

Weighted Correlation Network Analysis (WGCNA)

SIB course, Nov. 15-17, 2016

Introduction to Biological Network Analysis

Leonore Wigger

with Frédéric Burdet and Mark Ibberson

Agenda

Overview of WGCNA

Theory 1: Weighted correlation network, split into modules

Practical

LUNCH BREAK

Theory 2: Identify modules and genes of interest

Practical

Data for WGCNA

Gene expression data (microarray or RNA-Seq)

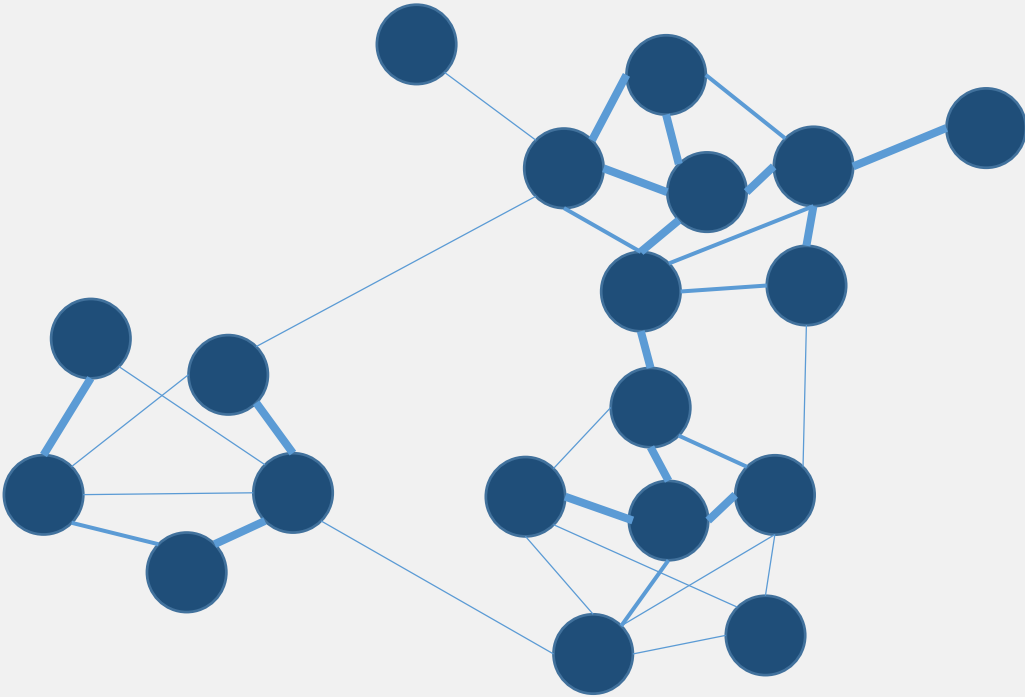
Recommendation: at least 20 individuals

Clinical/phenotypical traits from the same individuals (optional)

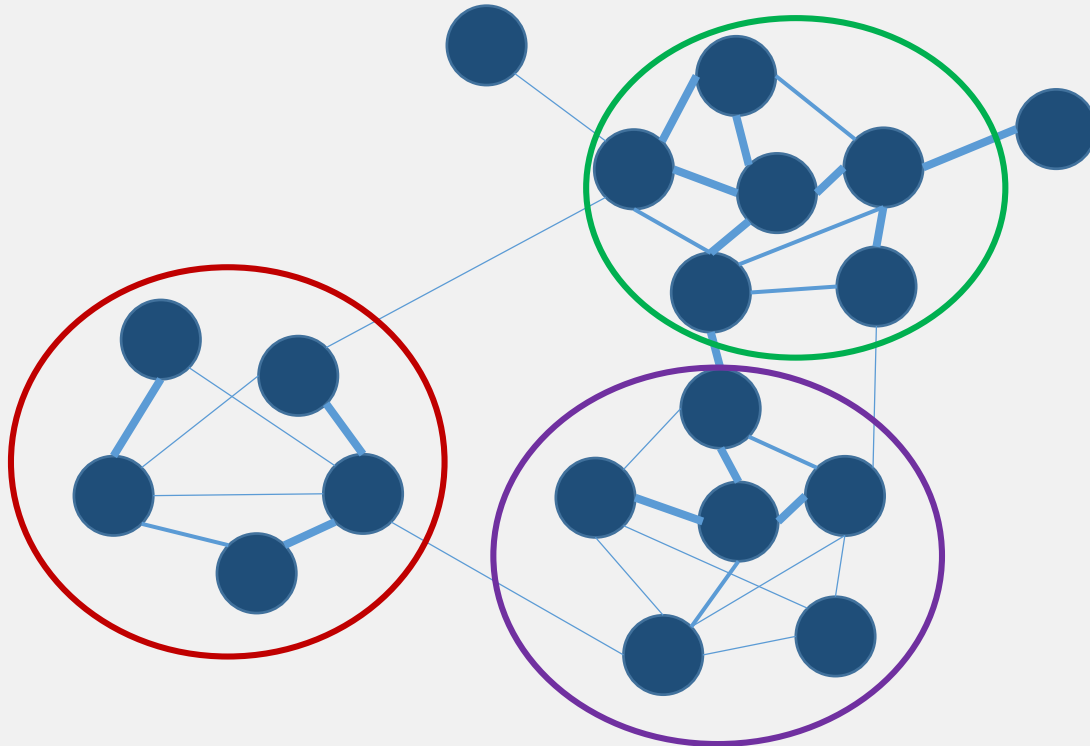
e.g. weight, insulin level, glucose level

Aims of WGCNA

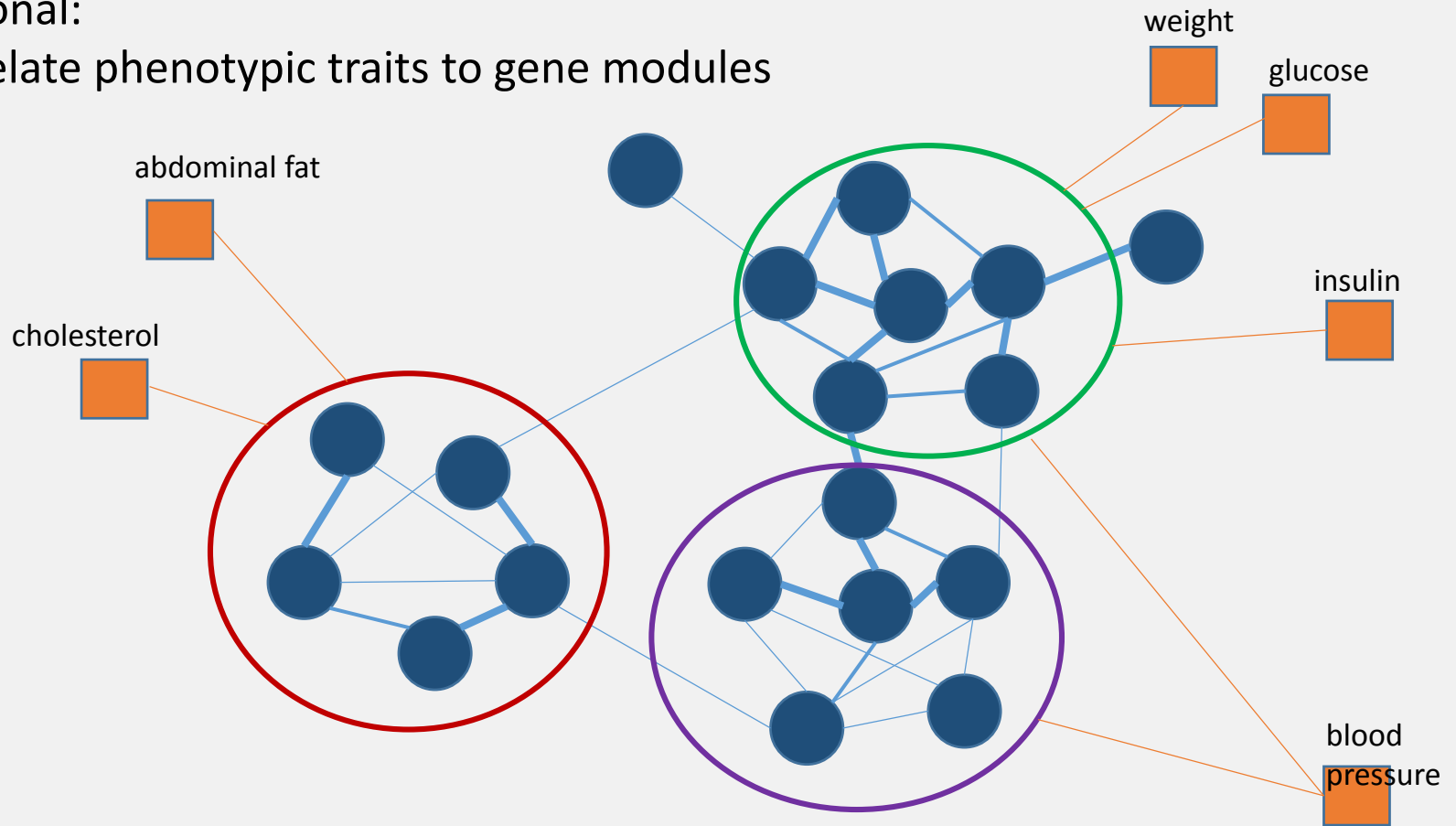
Construct a gene-gene similarity network



