

Reproduce Part of **WGNA** in Finding Weight-related Gene Modules

Bioinformatics Final Project

12212859 Sijie Li

June 3th, 2025

Curious When Using WGCNA...

- Network construction (esp. soft threshold)?

BMC Bioinformatics

 BioMed Central

Software

WGCNA: an R package for weighted correlation network analysis

Peter Langfelder¹ and Steve Horvath^{*2}

Open Access

2008, Cited: 22,262

A General Framework for Weighted Gene Co-Expression Network Analysis

Bin Zhang, Departments of Human Genetics and Biostatistics, University of California at Los Angeles
Steve Horvath, Departments of Human Genetics and Biostatistics, University of California at Los Angeles

2005, Cited: 5,970

- Dynamic Tree cut?

BIOINFORMATICS APPLICATIONS NOTE

Vol. 24 no. 5 2008, pages 719–720
doi:10.1093/bioinformatics/btm563

Gene expression

Defining clusters from a hierarchical cluster tree: the Dynamic Tree Cut package for R

Peter Langfelder^{1,†}, Bin Zhang^{2,†} and Steve Horvath^{1,*}

2008, Cited: 2,2136

Dynamic Tree Cut: in-depth description, tests and applications

Peter Langfelder ^{a*}; Bin Zhang ^{b*}; Steve Horvath ^{a†}

^aDept. of Human Genetics, University of California at Los Angeles, CA 90095-7088

^bRosetta Inpharmatics-Merck Research Laboratories, Seattle, WA

September 12, 2007

Apply to a scenario

OPEN ACCESS Freely available online

PLOS GENETICS

Integrating Genetic and Network Analysis to Characterize Genes Related to Mouse Weight

Anatole Ghazalpour¹✉, Sudheer Doss²✉, Bin Zhang^{2,3}, Susanna Wang², Christopher Plaisier², Ruth Castellanos¹, Alec Brozell¹, Eric E. Schadt⁴, Thomas A. Drake⁵, Aldons J. Lusis^{1,2,6,7}✉, Steve Horvath^{2,3}*

2006, Cited: 544

➤ Data

No WT!

- female F2 ApoE -/- high fat Western diet mice livers
- dye-switch microarray
- gene expression “variation”

➤ Result

- weight-related gene module
- potential regulatory locus of the module

Reproduce

- Construct Network
- Module Detection
- mQTL

(Module Quantitative Trait Loci)

\log_{10}

Expression of gene A in sample1
—
average expression of gene A in
150 random selected female
samples

Data Obtain

➤ GSE2814

The screenshot shows the NCBI GEO Accession Display page for dataset GSE2814. The page header includes the NCBI logo and the GEO Gene Expression Omnibus logo. The main content area displays various details about the dataset, such as its status (Public on Jul 27, 2006), title (Expression profiling of liver tissue from (C57BL/6J X C3H/He)F2 mice on ApoE null backgrounds), organism (Mus musculus), experiment type (Expression profiling by array), summary (The (C57BL/6J X C3H/He)F2 intercross consists of 334 animals of both sexes. All are ApoE null and received a high fat Western diet from 8-24 weeks of age.), overall design (Livers from 311 F2 female and male mice (animals fed a high fat "Western" diet from 8-24 weeks of age.) derived from C57BL/6J and C3H/He parental strains with both on ApoE null backgrounds. All samples were compared to a common pool created from equal portions of RNA from each of the samples. Keywords=Genetics of Gene Expression Keywords=C57BL/6J Keywords=C3H/He), and citation(s) (Wang S, Yehya N, Schadt EE, Wang H et al. Genetic and genomic analysis of a fat mass trait with complex inheritance reveals marked sex specificity. *PLoS Genet* 2006 Feb;2(2):e15. PMID: 16462940; Yang X, Schadt EE, Wang S, Wang H et al. Tissue-specific expression and regulation of sexually dimorphic genes in mice. *Genome Res* 2006 Aug;16(8):995-1004. PMID: 16825664; Ghazalpour A, Doss S, Zhang B, Wang S et al. Integrating genetic and network analysis to characterize genes related to mouse weight. *PLoS Genet* 2006 Aug 18;2(8):e130. PMID: 16934000; Chen Y, Zhu J, Lunn PY, Yang X et al. Variations in DNA elucidate molecular networks that cause disease. *Nature* 2008 Mar 27;452(7186):429-35. PMID: 18344982; Wu S, Mar-Heyming R, Dugum EZ, Kolaitis NA et al. Upstream transcription factor 1 influences plasma lipid and metabolic traits in mice. *Hum Mol Genet* 2010 Feb 15;19(4):597-608. PMID: 19995791; van Nas A, Ingram-Drake L, Sinsheimer JS, Wang SS et al. Expression quantitative trait loci: replication, tissue- and sex-specificity in mice. *Genetics* 2010 Jul;185(3):1059-68. PMID: 20439777).

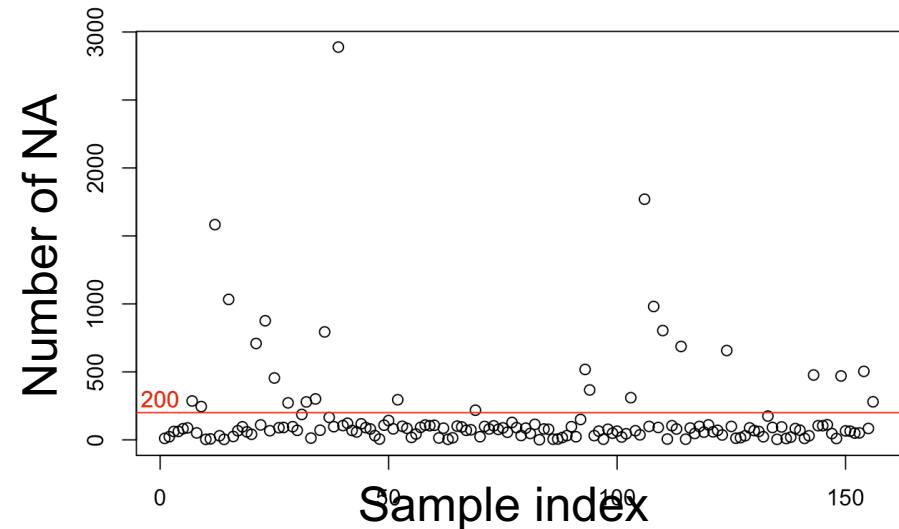
➤ Tutorial dataset

The screenshot shows a tutorial page titled "Tutorials for the WGCNA package" by Peter Langfelder and Steve Horvath. The page is hosted on a GitHub page (edo98811.github.io). It provides an overview of the WGCNA package, mentioning its compatibility with WGCNA version 1.13 and higher. The page includes a "Background and glossary" section and a "Network analysis of liver expression data from female mice: finding modules related to body weight" section. The "Background and glossary" section is described as a short text containing background information, an overview figure, and a short glossary of network analysis terms and concepts. The "Network analysis" section is described as a tutorial guiding the reader through the analysis of an empirical data set from livers of female mice of a specific F2 intercross.

- 3600 genes
- 135 samples

Process Data--Sample Reduction

1. Trim annotation → 23,388 genes × 311 sample matrix
2. Extract **female** sample
3. Filter NA, remove sample with
 - > 200 NA expression
 - any NA trait

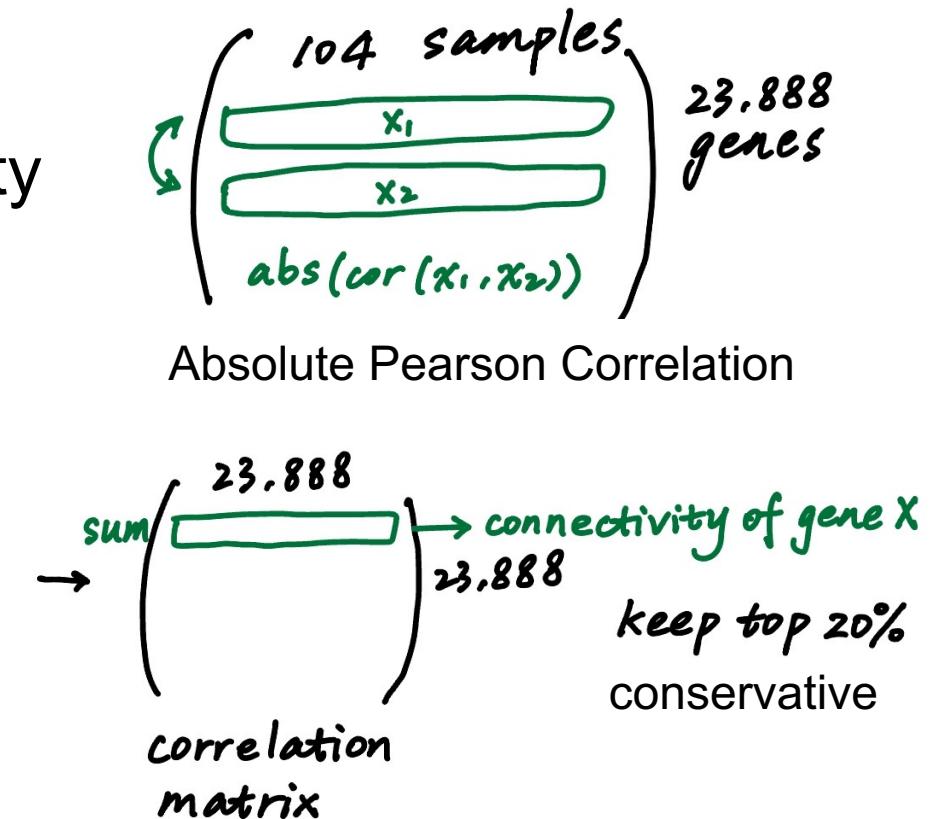
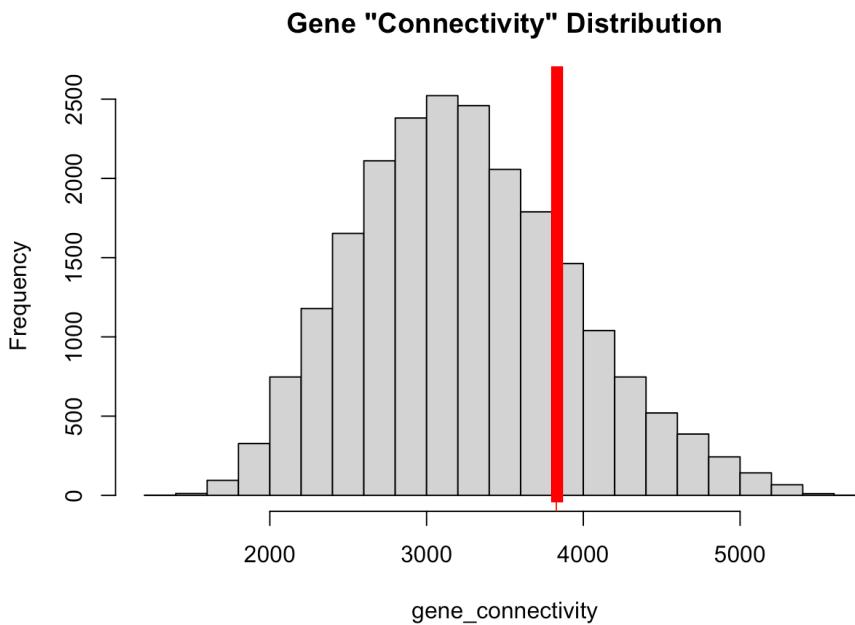


and comparison of the male network will be reported elsewhere. Only those mice with complete phenotype, genotype, and array data were used. This gave a final experimental sample of 135 female mice used

