

Integration of Co-expression Networks for Gene Clustering

Malay Bhattacharyya

&

Sanghamitra Bandyopadhyay

Machine Intelligence Unit,

Indian Statistical Institute, Kolkata

Year – 2009

February 6, 2009

ICAPR 2009

1

Outline of the presentation

- Motivation
- Related works
- Contribution
- Results

Motivation

Microarray experimentation

- Microarray profiling is prone to noise
- Multi-experimental data integration
- Avoiding missing value estimation

Related works

Probabilistic integration

- Log-likelihood score [Lee *et al.*, 2004]

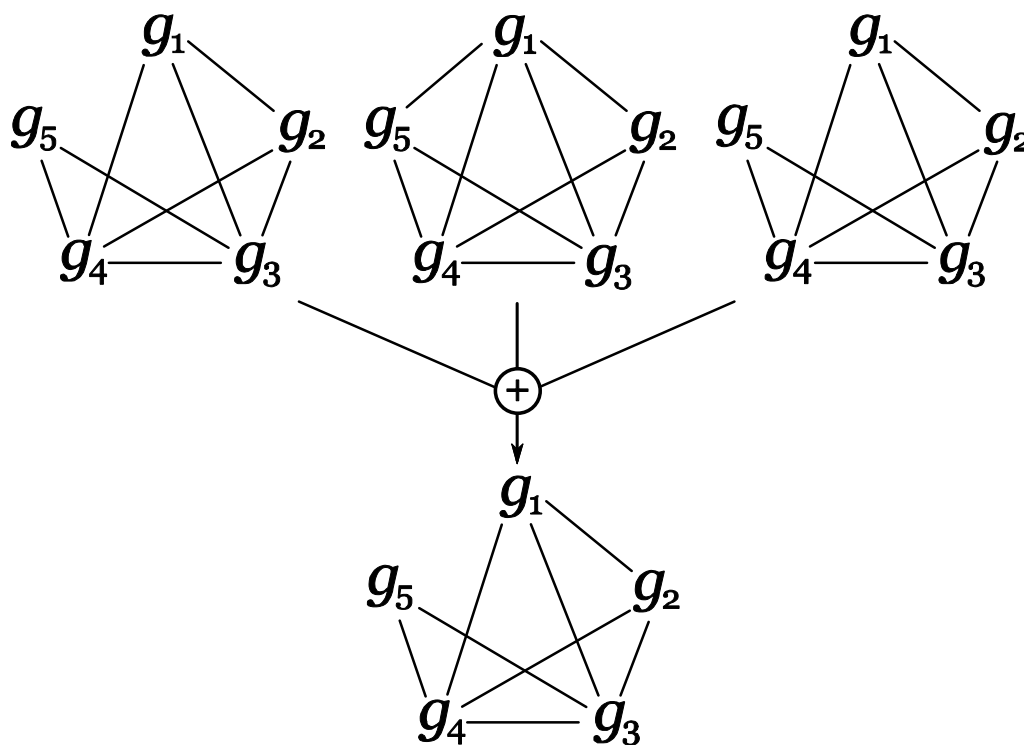
$$LLS = \ln \left(\frac{P(L | E) / \sim P(L | E)}{P(L) / \sim P(L)} \right)$$