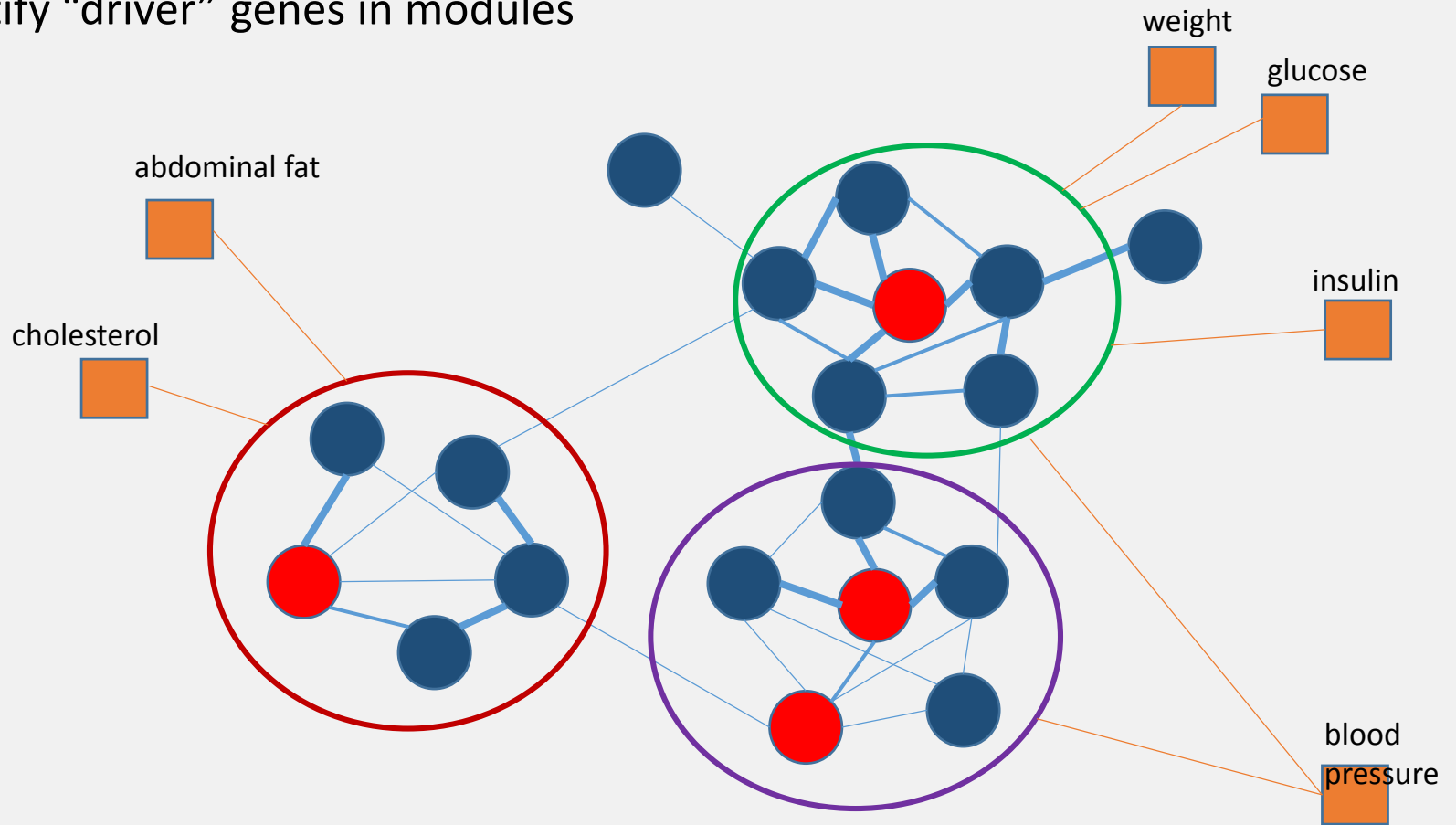
Divide network into modules
(Group genes with similar expression)



Optional:
Correlate phenotypic traits to gene modules

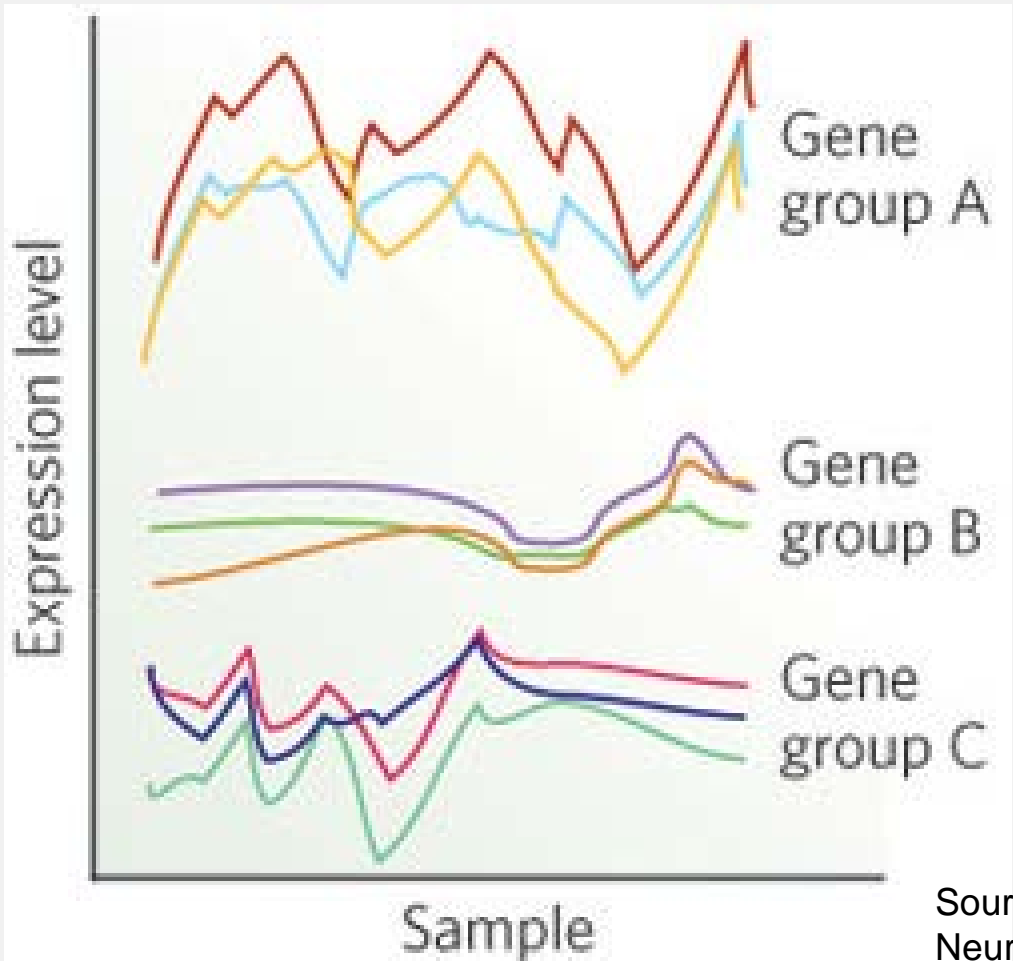


Identify “driver” genes in modules



Modules found in WGCNA:

Groups of co-expressed genes (with similar expression profiles over a large group of individuals)



Source: Daniel H. Geschwind & Genevieve Konopka. Neuroscience in the era of functional genomics and systems biology, Nature 461, 908-915

Central Hypothesis:

Genes with **similar expression patterns** are of interest because they may be

- tightly co-regulated
- functionally related
- members of the same pathway

WGCNA is a **guilt-by-association** approach:

Encourages hypotheses about genes based on their close network neighbors.

Workflow

Gene expression data



Weighted correlation network of genes



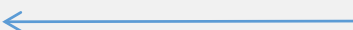
Modules of genes with similar profiles



Pathway analysis



Correlation of modules with clinical traits



Trait data



Identification of potential "driver" genes

Construct weighted correlation network

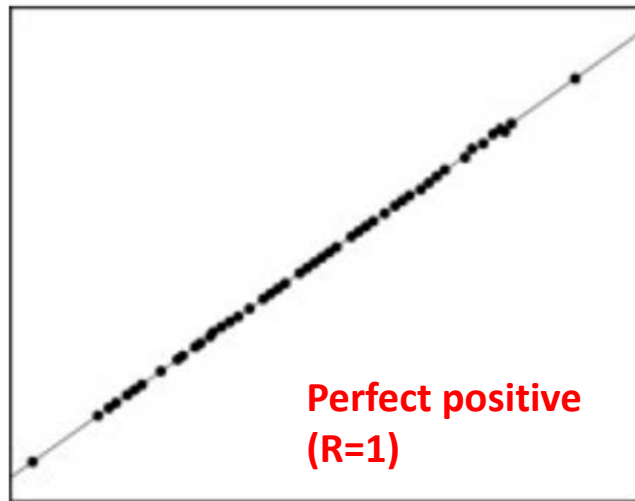
Correlation:

A statistical measure for the extent to which two variables fluctuate together.

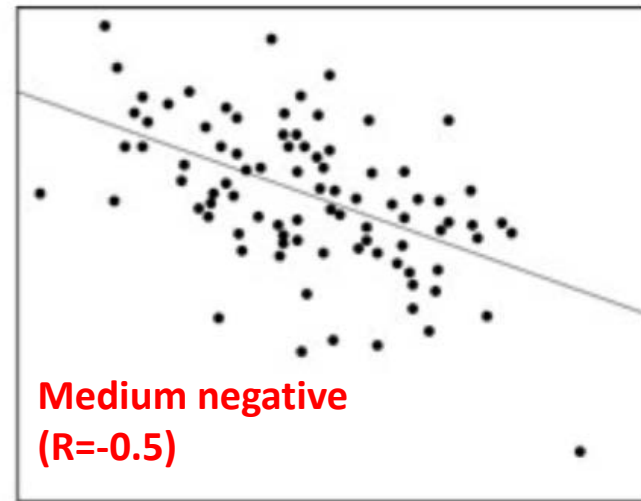
Positive correlation: variables increase/decrease together

Negative correlation: variables increase/decrease in opposing direction

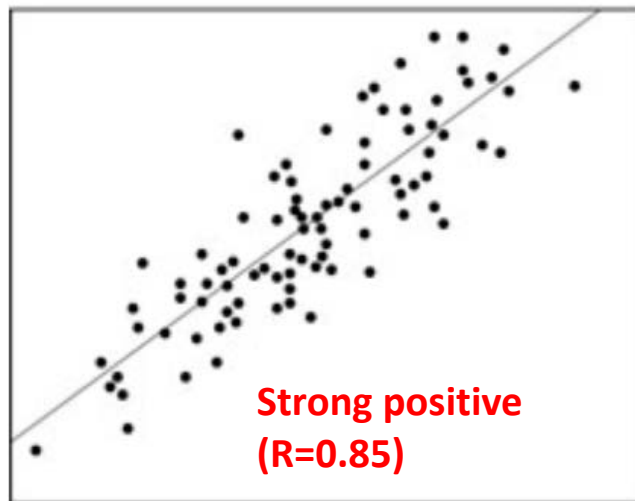
Correlation examples (Pearson correlation coefficient R)



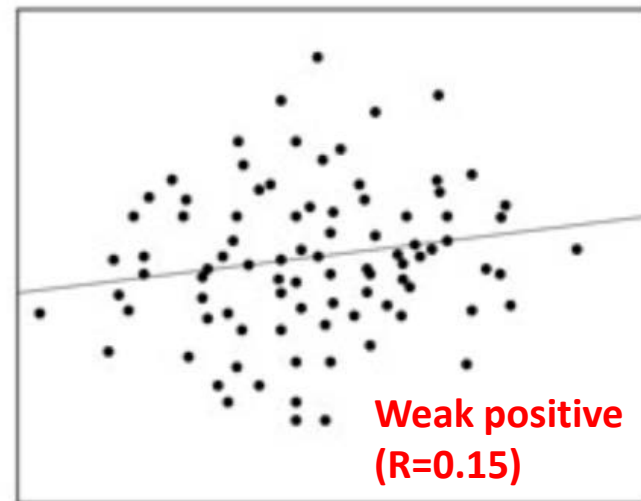
a



b



c



d

Scatterplots with correlations of a) +1.00; b) -0.50; c) +0.85; and d) +0.15.