104 samples left

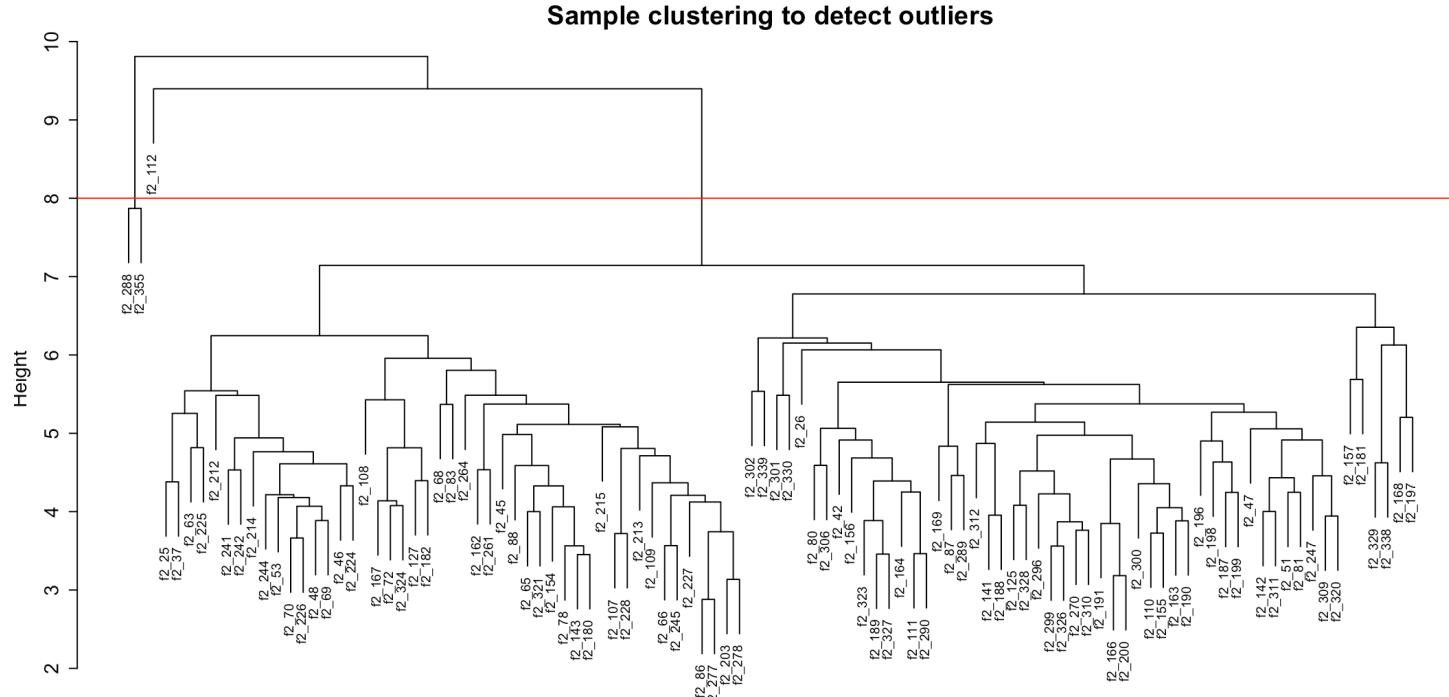
Process Data--Gene Reduction

4. Remove genes with NA records
5. Select genes with high connectivity



$$23,388 \times 311 \rightarrow 4,319 \times 104 \text{ VS } 3,600 \times 135$$

Check by Sample Clustering



Cluster the downstream analyze data

$23,388 \times 311 \rightarrow 4,319 \times 104 \text{ VS } 3,600 \times 135$

Error source 1: Data

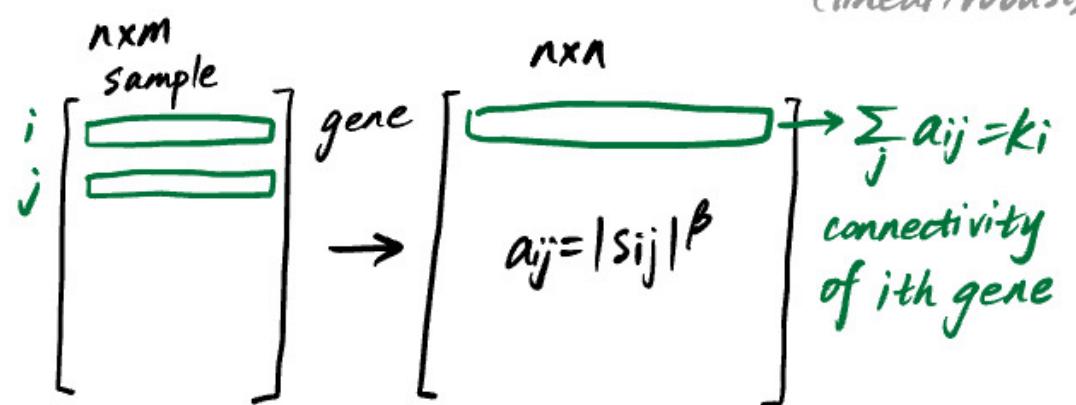
- Do not remove outliers
- Differ in sample filtration
- Differ in gene selection

💡 Use WGCNA to analyze my data to evaluate the effect of different data.

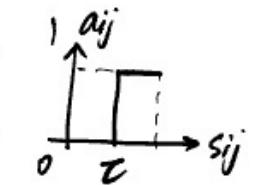
Construct Network--Overview

- Node (biological objects)
- Edge (**relationship** between biological quantitives of the objects)

► Co-expression Similarity ($\in [0,1]$) $\xrightarrow{A=f(s)}$ ► Adjacency Function (①) ② $f: [0,1] \rightarrow [0,1]$)
 $s_{ij} \equiv |\text{cor}(i,j)|$ → Pearson
 (linear, robust)

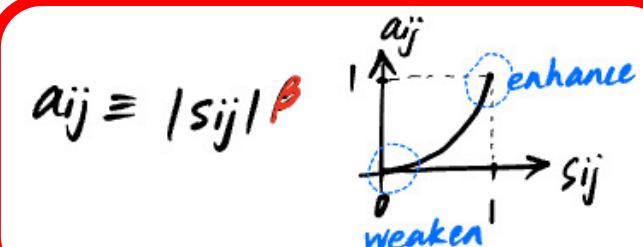


"Distance" of network (similarity)
 ① Hard $a_{ij} \equiv \begin{cases} 1 & \text{if } s_{ij} \geq \tau \\ 0 & \text{otherwise} \end{cases}$

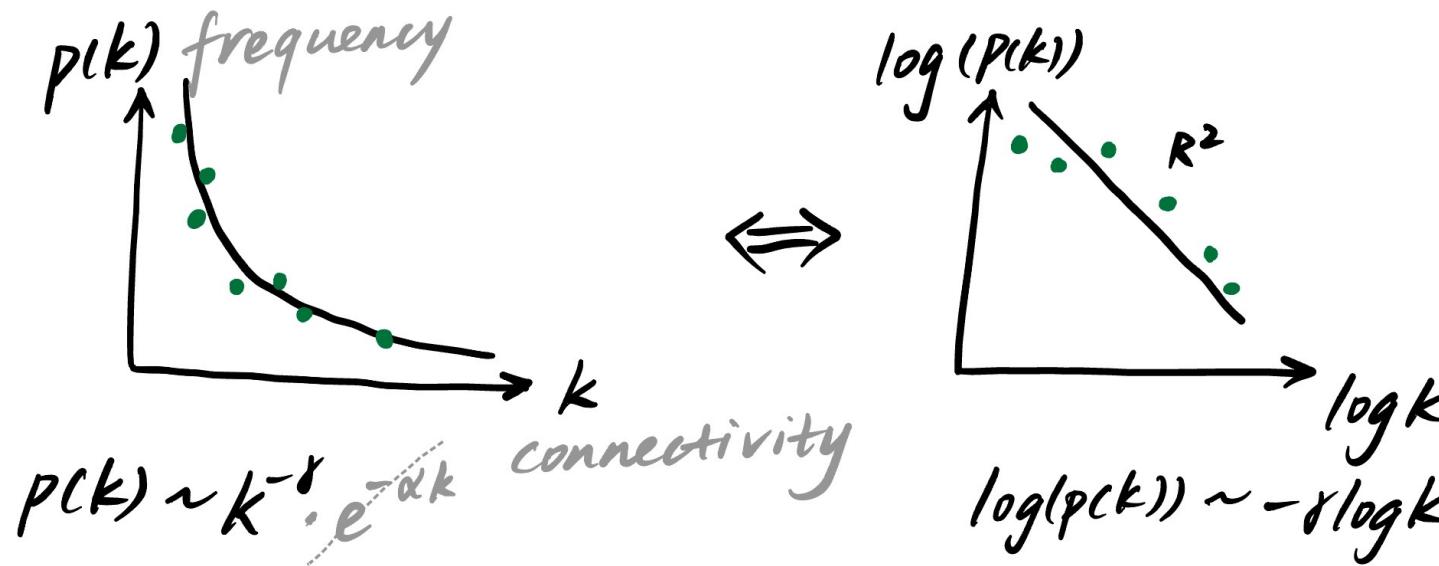
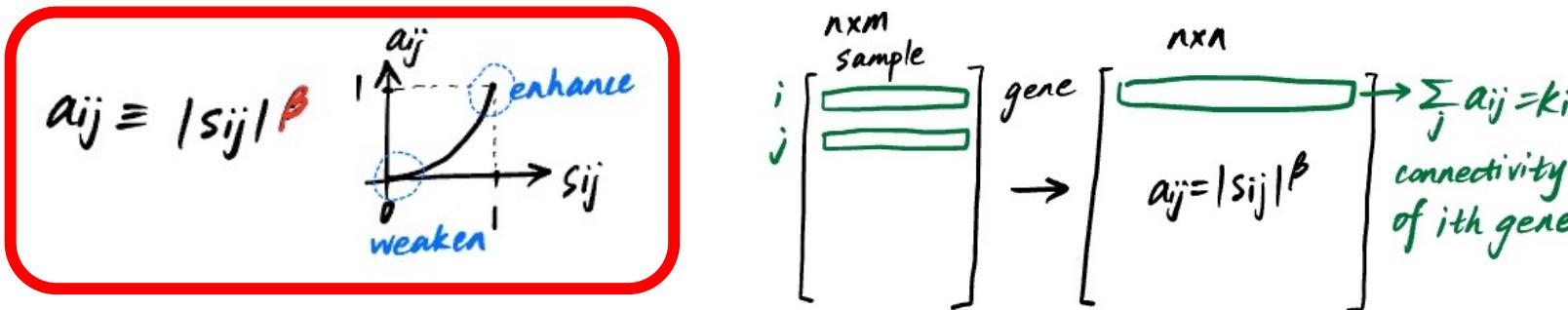


② Soft (continuous nature)

$$a_{ij} \equiv \frac{1}{1 + e^{-\alpha(s_{ij} - \tau_0)}}$$



Network Property

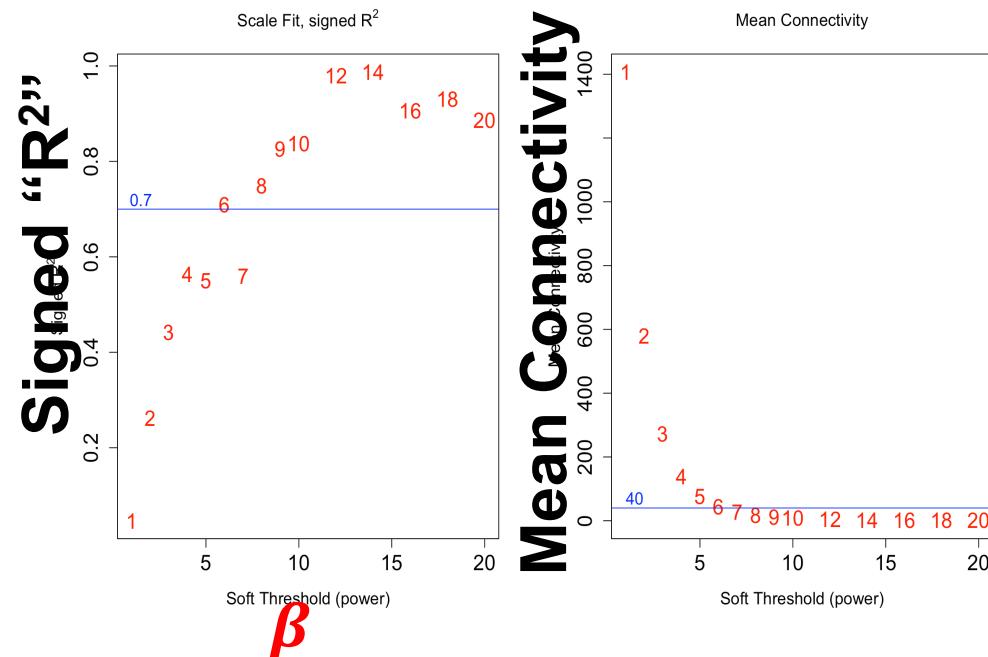
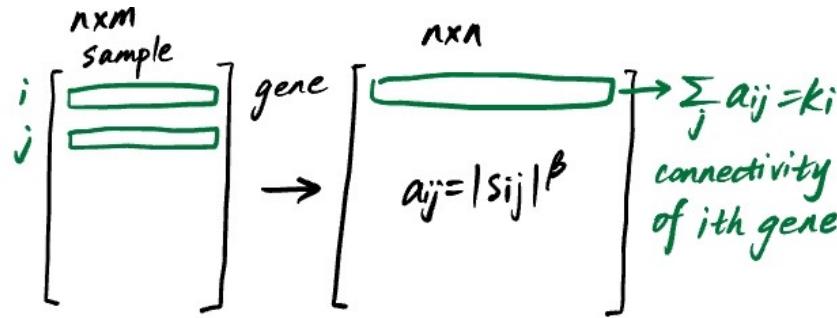
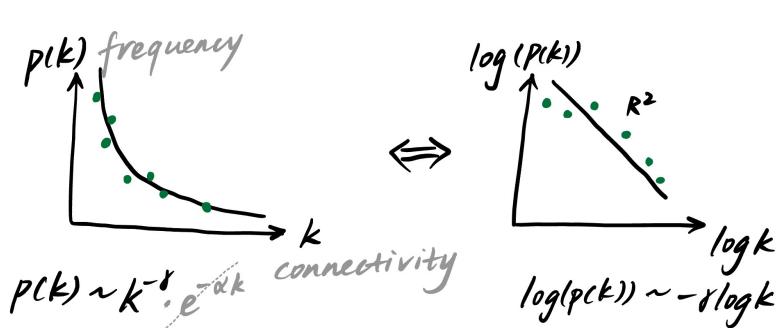


Scale-free

Property of many biological networks

- ✓ $R^2 >$ threshold
- ✓ High mean connectivity
- ✓ Signed “ R^2 ” (slope)

Soft Threshold β



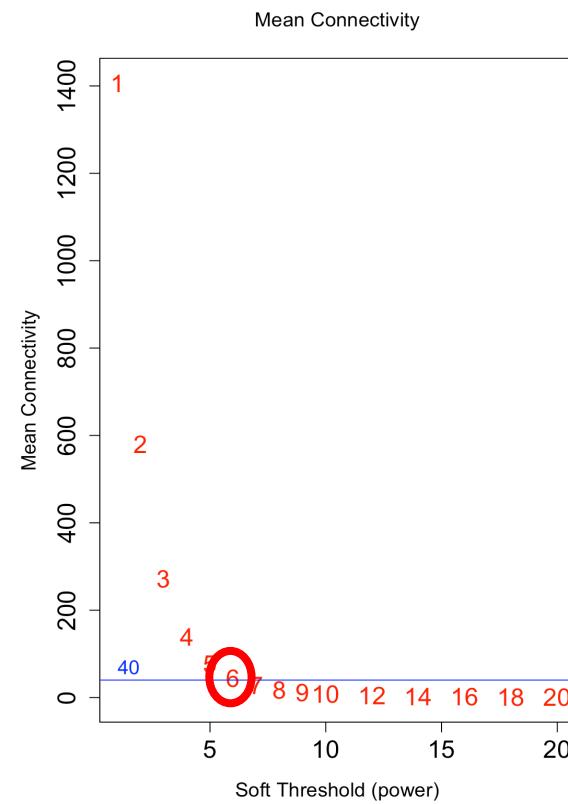
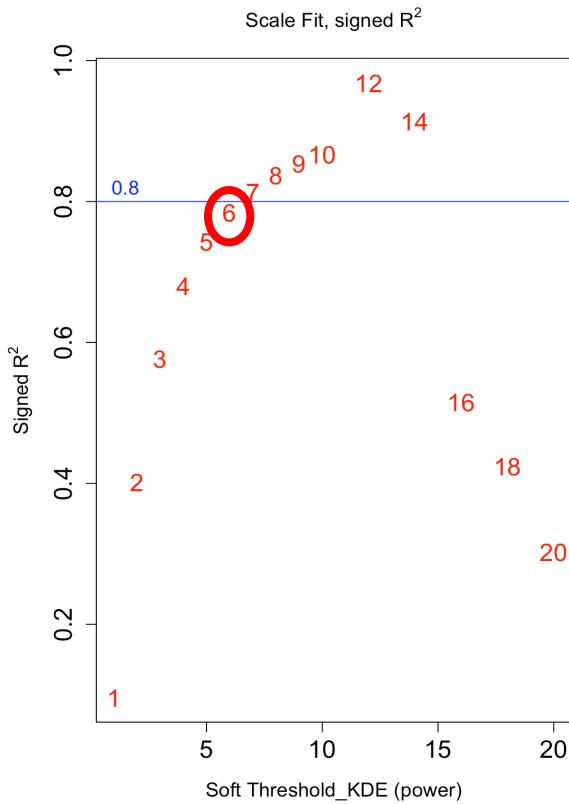
$p(k)$ for continuous k ?