- Depends on the ground truth knowledge

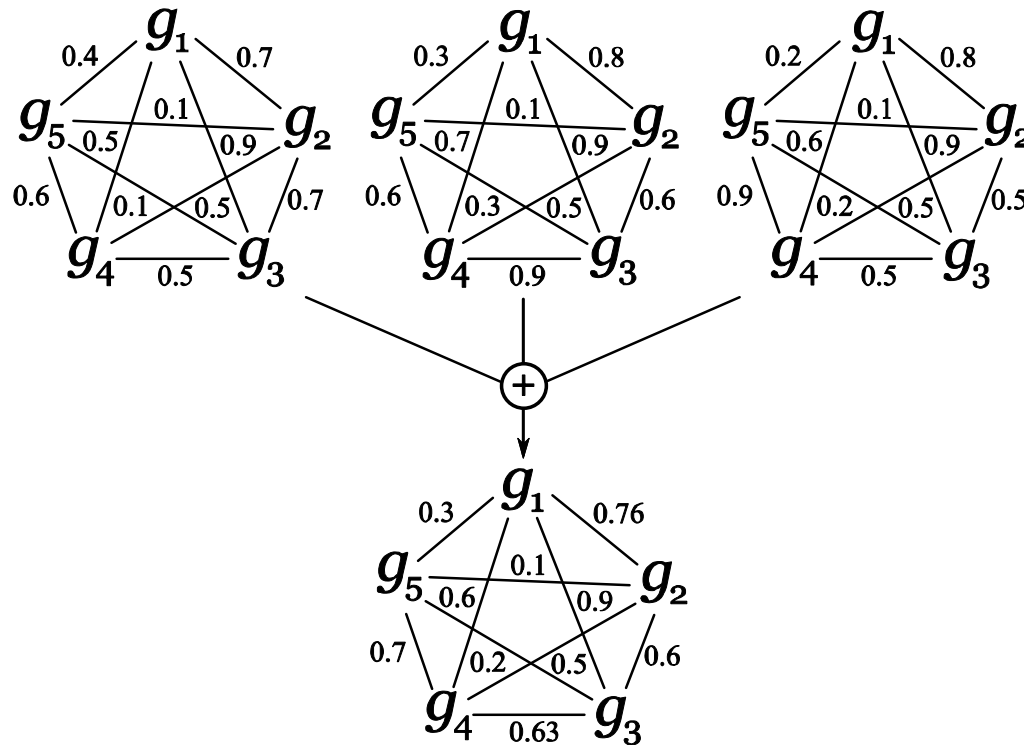
Summary graph approach

- On unweighted graphs [Hu *et al.*, 2005]
 - Frequency count of the edges
- On weighted graphs [Yan *et al.*, 2007]
 - Cutoff threshold on the aggregated weights

Consensus unweighted graph



Consensus weighted graph



Cross platform normalization

- On gene expression studies [Shablin *et al.*, 2008]
 - Merging expression studies

Contribution

Co-expression networks

$$N' = (N, A, W)$$

$$N = \{n_1, n_2, \dots, n_{|N|}\}$$

$$A \subseteq N \times N - \bigcup_{i=1}^{|N|} (n_i, n_i)$$

$$W : A \rightarrow [0,1]$$

Gene co-expression network similarity

$$N'_1 = (N, A, W_1) \quad N'_2 = (N, A, W_2)$$

$$S(N'_1, N'_2) = \frac{1}{|A|} \sum_{\forall i \in N \forall j \in N, i \neq j} 1 - |W_1(i, j) - W_2(i, j)|$$

Characteristics of the similarity measure

$$S(N'_1, N'_2) \in [0,1]$$

$$S(N', N') \in 1$$

$$S(N'_1, N'_2) = S(N'_2, N'_1)$$

Consensus gene co-expression network

A consensus gene co-expression network, $N'_c = (N, A, W_c)$, of a set of n networks $\{N'_1 = (N, A, W_1), N'_2 = (N, A, W_2), \dots, N'_n = (N, A, W_n)\}$, is defined to be a network having the maximum similarity with the given set of n networks, i.e.,

