

# Determining differentially expressed genes in a microarray expression dataset based on the global connectivity structure of pathway information

Tae Su Chung<sup>1,2</sup>, Keewon Kim<sup>1</sup>, Hye Won Lee<sup>1</sup>, Ju Han Kim<sup>1,2\*</sup>

<sup>1</sup> Seoul National University Biomedical Informatics (SNUBI), Seoul National University College of Medicine, Seoul 110-799, Korea

<sup>2</sup> Human Genome Research Institute, Seoul National University College of Medicine, Seoul, Korea

\*To whom correspondence should be addressed. E -mail: juhan@snu.ac.kr

---

## Abstract

Microarray expression datasets are incessantly cumulated with the aid of recent technological advances. One of the first steps for analyzing these data under various experimental conditions is determining differentially expressed genes (DEGs) in each condition. Reasonable choices of thresholds for determining differentially expressed genes are used for the next-step-analysis with suitable statistical significances. We present a model for identifying DEGs using pathway information based on the global connectivity structure. Pathway information can be regarded as a collection of biological knowledge, thus we are trying to determine the optimal threshold so that the consequential connectivity structure can be the most compatible with the existing pathway information. The significant feature of our model is that it uses established knowledge as a reference to determine the direction of analyzing microarray dataset. In the most of previous work, only intrinsic information in the microarray is used for the identifying DEGs. We hope that our proposed method could contribute to construct biologically meaningful network structure from microarray datasets.

## Introduction

Microarray technology makes it possible to measure the expressions of thousands of genes simultaneously under various experimental conditions. Identifying differentially expressed genes (or DEGs) in each condition is common first step for the DNA microarray data analysis.

Widely used methods for identifying DEGs include qualitative observation, heuristic rules and model-based probability analyses. Iyer *et al.* (1999) and DeRisi *et al.* (1997) discussed approaches choosing genes that have big fold-changes, with suitable thresholds, over their base

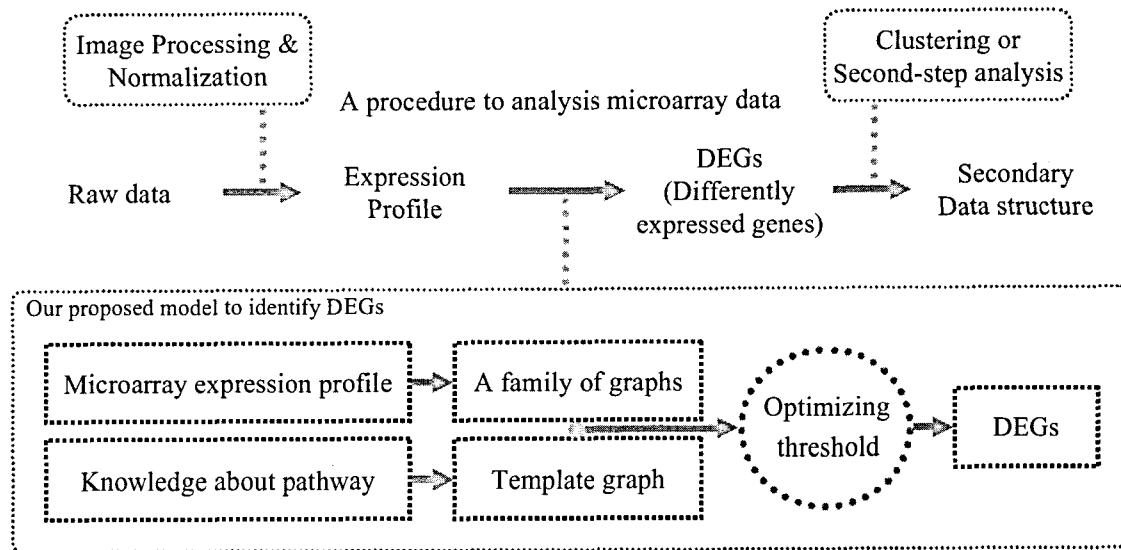


Fig. 1 Overview of our model

line expressions. Selecting meaningful thresholds are based on the heuristic rule focused on the absolute expression changes of genes.