➤ Bin, size?

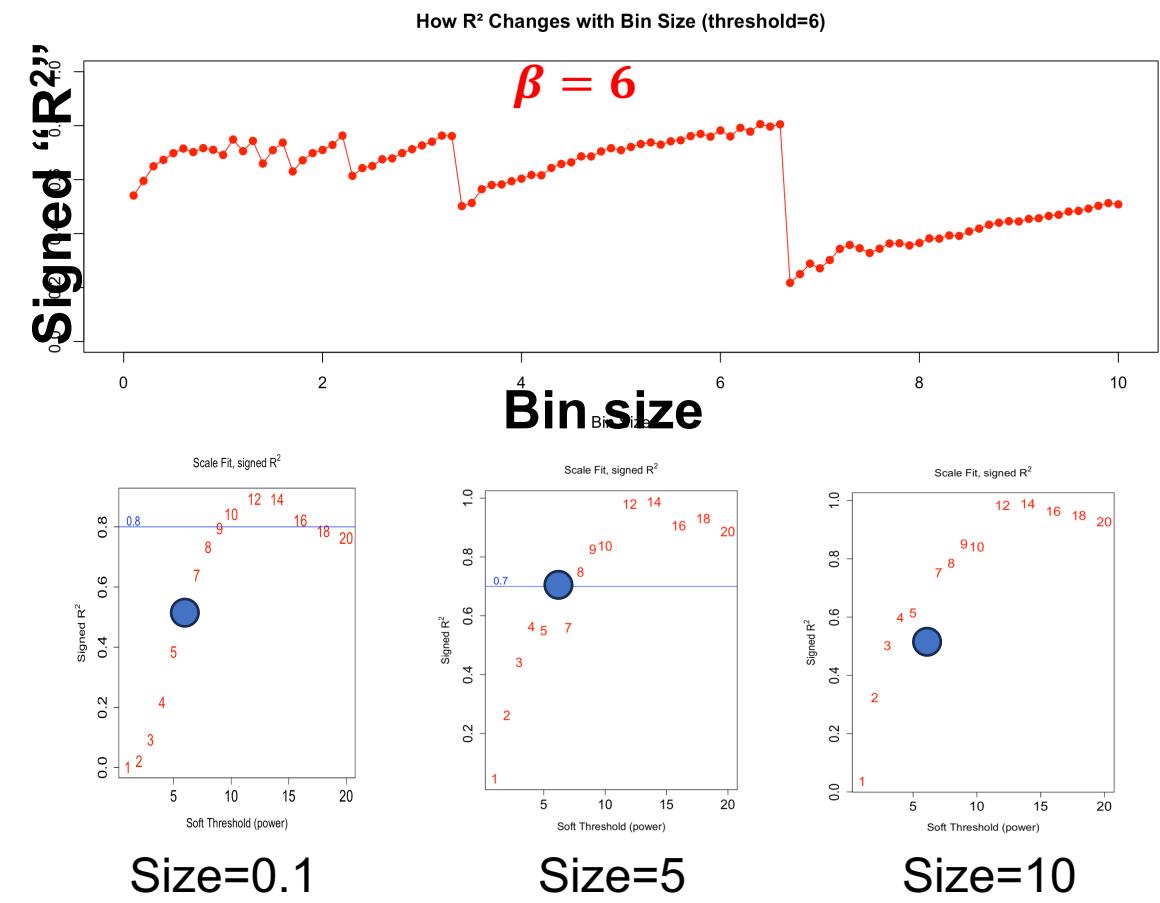
➤ Kernel Density Estimation (KDE)

Soft Threshold $\beta = 6$ Also indicated by the paper

➤ KDE

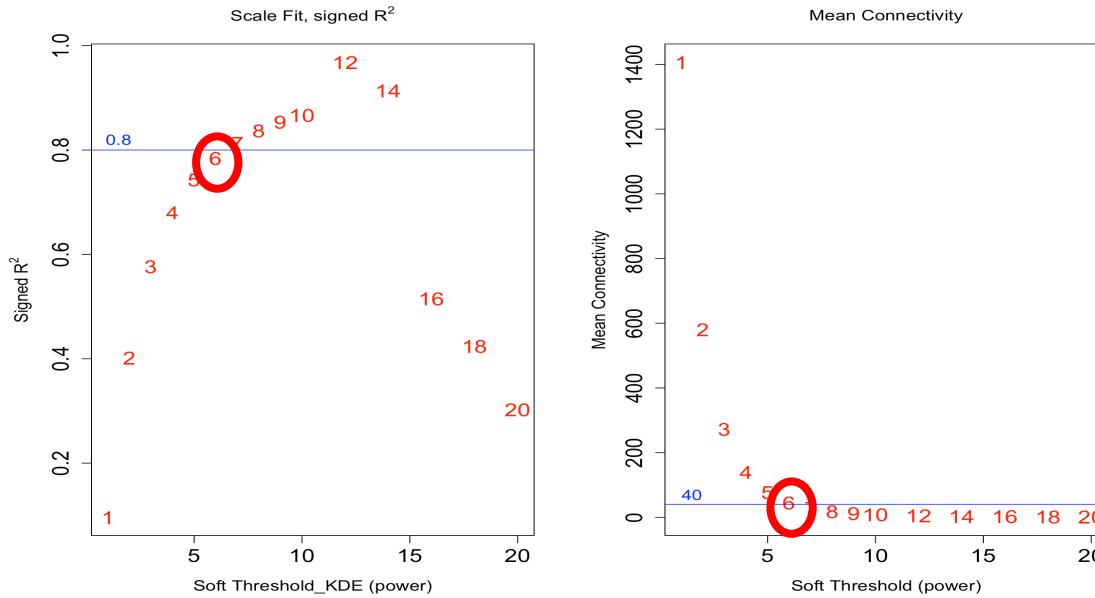


➤ Bin—size matters, a lot!



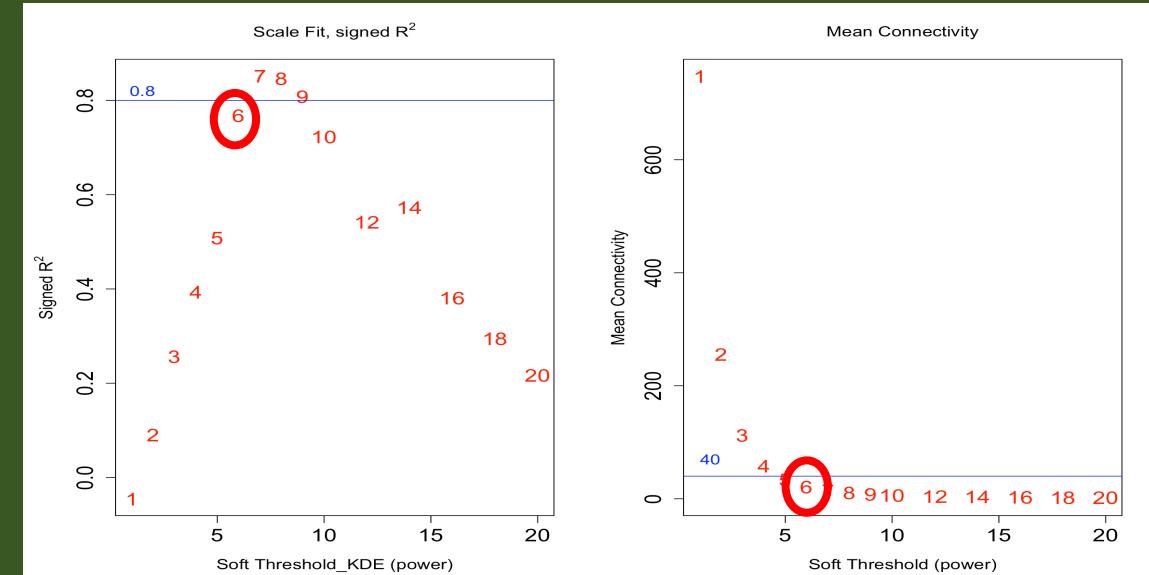
Soft Threshold $\beta = 6$

➤ KDE

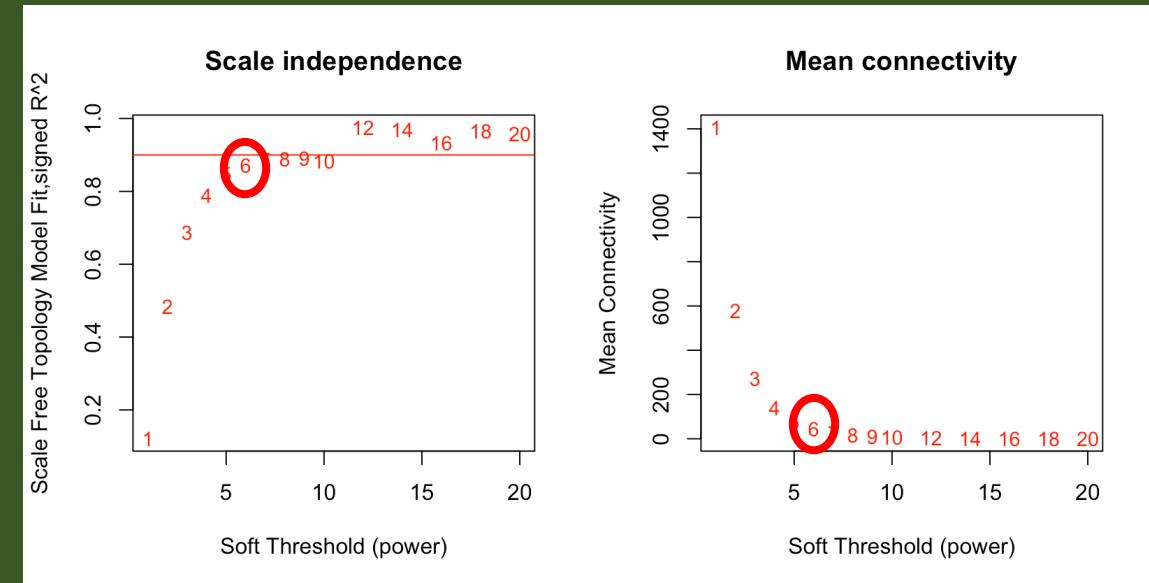


🤔 Though different, 6 can be pick out by all of it.
And my data seems to work, for the same soft threshold recommended by WGCNA

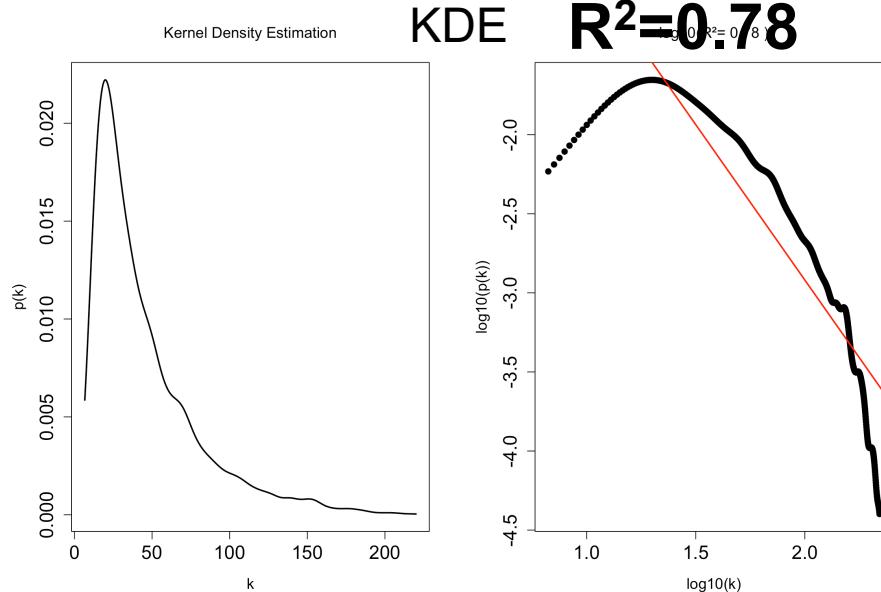
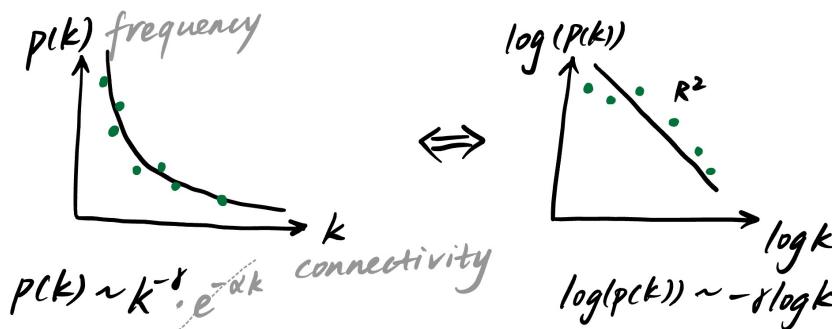
Tutorial Data, My Method (KDE)



