

A Geometric Interpretation of Gene Co-Expression Network Analysis

Steve Horvath, Jun Dong

Outline

- Network and network concepts
- Approximately factorizable networks
- Gene Co-expression Network
 - Eigengene Factorizability, Eigengene Conformity
 - Eigengene-based network concepts
- What can we learn from the geometric interpretation?

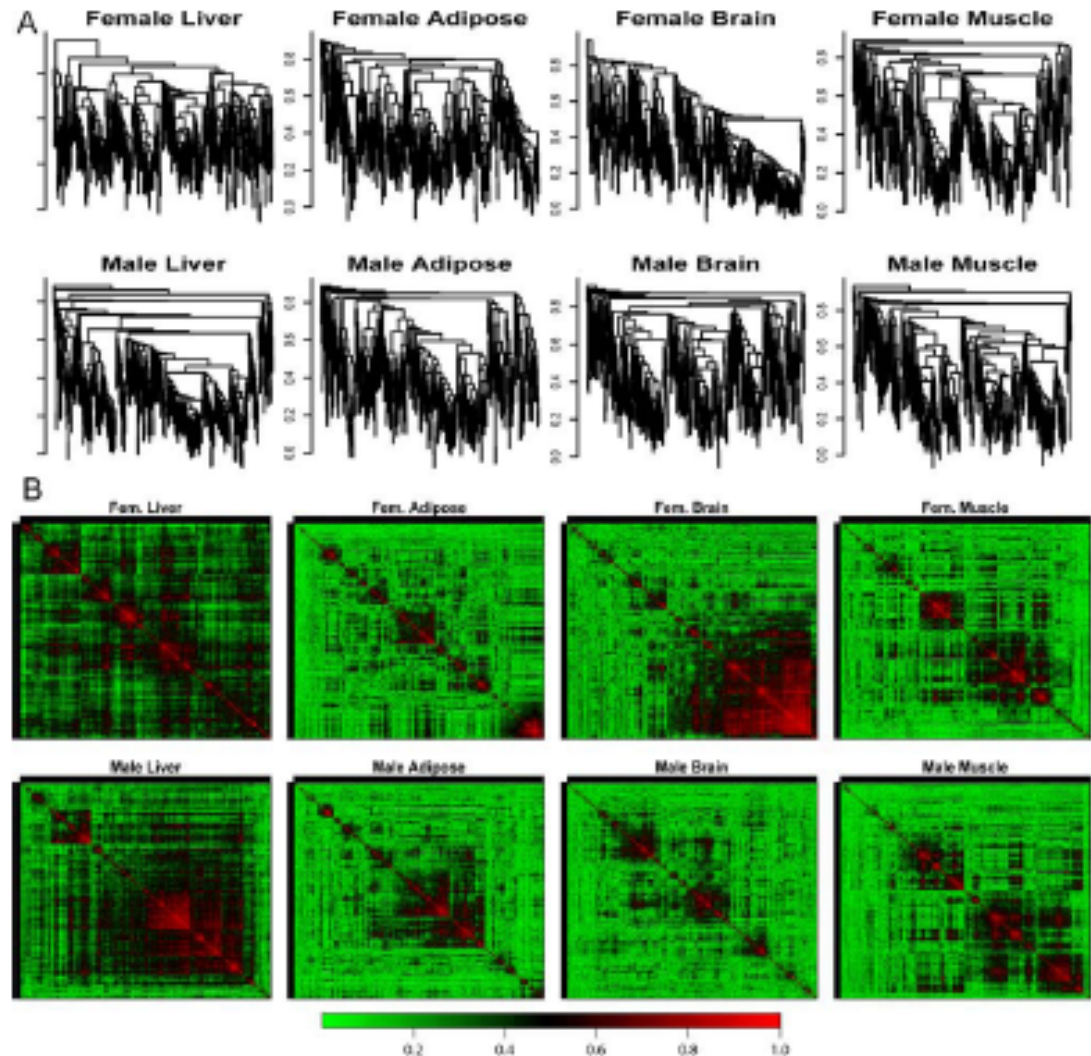
Network=Adjacency Matrix

- A network can be represented by an adjacency matrix, $A=[a_{ij}]$, that encodes whether/how a pair of nodes is connected.
 - A is a symmetric matrix with entries in $[0,1]$
 - For unweighted network, entries are 1 or 0 depending on whether or not 2 nodes are adjacent (connected)
 - For weighted networks, the adjacency matrix reports the connection strength between node pairs
 - Our convention: diagonal elements of A are all 1.

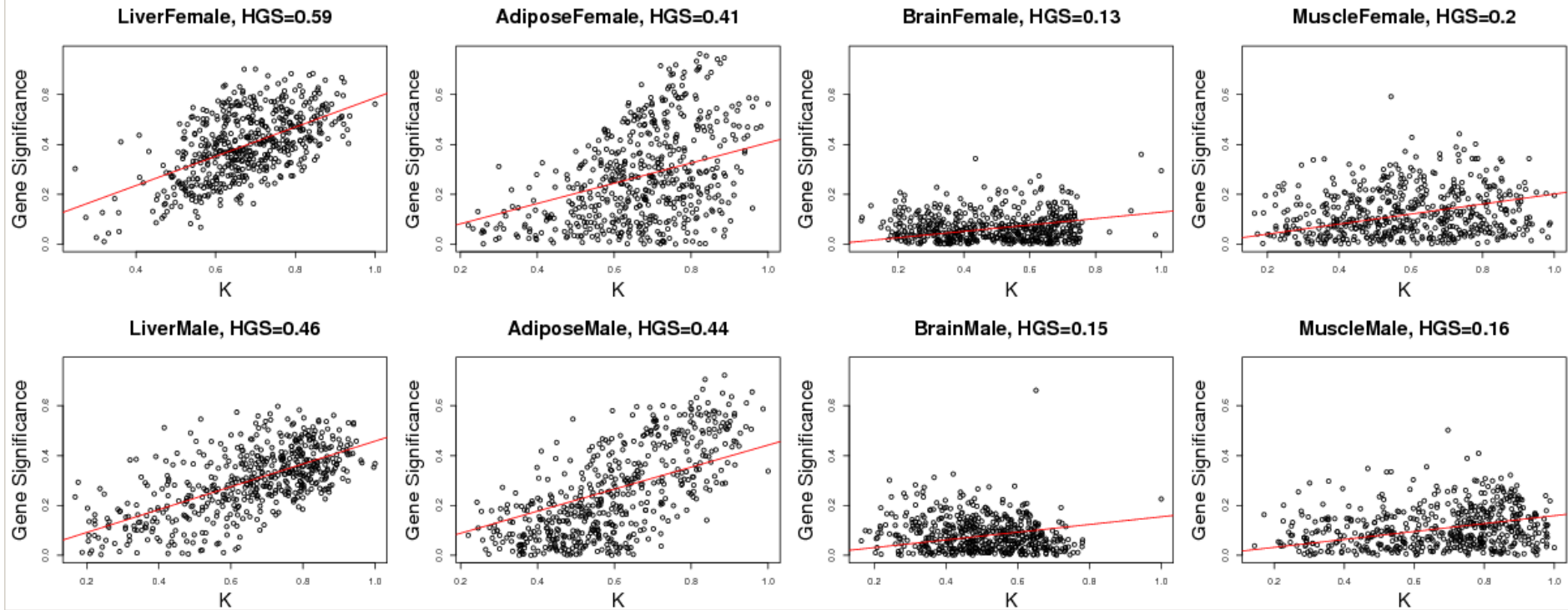
Motivational example I: Pair-wise relationships between genes across different mouse tissues and genders

Challenge:
Develop simple
descriptive measures that
describe the patterns.

Solution:
The following network
concepts are useful:
density, centralization,
clustering coefficient,
heterogeneity



Motivational example (continued)



Challenge: Find a simple measure for describing the relationship between gene significance and connectivity

Solution: network concept called hub gene significance

Backgrounds

- Network concepts are also known as network statistics or network indices
 - Examples: connectivity (degree), clustering coefficient, topological overlap, etc
- Network concepts underlie network language and systems biological modeling.
- Dozens of potentially useful network concepts are known from graph theory.

