



Combining multiple microarray studies and modeling interstudy variation

Jung Kyoong Choi^{1,2}, Ungsik Yu², Sangsoo Kim² and Ook Joon Yoo^{1,*}

¹Department of Biological Sciences, Korea Advanced Institute of Science and Technology, 371-1 Guseong-dong Yuseong-gu, Daejeon 305-701, Korea and

²Genome Research Center, Korea Research Institute of Bioscience and Biotechnology, Oun-dong 52 Yuseong-gu, Daejeon 305-333, Korea

Received on January 6, 2003; accepted on February 20, 2003

ABSTRACT

We have established a method for systematic integration of multiple microarray datasets. The method was applied to two different sets of cancer profiling studies. The change of gene expression in cancer was expressed as 'effect size', a standardized index measuring the magnitude of a treatment or covariate effect. The effect sizes were combined to obtain the estimate of the overall mean. The statistical significance was determined by a permutation test extended to multiple datasets. It was shown that the data integration promotes the discovery of small but consistent expression changes with increased sensitivity and reliability. The effect size methods provided the efficient modeling framework for addressing interstudy variation as well. Based on the result of homogeneity tests, a fixed effects model was adopted for one set of datasets that had been created in controlled experimental conditions. By contrast, a random effects model was shown to be appropriate for the other set of datasets that had been published by independent groups. We also developed an alternative modeling procedure based on a Bayesian approach, which would offer flexibility and robustness compared to the classical procedure.

Contact: jkchoi@kaist.ac.kr

Keywords: microarray, meta-analysis, effect size, Bayesian meta-analysis

INTRODUCTION

One of the most challenging tasks microarray analysts face is how to extract, compare, and integrate information from an enormous amount of accumulating data. However, complicated experimental variables embedded in microarray data act as an obstacle to this end. The lack of standards for microarray experiments generates heterogeneous datasets of which direct comparison is not possible. The reasonable approach, in this situation, is to

combine results of individual studies. For this purpose, Rhodes *et al.* (2002) recently applied a meta-analytic approach to microarrays. Here we introduce the most developed meta-analytic method, which expresses the result in the form of effect size.

An effect size approach has desirable features to be applied to microarray data:

1. It provides a standardized index. At present, the measure of expression levels is not interchangeable, in particular between oligo chips and cDNA chips. cDNA microarrays report only the relative change compared to a reference, which is rarely standardized. Obtaining effect sizes offers the direct comparison between the results from different measures.
2. It is based on a well-established statistical framework for combining different results. The main purpose of calculating effect sizes rather than traditional statistics is to draw a synthetic conclusion from multiple studies. We can efficiently integrate microarray data scattered across a multitude of applications.
3. It is superior to other meta-analytic methods in that it has the ability to handle the variability between studies. Appropriate modeling of the interstudy variation is a key factor for successful meta-analysis, especially in an area such as microarray analysis which is prone to study-to-study differences.

By virtue of a wide utility of effect size indices, the method is generally applicable to many types of microarray studies. Here we applied it to differential gene expression studies of clinical tumors where multiple datasets were available. They had been generated with a common objective to make comparisons between the gene expression profiles of tumor and non-tumor tissues. The goal of our meta-analysis was to draw a consensus among the datasets taking into account interstudy variation. With

*To whom correspondence should be addressed.

appropriate modeling, it is expected that the increased sample size enhances the statistical power. To this end, the excellent features of the effect size model will be explicitly presented throughout this paper.

SETS OF DATASETS

Two different sets of microarray datasets were used in this study. One set contains our own four cDNA microarray datasets on hepatocellular carcinoma, which were generated independently but in controlled experimental conditions. They share a common clone set and the same reference except for one experiment. The other set is comprised of four publicly available prostate cancer datasets, which are completely independent. Two datasets are from cDNA technology (Dhanasekaran *et al.*, 2001; Luo *et al.*, 2001), while the others are from oligo-based technology (Magee *et al.*, 2001; Welsh *et al.*, 2001). This set was used in the previous meta-analysis study which ignored inter-study variation (Rhodes *et al.*, 2002). We denote the former LC (liver cancer) datasets and the latter PC (prostate cancer) datasets.

MEASURING EFFECT SIZE

Meta-analysis can be adapted to various types of microarray analyses as there are many kinds of metrics that can be used to measure effect sizes. Discovering differentially expressed genes as in this study might be the most frequent application. If one needs to identify genes whose expression correlates with a quantitative parameter, such as the drug dose in a series of drug treatment experiments, effect size can be defined in terms of the Pearson correlation coefficient. Another no less important application using the correlation coefficient is measuring the gene-to-gene correlation. In contrast to the former cases, it does not require the condition of sharing a common study design. Furthermore, conventional t , F , and χ^2 statistics can be easily converted to effect size indices.