WGCNA Package, My Data (recommend 6)

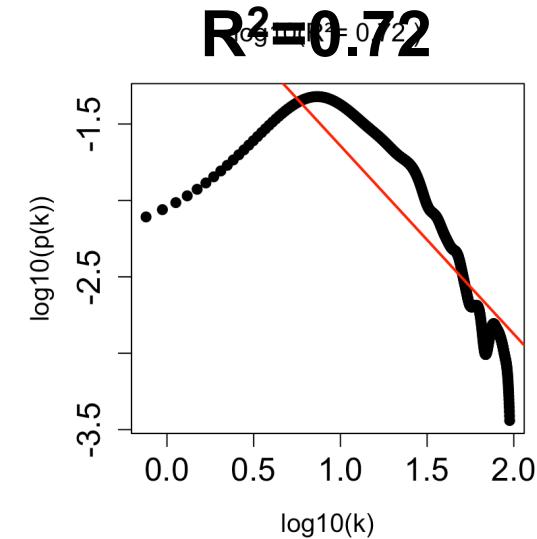
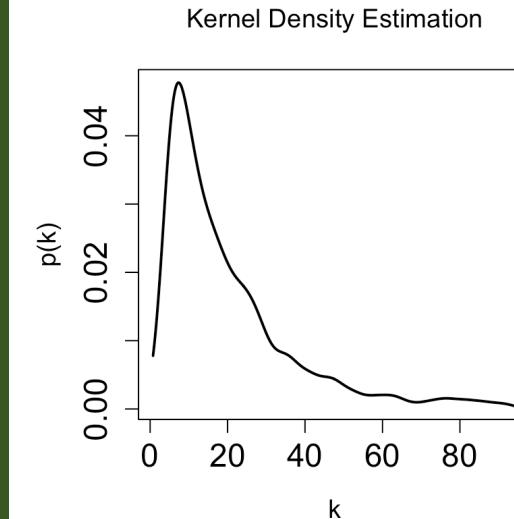


Validation: KDE $\beta = 6$

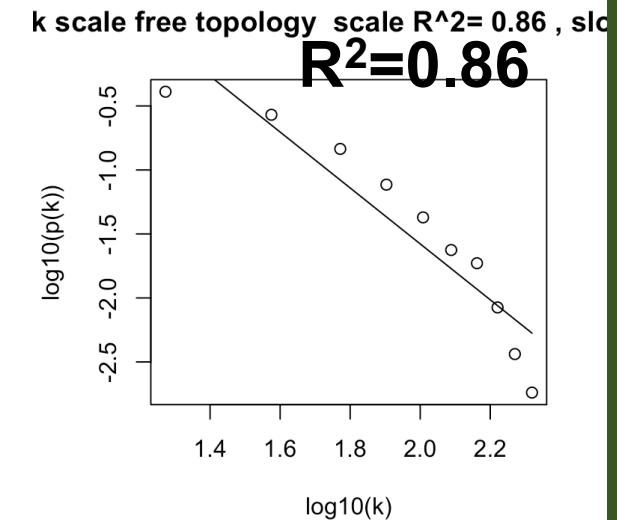
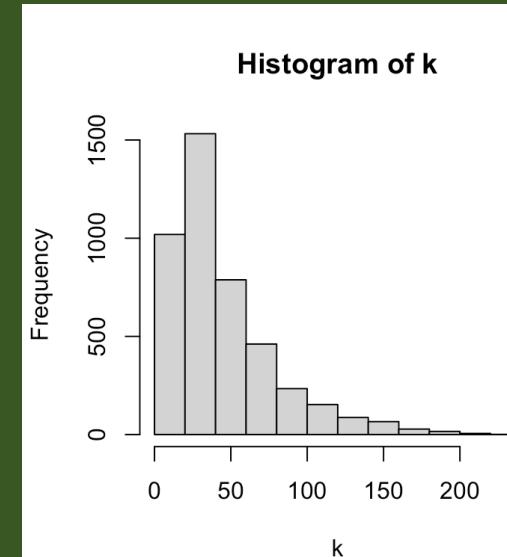


🤔 Both Method and Data matters
But do not vary a lot (I think)

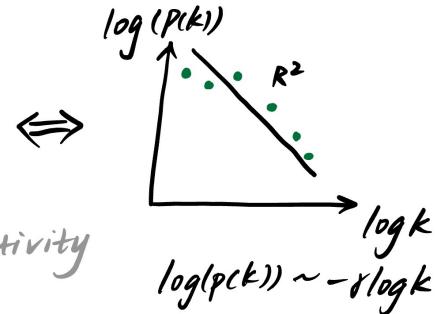
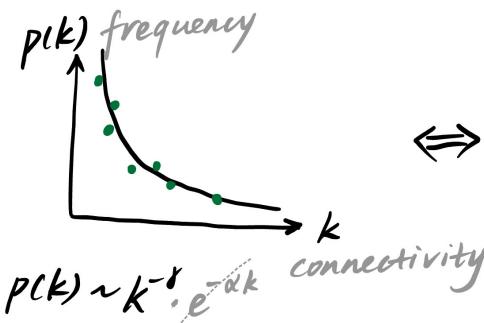
Tutorial Data, My Method (KDE)



WGCNA Package, My Data

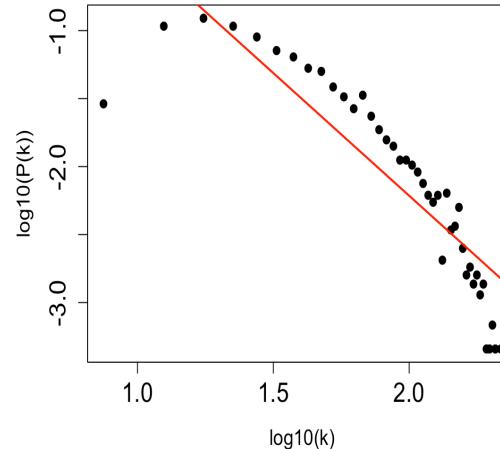
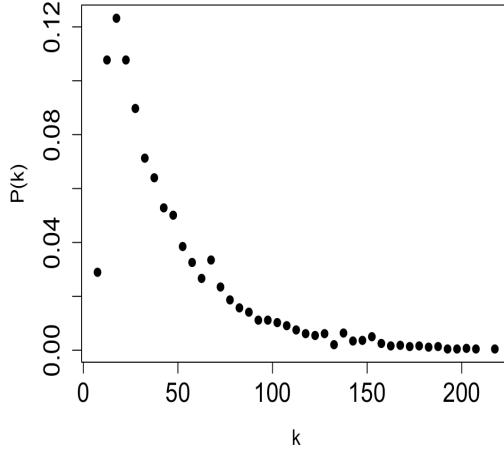


Validation: bin $\beta = 6$



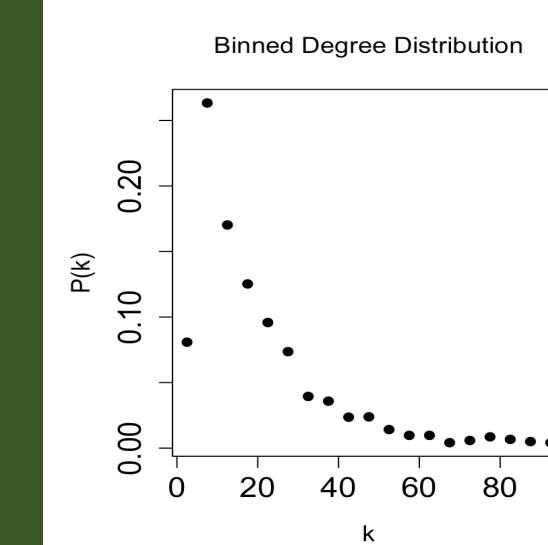
Binned Degree Distribution

$R^2 = 0.76$

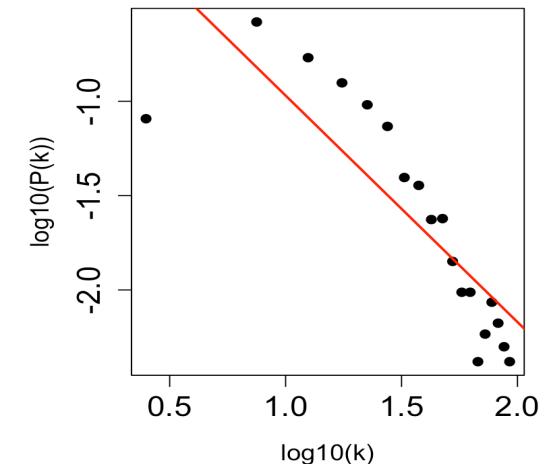


🤔 Both Method and Data matters
But do not vary a lot (I think)

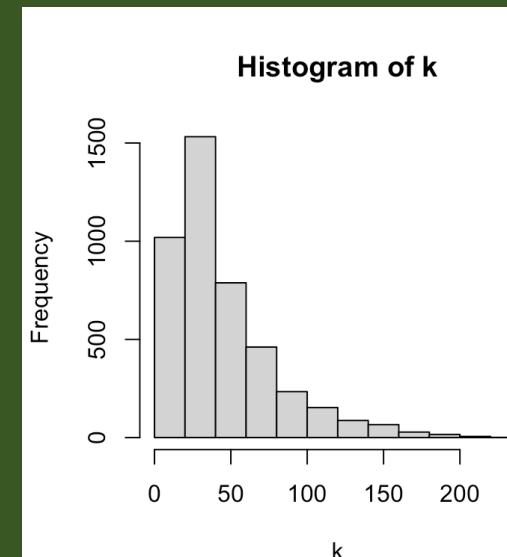
Tutorial Data, My Method (bin-size=5)



$R^2 = 0.71$

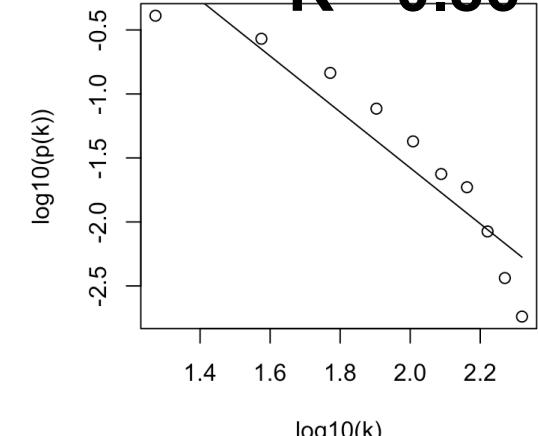


WGCNA Package, My Data



k scale free topology scale $R^2 = 0.86$, slope

$R^2 = 0.86$



Detect Module--overview

- Module **definition**: clusters of densely interconnected genes
- Measure interconnectivity: **dTOM**
- Detect **method**: **clustering**
- ❖ **Topological Overlap Matrix (TOM)**
consider **topological structure** in similarity measurement

$$\omega_{ij} = \frac{l_{ij} + a_{ij}}{\min\{k_i, k_j\} + 1 - a_{ij}}$$

(4) **Normalization s.t. $\omega_{ij} \in [0, 1]$**

Common Neighbor **Direct connect**

where $l_{ij} = \sum_u a_{iu}a_{uj}$, and $k_i = \sum_u a_{iu}$ is the node connectivity, see equation

- ❖ **Dissimilarity TOM (dTOM)** $d_{ij}^\omega = 1 - \omega_{ij}$.

Larger value, more similar

Tree Cut After Clustering

- ✗ Fix tree height
- ✗ Fix cluster number

“The first variant, called ‘Dynamic Tree’ cut...”