Review of *some*
fundamental network concepts
which are defined for all networks
(not just co-expression networks)

Connectivity

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

$$\text{Connectivity}_i = k_i = \sum_{j \neq i} a_{ij}$$

$$\text{Scaled connectivity} = K_i = \frac{k_i}{\max(k)}$$

Density

- Density= mean adjacency
- Highly related to mean connectivity

$$\text{Density} = \frac{\sum_i \sum_{j \neq i} a_{ij}}{n(n-1)} = \frac{\text{mean}(k)}{n-1}$$

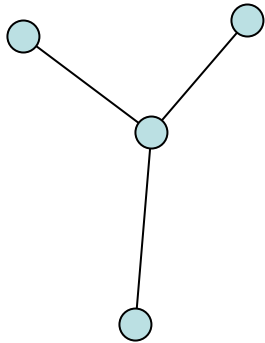
where n is the number of network nodes.

Centralization

$$Centralization = \frac{n}{n-2} \left(\frac{\max(k)}{n-1} - Density \right) \approx \frac{\max(k)}{n-1} - Density$$

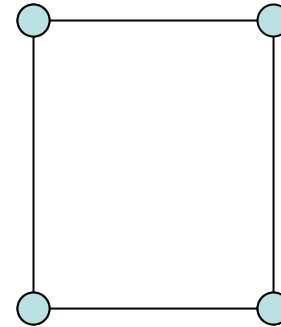
= 1 if the network has a star topology

= 0 if all nodes have the same connectivity



Centralization = 1

because it has a star topology



Centralization = 0

because all nodes have the same connectivity of 2

Heterogeneity

- Heterogeneity: coefficient of variation of the connectivity
- Highly heterogeneous networks exhibit hubs

$$\textit{Heterogeneity} = \frac{\sqrt{\textit{variance}(k)}}{\textit{mean}(k)}$$

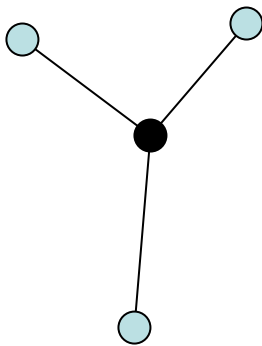
Clustering Coefficient

Measures the cliquishness of a particular node

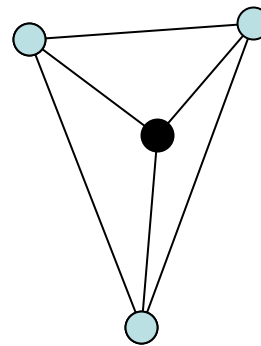
« A node is cliquish if its neighbors know each other »

$$ClusterCoef_i = \frac{\sum_{l \neq i} \sum_{m \neq i, l} a_{il} a_{lm} a_{mi}}{\left(\sum_{l \neq i} a_{il} \right)^2 - \sum_{l \neq i} a_{il}^2}$$

This generalizes directly to weighted networks (Zhang and Horvath 2005)



Clustering Coef of
the black node = 0



Clustering Coef = 1

The topological overlap dissimilarity is used as input of hierarchical clustering

$$TOM_{ij} = \frac{\sum_{u \neq i, j} a_{iu} a_{uj} + a_{ij}}{\min(k_i, k_j) + 1 - a_{ij}}$$

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

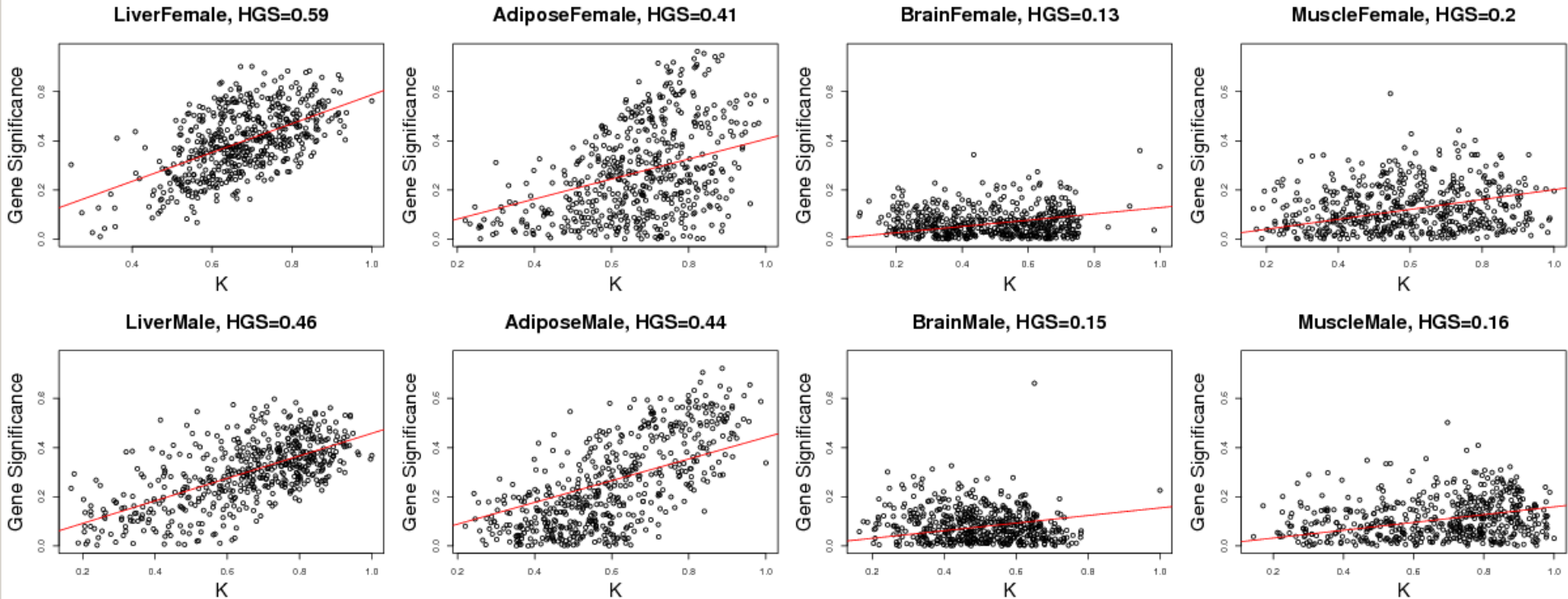
- Generalized in Zhang and Horvath (2005) to the case of weighted networks
- Generalized in Li and Horvath (2006) to multiple nodes
- Generalized in Yip and Horvath (2007) to higher order interactions

Network Significance

- Defined as average gene significance
- We often refer to the network significance of a module network as module significance.

$$\textit{NetworkSignif} = \frac{\sum GS_i}{n}$$

Hub Gene Significance= slope of the regression line (intercept=0)



$$HubGeneSignif = \frac{\sum GS_i K_i}{\sum (K_i)^2}$$

Q: What do all of these fundamental network concepts have in common?

They are functions of the adjacency matrix A and/or a gene significance measure GS .

CHALLENGE

Find relationships between these and other seemingly disparate network concepts.

- For general networks, this is a difficult problem.
- But a solution exists for a special subclass of networks: approximately factorizable networks

Definition of an approximately factorizable network

Definitions:

The adjacency matrix A is **approximately factorizable** if there exists a vector CF with non-negative elements such that

$$a_{ij} \approx CF_i CF_j \quad \text{for all } i \neq j$$

CF_i is referred to as the **conformity** of the i -th node

Why is this relevant?

Answer: Because modules are often approximately factorizable

Algorithmic definition of the conformity and a measure of factorizability

We define the conformity as a maximizer of the factorizability function

$$F_A(v) = 1 - \frac{\sum_i \sum_{j \neq i} (a_{ij} - v_i v_j)^2}{\sum_i \sum_{j \neq i} (a_{ij})^2}$$

We use an iterative algorithm to approximate the conformity vector CF .

A measure of factorizability $F(A)$ is defined as $F_A(CF)$.

Conceptually related to a factor analysis of A .