For the measure of differential expression of a gene, a standardized mean difference was obtained as an effect size index. A well-established estimator for the standardized mean difference, found in Hedges and Olkin's (1985) work, was used as

$$d = \frac{\bar{X}_t - \bar{X}_n}{S_p}$$

\bar{X}_t and \bar{X}_n represent the means of the tumor and normal group, respectively and S_p indicates an estimate of the pooled standard deviation. When a study consists of n samples, the unbiased estimate is obtained as $d' = d - 3d/(4(n-2) - 1)$, which indicates the correction for sample size bias. Meanwhile, the estimated variance of the unbiased effect size is given as

$$\hat{\sigma}_d^2 = (n_t^{-1} + n_n^{-1}) + d^2(2(n_t + n_n))^{-1},$$

where n_t and n_n are the sample sizes of each group and d is the unbiased effect size (Hedges and Olkin 1985). It indicates the precision of the measure each study provides.

EFFECT SIZE MODELS

Let μ be the overall mean, typically the parameter of interest, and y_i be the observed effect size for independent studies $i = 1, 2, \dots, k$. The general model is given hierarchically as

$$\begin{aligned} y_i &= \theta_i + \varepsilon_i, & \varepsilon_i &\sim N(0, s_i^2), \\ \theta_i &= \mu + \delta_i, & \delta_i &\sim N(0, \tau^2), \end{aligned}$$

where between-study variance τ^2 represents the variability between studies while within-study variance s_i^2 represents the sampling error conditioned on the i th study. In this application, y_i and s_i^2 are given by d and $\hat{\sigma}_d^2$ described above. Therefore, μ means the average measure of differential expression across the datasets for each gene.

A fixed-effects model (FEM) assumes that the differences of observed effect sizes are from sampling error alone. Thus $\tau^2 = 0$, and consequently $y_i \sim N(\mu, s_i^2)$. On the other hand, a random-effects model (REM) postulates that each effect size is a draw from a distribution with a study-specific mean θ_i and variance s_i^2 . Furthermore, each θ_i is assumed to be a draw from some superpopulation with the overall mean μ and variance τ^2 . Thus $y_i \sim N(\theta_i, s_i^2)$ and $\theta_i \sim N(\mu, \tau^2)$. Therefore, the FEM can be considered as a special case of the REM.

The question of which model is appropriate for given studies can be addressed by testing for the homogeneity of study effects. It is equivalent to assessing the hypothesis that τ^2 is actually zero. Cochran (1954) proposed a now widely-used test of homogeneity based on the statistic

$$Q = \sum w_i (y_i - \hat{\mu})^2,$$

where $w_i = s_i^{-2}$, and $\hat{\mu} = (\sum w_i y_i) / \sum w_i$ is the weighted least squares estimator which ignores between-study variance. Under the hypothesis of homogeneity, it follows a χ_{k-1}^2 distribution. A large observed value of the statistic Q relative to this distribution indicates rejection of the hypothesis of homogeneity, which should indicate the appropriateness of the REM. For microarray studies, the genes can be treated as independent samplings and the homogeneity can be explored over all the genes. The quantile-quantile plots of the observed vs. expected Q values are shown in Figure 1. They strongly suggest that the PC datasets, created independently without any experimental control, generate significantly variable results. By contrast, the effect sizes from the LC datasets can be considered as samples from a common population with the sampling error alone.

