# GENE CO-EXPRESSION NETWORK DISCOVERY WITH CONTROLLED STATISTICAL AND BIOLOGICAL SIGNIFICANCE

Dongxiao Zhu<sup>a,b</sup> and Alfred O Hero<sup>b</sup>

<sup>a</sup>Bioinformatics Program, <sup>b</sup>Depatments of EECS, Biomedical Engineering and Statistics University of Michigan, Ann Arbor, MI 48105

### ABSTRACT

Many biological functions are executed as a module of coexpressed genes which can be conveniently viewed as a coexpression network. Genes are network vertices and significant pairwise co-expressions are network edges. Traditional network discovery methods controls either statistical significance or biological significance, but not both. We have designed and implemented a two-stage algorithm that controls both statistical significance (False Discovery Rate, FDR) and biological significance (Minimium Acceptable Strength, MAS) of the discovered network. Based on the estimation of pairwise gene profile correlation, the algorithm provides an initial network discovery that controls only FDR, which is then followed by a second network discovery which controls both FDR and MAS. We illustrate the algorithm for discovery of co-expression networks for yeast galactose metabolism with controlled FDR and MAS.

## 1. INTRODUCTION

Microarray gene expression data enable researchers to interrogate gene expression levels simultaneously on the genome scale. Detection of co-expressed genes from microarray data has attracted much attention since many co-expressed genes are found to have functional relationships, e.g. lying in the same signal transduction pathway. Many coexpression detection techniques such as relevance network and hierarchical clustering rely on the quantitative or qualitive assessment of similarities between the expression profiles of gene pairs, which is one of the fundamental objectives in functional genomics and system biology. Traditional methods either screen statistically significant or biologically significant co-expressed gene pairs. The former does not control error rate, and the latter leads to screeningin some weakly correlated gene pairs that are difficult to verify by follow-up experiments such as real time RT-PCR.

In this paper, we present a two-stage algorithm that simultaneously controls statistical and biological significance of the discovered co-expression network. The algorithm implements Pearson correlation coefficients and Kendall correlation coefficients in order to capture both linear and nonlinear types of dependencies between all pairs of gene expression profiles. A two-stage error control procedure is then implemented through which a number of gene pairs are declared to be both statistically and biologically significant as measured by FDR and MAS of association. These gene pairs form the edges of the relevance network that represents the complicated web of gene co-expression among all pairs of genes.

We demonstrate the application of our two-stage algorithm by constructing relevance networks from yeast galactose metabolism data [1]. This data represents approximately 6200 gene expression levels on two-color cDNA microarrays over 20 physiological/genetic conditions (nine mutants and one wild type strains incubated in either GALinducing or non-inducing media) with four replicates in each condition.

The paper is organized into five parts: Introduction of Kendall and Pearson statistics for strength of association (Sec. 2); Formulation of the problem of network discovery as a composite hypothesis test with multiple comparisons (Sec 3); Introduction of two-stage procedure for testing these hypotheses (Sec 4); Validation of the two-stage algorithm and application to yeast data (Sec 5).

## 2. MEASURING THE STRENGTH OF ASSOCIATION

We use  $\Gamma$  to denote the true strength of association between a pair of gene expression profiles. Under a Gaussian linear hypothesis, the sample Pearson correlation coefficient  $\hat{\rho}$ is an appropriate metric. A robust distribution-free alternative is the sample Kendall rank correlation coefficient  $\hat{\tau}[2]$ . The Pearson and Kendall correlation coefficients are special cases of the generalized correlation coefficient  $\Gamma$ . We define  $\{g_p\}_{p=1}^G$  as the indices of G gene probes on the microarray;  $\{X_{g_p}\}_{p=1}^G$  as normalized probe responses (random variables); and  $\{\{x_{g_p(n)}\}_{p=1}^G\}_{n=1}^N$  as realizations of  $\{X_{g_p}\}_{p=1}^G$  under N i.i.d. microarray experiments.