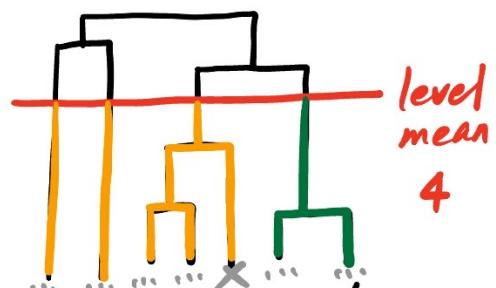
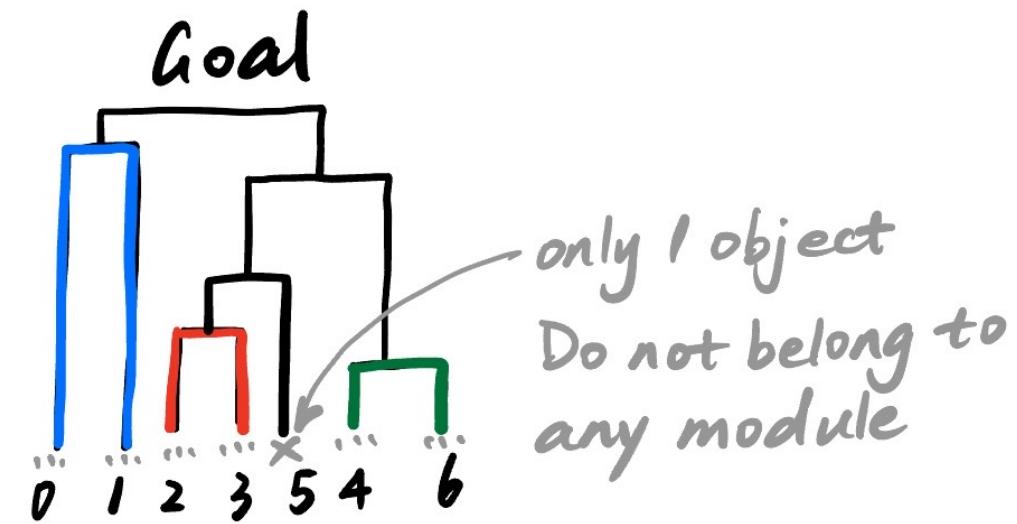
on it. This method has been utilized for identifying biologically meaningful gene modules in gene co-expression networks from several species such as yeast [2], [3], mouse [5] and human cell lines [4], but has not previously been systematically described nor m... Hybrid” cut, builds the clusters from dissimilarity information among the of each cluster. We describe each of

[5] A. Ghazalpour, S. Doss, B. Zhang, C. Plaisier, S. Wang, E.E. Schadt, A. Thomas, T.A. Drake, A.J. Lusis, and S. Horvath. Integrating genetics and network analysis to characterize genes related to mouse weight. PloS Genetics, 2(8):e130, 2006.

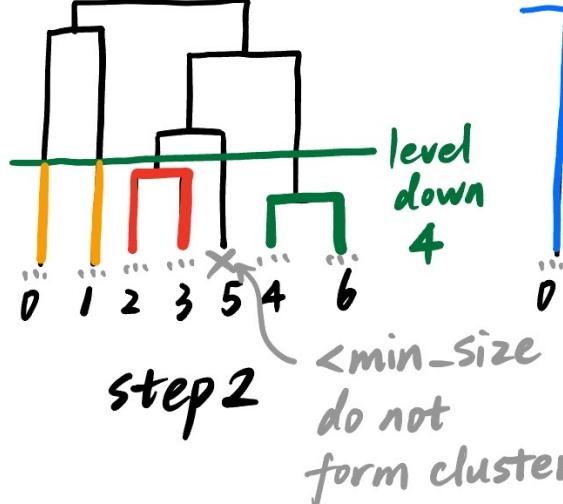
(Langfelder et al., 2008)

But existing package is the **second variant** → reproduce by myself...

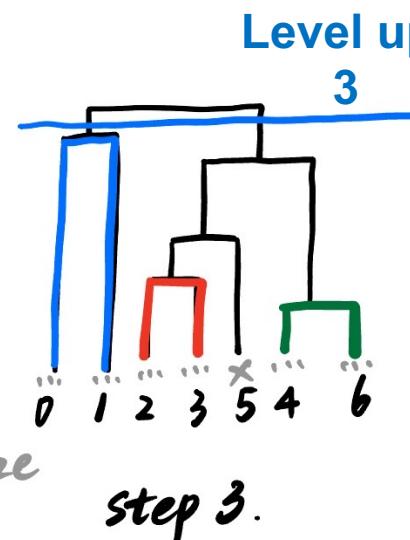
DynamicTreeCut Principle



step 1



step 2



step 3.

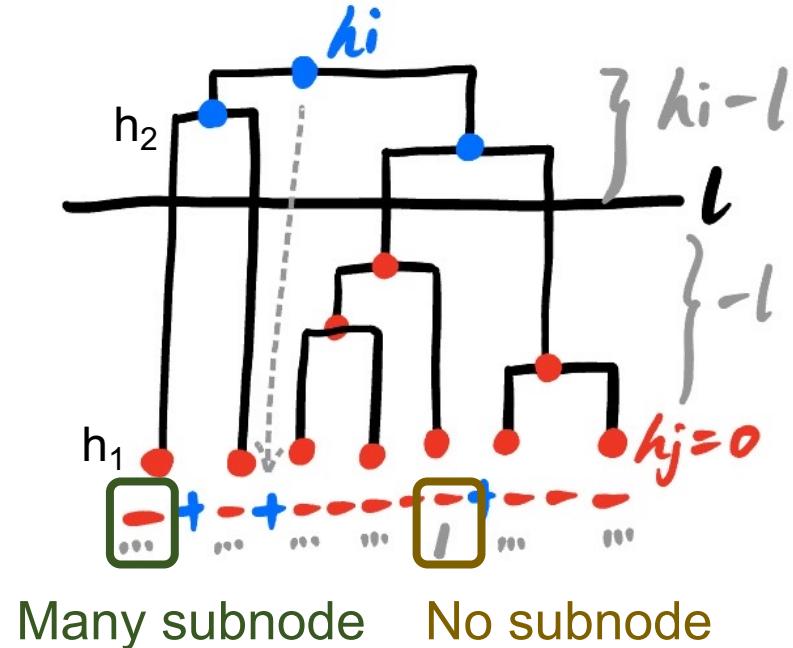
Error source 2: TreeCut Method

- 💡 Use my method to analyze tutorial data to evaluate.

$$\begin{aligned}
 l_m &= \frac{1}{n} \sum_{i=1}^n h_i, \\
 l_u &= \frac{1}{2}(l_m + \max\{h_1, h_2, \dots, h_n\}), \\
 l_d &= \frac{1}{2}(l_m + \min\{h_1, h_2, \dots, h_n\}).
 \end{aligned}$$

Reproduce DynamicTreeCut

`def TreeCut(H,l,min_size):` My `min_size` is from the paper: 34 +33



- Nodes here including merging nodes
- # Nodes in the cluster = # **continuous** **minus**
- If # > min_size → form a cluster

💡 Order of H = $[h_1, h_2, \dots]$?
• `hclust()` and `linkage()` only do return the merging process
• Order is the **same as dendrogram!**

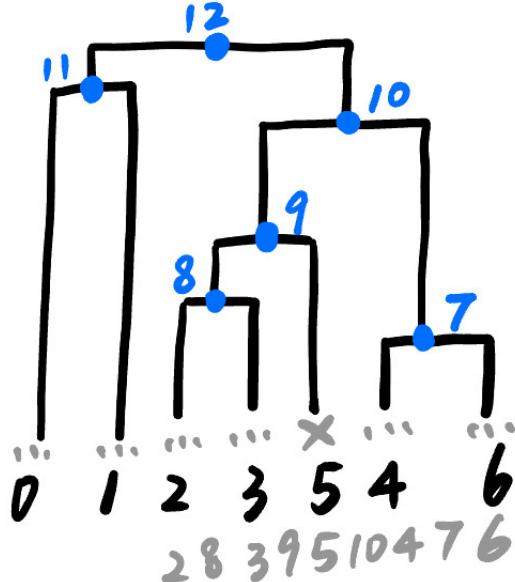
2 🤔 You're right to be confused. The only explanation that makes sense is that "dendrogram height sequence" has an unconventional meaning the authors never managed to mention. There's a possible clue in an implementation comment that says it's "the sum of the dissimilarity measure from the root to the node" <http://www.bioinformatics.org/cgi-bin/viewvc.cgi/catch/branches/catch-engine/splitter/splitter.c?r1=482&r2=483&view=patch> - Gene

Arduous!
Confusions in the paper!

(From stack overflow)

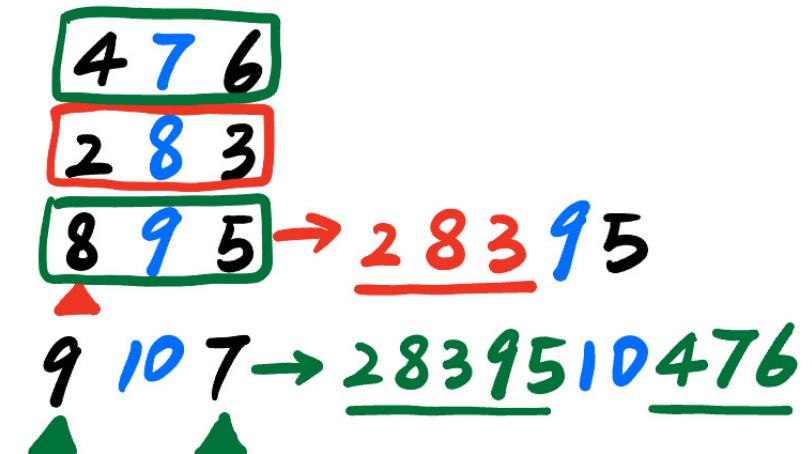
Reproduce DynamicTreeCut

```
def node_order(link):
```



merging		height	# obj	new node
node1	node2			
4	6	0.1	2	→ 7
2	3	0.14	2	→ 8
8	5	0.2	3	→ 9
9	7	0.3	5	→ 10
...				

link

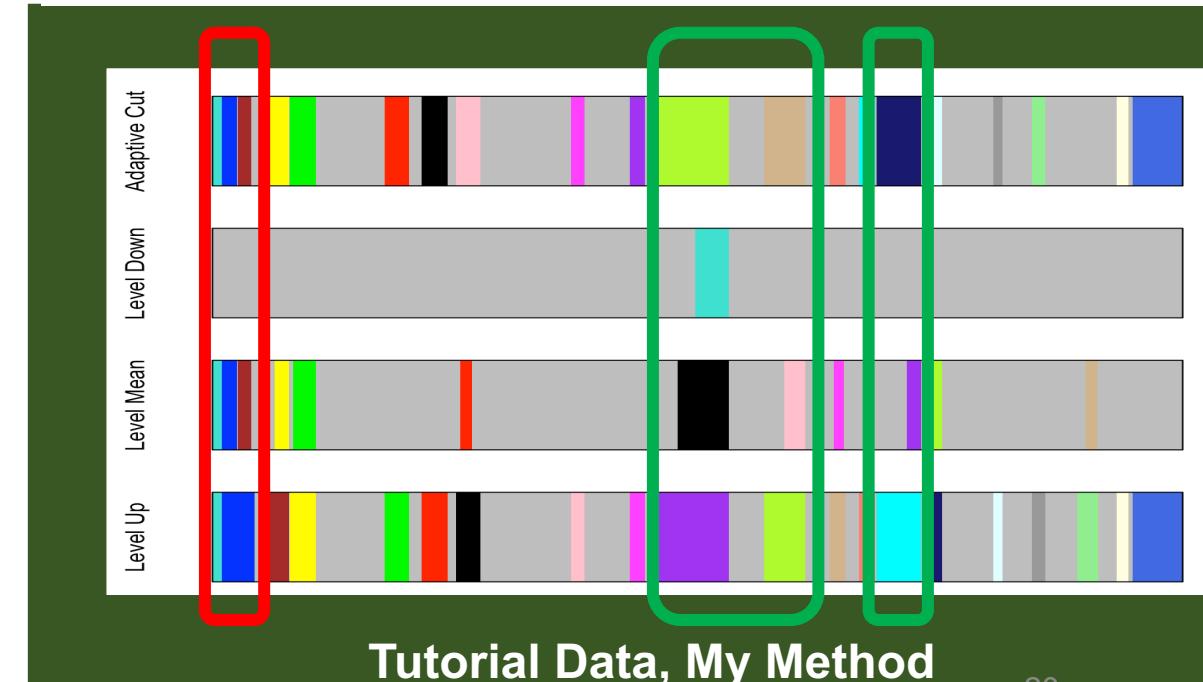
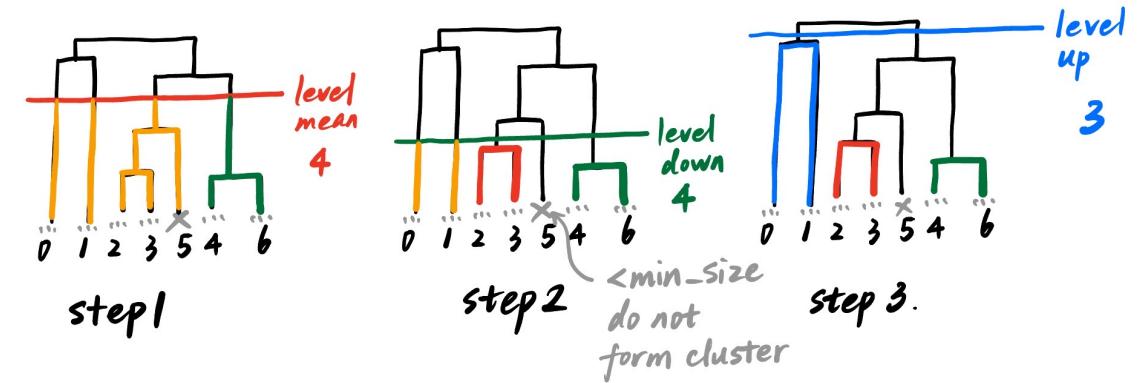
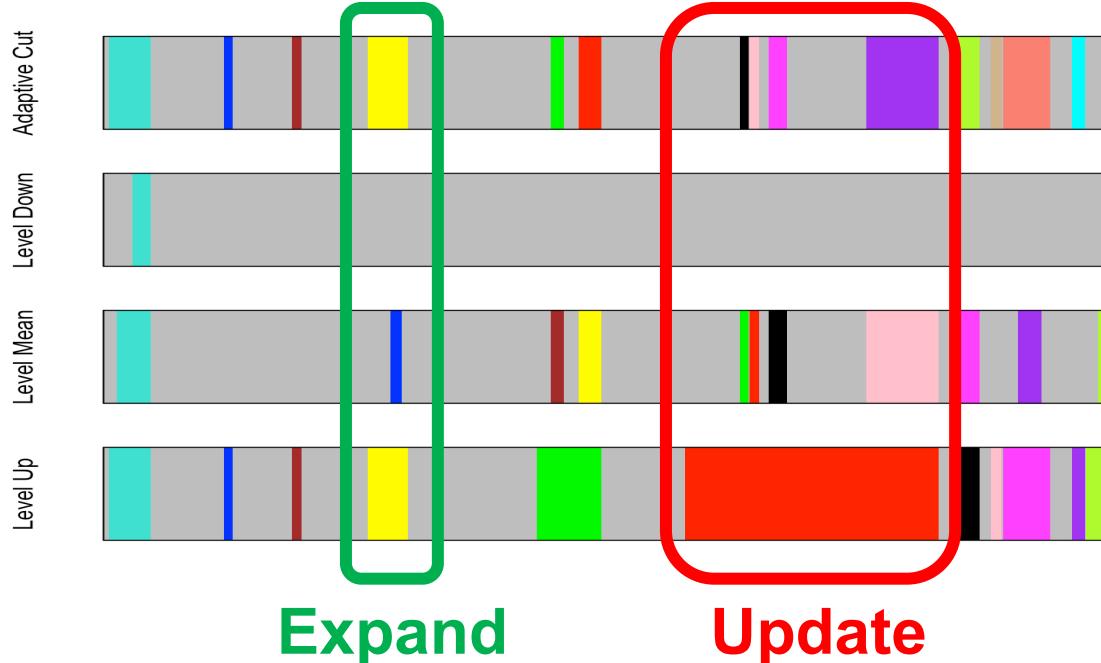