Empirical Observation 1

- Sub-networks comprised of module genes tend to be approximately factorizable, i.e.

$$a_{ij} \approx CF_i CF_j \quad \text{for all } i \neq j$$

Empirical evidence is provided in the following article:

Dong J, Horvath S (2007) Understanding Network Concepts in Modules BMC Systems Biology 2007, 1:24

This observation implies the following observation 2...

Observation 2: Approximate relationships among network concepts in approximately factorizable networks

$$\text{mean}(\text{ClusterCoef}) \approx (1 + \text{Heterogeneity}^2)^2 \times \text{Density}$$

$$\text{TopOverlap}_{ij} \approx \frac{\max(k_i, k_j)}{n} \times (1 + \text{Heterogeneity}^2)$$

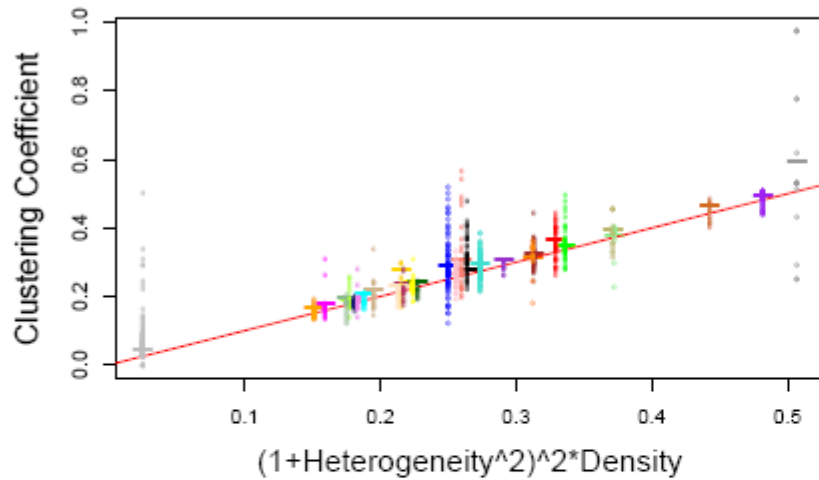
$$\text{TopOverlap}_{[1]j} \approx (\text{Centralization} + \text{Density}) \times (1 + \text{Heterogeneity}^2)$$

where [1] denotes the index of the most highly connected hub

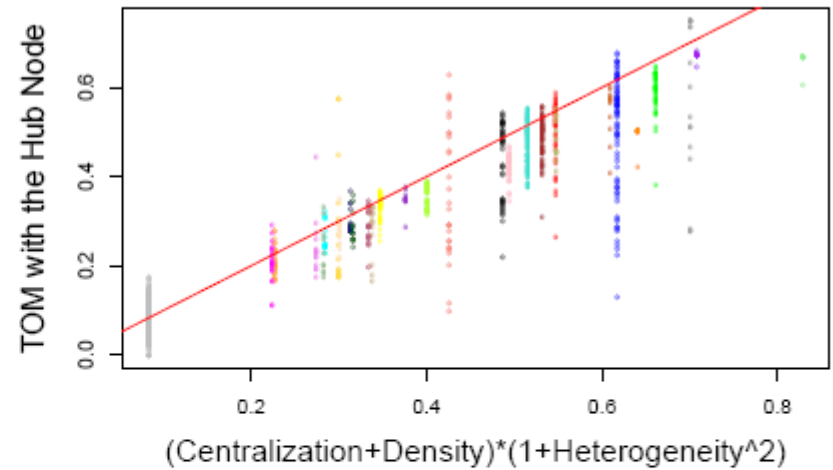
Drosophila PPI module networks: the relationship between fundamental network concepts.

A

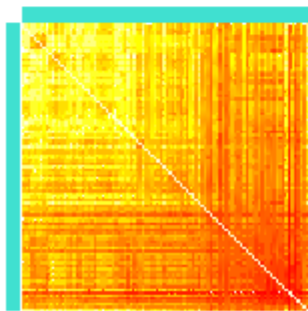
$R^2=0.87$



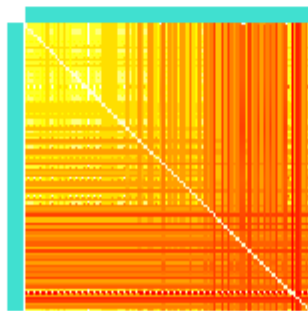
$R^2=0.91$



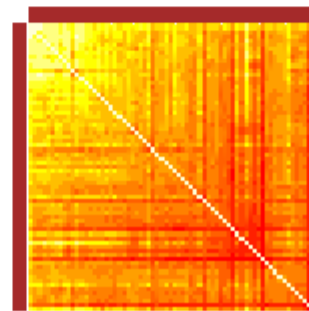
C



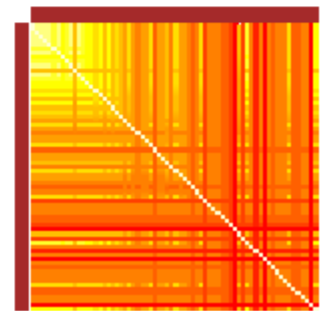
D



E



F



What if we focus on gene co-expression network?

Weighted Gene Co-expression Network

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

where x_i is the expression profile for gene i ,
and mathematically a vector of expression values
across multiple samples.

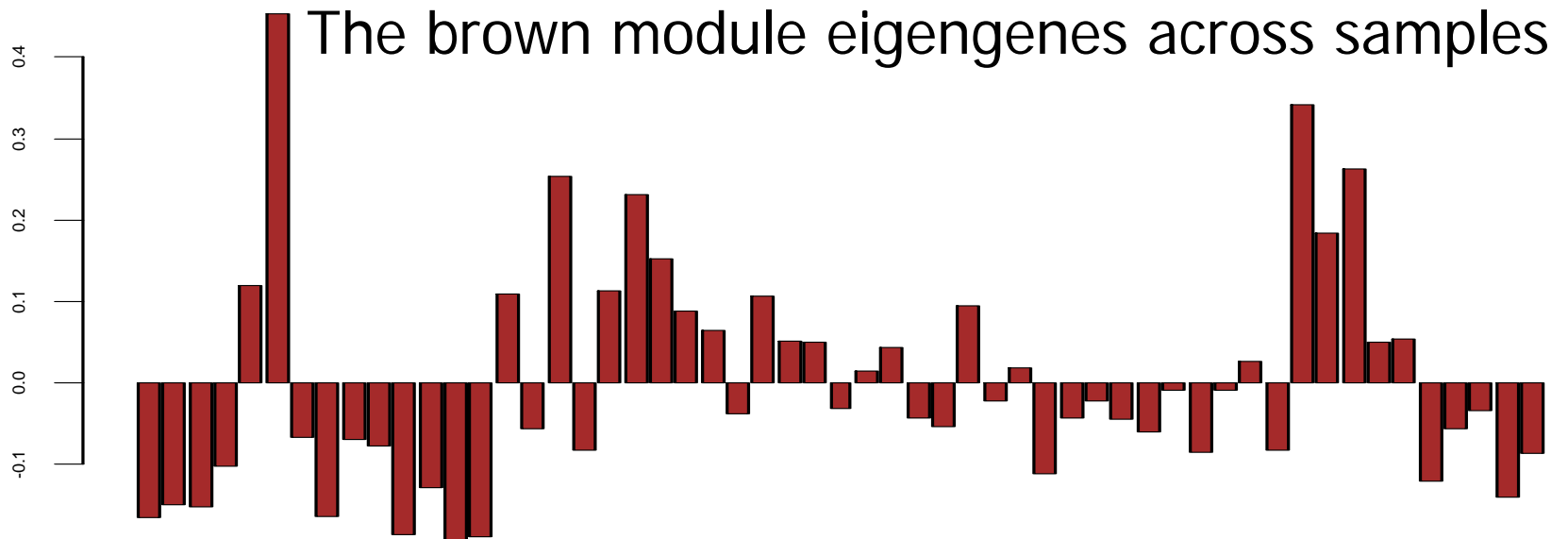
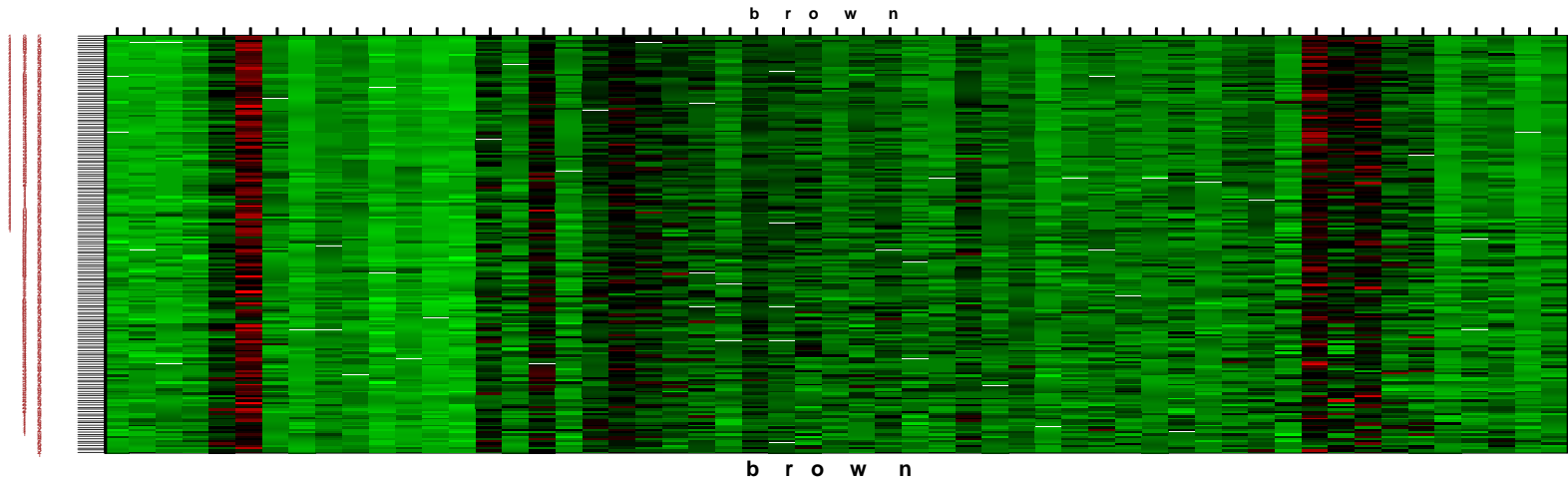
Note: Unweighted Network is