Several probability approaches have been proposed to detect differentially expressed genes. Tusher *et al.* (2001) proposed a t-statistics to select genes that have significantly different means between conditions. They overcome the traditional problems that some genes with small differences between conditions may be selected because of their very small group -variation. But they still have normality assumption of expression measurements that often not adequate (Hunter *et al.*, 2001). Dudoit *et al.* (2002) used a non-parametric t-test and Efron and Tibshirani (2002) considered a Wilcoxon statistic and estimated the associated distributions using an empirical Bayes approach. Zhao and Pan *et al.* (2002) proposed a modified statistic which overcomes the disagreement of the null statistic and test statistic under the null hypothesis.

In this paper, we present a model for

identifying differentially expressed genes in each condition using pathway information based on the global connectivity structure. Most of previous models that share our goal used only intrinsic information in the microarray expression data. The significant feature of our model is that it uses established knowledge as a reference to determine the direction of analyzing microarray dataset. We hope that our proposed method could contribute to construct biologically meaningful network structure from microarray data sets.

## Materials and Methods

### Construction graphs from microarray data and biological pathway information

As a source of microarray expression data to analysis, we use the Rosetta compendium dataset (Hugh *et al.*, 2000), which is hitherto the most systematic approach to profile yeast genes. The dataset is consisted of 300 microarray experiment results, which contain 287 diverse gene mutations and 13 chemical treatments. They all cover 6,153

genes in each microarray data. The log-expression ratio values are used as entries of expression matrix, and these values are normalized so that mean and standard deviation of each column are 0 and 1, respectively.

In each experiment, we are trying to identify differentially expressed genes, which are usually determined by genes whose expression levels (or its absolute values) exceed some threshold. We here note that once the thresholds are determined, the graph structure on the whole genes is naturally introduced by linking co-differentially expressed genes. Here by co-differently expressed genes we mean that they are DEGs under same experimental condition. In this paper, we find optimal thresholds in the sense that the resulting secondary graph structure is most similar to the graph constructed from pathway knowledge.

The source of pathway knowledge is KEGG (Kyoto Encyclopedia of genes and genomes) database (Kanehisa, 1996), which provides 88 biological pathways- 84 metabolic pathways and 4 regulatory pathways. Among these 88 pathways, we select 43 pathways that include 12 or more genes to avoid the perturbation caused by scarcity of basis knowledge. KEGG database presents a pathway with participating genes and relations between them. The relations are divided into three categories: EC relation, PP relations and GE relation. EC relation stands for relations between two genes whose protein products share the same metabolite in a metabolic pathway. When two proteins interact directly, genes coding them are said to have PP relation. GE relation means that one gene or its product regulates the expression of the other gene. In addition, we define co-member

relation that refers to the genes assigned to the same function in a pathway.

In constructing 43 pathway graphs from these information, we make a node for each gene and link a pair of nodes when they are assigned one of the relations listed above. Merging these 43 pathways graphs, we build the single pathway graph of 570 genes as nodes, which will be used as a template to determine DEGs.

### The global compatibility between two graphs

Here we introduce the notion of compatibility between two graphs, in the general context of graph theory. The geodesic distance  $d(g, h; G)$  between nodes  $g$  and  $h$  is defined by the length of shortest path from nodes  $g$  to  $h$  in the graph  $G$ . This distance represents the global structure of graph. (Chatrand, *et al.*, 1998) If two graphs  $G_1$  and  $G_2$  are constructed on the same set of nodes, then the geodesic distance of two graphs can be easily extended by the average of differences of all geodesic distances of all pairs of nodes in each graph, i.e.

$$dist(G_1, G_2) = \frac{\sum |d(g, h; G_1) - d(g, h; G_2)|}{n(n-1)/2},$$

where the summation is taken over all (ordered) pairs  $(g, h)$  of nodes and  $n$  is the number of common set of nodes. We here note that it is symmetric and satisfies triangle inequality. (Chatrand, *et al.*, 1998)

Let  $G_P$  be the pathway graph obtained from KEGG pathway database. And let  $G_M = G_M(\theta)$  be a graph constructed from the Rosetta compendium dataset by linking co-differentially expressed genes. The threshold  $\theta$  is used to