Multiple measures for correlation exist

Implemented in WGCNA:

Pearson	(function <code>cor</code>)
Spearman, Kendall	(function <code>cor</code>)
biweight midcorrelation	(function <code>bicor</code>)

The Basis of WGCNA: Weighted Correlation Network of Genes

Adjacencies

Compute a correlation raised to a power between every pair of genes (i, j)

$$a_{i,j} = |\text{cor}(i, j)|^\beta$$

Effect of raising correlation to a power:

Amplifies disparity between strong and weak correlations

Example: Power term $\beta = 4$

Correlations

$$\mathit{cor}(i, j) = 0.8$$

$$\mathit{cor}(k, l) = 0.2$$

0.8/0.2:

4-fold difference

→

→

→

Adjacencies

$$|0.8|^4 = 0.4096$$

$$|0.2|^4 = 0.0016$$

0.4096/0.0016:

256-fold difference

Strong corr.

Weak corr.

The Basis of WGCNA: Weighted Correlation Network of Genes

Adjacencies

Compute a **correlation raised to a power** between every pair of genes (i, j)

$$a_{i,j} = |\text{cor}(i, j)|^\beta$$

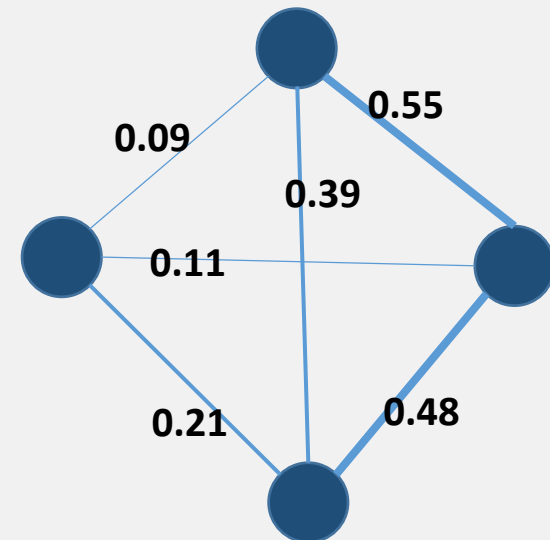
Network

Construct a **fully connected** network;
Genes as nodes, $a_{i,j}$ as edge weights.

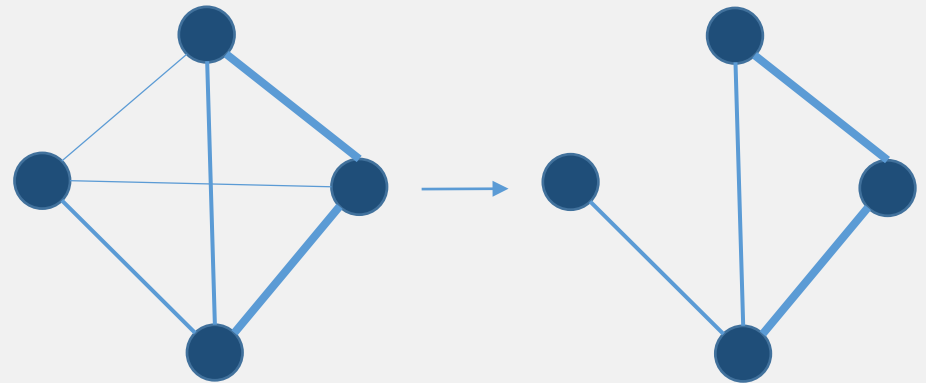
high correlation – strong connection
low correlation – weak connection

Adjacency matrix of 4 genes

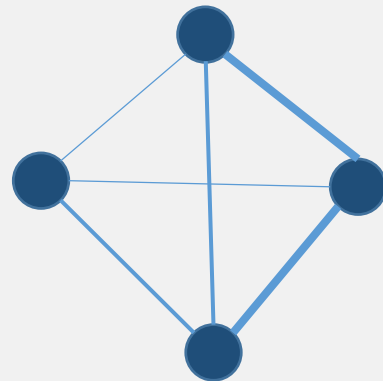
$a_{i,j}$	gene1	gene2	gene3	gene4
gene1	1	0.55	0.39	0.09
gene2	0.55	1	0.48	0.11
gene3	0.39	0.48	1	0.21
gene4	0.09	0.11	0.21	1



For visualizations, set a threshold on edge weight and **remove the weakest links.**



In most computations, work with all edges of the **fully connected network.**



Connectivity (degree) in a weighted network

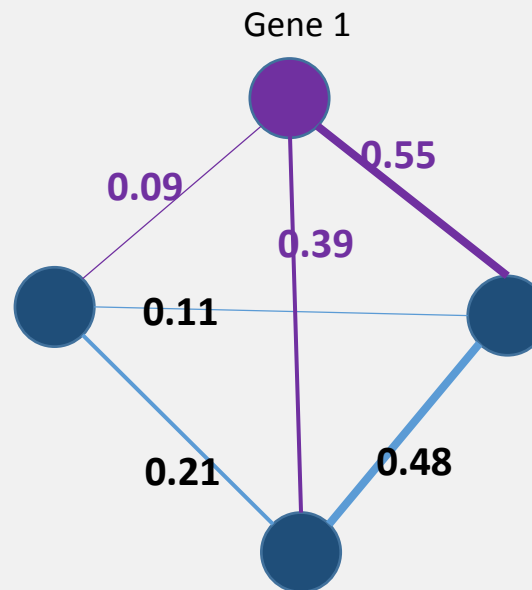
Connectivity of a gene:

Sum of the weights of all edges connecting to this gene

Example:

Connectivity of gene 1:

$$0.55 + 0.39 + 0.09 = 1.03$$



Weighted Correlation: Decisions to make

Selecting a network type

unsigned, signed etc.

Choosing a correlation method

Pearson, Spearman, biweight midcorr. etc.

Picking a Power term

1, ..., 20 etc.

Selecting a network type

Unsigned Network (Default)

No differentiation between **positive** and **negative** correlations.

Use this if negative correlation are of interest

Signed Hybrid Network

Only **positive** correlations are taken into account. Negative correlations are set to 0.

Use this if negative correlations are NOT of interest.

Signed Network

[not covered in course, use as an alternative to signed hybrid network]

Choosing a correlation method

$$a_{i,j} = |\text{cor}(i, j)|^\beta$$

Fastest, but sensitive to outliers:

Pearson correlation `cor(x)`
“standard” measure of linear correlation

Less sensitive to outliers but much slower:

Biweight mid-correlation `bicor(x)`
robust, recommended by the authors for most situations
[needs modification for correlations involving binary/categorical variables]

Spearman correlation `cor(x, method=“spearman”)`
rank-based, works even if relationship is not linear
less sensitive to gene expression differences
[can be used as-is for correlations involving binary/categorical variables]

Default correlation method in WGCNA: `cor` (Pearson).

Caveat: use it only if there are no outliers, or for exercises/tutorials.

Picking a power term

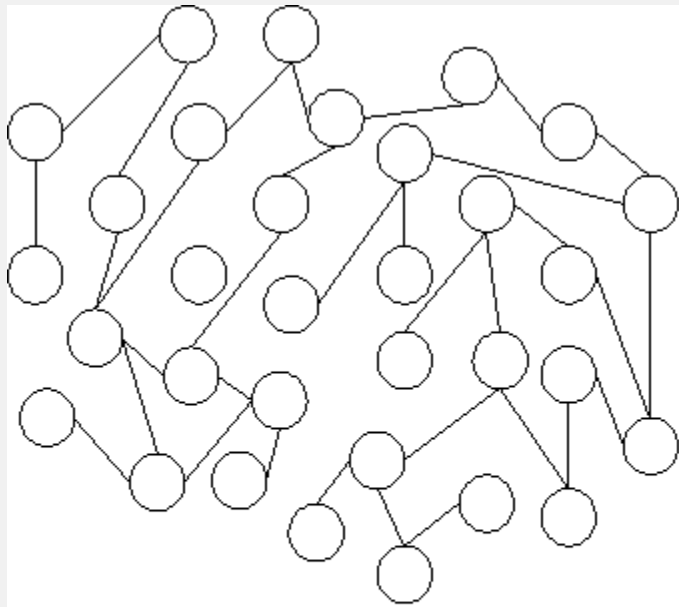
$$a_{i,j} = |\text{cor}(i, j)|^\beta$$

Selection criterion: Pick lowest possible β that leads to an approximately **scale-free network topology**

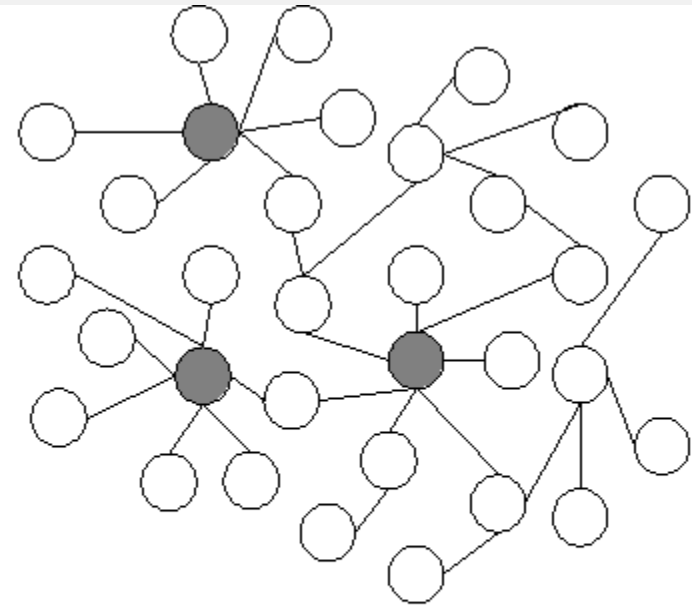
- few nodes with many connections ("hubs")
- many nodes with few connections

Degree distribution follows a power law:

the probability for a node of having k connections is $k^{-\gamma}$



(a) Random network



(b) Scale-free network

Source: Carlos Castillo: Effective Web Crawling,
PhD Thesis, University of Chile, 2004.
(obtained from Wikimedia Commons)

Why scale-free network topology?