$$A = [a_{ij}] = [I(| \text{cor}(x_i, x_j) | \geq \tau)]$$

where $I(.)$ is an indicator function.

Module Eigengene= measure of over-expression=average redness

Rows, =genes, Columns=microarray



Recall that the module eigengene is defined by the singular value decomposition of X

- X = gene expression data of a module
- Aside: gene expressions (rows) have been standardized across samples (columns)

$$X = UDV^T$$

$$U = (u_1 \quad u_2 \quad \cdots \quad u_m)$$

$$\tilde{V} = (v_1 \quad v_2 \quad \cdots \quad v_m)$$

$$D = \text{diag}(|d_1|, |d_2|, \dots, |d_m|)$$

Message: v_1 is the module eigengene E

Question: When are co-expression modules factorizable?

Question: Characterize gene expression data X that lead to an approximately factorizable correlation matrix

Solution:

Define the eigengene based factorizability as follows

$$EF(X) = \frac{|d_1|^4}{\sum_j |d_j|^4} = 1 - \frac{\|cor(X) - C(C)^\tau\|_F^2}{\|cor(X)\|_F^2}$$

where $C_i = cor(x_i, E)$.

Thus, $cor(X)$ is approximately factorizable if $EF(X) \approx 1$.

Note that a factorizable correlation matrix implies a factorizable weighted co-expression network

$$\begin{aligned} a_{i,j} &= | \text{cor}(x_i, x_j) |^\beta \\ &\approx | \text{cor}(x_i, E) |^\beta | \text{cor}(x_j, E) |^\beta = a_{e,i} a_{e,j} \end{aligned}$$

We refer to the following as weighted eigengene conformity

$$a_{e,i} = | \text{cor}(x_i, E) |^\beta$$

If $EF(X) \approx 1$

Weighted gene co-expression network and its eigengene-based approximation if $EF(X^{(q)}) \approx 1$

Co-Expression Network

Eigengene-based counterpart

network

$$A = |cor(X)|^\beta$$

$$A_E = a_e a_e^T$$

gene significance

$$GS_i = |cor(x_i, T)|^\beta$$

$$GS_{E,i} = a_{e,i} a_{e,t}$$

Connectivity(i)

$$k_i = \sum_{j \neq i} a_{ij}$$

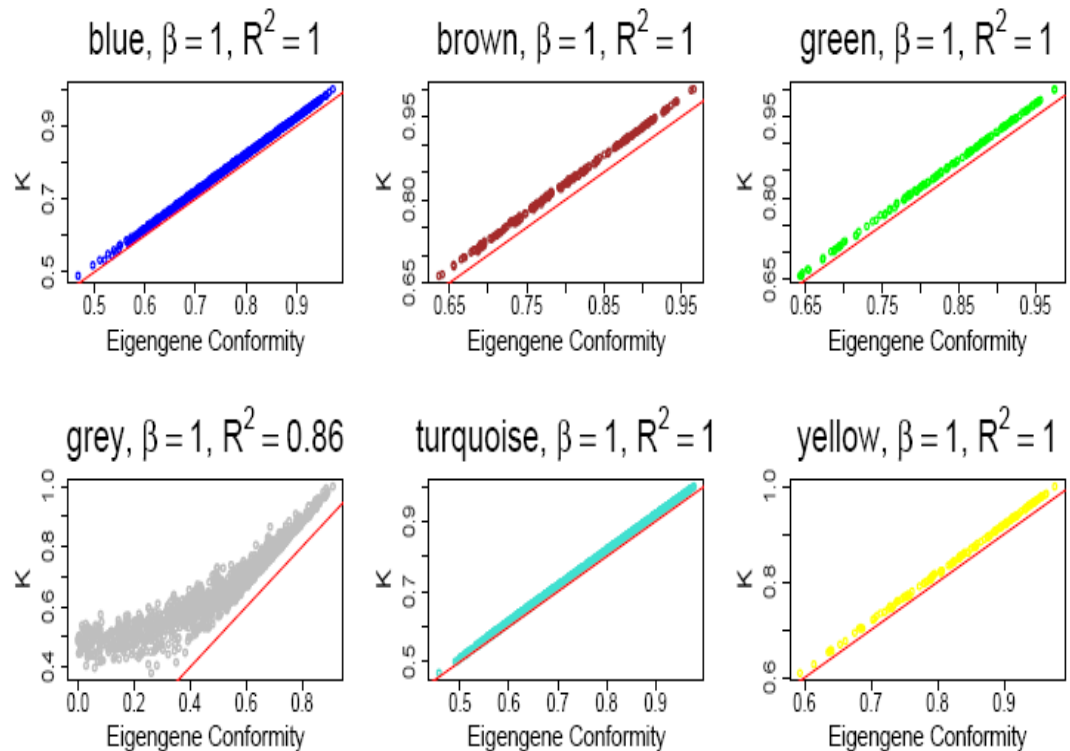
$$k_{E,i} = a_{e,i} \sum_j a_{e,j}$$

Theoretical relationships in co-expression modules with high eigengene factorizability

Result "Group conform behavior leads to a lot of friends."

More precisely, the scaled intramodular connectivity K_i approximates the eigengene conformity, i.e. $K_i \approx a_{e,i} = |cor(x_i, E)|^\beta$.