Then map the height to the node with link get height_order

Reproduce DynamicTreeCut

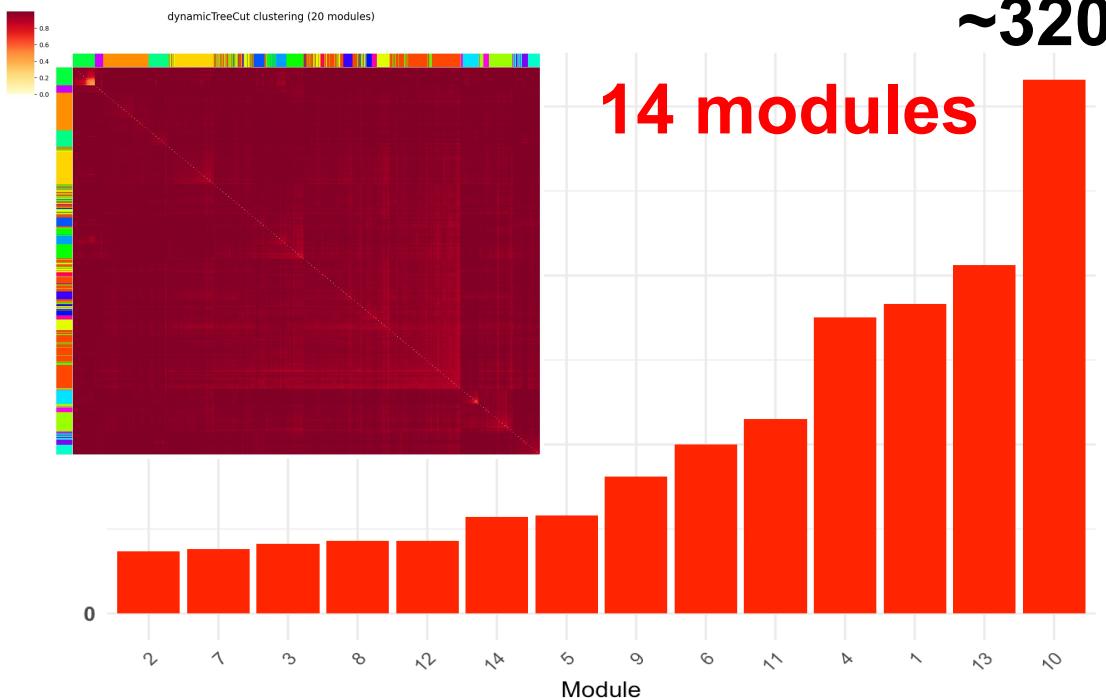
Def AdaptiveTreeCut(H):

- If ≥ 2 sub-clusters in the lower cut level result, **update**
- Or else, **expand** the cluster range

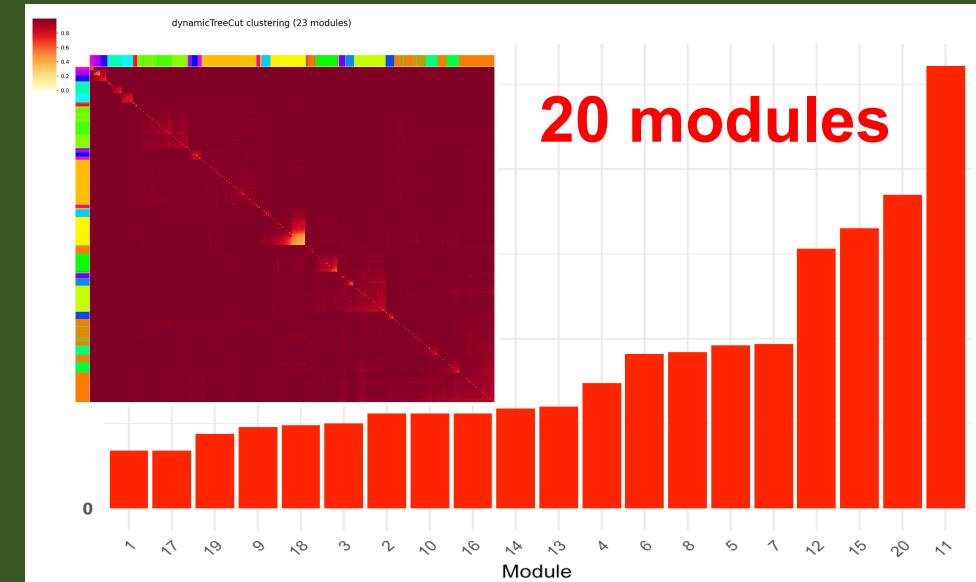


Number of Clusters

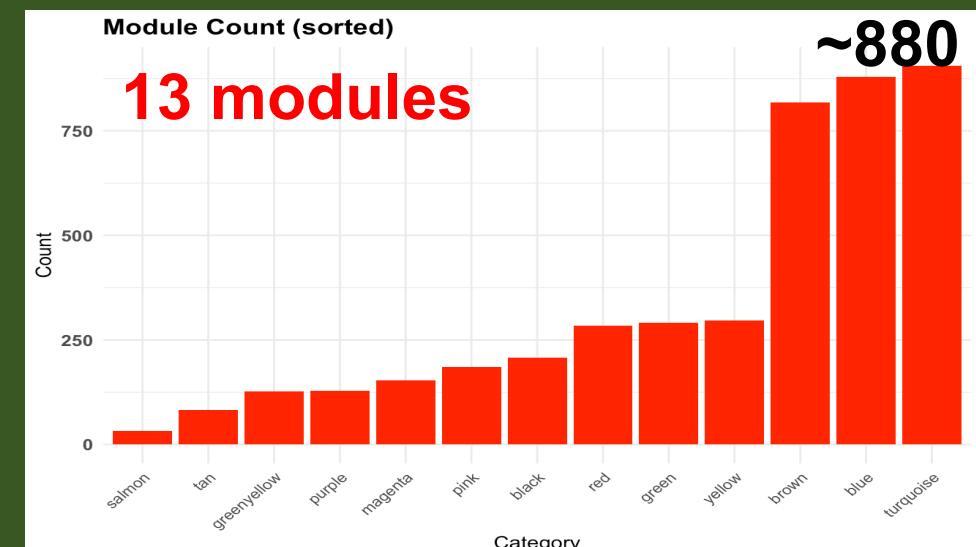
dTOM plot is better in determining the number of



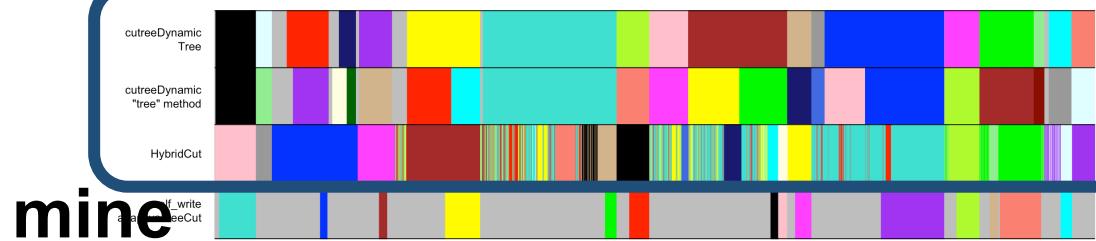
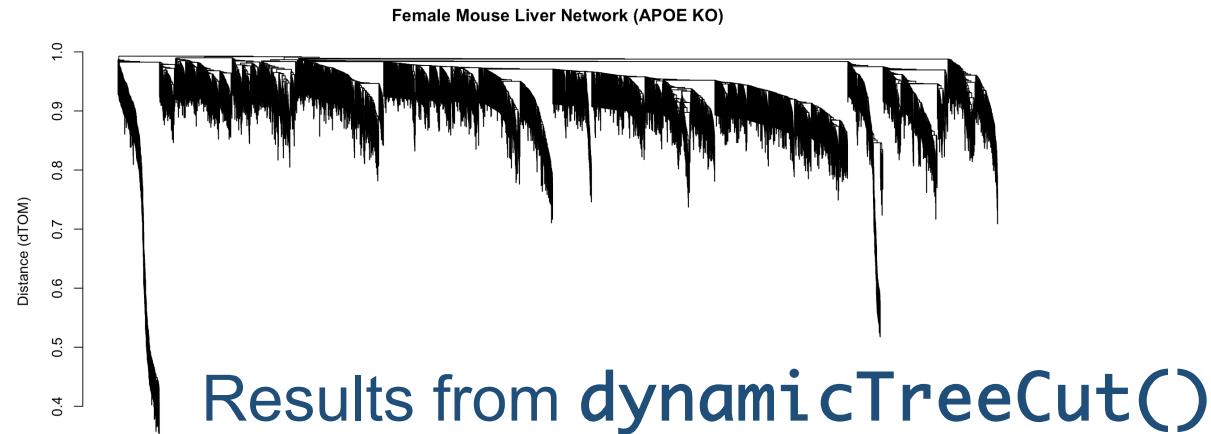
Tutorial Data, My Method



WGCNA Package, My Data

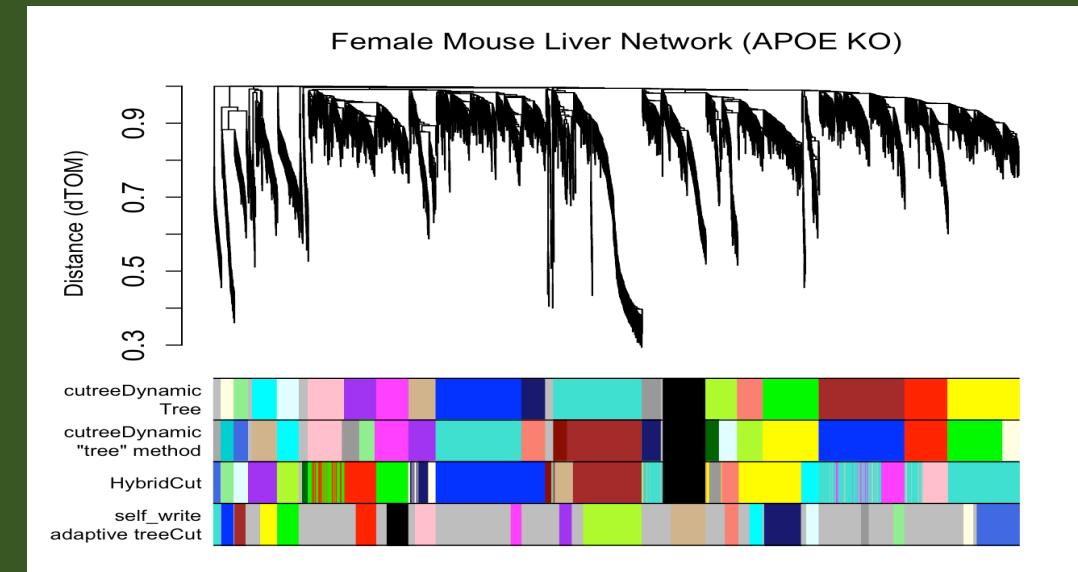


DynamicTreeCut

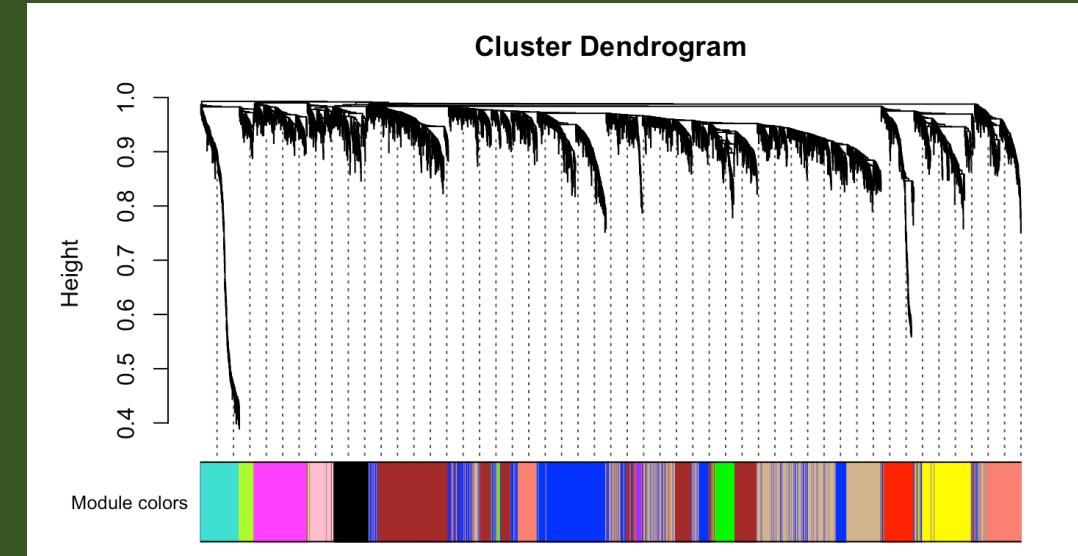


🤔 my results are much sparser because do not do assignment (variant 2) —hard to evaluate...