Background:

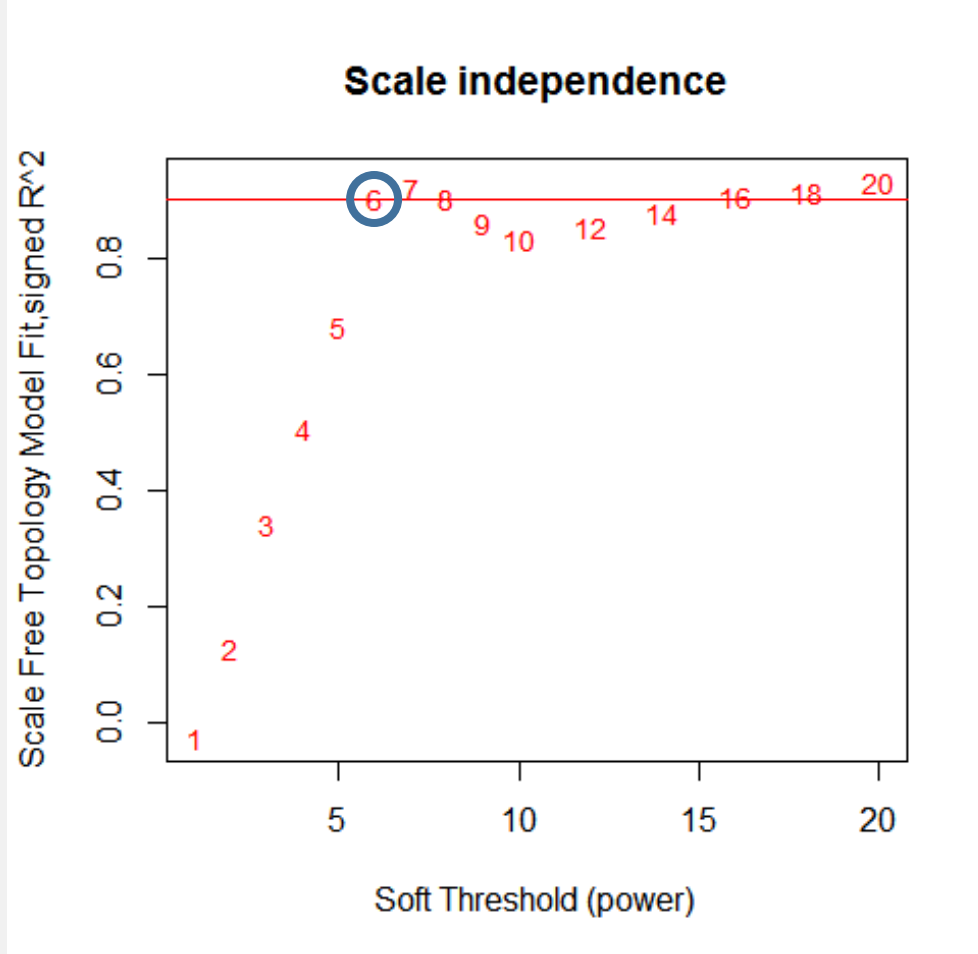
Barabási et al. found many types of network in many domains to be approximately scale-free, including metabolic and protein interaction

Aim in WGCNA:

Build a biologically “realistic” network.

Barabási, Albert-László; Bonabeau, Eric (May 2003). "[Scale-Free Networks](#)" (PDF). *Scientific American*. **288** (5): 50–9.

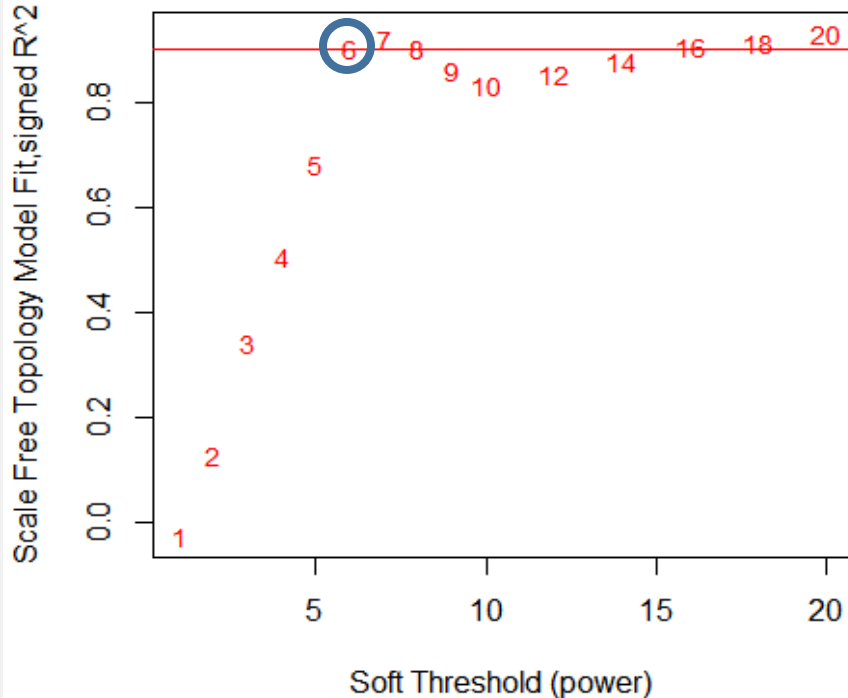
Pick a power term: Visual Aid in WGCNA



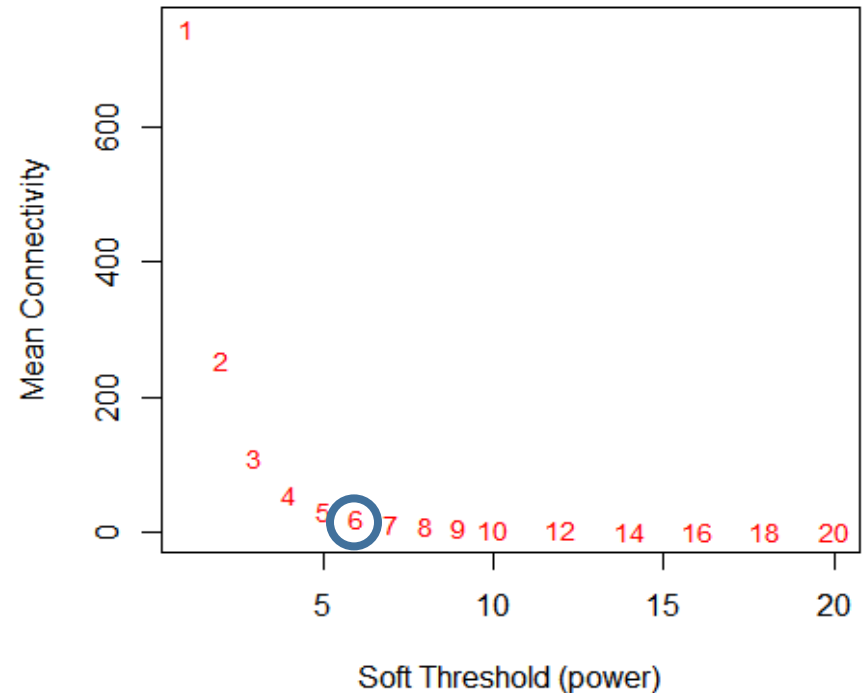
Choose power 6. Lowest possible power term where topology approximately fits a scale free network (on or above red horizontal line).

Pick a power term: Visual Aid in WGCNA

Scale independence

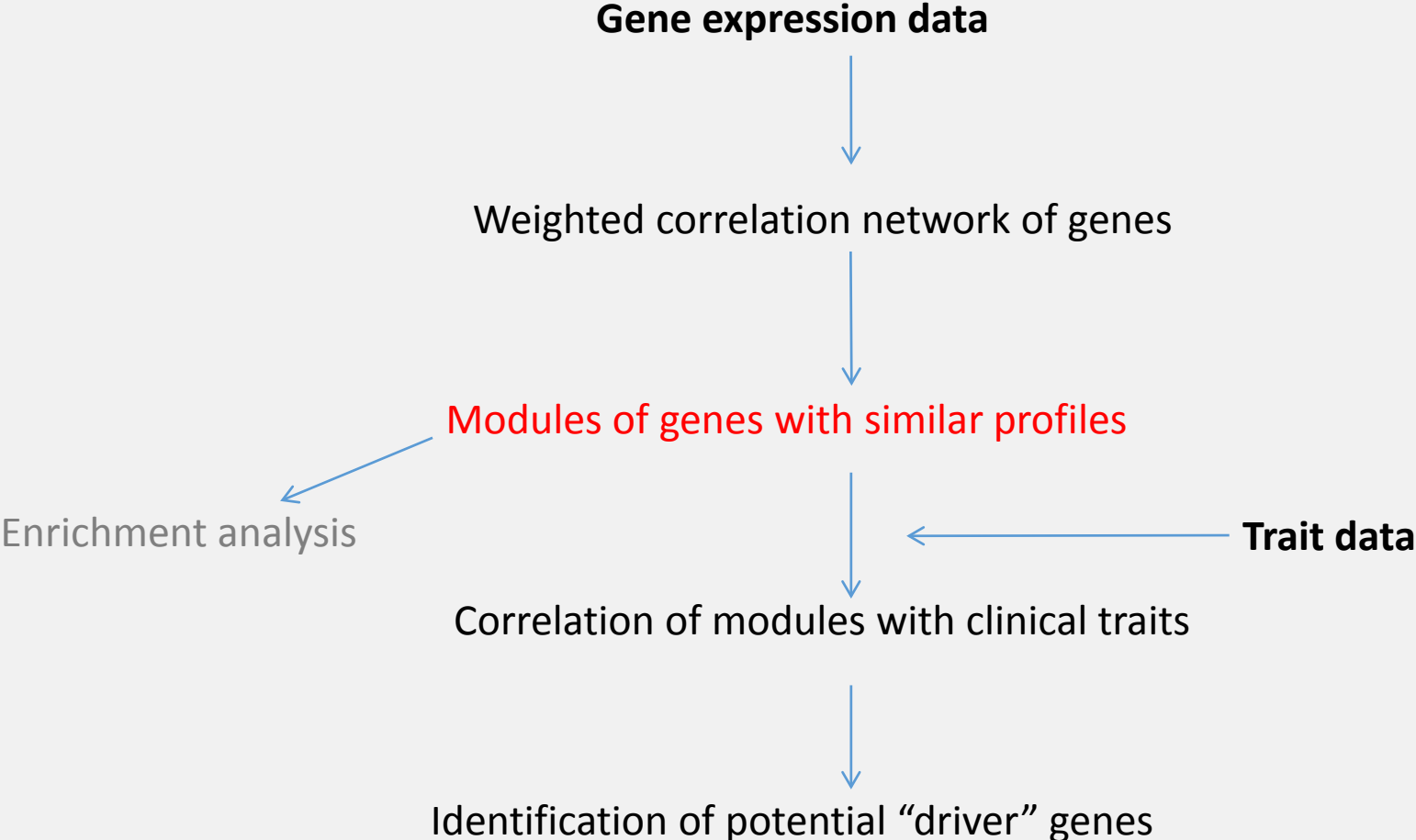


Mean connectivity



- Left plot: Choose power 6. Lowest possible power term where topology approximately fits a scale free network (on or above red horizontal line).
- Right plot: mean connectivity drops as power goes up. Must not drop too low.

Workflow



Detect modules of co-expressed genes

4 steps to get from network to modules

1. Compute dissimilarity between genes
(topological overlap measure dissimilarity)
2. Perform hierarchical clustering of genes
(obtain tree structure)
3. Divide clustered genes into modules
(cut tree branches)
4. Merge very similar modules
(use module “eigengenes”)

Step 1: Compute dissimilarity between genes

Similarity/dissimilarity between genes

Topological Overlap Measure (TOM):

- is a pairwise similarity measure between network nodes (genes)
- $TOM(i,j)$ is **high** if genes i,j have **many shared neighbors**
(overlap of their network neighbors is large)
- A high $TOM(i,j)$ implies that genes have similar expression patterns

Calculating TOM

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

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

Originally defined for unweighted networks

TOM similarity between two nodes:

1. Count number of shared neighbors
("agreement" of the set of neighboring nodes)

2. Normalize to [0,1]

TOM(i,j) = 0: no overlap of network neighbors

TOM(i,j) = 1: identical set of network neighbors

***Generalized to the case of weighted networks in Zhang and Horvath (2005),
first WGCNA paper***

All nodes are neighbors; counting them is not informative.