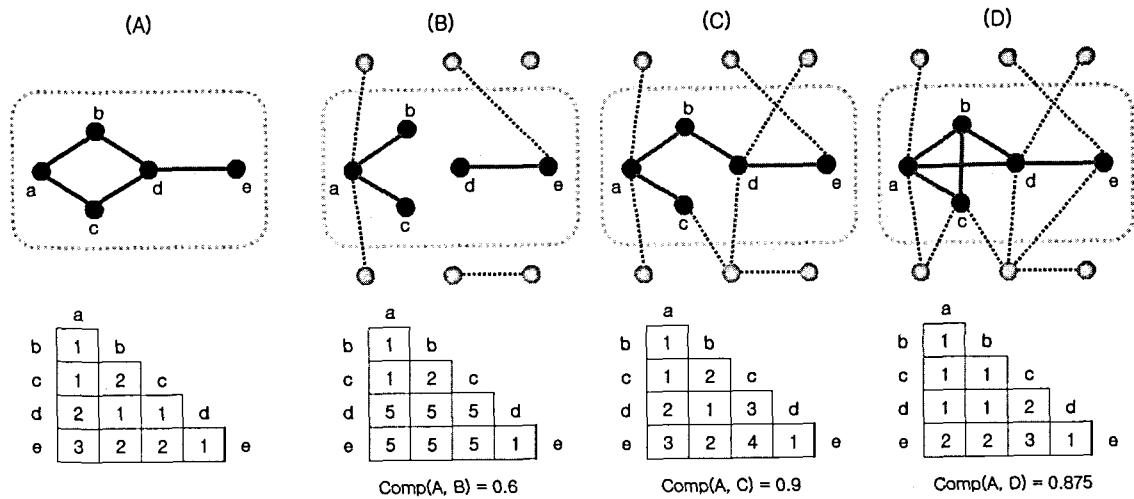


Fig. 2 Compatibilities between pathway graph and various graphs from microarray data: (A) is the pathway graph and (B)-(D) are graphs from microarray data with various thresholds. Triangular tables represent geodesic distances in graph (A) and subgraph of (B)-(D), respectively.

determine DEGs in each condition, i.e., a gene pair  $(g, h)$  is linked in the graph  $G_M(\theta)$  if the absolute value of normalized log-ratios of expressions in  $g$  and  $h$  are both greater than  $\theta$ . Then the compatibility  $Comp(G_P, G_M)$  is obtained by

$$Comp(G_P, G_M) = 1 - \text{dist}(G_P, G_M | G_P) / (n - 1).$$

Here  $n$  is the number of genes in  $G_P$ , and  $G_M | G_P$  is the relative subgraph of  $G_M$  to  $G_P$ . Since the pathway graph  $G_P$  contains only subset of genes that are described in graph  $G_M$ , it is natural to compare  $G_P$  and subgraph of  $G_M$ . It is clear that the compatibility lies between 0 and 1 and it becomes 1 only when the graph  $G_M$  includes exactly same structure of  $G_P$ .

The effect of the optimization on compatibility is illustrated in Fig. 2. In the figure, (B), (C) and (D) are possible graphs from microarray dataset by taking different threshold determining DEGs. And (A) is the pathway graph which will be used as a reference to select one

graph from (B), (C) and (D). By calculating compatibilities, which represents the global similarity of graph structures, we can take the threshold used in constructing graph (B) as the optimal one.

## Results

### Correlation of microarray data and pathway knowledge

We try to determine DEGs in each experimental condition of microarray expression data based on the template structure of pathway knowledge. The underlying assumption of our approach is that pathway knowledge can be conjectured from microarray data, or that microarray data reflect the biological pathway knowledge. We investigate the validity of this assumption by investigating the correlation between co-degree in microarray expression data and relations within pathways.

The co-degree of gene pair  $(g, h)$  with parameter  $\theta$  is defined by the number of

