a hierarchical cluster tree: the Dynamic Tree Cut package for R. Bioinformatics.;24(5):719-20. PMID: 18024473
- Gene screening based on intramodular connectivity identifies brain cancer genes that validate.
 - Horvath S, Zhang B, Carlson M, Lu KV, Zhu S, Felciano RM, Laurance MF, Zhao W, Shu, Q, Lee Y, Scheck AC, Liau LM, Wu H, Geschwind DH, Febbo PG, Kornblum HI, Cloughesy TF, Nelson SF, Mischel PS (2006) "Analysis of Oncogenic Signaling Networks in Glioblastoma Identifies ASPM as a Novel Molecular Target", PNAS | November 14, 2006 | vol. 103 | no. 46 | 17402-17407
- How to integrate SNP markers into weighted gene co-expression network analysis?
 - Plaisier CL et al Pajukanta P (2009) A systems genetics approach implicates USF1, FADS3 and other causal candidate genes for familial combined hyperlipidemia. PLoS Genetics;5(9):e1000642 PMID: 19750004
- Differential network analysis:
 - Fuller TF, Ghazalpour A, Aten JE, Drake TA, Lusis AJ, Horvath S (2007) "Weighted Gene Co-expression Network Analysis Strategies Applied to Mouse Weight", Mammalian Genome.

More WGCNA R-package tutorials:

<http://www.genetics.ucla.edu/labs/horvath/CoexpressionNetwork>

My E-mail address:

ksuhecheng90@gmail.com

Thank you