Kendall's  $\tau$  statistic is a measure of correlation that captures both linear and non-linear associations[2]. The  $\tau$  is defined as:  $\tau = P_+ - P_-$ , where, for any two independent pairs of observations  $(x_{g_i(n)}, x_{g_j(n)})$ ,  $(x_{g_i(m)}, x_{g_j(m)})$ from the population:  $P_+ = P[(x_{g_i(n)} - x_{g_i(m)})(x_{g_j(n)} - x_{g_j(m)}) > 0]$  and  $P_- = P[(x_{g_i(n)} - x_{g_i(m)})(x_{g_j(n)} - x_{g_j(m)}) < 0]$ . An unbiased estimator of  $\tau$  is given by the Kendall  $\tau$  statistic:  $\hat{\tau} = 2 \sum \sum_{1 \le i \le j \le N} \frac{K_{ij}}{N(N-1)}$ . Here  $K_{ij}$  is a indicator variable defined as  $K_{ij} = \operatorname{sgn}(x_{g_i(n)} - x_{g_i(m)})\operatorname{sgn}(x_{g_j(n)} - x_{g_j(m)})$  for each set of pairs  $\{X_{g_i}\}_{i=1}^G, \{X_{g_j}\}_{j=1}^G$  of sample observations.

To make the estimated correlation robust against spurious outliers yet sensitive to strong similarities in expression patterns, we adopted a leave-one-out cross-validation technique, using the median estimate as a robust estimator of the correlation.

### 3. HYPOTHESIS TESTING SCHEME

For G genes on each microarray, we need to simultaneously test  $\mathcal{G} = \begin{pmatrix} G \\ 2 \end{pmatrix}$  pairs of two-sided hypotheses:

$$H_0: \Gamma_{g_i,g_j} \le cormin \text{ versus } H_\alpha: \Gamma_{g_i,g_j} > cormin,$$
  
for  $g_i \ne g_j$ , and  $g_i, g_j \in (1, 2, ...G)$  (1)

where *cormin* is a minimium acceptable strength of correlation. The sample correlation coefficient  $\hat{\Gamma}$  ( $\hat{\rho}$  or  $\hat{\tau}$ ) is used as a decision statistic to decide on pairwise dependency of two genes in the sample. For N realizations of any pair of gene probe responses,  $(x_{gi(n)}, x_{gj(n)})$ , we first calculate  $\hat{\tau}$ or  $\hat{\rho}$ . For large N, the Per Comparison Error Rate (PCER) p-values for  $\rho$  or  $\tau$  are:

$$p_{\rho} = 2\left(1 - \Phi\left(\frac{\tanh^{-1}(\hat{\rho})}{(N-3)^{-1/2}}\right)\right)$$
$$p_{\tau} = 2\left(1 - \Phi\left(\frac{K}{N(N-1)(2N+5)/18^{1/2}}\right)\right)$$

where  $\Phi$  is the standard Gaussian cumulative density function, and  $K = \sum \sum_{1 \le i \le j \le N} K_{ij}$ . The above expressions are based on asymptotic Gaussian approximations[2].

The PCER p-value refers to the probability of Type I error rate incurred in testing a single pair of hypothesis for a single pair of genes  $g_i, g_j$ . It is the probability that purely random effects would have caused  $g_i, g_j$  to be erroneously selected based on observing correlation between this pair of genes only. When considering the  $\mathcal{G}$  multiple hypotheses for all possible pairs, two adjusted error rates have frequently been considered in microarray studies. These are familywise error rate (FWER) and false discovery rate (FDR). The FWER is the probability that the test of all  $\mathcal{G}$  pairs of hypotheses yields at least one false positive in the set of declared positive responses. In contrast, the FDR is the average proportion of false positives in the set of declared positive responses. The FDR is dominated by the FWER and is therefore a less stringent measure of significance. As in previous studies, we adopt the FDR to control statistical significance of the selected gene pair correlations in our screening procedure[3].

### 4. TWO-STAGE SCREENING PROCEDURE

Select a level  $\alpha$  of FDR and a level *cormin* of MAS significance levels. We use a modified version of the two-stage screening procedure applied to gene screening [3]. This procedure consists of:

Stage I. Test the simple null hypothesis.

$$H_0: \Gamma_{g_i, g_i} = 0$$
 versus  $H_\alpha: \Gamma_{g_i, g_i} \neq 0$ 

at FDR level  $\alpha$ . The step-down procedure of Benjamini and Hochberg [4] is used.

Stage II. Suppose  $G_1$  pairs of genes pass the stage I procedure. In stage II, we first construct asymptotic PCER Confidence Intervals (PCER-CI's) : $I^g(\alpha)$  for each  $\Gamma$  ( $\rho$  or  $\tau$ ) in subset  $G_1$ , and convert into FDR Confidence Intervals (FDR-CI's) : $I^g(G1\alpha/\mathcal{G})$ [5]. A gene pair in subset  $G_1$  is declared to be both statistically significant and biologically significant if its FDR-CI does not intersect the MAS interval [-cormin, cormin] (see Fig 3).