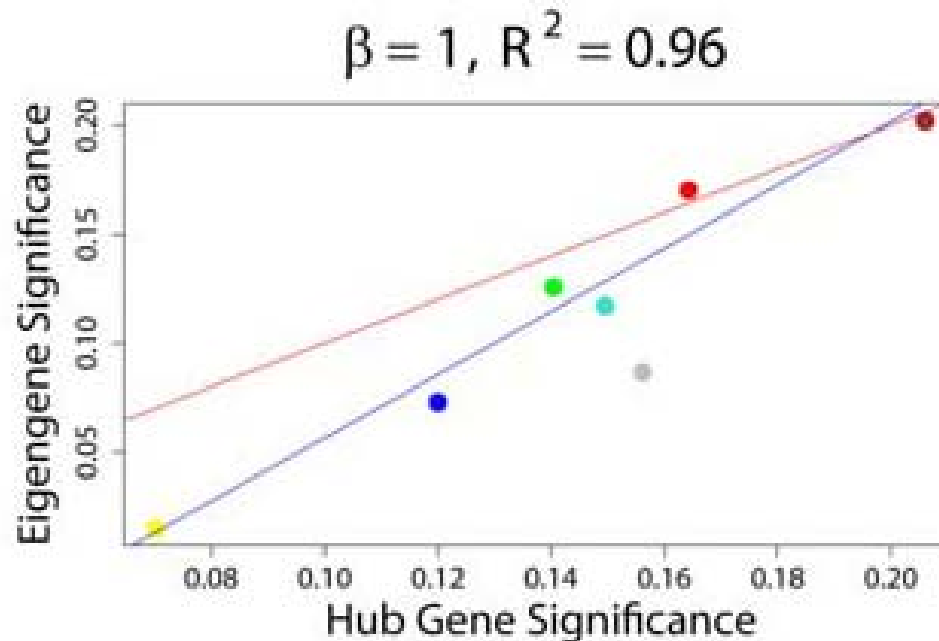
Message: the smaller the angle between x_i and E , the more connected is the i-th gene.



Result about hub gene significance:

Given a trait based gene significance measure $GS_T(i) = |cor(x_i, T)|^\beta$, the hub gene significance approximates the eigengene significance, $HGS \approx |cor(E, T)|^\beta$.

Message: the smaller the angle between E and T , the higher is the trait-significance of intramodular hubs and the higher is the module significance (average GS).

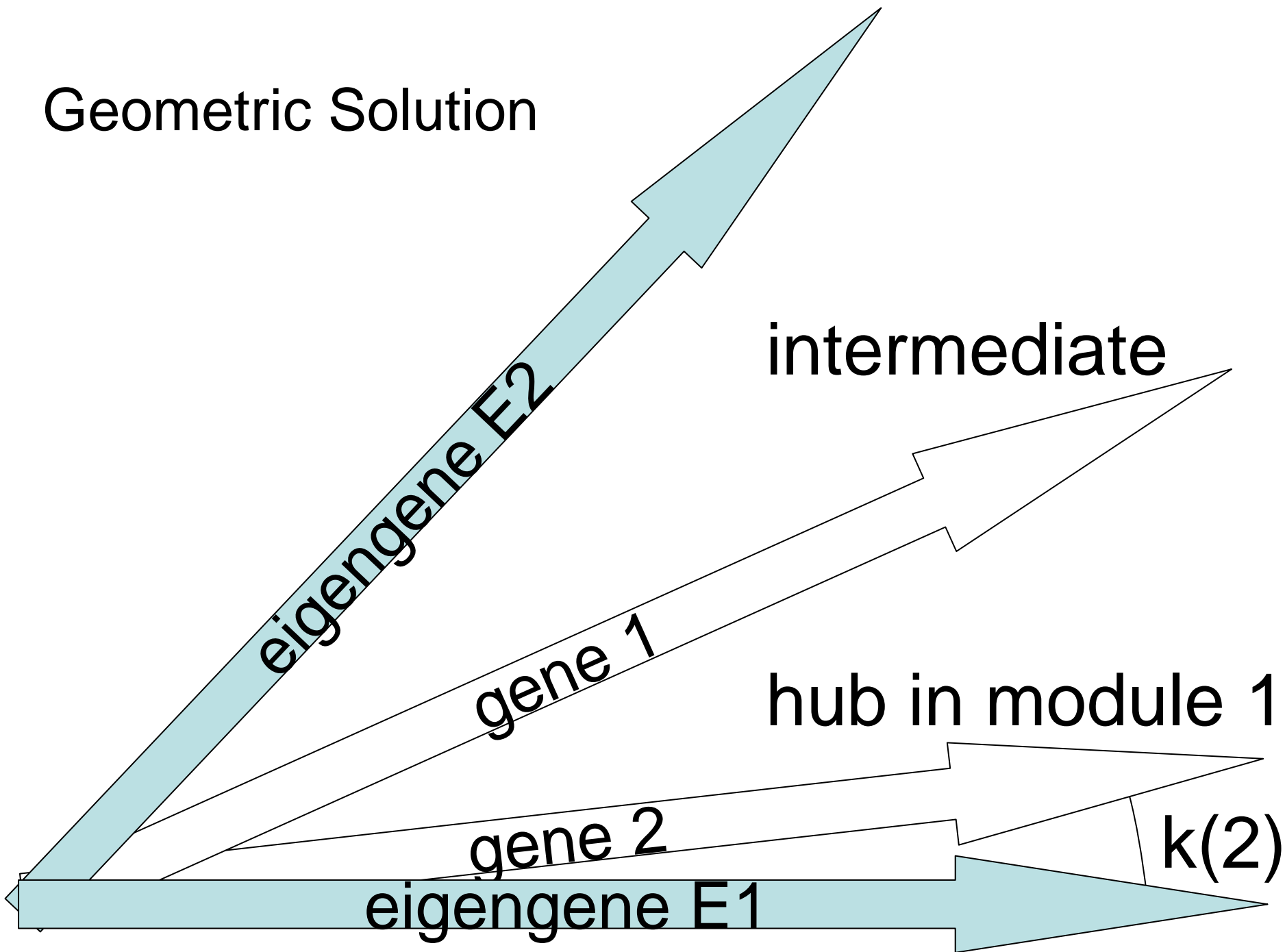


**What can network theorists learn
from the geometric interpretation?
Some examples...**

Problem

- Show that genes that lie intermediate between two distinct co-expression modules cannot be hub genes in these modules.