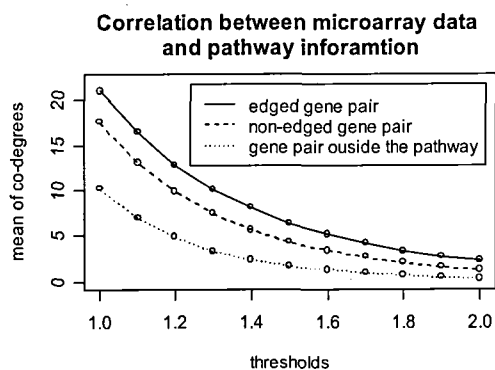


Fig. 3 Correlation between microarray data and pathway information:

conditions satisfying that both  $g$  and  $h$  become DEGs simultaneously under the condition with respect to the threshold  $\theta$ . Based on the pathway information, however, we can categorize gene pairs into three classes: The first class comprises gene pairs in which two genes belong to the pathway and are connected by a specific biological relation. The second class comprises gene pairs whose genes belong to the pathway but are not related by any biological relation. The third one is comprises gene pairs whose genes are outside the pathway. Among these three classes, we calculate the average of co-degrees as a correlation of microarray data and pathway knowledge.

We naturally expected that the average co-degrees on gene pairs in the first class is bigger than that of the second class, and that of the third class places smallest value. Fig. 3 shows the result which agrees with our expectation and so we validate that there is a correlation between the relational information of pathway information and microarray data.

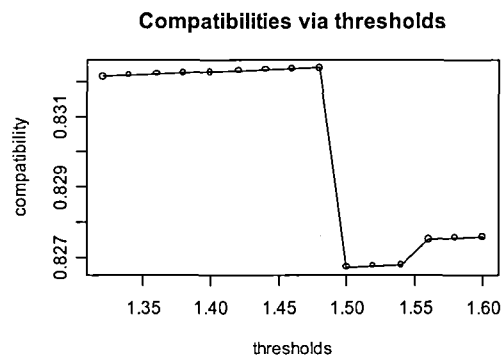


Fig. 4 Compatibilities via thresholds:

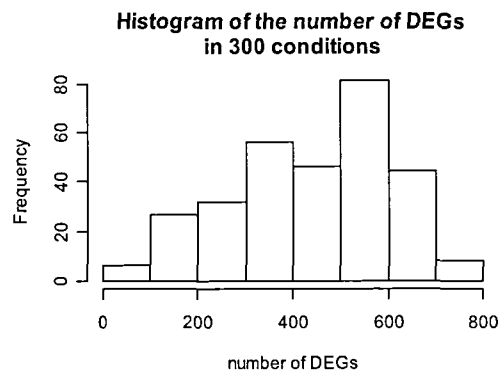


Fig. 5 Histogram of the number of DEGs in 300 conditions of the rosetta compendium dataset.

### Differentially expressed genes in the Rosetta compendium dataset

Applying our model to rosetta compendium dataset, we find the optimal threshold determining DEGs in whole 300 conditions. The compatibility plot via various thresholds is provided in Fig. 4. In the figure, the shape of curve is nearly unimodal, so there is no doubt about selecting threshold in the peak as the optimal one. Actually, we just adapt greed algorithm to find optimal threshold 1.4926 under the restriction that the threshold is independent from conditions for the

simplicity of algorithm.

Fig. 5 shows the distribution of the number of DEGs in 300 conditions applying the optimal threshold. Among whole 6,153 yeast genes, we see that our model select about 10% of genes as DEGs in each condition. This result agrees with the common criterion which can be found in many biological literatures.

## Discussion

In this paper, we suggest a novel model to identify differentially expressed genes in each conditions of a microarray dataset. The procedure is one of the first steps of analyzing expression profiles. The significant feature of our model is that it uses established knowledge as a reference to determine the direction of analyzing microarray dataset and it does not consider the individual structure but consider the global network structure to determine DEGs in each single condition. And because of using information outside array, it does not need any assumptions on the distributions of expression profiles. We also investigate the validity of our basic assumption that relational information in existing pathway knowledge reflects the expression level, or moreover network structure of microarray data.

Actually, knowledge about pathways covers only tiny portion of genes that exist. In the case of yeast, less than 10% of whole genes are revealed to participate in some known pathway among 6,000 genes as a whole. And graph structure is too simple to use as a template structure to compare the pathway and microarray data, because it ignore the orders and types of biological relations.

These incompleteness and simplicity may give insignificant results of our model. But applying other biological information, such as protein interaction data and relational information of genes or proteins from biological literature, and modifying graph structure into more complicated structure, we can overcome the problem of our current model.

## Acknowledgements

This study was supported by a grant from Korea Health 21 R&D Project, Ministry of Health & Welfare, Republic of Korea (0405-BC02-0604-0004).