**Fig. 1**. Verification of null sampling distribution (a) and variance approximation (b). (a) QQ plot of transformed sampling distribution of Pearson correlation coefficient versus normal distribution. (b) Variance approximation of transformed sampling distribution of Pearson correlation coefficient.

### 5. VALIDATION OF TWO-STAGE ALGORITHM

#### 5.1. Validating asymptotic null distribution

Here we verify that the two-stage algorithm controls FDR at a specified MAS level using simulated data. Since the pvalues are based on asymptotic distribution approximations, we verify in Fig 1 that the sampling distribution is Gaussian distributed using QQ plot. Moreover, since the construction of confidence intervals requires estimation of sampling distribution variance, the accuracy of variance approximation is vital, which can be accessed by calculating squared error:(*s.e.* denotes standard error, and  $F_X$  denotes sampling distribution)

$$\hat{\sigma}_{\rho} = (s.e.(\tanh^{-1}(F_{\hat{\rho}})) - (N-3)^{-1/2})^2$$
$$\hat{\sigma}_{\tau} = (s.e.(F_{\hat{\tau}}) - (\frac{2}{N(N-1)} \frac{2(N-2)}{N(N-1)^2} \sum_{i=1}^{N} (C_i - \overline{C}) + 1 - \hat{\tau}))^2$$

Fig1b shows that the variance of sampling distribution is close to its approximation value even for small sample size (N < 20).

#### 5.2. Validating error control procedure

In order to validate our FDR and MAS error control procedure, we simulated pairwise gene expression data based on pre-specified population covariances. The actual FDR at a MAS level is calculated as a ratio of the number of screened gene pairs whose corresponding population correlation parameters  $\Gamma$ 's are less than the MAS level specified, divided by the total number of screened gene pairs. The actual MAS is the minimium discovered population correlation  $\Gamma$  among the screened pairs. We pre-specified 16 pairs of (FDR,MAS) criteria (Four FDR levels: 0.2, 0.4, 0.6, 0.8; Four MAS levels: 0.2, 0.4, 0.6, 0.8), and each is plotted as a different point character (red) in Fig 2. The 16 corresponding pairs of actual (FDR,MAS) criteria are also shown in Fig 2 using the same set of point characters (Blue). It can be observed that the actual FDR's (blue points) fall below the prespecified constraint (red points) and the actual MAS's (blue points) fall above the pre-specified constraint (red points). The deviations of actual FDR's and MAS's from their prespecified levels are due to the conservative asymptotic approximation. This will translate into a reduction of power in discovering co-expressed pairs at the specified levels.

## 6. CONSTRUCTING A RELEVANCE NETWORK WITH CONTROLLED FDR AND MAS

Relevance networks are implemented as a graph where n nodes (genes) are connected by p sets of edges (co-express ions). Each of the p edges represents the similarity measure between pairs of nodes[6].

For the yeast galactose metabolism dataset, a subset of 997 genes were identified by Ideker et al using generalized likelihood ratio test [1]. Genes having a likelihood statistic  $\lambda \leq 45$  were selected as differentially expressed, whose mRNA levels differed significantly from reference under one or more perturbations. We used the average expression profiles over four replicates for subsequent analysis, which



**Fig. 2**. Verification of two-stage error control procedure based on Pearson correlation coefficient(a) and Kendall correlation coefficient(b). Sample size N = 20.

implicitly assumes that the between-replicates variances for a gene over different experimental conditions are equal.



**Fig. 3**. Segments of lower bounds (a) and upper bounds (b) specifying the 5% FDR-CI's on the positive Pearson correlation coefficients (a) and negative Pearson correlation coefficients (b) for the galactose metabolism study. Only those gene pairs whose FDR-CI's do not intersect [-cormin, cormin] are selected by the second stage of screening. When the MAS strength of association criterion is cormin = 0.5, these gene pairs are obtained by thresholding the curves as indicated.

Fig 3a and Fig3b illustrate the direct implementation of the two-stage procedure to screen positively or negatively correlated gene pairs based on the Pearson correlation coefficient. See [3] for more details on how to intepret these plots. The direct screening procedure is constrained by FDR criterion  $\alpha = 0.05$  and MAS criterion cormin = 0.5.