$$\prod_{i=1}^n S(N'_i, N'_c)$$

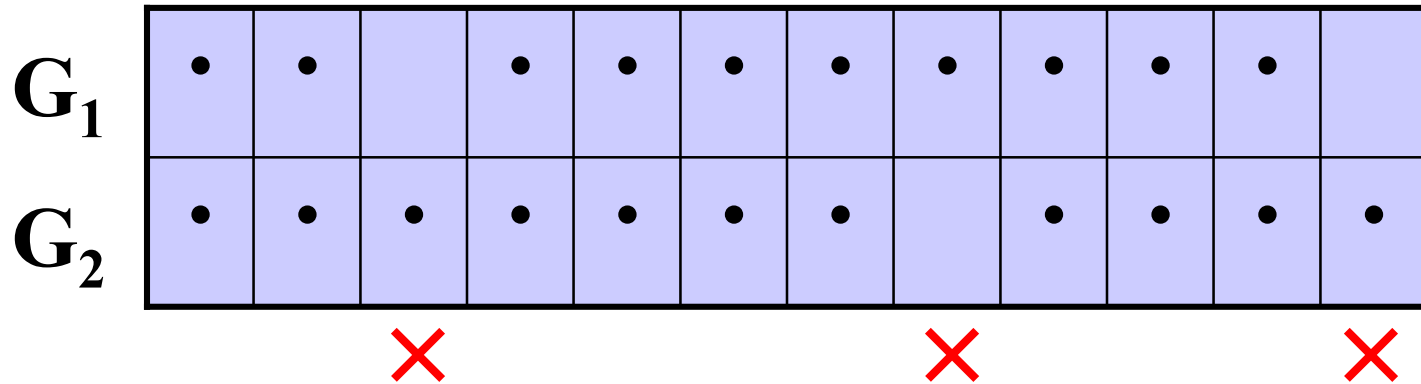
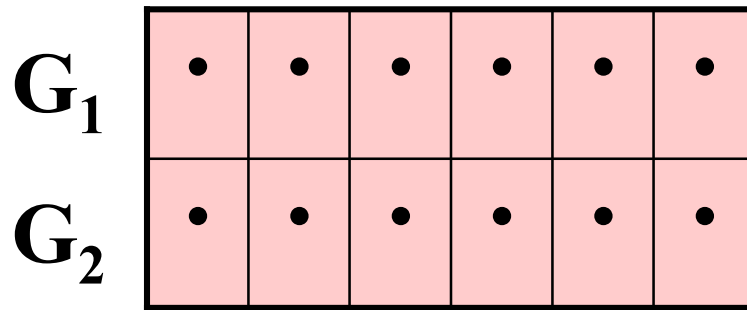
becomes maximum.

Normalized fuzzy integration (NFA)

$$W_c(i, j) = \sqrt[\alpha]{\sum_{k=1}^n \xi_k(i, j) W_k(i, j)^\alpha}$$

$$\xi_k(i, j) = \frac{\# \text{Condition}_k}{\sum_{k=1}^n \# \text{Condition}_k}$$

Avoiding missing value estimation



Results

Comparative results on the Eisen dataset [Eisen *et al.*, 1998]

Integration Method	#Clusters	<i>SI</i>	z-score
NIL	10	0.035	4.56
AVERAGE	10	-0.057	8.56
NFA ($\alpha = 2$)	10	-0.023	11

Comparative results on the Gasch dataset [Gasch *et al.*, 2001]

Integration Method	#Clusters	<i>SI</i>	z-score
NIL	10	-0.154	0.626
AVERAGE	10	-0.201	1.2
NFA ($\alpha = 2$)	10	-0.19	1.62

Major references

1. I. Lee *et al.* A probabilistic functional network of yeast genes. *Science*, 306:1555–1558, 2004.
2. H. Hu *et al.* Mining coherent dense subgraphs across massive biological networks for functional discovery. *Bioinformatics*, 21:i213–i221, 2005.
3. X. Yan *et al.* A graph-based approach to systematically reconstruct human transcriptional regulatory modules. *Bioinformatics*, 23:i577–i586, 2007.
4. A. Shabalin *et al.* Merging two gene-expression studies via cross-platform normalization. *Bioinformatics*, 24(9):1154–1160, 2008.
5. F. D. Gibbons and F. P. Roth. Judging the quality of gene expression-based clustering methods using gene annotation. *Genome Research*, 12:1574–1581, 2002.
6. M. B. Eisen *et al.* Cluster analysis and display of genome-wide expression patterns. *Proceedings of National Academy of Sciences*, 95:14863–14868, 1998.
7. A. P. Gasch *et al.* Genomic expression responses to dna-damaging agents and the regulatory role of the yeast atr homolog mec1p. *Molecular Biology of the Cell*, 12:2987–3003, 2001.

Thank you