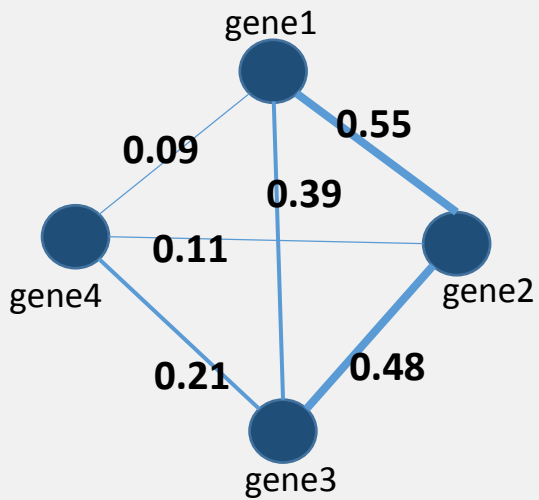
Compute agreement of the set of neighboring nodes based on edge strengths.

TOM as a **similarity** measure can be transformed into a **dissimilarity** measure $\text{distTOM} = 1 - \text{TOM}$.

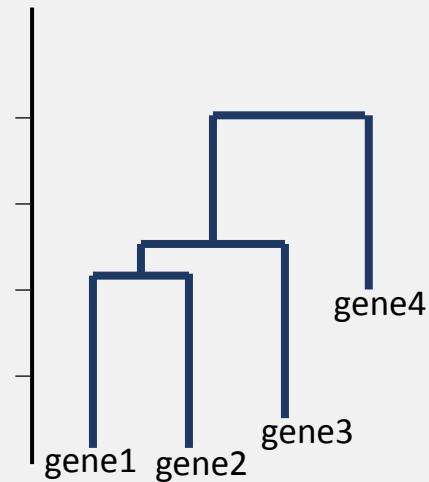
Step 2: Perform hierarchical clustering of genes

Compute gene dendrogram

Weighted correlation network
from gene expression data



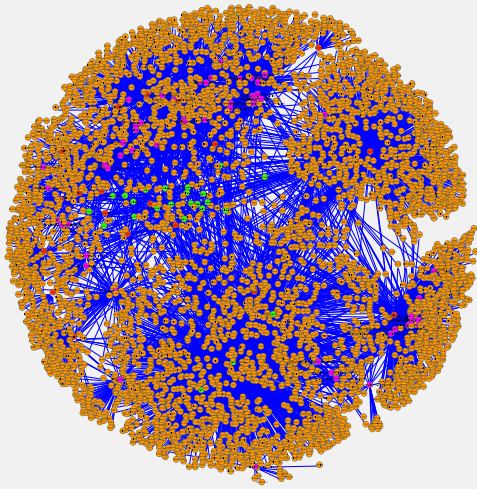
Clustering dendrogram



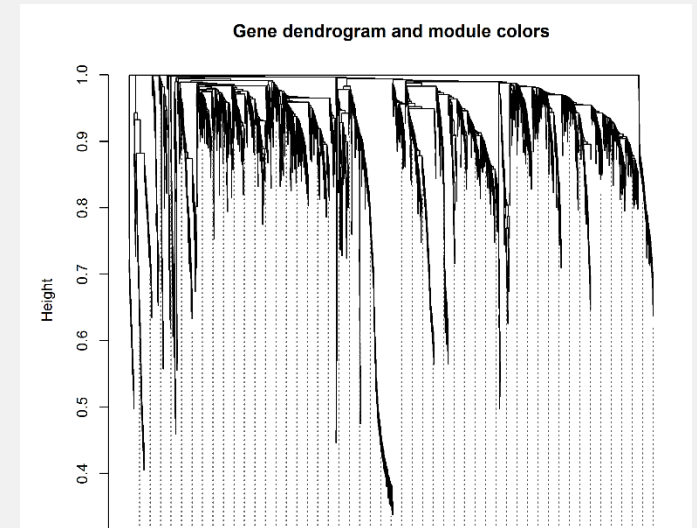
***(dis)similarity between genes:
Topological Overlap Measure TOM***

Compute gene dendrogram

Weighted correlation network
from gene expression data



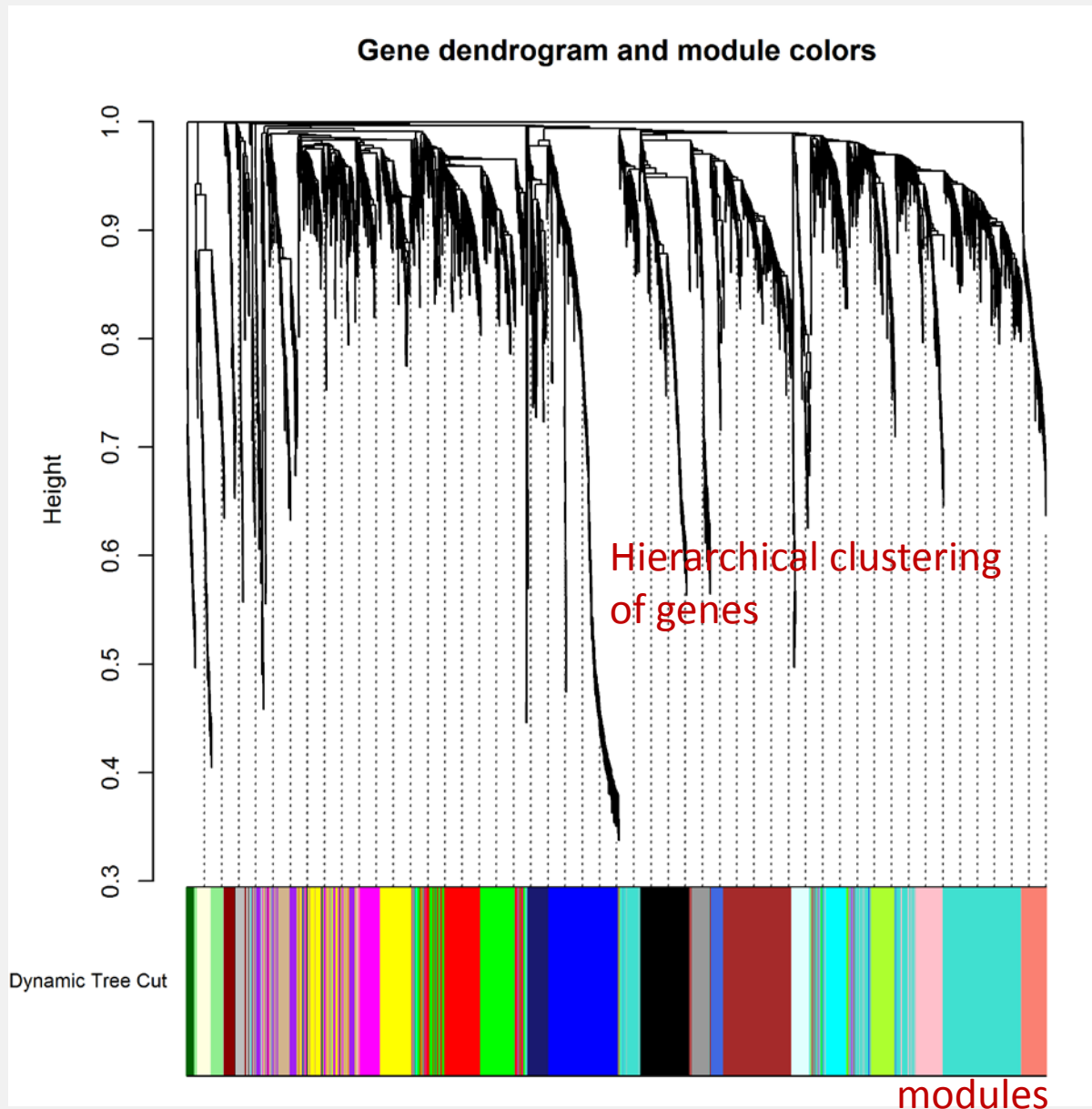
Clustering dendrogram



*(dis)similarity between genes:
Topological Overlap Measure TOM*

Step 3: Divide clustered genes into modules

Gene dendrogram and detected modules



Dynamic tree cut algorithm groups genes into modules

corFnc="pearson"; power=6; min. module size=30

Step 4: Merge very similar modules

Module eigengenes

A module eigengene is

a 1-dimensional data vector,
summarizing the expression data of the genes that form a module

How it is computed:

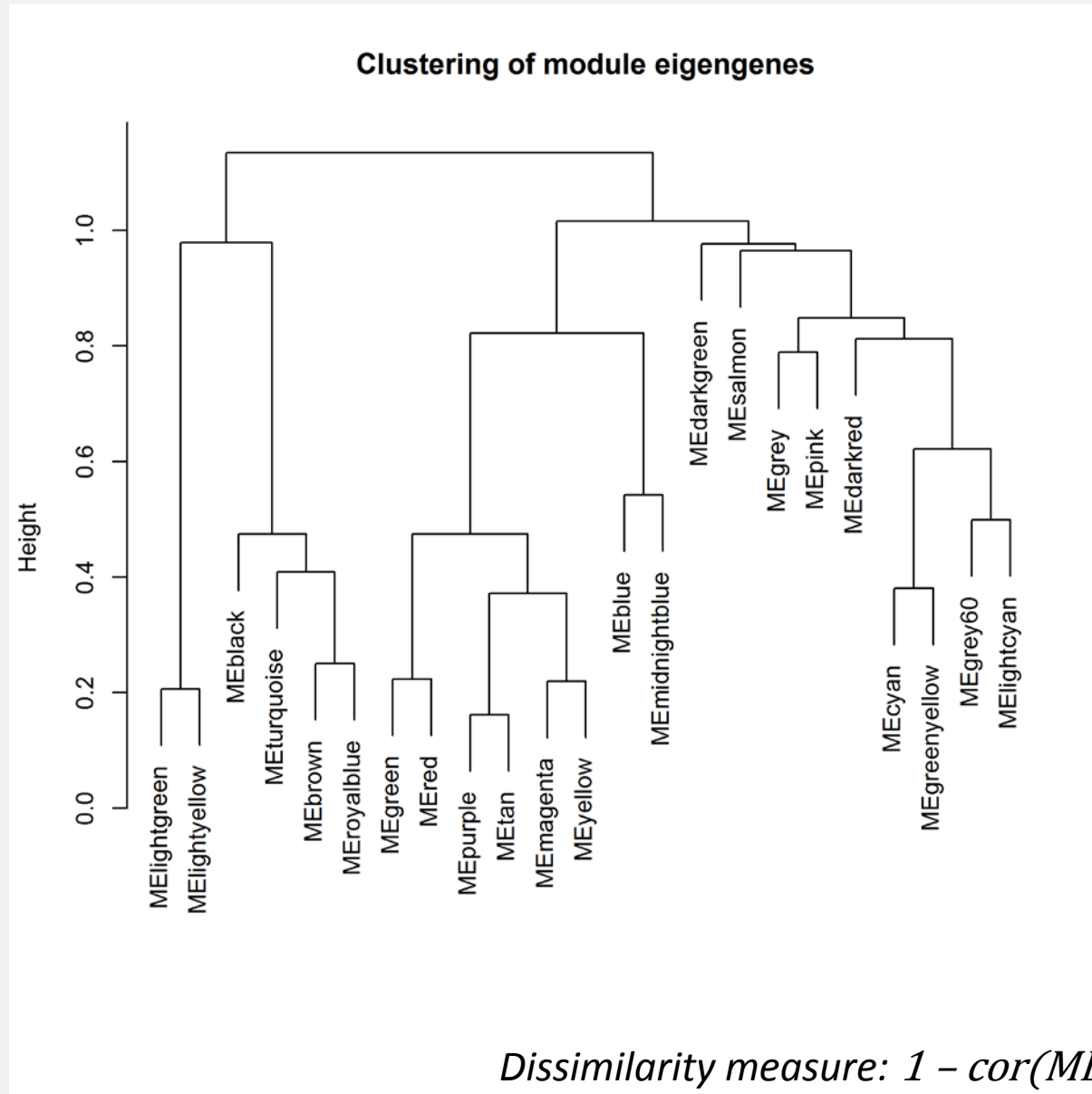
1st principal component of the expression data