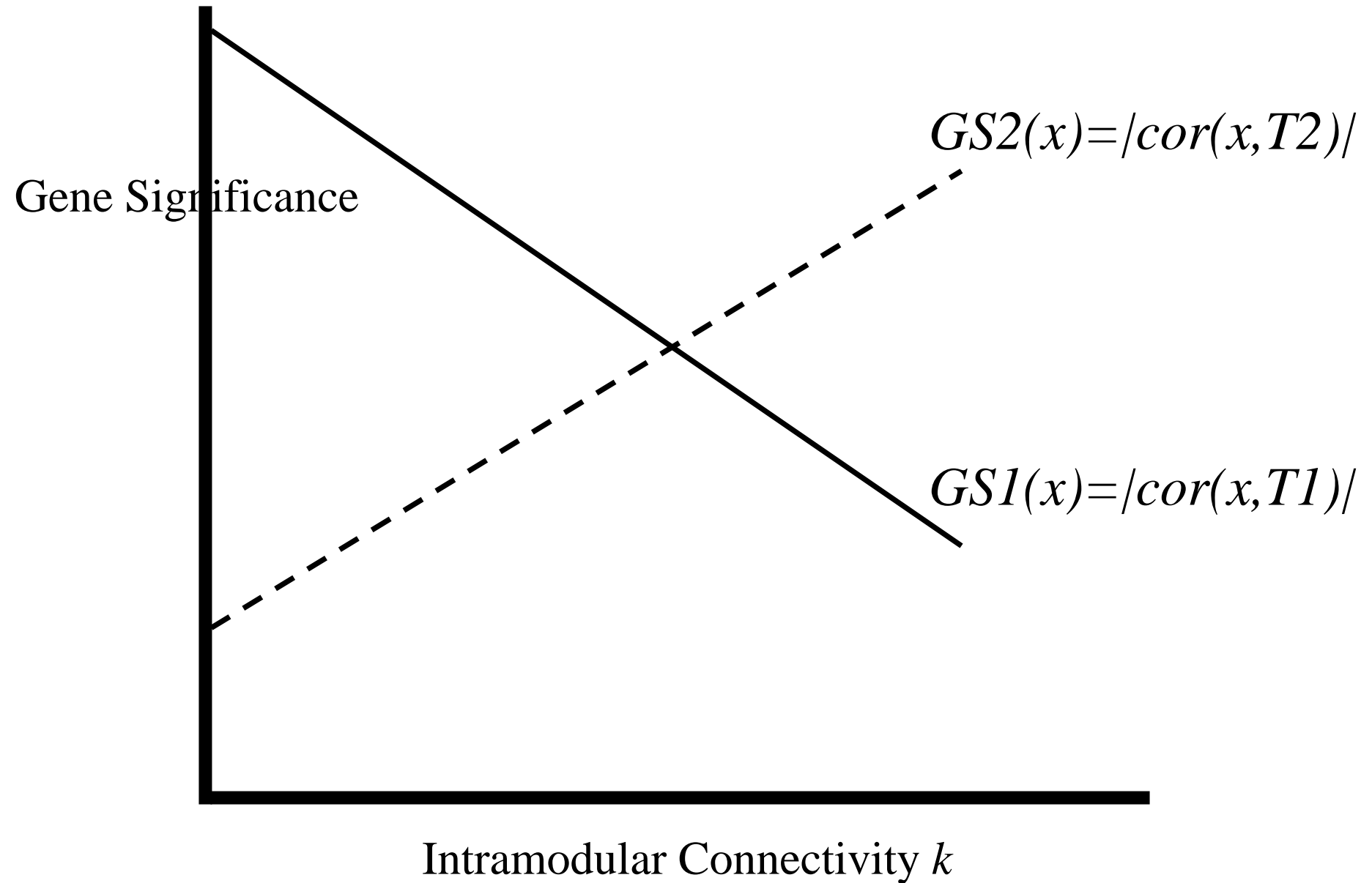
Geometric Solution

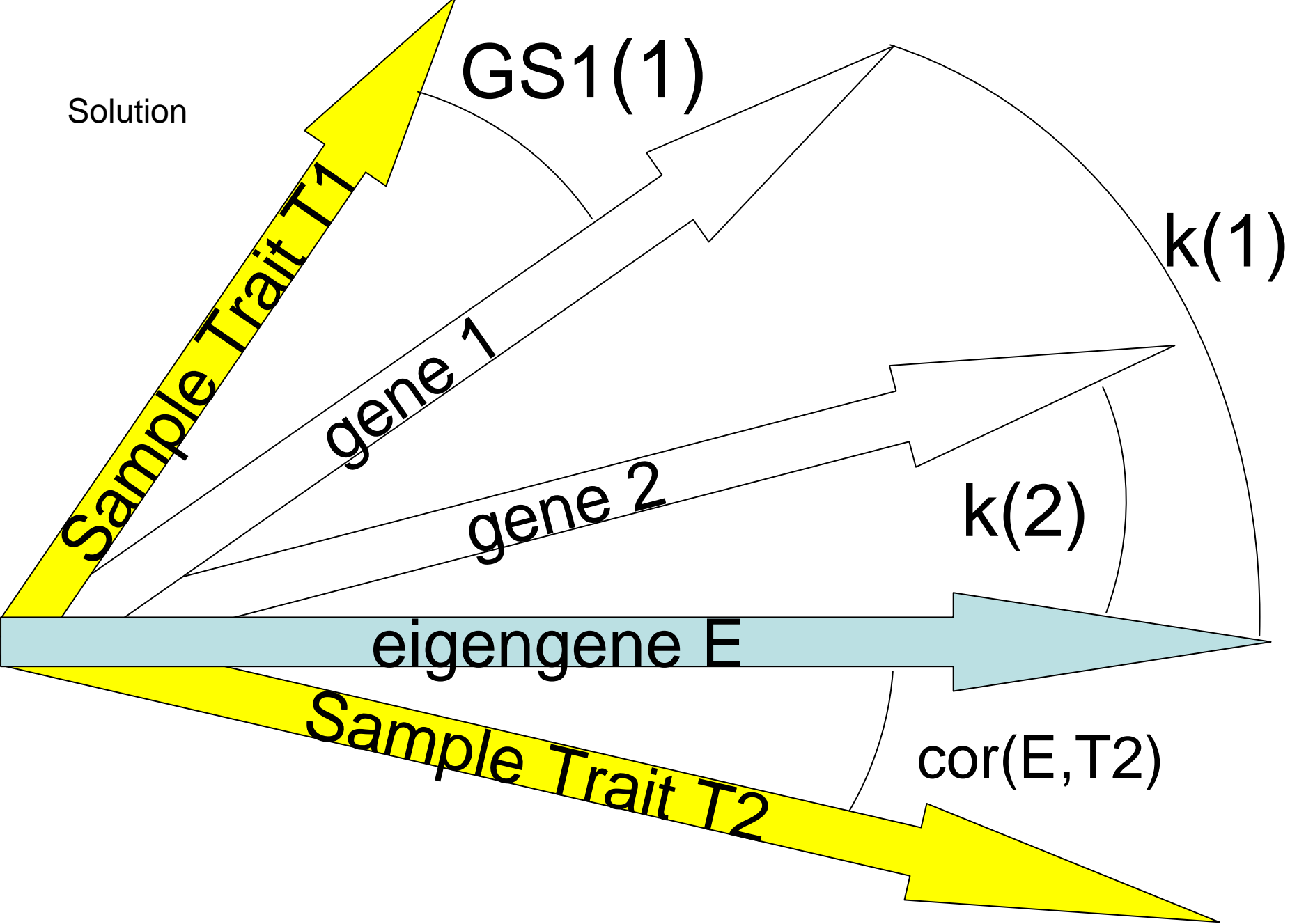


Problem

- Setting: a co-expression network and a trait based gene significance measure
 $GS(i) = |\text{cor}(x(i), T)|$
- Describe a situation when the sample trait (T1) leads to a trait-based gene significance measure with low hub gene significance
- Describe a situation when the sample trait (T2) leads to a trait-based gene significance measure with high hub gene significance