Fig 4 presents the discovered network topology with a FDR level of 0.05 (5% discovered edges are expected to be false positive) at the MAS level of 0.9 (*cormin* = 0.9). The network is composed of 91 connected vertices and 138 edges. Similiar to some other biological networks, the network marginal degrees appear power-law distributed, which is tested by verifying goodness of fit to the log-transformed



**Fig. 4**. Network topology visualization. The network is discovered by constraining  $FDR \le 5\%$  at a MAS level of 0.9. No significant negative correlation is discovered at this level. The graph is drawn using Pajek [9].

power-law model, (goodness of fit criterion  $R^2 = 0.95$ ) i.e.,  $\log P(K_i) = -\gamma \log K_i + \log \alpha + \varepsilon_i$ , (i = 1, 2, n), here  $\gamma$ and  $\alpha$  are parameters,  $\varepsilon_i$  is a residual fitting error.  $K_i$  is the degree and  $P(K_i)$  is the corresponding probability.

Genes that are of considerable interest to the biologist are the highly connected genes that dominates the network topology. These are called "hub genes", such as RPL33A and RPS4A in Fig 4 and are minimally sensitive to the network discovery criteria. Most of the "hub genes" in each discovered network fall into two categories: "RPL" and "RP S". The former encodes "Ribosome Protein Large (60S) subunit," and the latter encodes "Ribosome Protein Small (40S) subunit". Both of which are structural components of the ribosome that is responsible for protein biosynthesis. Protein biosynthesis plays the central role in galactose metabolism because galactose is not a primary carbon source for yeast, and different types of proteins including transporters, enzymes, and regulators have to be synthesized upon induction [7]. Interestingly, the list of "hub genes" contains many hypothetical Open Reading Frames (ORFs)(data not shown), which are presumably indispensable for galactose metabolism [8].

## 7. CONCLUSION

We have introduced a method to construct gene co-expression n networks with controlled FDR at different levels of MAS. By replacing correlation coefficient with partial correlation coefficient, the method can be naturally extended to the Gaussian Graphic Model framework.

**Table 1.** Top ten "hub genes". The rank of each gene is the average rank over five networks. Each of five networks is constraint by a different pair of (FDR,MAS) criteria. Highest rank is the most connected and most stable gene under varying constraints of (FDR,MAS)

Gene Name	Average Rank
RPL42B	4.2
RPS3	5.8
RPL14A	7.0
RPS16B	7.6
GTT2	8.4
RPS4A	9.8
RPL33A	11.8
RPL23B	15.8
RPS7A	16
RPL27A	17.4

#### 8. REFERENCES

[1] Ideker, T., V. Thorsson, et al. (2000) Testing for differen tially-expressed genes by maximum-likelihood analysis of microarray data. J Comput Biol 7(6): 805-17.

[2] Hollander A and Wolfe D. (1999) Nonparametric Statistical Methods. 2nd Edition Wiley-Interscience.

[3] Hero, G. Fleury, A et al. (2004) Multicriteria gene screening for analysis of differential expression with DNA microarrays, EURASIP Journal on Applied Signal Processing, 2004(1):43-52.

[4] Benjamini Y and Hochberg Y. (1995) Controlling the false discovery rate - a practical and powerful approach to multiple testing. J Roy Stat Soc B Met 57 (1): 289-300.

[5] Benjamini Y and Yekutieli D. (2004) False Discovery Rate adjusted multiple confidence intervals for selected parameters. Submitted to Journal of American Statistical Association.

[6] Butte, A. J., P. Tamayo, et al. (2000) Discovering functional relationships between RNA expression and chemotherapeutic susceptibility using relevance networks. Proc Natl Acad Sci U S A 97(22): 12182-6.

[7] Wieczorke, R., S. Krampe, et al. (1999) Concurrent knock-out of at least 20 transporter genes is required to block uptake of hexoses in Saccharomyces cerevisiae. FEBS Lett 464(3): 123-8.

[8] Jeong, H., S. Mason, A.-L. Barabsi and Z. N. Oltvai. (2001) Lethality and centrality in protein networks. Nature 411: 41-42.

[9] Batagelj A., A. Mrvar. (1998) Pajek - Program for Large Network Analysis. Connections 21(2): 47-57.