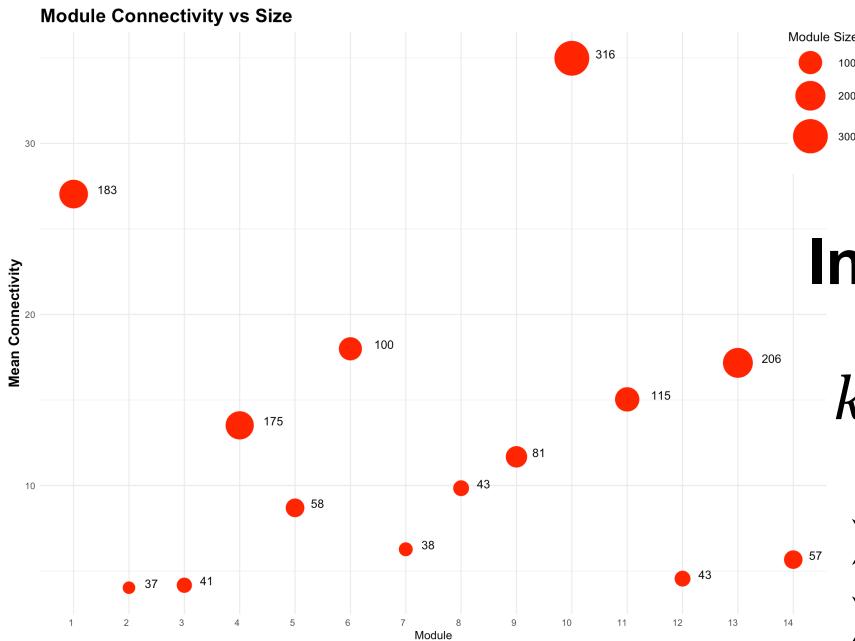
Tutorial Data, My Method



WGCNA Package, My Data



Cluster Evaluation



$$a_{ij} = |s_{ij}|^\beta$$

$n \times n$

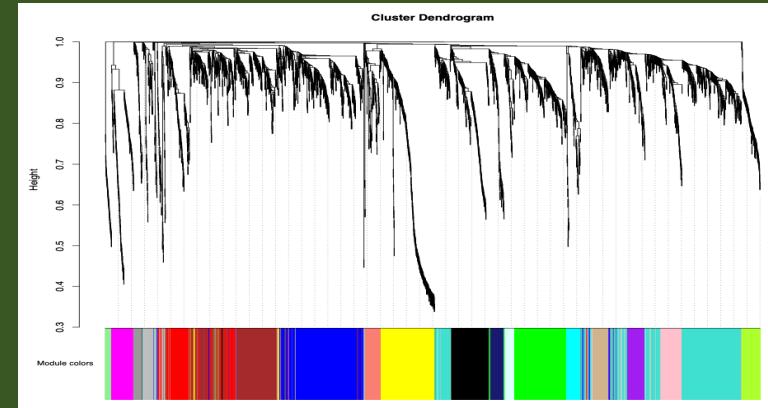
$\sum_j a_{ij} = k_i$
 connectivity
 of i th gene

Intramodular Connectivity

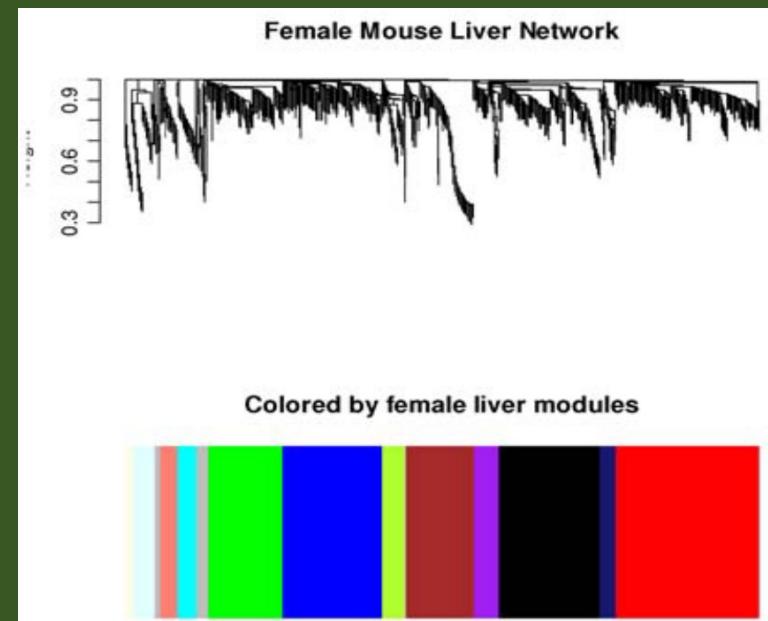
$$k_{in} = \frac{\sum_i k_i}{\# genes in module}$$

- Larger module size
- Generally higher k_{in}

🤔 Even tutorial data processed with WGCNA package **do not fully align** to the paper result! Maybe should focus on the **biological interpretation** of the result to **evaluate...**

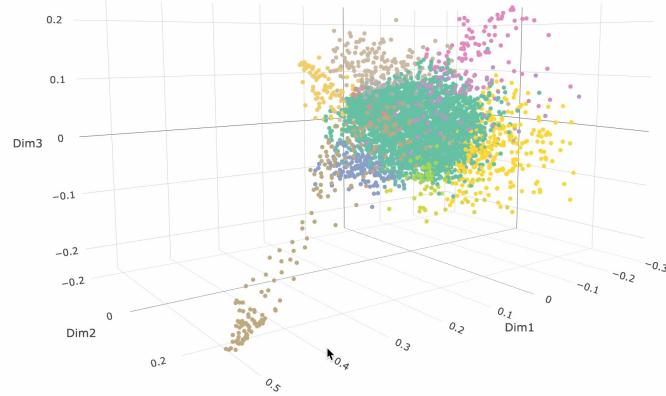


Paper Figure

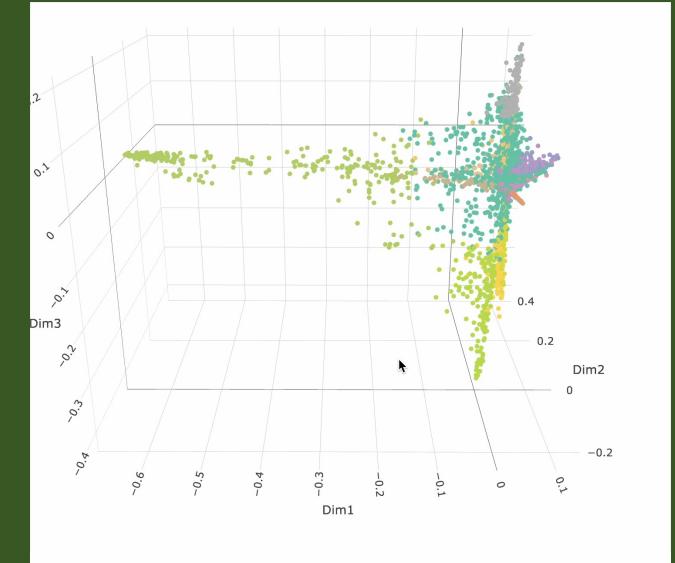


Evaluation--MDS Plot

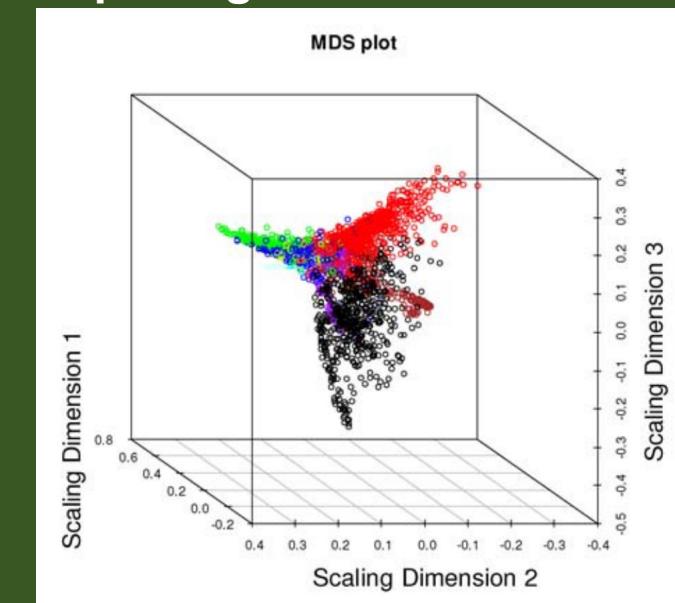
3D



2D



Paper Figure



🤔 Similar shape → my data may work
unassigned genes in the center →
clustering performs well

Module Significance (MS)

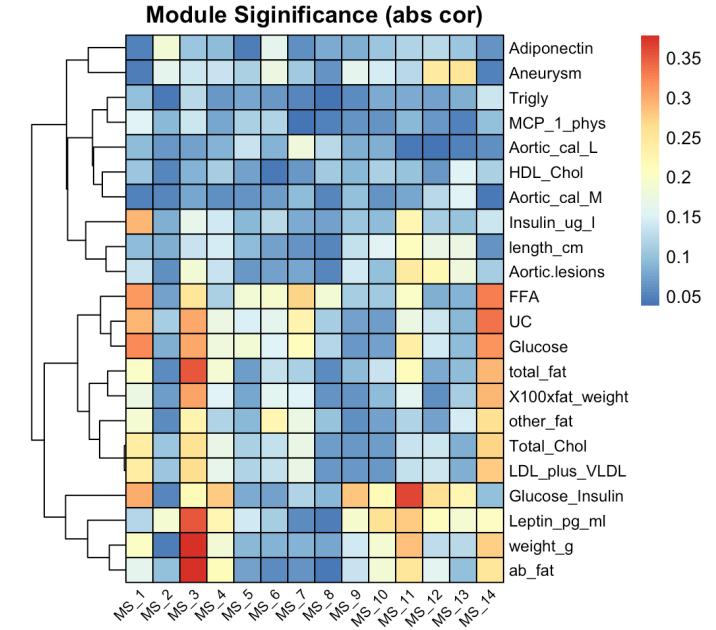
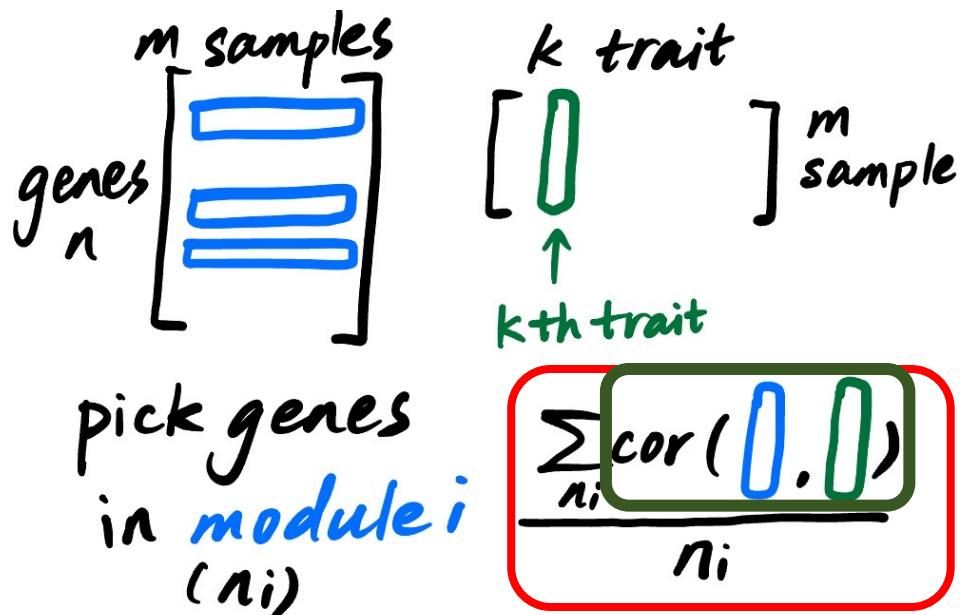
❖ Gene significance (GS_i)

the relationship strength between ith gene and kth trait

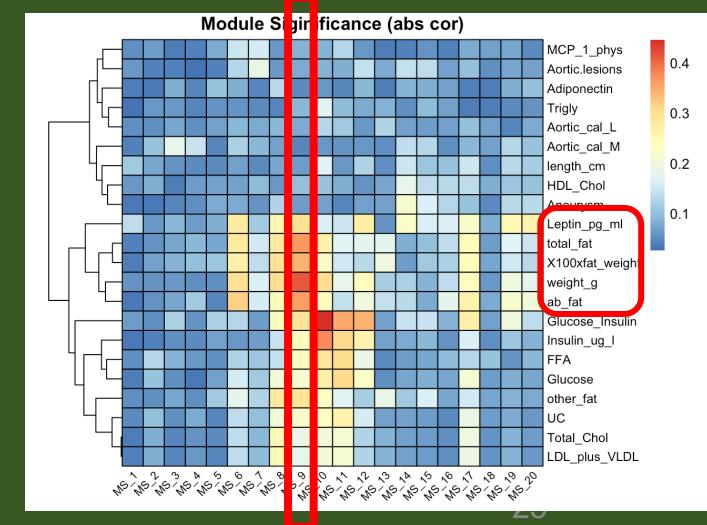
$$GS_i = \text{cor}(x_i, T_k)$$

❖ Module significance (MS)