## References

- [1] Chartrand, G., Kubicki, G. and Schultz, M. (1998), Graph similarity and distance in graphs, *Aequationes Math.* **55**, 129-145.
- [2] DeRisi, J. L., Iyer, V. R. and Brown, P. O. (1997) Exploring the metabolic and genetic control of gene expression as a genomic scale, *Science*, **278**, 680-686.
- [3] Dudoit, S., Yang, Y. H., Gallow, M. J. and Speed, T. P. (2002) Statistical methods for identifying differentially expressed genes in replicated cDNA microarray experiments, *Statistica Sinica*, **12**, 111-139.
- [4] Efron, B. and Tibshirani, R. (2002) Empirical Bayes methods and false discovery rates for microarrays, *Genetic Epidemiology*, **23**, 70-86.
- [5] Farkas, I., Jeong, H., Vicsek, T., A.-L. Barabási, A.-L. and Oltvai, Z.N. (2003) The topology of the transcription regulatory network in the yeast, *Saccharomyces*

- cerevisiae*, *Physica A* **318**, 601-612.
- [6] Holme, P., Huss, M. and Jeong, H. (2003), Subnetwork hierarchies of biochemical pathways, *Bioinformatics*, **19**(4), 532-538.
- [7] Hughes, T. R., Marton, M. J., Jones, A. R., Roberts, C. J., Stoughton, R., Armour, C. D., Bennett, H. A., Coffey, E., Dai, H., He, Y. D., Kidd, M. J., King, A. M., Meyer, M. R., Slade, D., Lum, P. Y., Stepaniants, S. B., Shoemaker, D. D., Gachotte, D., Chakraburty, K., Simon, J., Bard, M. and Friend, S. H. (2000) Functional Discovery via a Compendium of expression Profiles, *Cell*, **102**, 109-126.
- [8] Hunter, L., Taylor, R. C., Leach, S. M. and Simon, R. (2001) GEST: a gene expression search tool based on a novel Bayesian similarity metric. *Bioinformatics*, **17**(Suppl. 1), s115-s122.
- [9] Iyer, V. R., Eisen M. B., Ross D. T., Schuler G., Moore T., Lee J. C., Trent J. M., Staudt L. M., Hudson J. Jr, Boguski M. S., Lashkari D., Shalon D., Botstein D. and Brown P. O. (1999) The transcriptional program in the response of human fibroblasts to serum, *Science*, **283**, 83-87.
- [10] Kanehisa, M. (1996) Toward pathway engineering: a new database of genetic and molecular pathways. *Science & Technology Japan*, No. 59, pp. 34-38
- [11] Kanehisa, M. (2000) *Post-Genome Informatics*, Oxford University Press.
- [12] Newton, M. A., Kendzioriski, C. M., Richmond, C. S., Blattner, F. R. and Tsui, K. W. (2001) On differential variability of expression ratios: improving statistical inference about gene expression change from microarray data. *Journal of Computational Biology*, **8**, 37-52.
- [13] Qian, J., Lin, J., Luscombe, N.M., Yu, H. and Gerstein, M. (2003), Prediction of regulatory networks: genome-wide identification of transcription factor targets from gene expression data, *Bioinformatics*, **19**(15), 1917-1926.
- [14] Rung, J., Schlitt, T., Brazma, A., Freivalds, K. and Vilo, J. (2002) Building and analysing genome-wide gene disruption networks, *Bioinformatics*, **18**(Suppl. 2). s202-s210.
- [15] Shmulevich, I. and Zhang, W. (2002), Binary analysis and optimization-based normalization of gene expression data, *Bioinformatics*, **18**(4), 555-565.
- [16] Tusher, V. G., Tibshirani, R. and Chu, G. (2001) Significance analysis of microarrays applied to the ionizing radiation response, *Proc. Natl. Acad. Sci. USA*, **98**, 5116- 5121.
- [17] Wagner, A. (2001) How to reconstruct a large genetic network from n gene perturbations in fewer than n<sup>2</sup> easy steps, *Bioinformatics*, **17**(12), 1183-1197.
- [18] Zhao, Y. and Pan, W. (2003) Modified nonparametric approached to detecting differentially expressed genes in replicated microarray experiments, *Bioinformatics*, **19**, 1046-1054.