Another way of stating the problem: Find $T2$ and $T1$ such that



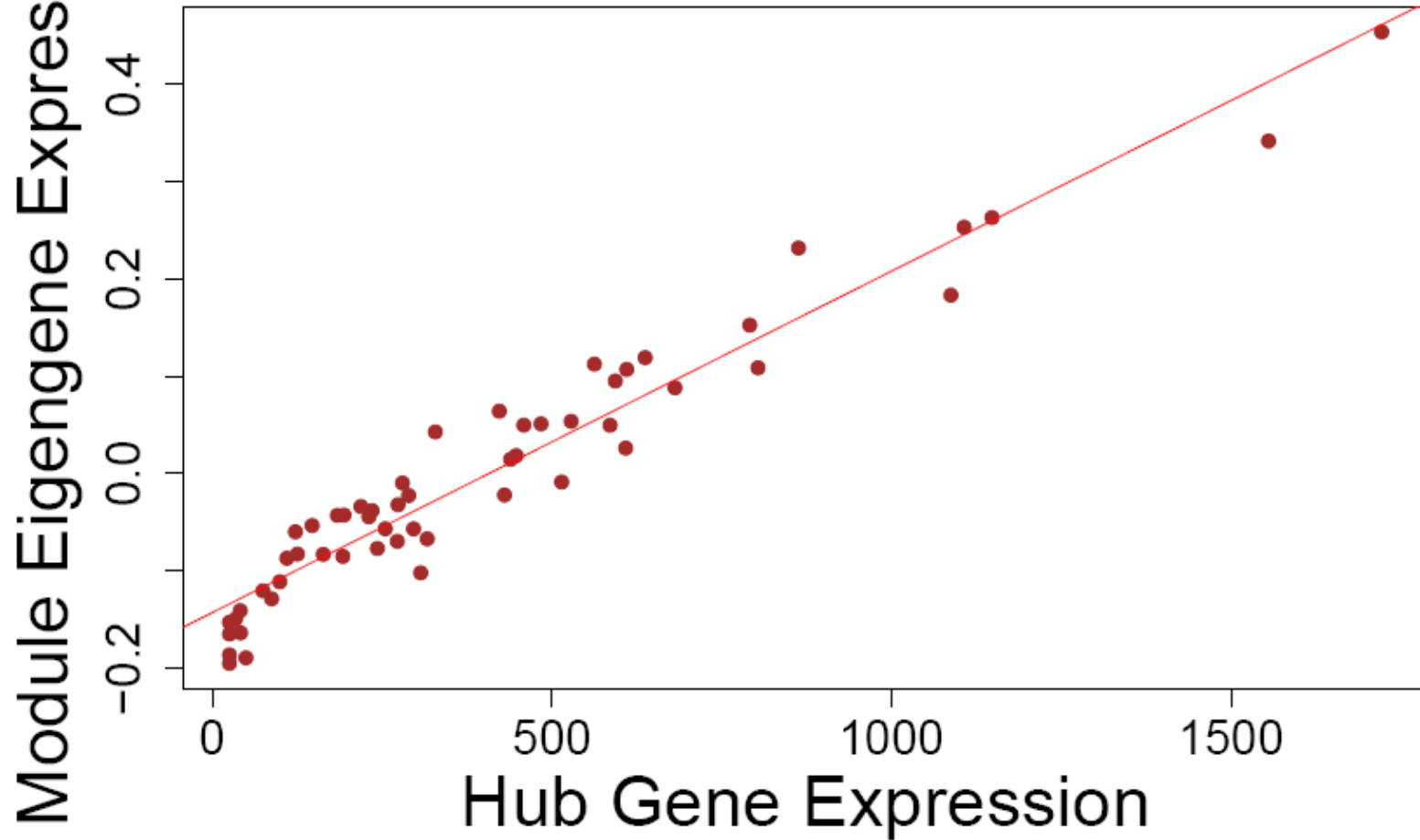


What can a microarray data analyst learn from the geometric interpretation?

Some insights

- Intramodular hub gene= a genes that is highly correlated with the module eigengene, i.e. it is a good representative of a module
- Gene screening strategies that use intramodular connectivity amount to path-way based gene screening methods
- Intramodular connectivity is a highly reproducible “fuzzy” measure of module membership.
- Network concepts are useful for describing pairwise interaction patterns.

brown, $\beta = 1$, $R^2 = 0.93$



The module eigengene is highly correlated with the most highly connected hub gene.

Dictionary for translating between general network terms and the eigengene-based counterparts.

Term	General network	Gene expression
	Adjacency matrix $A^{(q)} = [a_{ij}]$	Microarray data $X^{(q)}$
Decomposition	Factor analysis of A	Singular value decomposition of $X = UDV^T$
Centroid	Intramodular hub gene	Module eigengene E
Conformity(i)	CF_i defined as 1st factor of A	$a_{e,i} = \text{cor}(x_i, E) ^\beta$
Approximately factorizable means	$a_{ij} \approx CF_i CF_j$	$x_i \approx u_1(i) d_1 E$
Factorizability measure	$F(A) = 1 - \frac{\ (A - I) - (A_{CF} - I)\ _F^2}{\ A - I\ _F^2}$	$EF(X) = \frac{ d_1 ^4}{\sum_j d_j ^4}$
CentroidSignif(i)	$GS_{i,centroid}$	$a_{e,t} = \text{cor}(E, T) ^\beta$
CentroidConformity(i)	$a_{i,centroid,i}$	$a_{e,i} = \text{cor}(E, x_i) ^\beta$

Weighted gene coexpression network and its eigengene-based approximation if $EF(X^{(q)}) \approx 1$

	Coexpression network	Eigengene-based counterpart
Network	$A = \text{cor}(X) ^\beta$	$A_E = a_e a_e^T$
Gene significance(i)	$GS_i = \text{cor}(x_i, T) ^\beta$	$GS_{E,i} = a_{e,i} a_{e,t}$
Connectivity(i)	$k_i = \sum_{j \neq i} a_{ij}$	$k_{E,i} = a_{e,i} \sum_j a_{e,j}$

If also $\max_j (a_{e,j}) \approx 1$

Network concepts based on a network concept function $NCF(\cdot)$ if $EF(X^{(q)}) \approx 1$ and $\max_j (a_{e,j}) \approx 1$

Concept	intramodular $NCF(A, GS)$	eigengene-based $NCF(A_E, GS_E)$
Scaled Connectivity(i)	$K_i = \frac{k_i}{k_{max}}$	$K_{E,i} \approx a_{e,i}$
Density	$\frac{\sum_i \sum_{j \neq i} a_{ij}}{n(n-1)}$	$\left(\frac{S_1(a_e)}{n}\right)^2$
Centralization	$\frac{n}{n-2} \left(\frac{k_{max}}{n-1} - Density\right)$	$\sqrt{Density_E}(1 - \sqrt{Density_E})$
Heterogeneity	$\frac{\sqrt{variance(k)}}{mean(k)}$	$\frac{\sqrt{variance(a_e)}}{mean(a_e)}$
Clustering Coefficient(i)	$\frac{\sum_{l \neq i} \sum_{m \neq i, l} a_{il} a_{lm} a_{mi}}{(\sum_{l \neq i} a_{il})^2 - \sum_{l \neq i} a_{il}^2}$	$\left(\frac{S_2(a_e)}{S_1(a_e)}\right)^2$
Max. Adjacency Ratio(i)	$\frac{\sum_{j \neq i} a_{ij}}{\sum_{j \neq i} a_{ij}}$	$a_{e,i} \frac{S_2(a_e)}{S_1(a_e)}$
Hub Gene Significance	$\frac{\sum_i GS_i K_i}{\sum_i (K_i)^2}$	$a_{e,t}$
Module Significance	$\frac{\sum_{i \in I(\cdot)} GS_i}{ I(\cdot) }$	$\sqrt{Density_E} \times a_{e,t}$

Summary

- The unification of co-expression network methods with traditional data mining methods can inform the application and development of systems biologic methods.
- We study network concepts in special types of networks, which we refer to as approximately factorizable networks.
- We find that modules often are approximately factorizable
- We characterize co-expression modules that are approximately factorizable
- We provide a dictionary for relating fundamental network concepts to eigengene based concepts
- We characterize coexpression networks where hub genes are significant with respect to a microarray sample trait
- We show that intramodular connectivity can be interpreted as a fuzzy measure of module membership.

Summary Cont'd

- We provide a geometric interpretation of important network concepts (e.g. hub gene significance, module significance)
- These theoretical results have important applications for describing pathways of interacting genes
- They also inform novel module detection procedures and gene selection procedures.

Acknowledgement

Biostatistics/Bioinformatics

- Tova Fuller
- Peter Langfelder
- Ai Li
- Wen Lin
- Mike Mason
- Angela Presson
- Lin Wang
- Andy Yip
- Wei Zhao

Brain Cancer/Yeast

- Paul Mischel
- Stan Nelson
- Marc Carlson

Comparison Human-Chimp

Dan Geschwind
Mike Oldham
Giovanni

Mouse Data

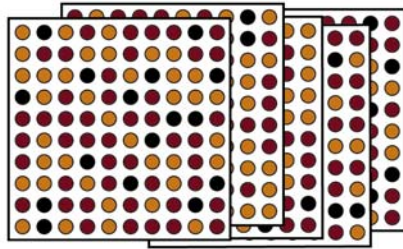
Jake Lusis
Tom Drake
Anatole Ghazalpour
Atila Van Nas

APPENDIX

(back up slides)