$$MS = \sum GS_i , GS_i \text{ of module}$$



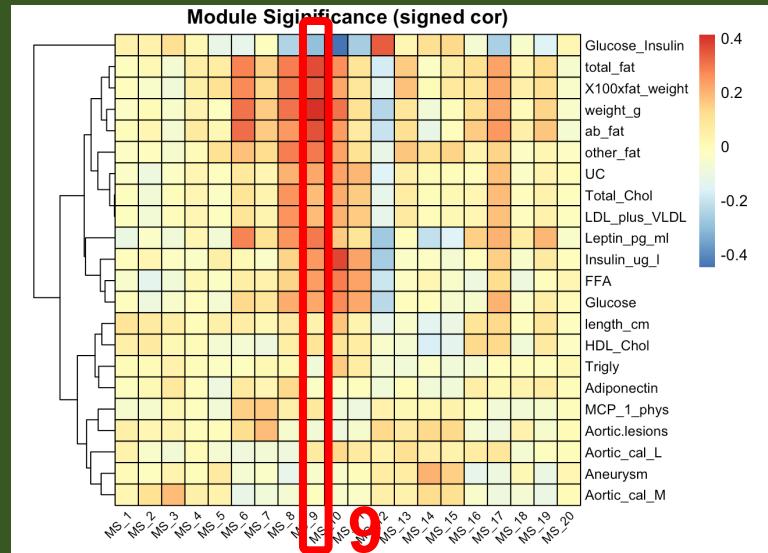
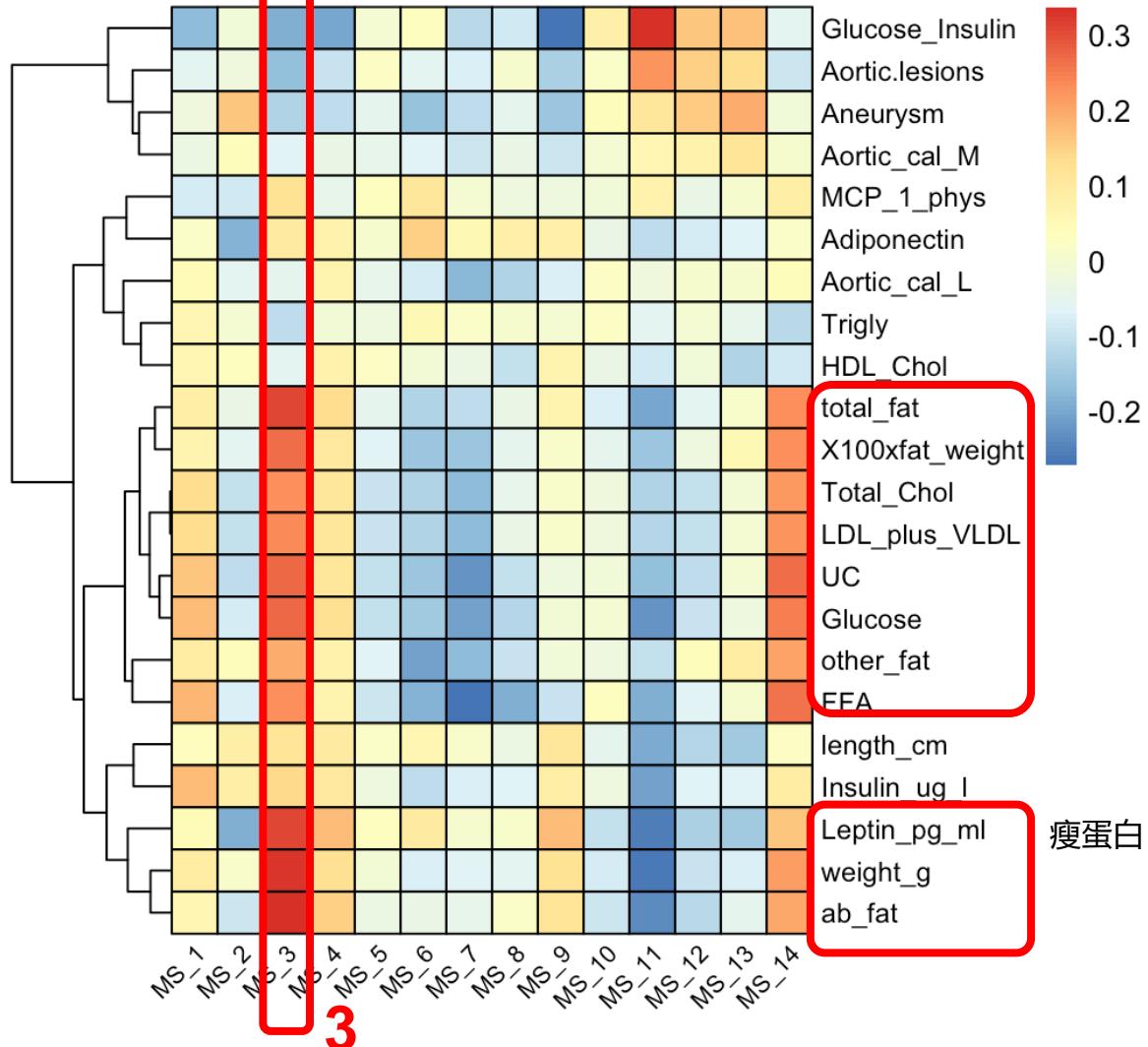
Tutorial Data, My Method
Absolute correlation result



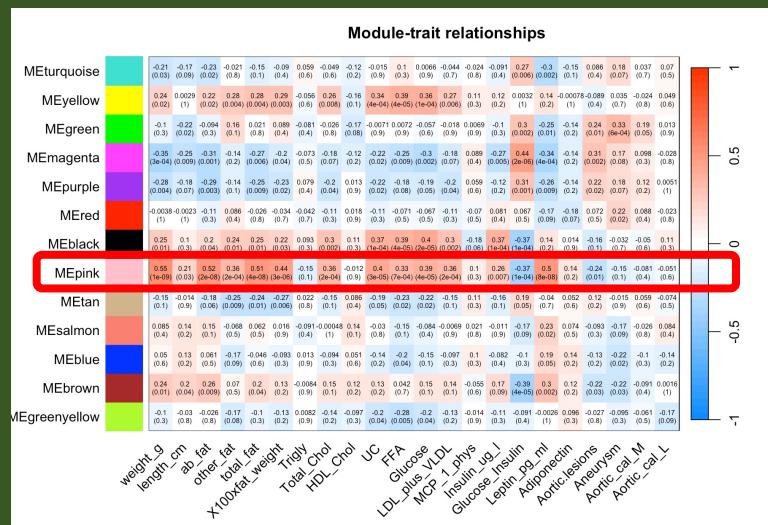
Tutorial Data, My Method

Module Significance (MS)

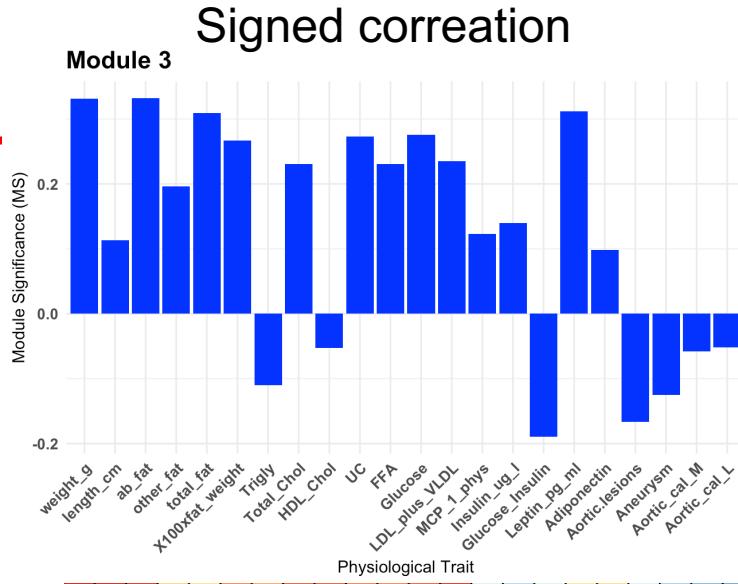
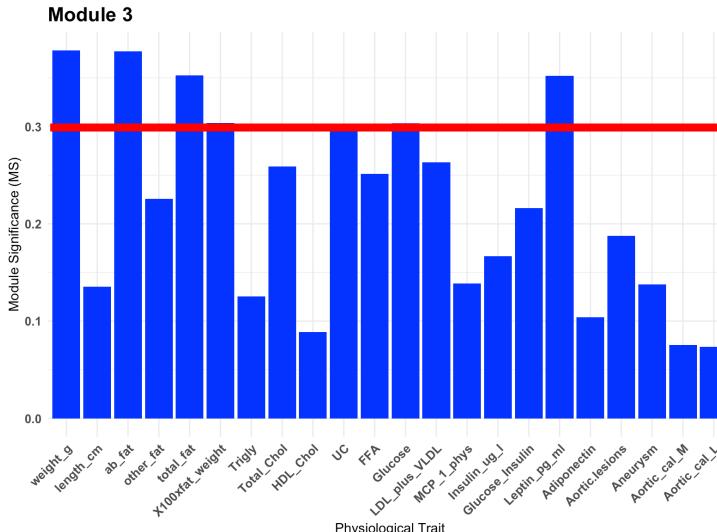
Signed correation



WGCNA Package, My Data

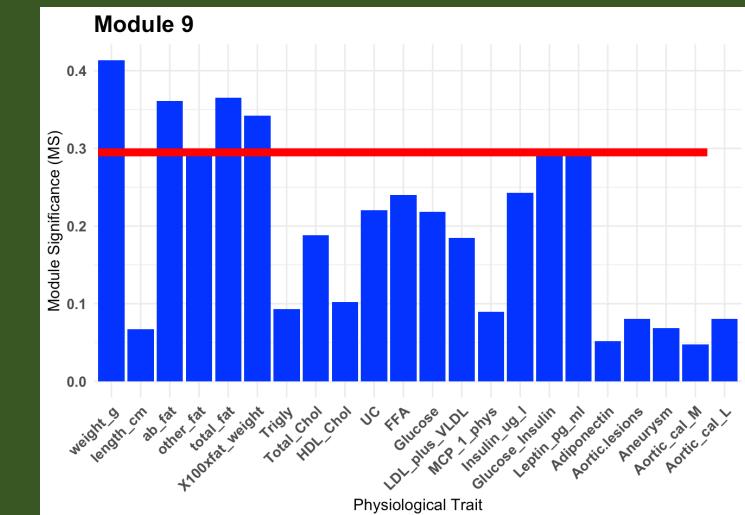


Weight-related Module

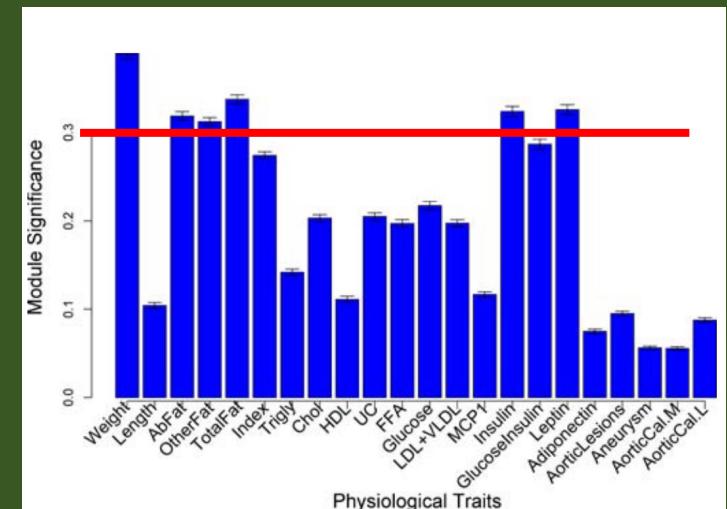


🤔 my **compact module** show **higher correlation** with the traits
Mostly the **same trend** as the paper result → my **method** seems to work!

Tutorial Data, My Method



Paper Result



Wrap Up

➤ Done

- Reproduce: **data processing, network construction, module detection** (with “my” clustering method), and **module significance evaluation**.
- From visualization, both **my data and method seem to work**. They can provide **biological insights** we want. **Do downstream analyses** are required to further validate.

➤ To be improved...

- Should **simplify codes**
- Should deepen **understanding** network property and similarity measurement

➤ Overall, thanks this project!

- I've become **stronger**—both mentally and technically
- I just learned that analysis can be done **without explicit control group!**

Reference List

1. Ghazalpour A, Doss S, Zhang B, Wang S, Plaisier C, et al. (2006) Integrating Genetic and Network Analysis to Characterize Genes Related to Mouse Weight. *PLOS Genetics* 2(8): e130. <https://doi.org/10.1371/journal.pgen.0020130>
2. Langfelder, P., & Horvath, S. (2008). WGCNA: An R package for weighted correlation network analysis. *BMC Bioinformatics*, 9(1), 559. <https://doi.org/10.1186/1471-2105-9-559>
3. Langfelder, P., Zhang, B., & Horvath, S. (2008). Defining clusters from a hierarchical cluster tree: The dynamic tree cut package for R. *Bioinformatics*, 24(5), 719–720. <https://doi.org/10.1093/bioinformatics/btm563>
4. Zhang, B., & Horvath, S. (2005). A general framework for weighted gene Co-expression network analysis. *Statistical Applications in Genetics and Molecular Biology*, 4(1). <https://doi.org/10.2202/1544-6115.1128>