What it is used for:

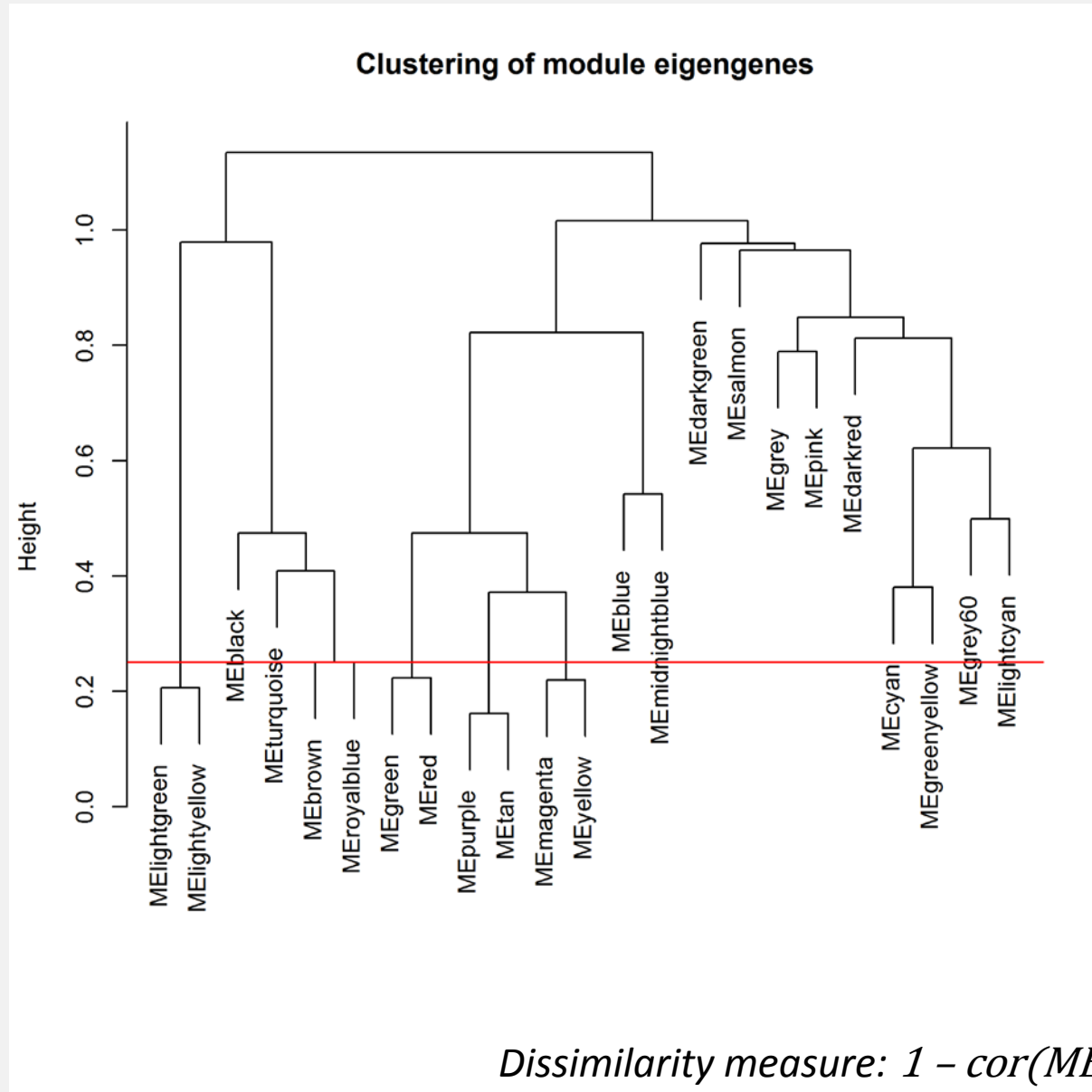
to represent the module in mathematical operations

- modules can be correlated with one another
- modules can be **clustered**
- modules can be correlated with external traits