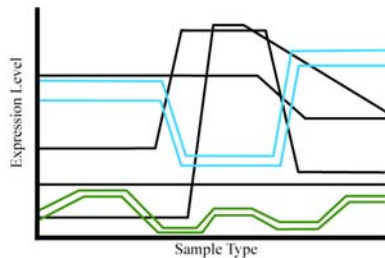
figure 1

A Array Data



Data contains correlations

B Correlation Analysis



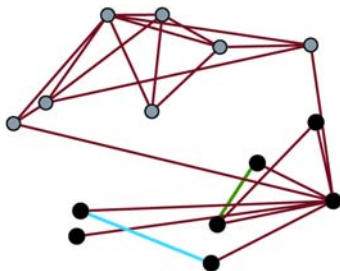
Correlation coefficients for all genes

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

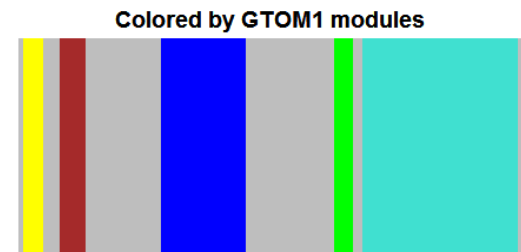
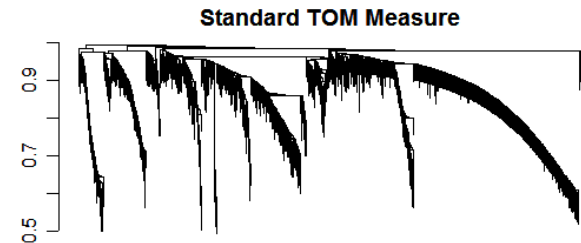


Steps for constructing a co-expression network

- Microarray gene expression data
- Measure concordance of gene expression with a Pearson correlation
- The Pearson correlation matrix is either dichotomized to arrive at an adjacency matrix \rightarrow unweighted network
Or transformed continuously with the power adjacency function \rightarrow weighted network

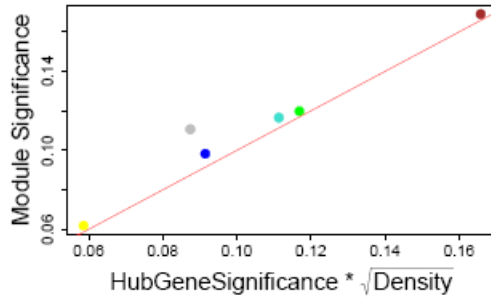
Definition of module (cluster)

- **Module=cluster of highly connected nodes**
 - Any clustering method that results in such sets is suitable
- We define modules as branches of a hierarchical clustering tree using the topological overlap matrix



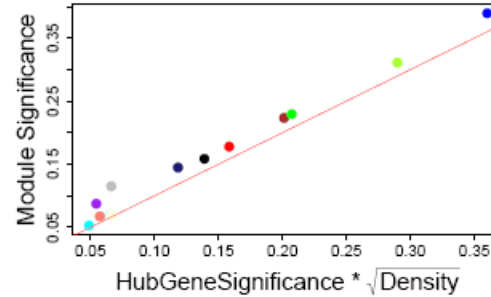
Relationship between Module significance and hub gene significance

$\beta = 1, R^2 = 1$



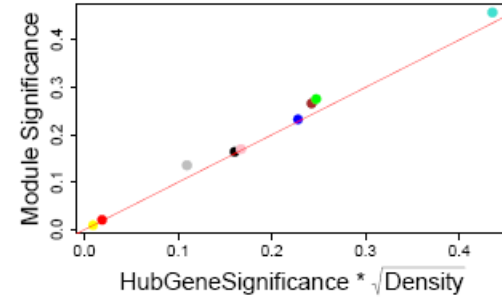
(a)

$\beta = 1, R^2 = 0.99$



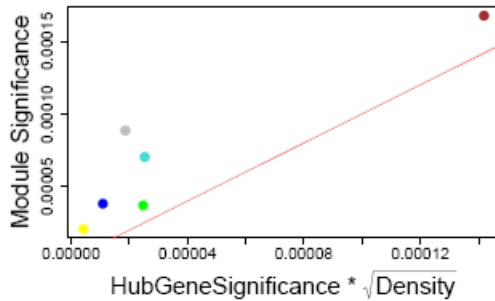
(b)

$\beta = 1, R^2 = 1$



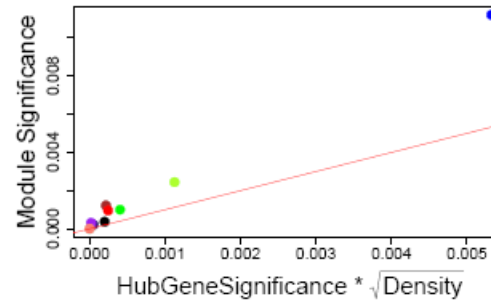
(c)

$\beta = 6, R^2 = 0.95$



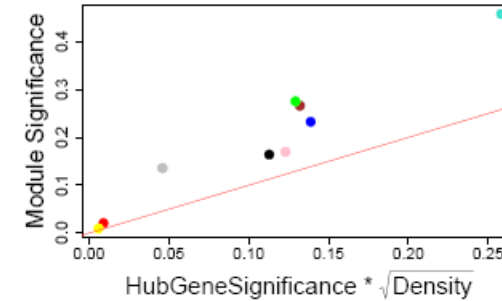
(d)

$\beta = 6, R^2 = 0.99$



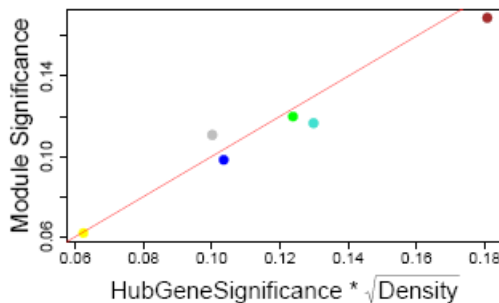
(e)

$\beta = 6, R^2 = 0.95$

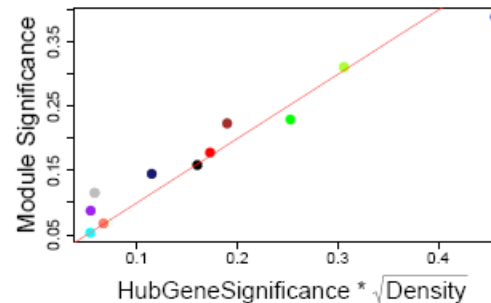


(f)

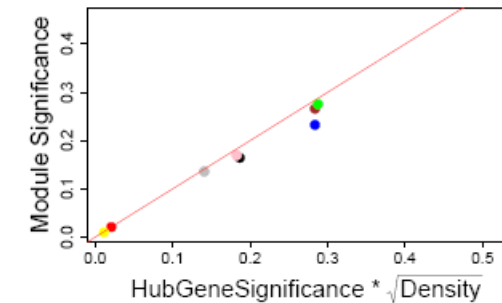
$\tau = 0.5, R^2 = 0.99$



$\tau = 0.5, R^2 = 0.96$



$\tau = 0.5, R^2 = 0.99$



Application: Brain Cancer Data

Module	blue	brown	green	grey	turquoise	yellow
Size ($n^{(q)}$)	606	185	136	1418	1112	143
Eigengene Factorizability ($EF(X^{(q)})$)	0.974	0.986	0.991	0.695	0.978	0.990
$VarExplained(\mathbf{E}^{(q)})$	0.591	0.655	0.702	0.288	0.570	0.710
$max(a_{e,i})$	0.971	0.966	0.975	0.910	0.979	0.976
<i>Density</i>	0.580	0.647	0.692	0.295	0.554	0.699
$Density_E$	0.581	0.652	0.699	0.242	0.554	0.706
<i>Centralization</i>	0.161	0.130	0.121	0.153	0.174	0.119
$Centralization_E$	0.160	0.132	0.121	0.206	0.175	0.119
<i>Heterogeneity</i>	0.137	0.100	0.105	0.189	0.170	0.110
$Heterogeneity_E$	0.138	0.101	0.105	0.433	0.171	0.111
$Mean(ClusterCoef)$	0.603	0.660	0.707	0.329	0.587	0.716
$ClusterCoef_E$	0.603	0.662	0.710	0.342	0.587	0.718
<i>ModuleSignif</i>	0.0983	0.169	0.120	0.111	0.117	0.0618
$ModuleSignif_E$	0.0554	0.163	0.105	0.0512	0.0871	0.0123
<i>HubGeneSignif</i>	0.120	0.206	0.141	0.161	0.150	0.070
$HubGeneSignif_E$	0.0707	0.196	0.123	0.0947	0.1145	0.0143
$EigengeneSignif = a_{e,t}^{(q)}$	0.0728	0.202	0.126	0.104	0.117	0.0147