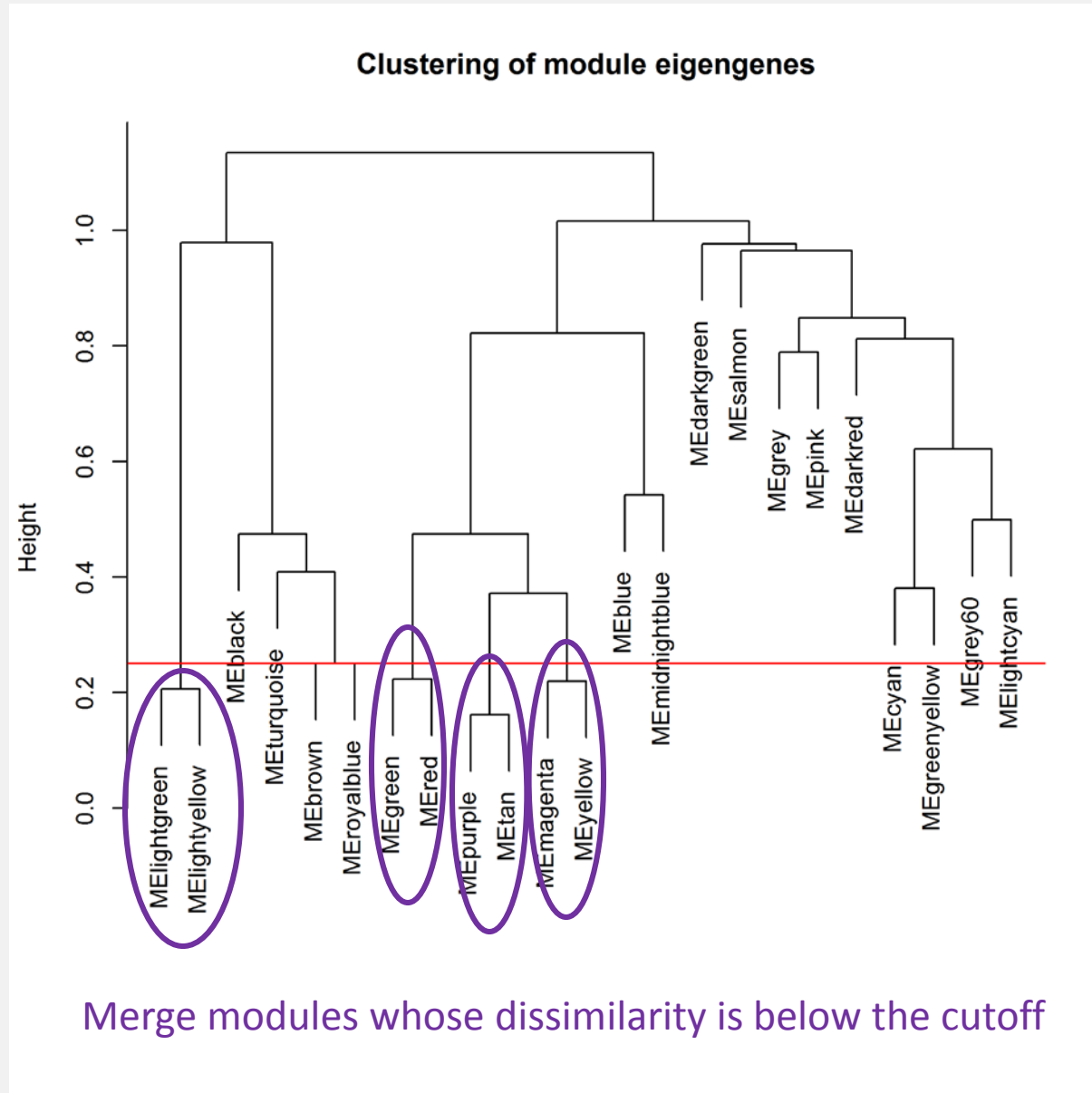
Clustering of module eigengenes



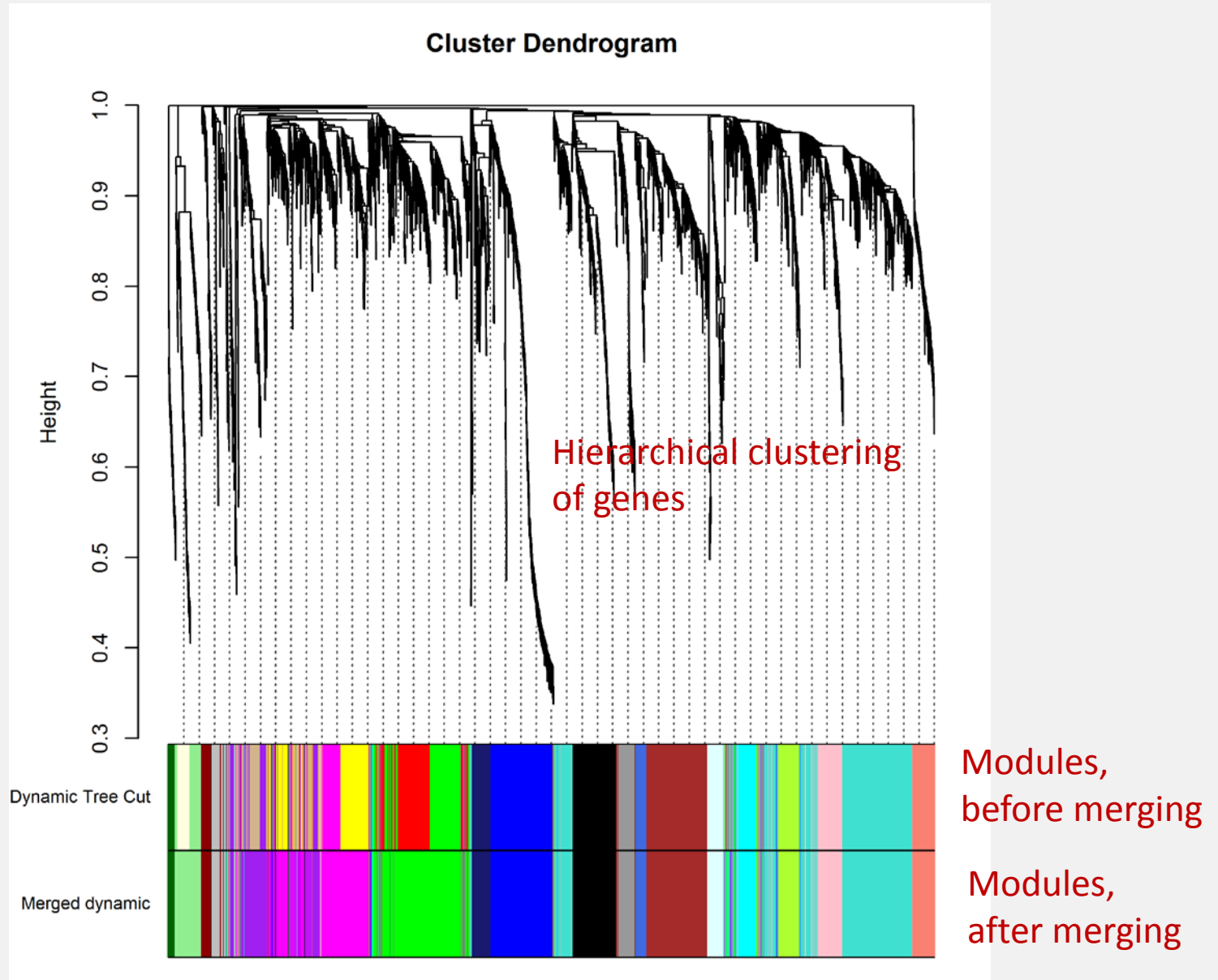
Clustering of module eigengenes



Clustering of module eigengenes



Gene dendrogram and detected modules, before and after merging



corFnc="pearson"; power=6; min. module size=30

Module Detection: Decisions to make

Dynamic Tree cut

Minimal module size

typically 20 or 30

Module Merging

Cutoff for module eigengene dendrogram

typically between 0.15 and 0.25

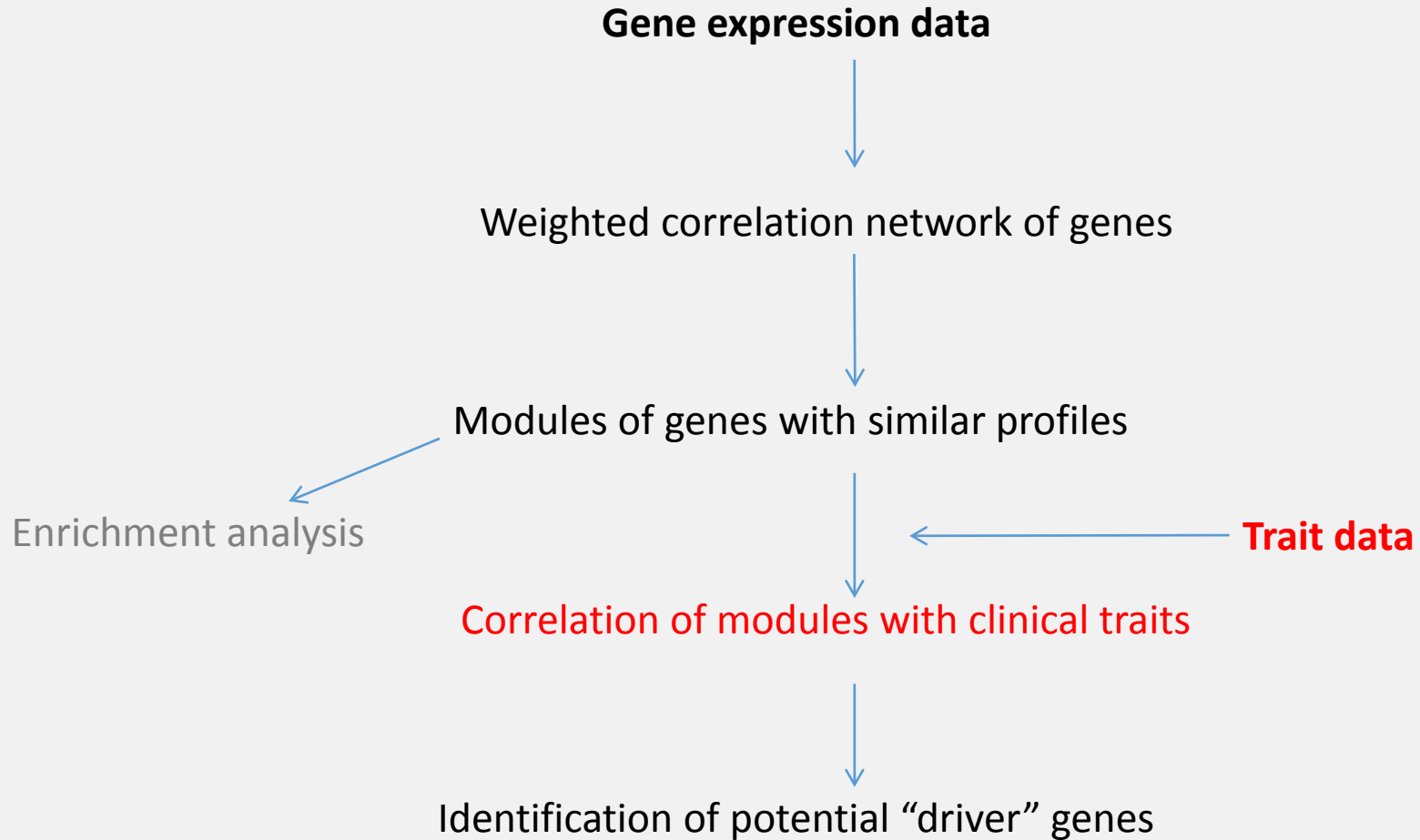
check if clusters look ok on dendrogram

Merge once or several times?

usually once, but merge step can be repeated

- *if some modules are very similar*
- *if we want larger modules*

Workflow



Correlate modules to external traits

Example: trait variables from Mouse Liver Data Set

weight_g
length_cm
ab_fat
other_fat
total_fat
Trigly
Total_Chol
HDL_Chol
UC
FFA
Glucose
LDL_plus_VLDL
MCP_1_phys

Insulin_ug_l
Glucose_Insulin
Leptin_pg_ml
Adiponectin
Aortic lesions
Aneurysm
Aortic_cal_M
Aortic_cal_L
CoronaryArtery_Cal
Myocardial_cal
BMD_all_limbs
BMD_femurs_only

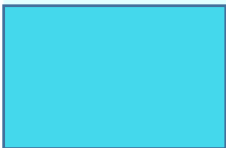



Compute correlations:
each module eigengene to each trait variable

cor(MEs, traitDat)

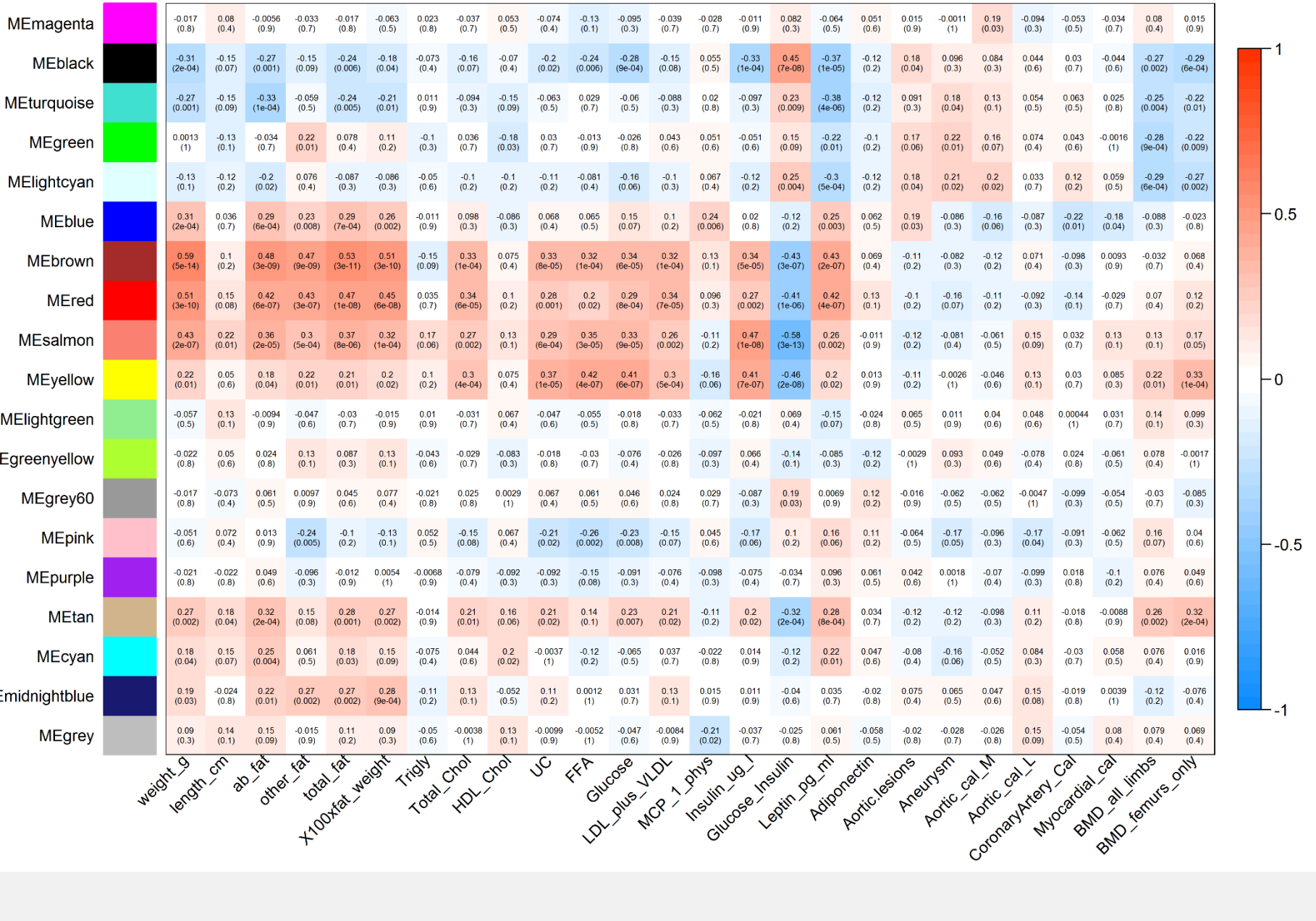
Table of module – trait correlations

- Identify modules highly correlated to traits of interest
- Identify traits highly correlated to multiple modules

Table excerpt:

		Weight (g)	Length (cm)	Ab_fat
MEcyan		-0.13 (0.1)	-0.12 (0.2)	-0.2 (0.02)
MEblue		0.31 (2e-04)	0.036 (0.7)	0.29 (6e-04)
MEbrown		0.59 (5e-14)	0.1 (0.2)	0.48 (3e-09)
MEred		0.51 (3e-10)	0.15 (0.08)	0.42 (6e-07)

Module-trait relationships



What to do if there are binary traits?

Examples of binary traits:

Sex (male/female)

Disease status (healthy/diseased)

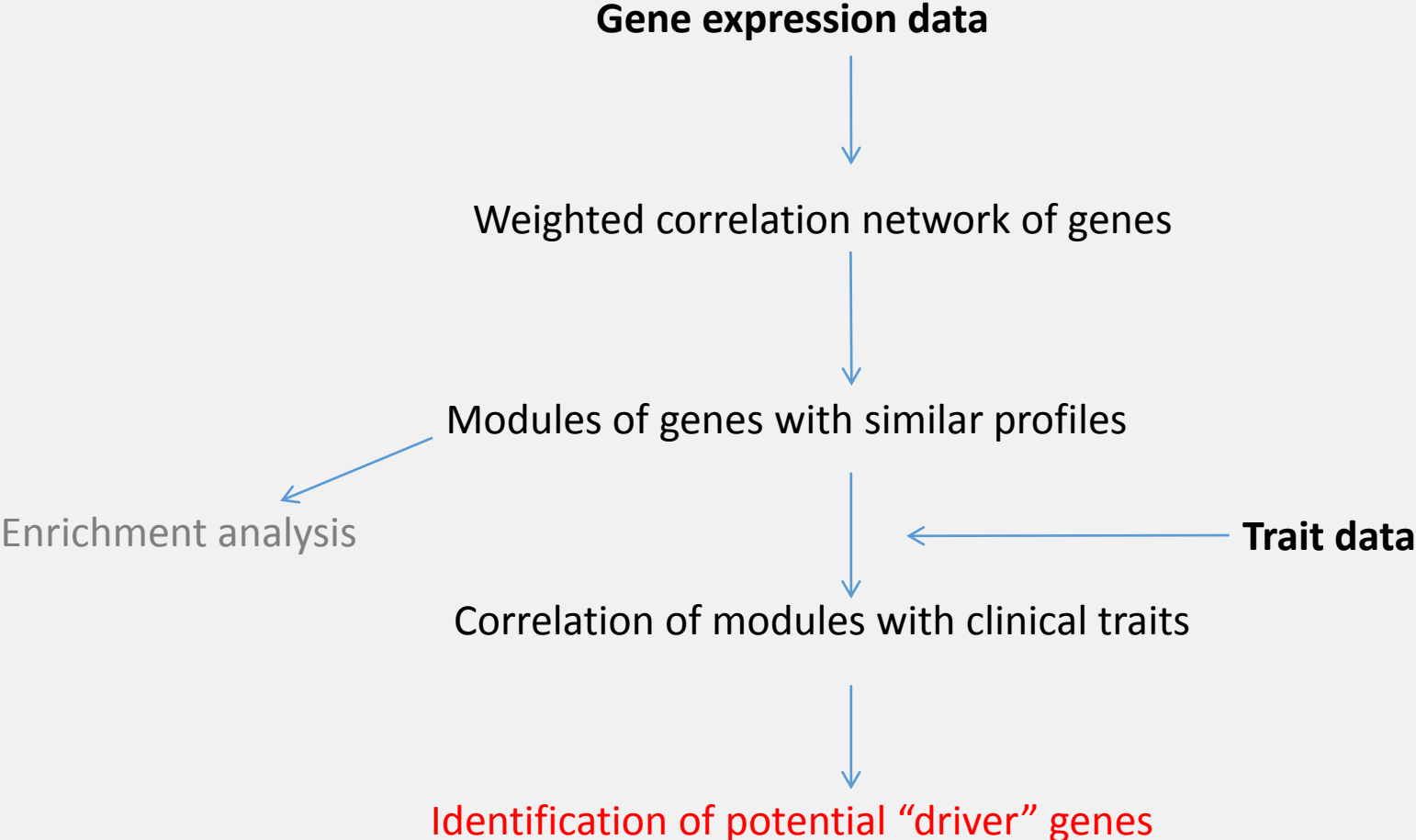
Pearson and **Spearman** correlation (function `cor`): will work

Biweight midcorrelation (function `bicor`) needs modification

- restrict number of values that will be treated as outliers
- turn off robust treatment for the binary variable

```
bicor(moduleEigengenes, datTraits, use="p",  
robustY=FALSE, maxPOutliers=0.1)
```

Workflow



Identify potential “driver” genes

Potential driver genes

Aim: Identify key genes that

- May **influence** the expression or function of **other genes**
- May be **causal** factors for a **trait** of interest

Caveat:

- WGCNA **cannot** show whether gene-gene or gene-trait relationships are **causal**.
- WGCNA **can** help find **candidate** genes for further study.

Potential driver genes

Strategy:

Identify those genes within a module that are

1) Highly connected within the module (hub genes)

AND

2) Most strongly correlated with a clinical/phenotypical trait of interest

How to detect hub genes inside a module?

The straightforward way:

gene(s) with highest intramodular connectivity
(=sum of in-module edge weights)

Alternative way proposed in WGCNA:

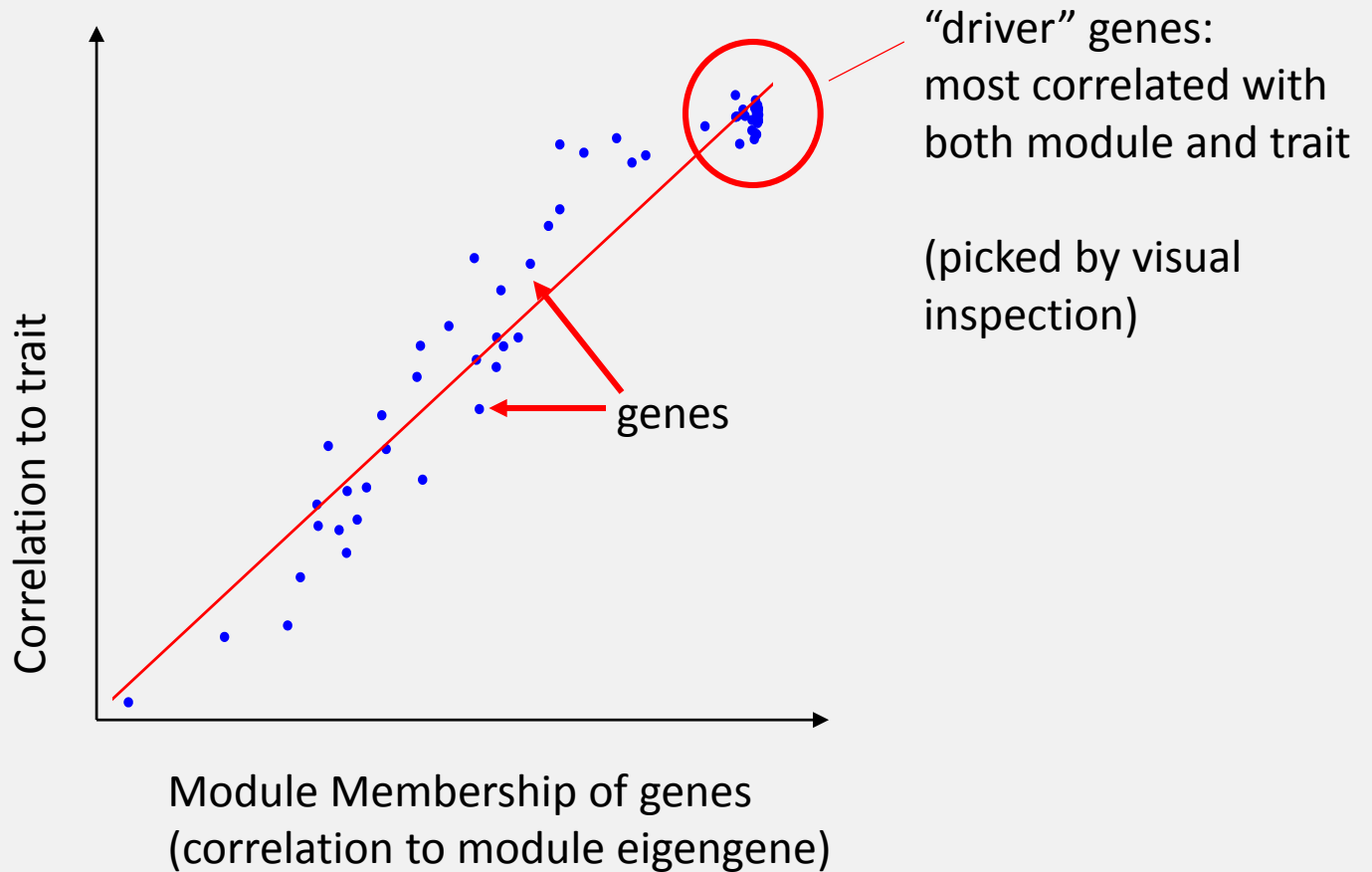
gene(s) with highest module membership

Module Membership of Genes

Module membership: **Correlation** of a **gene** to a **module eigengene**

- Genes with high module membership are good representatives of the overall expression profile in the module
- Genes with high module membership tend to be “hub” genes in the module (high intramodule connectivity)
- A gene can have high membership in several modules (not just the one to which it is assigned)

Potential driver genes



Evaluate module quality

Check module sizes

Potential issues:

Very large modules

- may make biological sense, but difficult to handle

Many similar modules

- level of merging not sufficient

Quality checks on modules

Are modules better than random groupings of genes?

Connectivity

mean intra-module connectivity

mean ratio of intra-module / total connectivity

Trait correlations

strong correlation between module eigengenes and traits of interest

strong correlation between gene module membership and gene-trait corr.

Functional enrichment

many functionally related genes in the same module →

Next up:
Presentation on
GO enrichment
analysis

Quality checks on modules can be simple or complex

Assessing modules by intramodular connectivity

- Simple ranking of modules
(highest to lowest mean connectivity)
- Statistical analysis
obtain p-values, e.g. via bootstrap
(is module internally more connected than would be expected by chance?)

Quality checks on modules can be simple or complex

Assessing modules by functional enrichment of GO terms

multiple methods exists

with or without taking GO hierarchy into account (next presentation)

Two purposes of module evaluation:

1) Identify modules of interest for further analysis

2) Assess overall quality of modules

if not satisfactory, re-do network construction and/or module detection
with changed parameters



Resources

WGCNA theory papers and applied papers, including PowerPoint presentations

<https://labs.genetics.ucla.edu/horvath/CoexpressionNetwork/>

Examples of topics

Constructing a gene coexpression network (2005)

Understanding the Topological Overlap Measure (2007)

Dynamic Tree Cut (2007)

Eigenegene networks (2007)

When Is Hub Gene Selection Better than Standard Meta-Analysis (2013)

WGCNA R package Website, including many tutorials

<https://labs.genetics.ucla.edu/horvath/CoexpressionNetwork/Rpackages/WGCNA/>