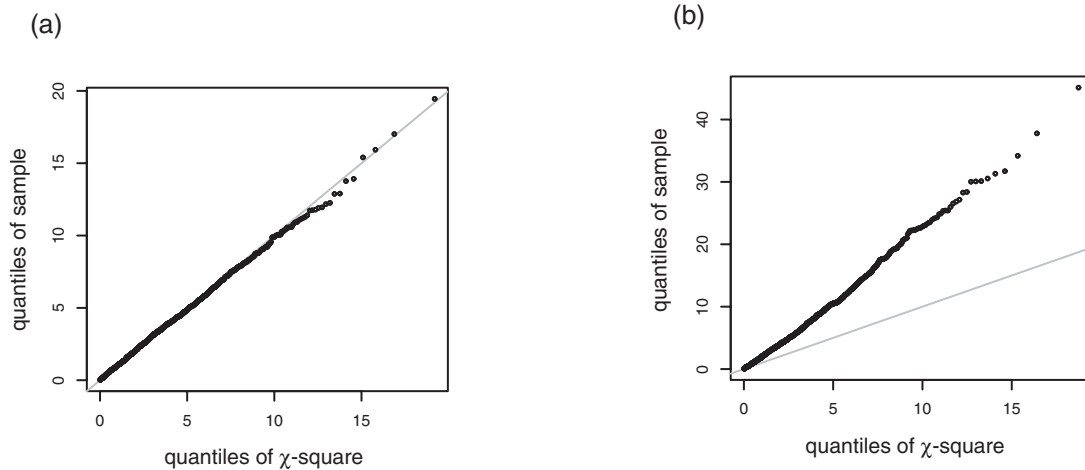


Fig. 1. Gene by gene testing for the homogeneity of study effects. Overall test results are shown by the plot of the observed vs. expected Q quantiles for the (a) LC datasets ($\bar{Q} = 2.91, s_Q^2 = 5.45$) (b) PC datasets ($\bar{Q} = 6.27, s_Q^2 = 30.59$). In the LC datasets, the sample mean and variance as well as the observed Q values are very close to their expected values. In both cases, the genes observed in all four datasets were used in the test. The expected Q values are from the χ_3^2 distribution.

It is standard (Cooper and Hedges, 1994) when estimating μ to use a point estimate for τ^2 in

$$\hat{\mu}(\tau^2) = \frac{\sum (s_i^2 + \tau^2)^{-1} y_i}{\sum (s_i^2 + \tau^2)^{-1}},$$

and

$$\text{Var}[\hat{\mu}(\tau^2)] = \frac{1}{\sum (s_i^2 + \tau^2)^{-1}}.$$

DerSimonian and Laird (1986) developed a method of moments estimator for τ^2 from the expected value of Q :

$$\hat{\tau}_{DL}^2 = \max \left\{ 0, \frac{Q - (k - 1)}{S_1 - (S_2/S_1)} \right\},$$

where $S_r = \sum w_i^r$. Note that it is given for the more general case where the errors ε_i and δ_i do not follow a normal distribution. If they do, $\hat{\mu}(\tau^2)$ is the minimum variance unbiased estimator of μ (Mood *et al.*, 1995). The z statistic was computed as a ratio of $\hat{\mu}(\tau^2)$ over its standard error. The FEM is treated the same way except that $\tau^2 = 0$. Statistically significant genes were chosen by comparing the z statistic assigned to each gene with a given threshold z_{th} . To assess the statistical significance not assuming a normal distribution, empirical distributions were generated by random permutations. In addition, we applied a Bayesian approach to the REM under normal assumption for ε_i and δ_i . The main goal of the Bayesian meta-analysis was to obtain posterior distributions for the overall mean μ and for the study-specific mean θ_i as well.

STATISTICAL SIGNIFICANCE OF COMBINED RESULTS

The effect sizes were combined under the FEM for the LC datasets and the REM using $\hat{\tau}_{DL}^2$ for the PC datasets. The z score of the average effect size, z_j , was obtained from y_{ij} and s_{ij} , for $i = 1, 2, \dots, k$, studies and $j = 1, 2, \dots, p$ genes. We could estimate the statistical significance addressing the multiple testing problem by adapting the core algorithm of SAM (Tusher *et al.*, 2001) to our meta-analysis. It introduced into microarray the concept of false discovery rate (FDR) recently proposed by Benjamini and Hochberg (1995). Column-wise permutations were performed within each dataset, not allowing the expression data to be mixed between the studies. For each permutation $b = 1, 2, \dots, B$, randomized data were created to generate y_{ij}^{*b}, s_{ij}^{*b} . From these values, the overall mean μ_j^{*b} and the variance were estimated to produce z_j^{*b} . The order statistics $z_{(j)}$ ($z_{(1)} \leq \dots \leq z_{(p)}$) and $z_{(j)}^{*b}$ ($z_{(1)}^{*b} \leq \dots \leq z_{(p)}^{*b}$) were obtained. FDR was estimated by

$$\text{FDR} = \frac{(1/B) \sum_b \sum_{(j)} I(|z_{(j)}^{*b}| \geq z_{th})}{\sum_{(j)} I(|z_{(j)}| \geq z_{th})},$$

where $I(\cdot)$ is the indicator function equaling 1 if the condition in parentheses is true, and 0 otherwise. The denominator represents the number of genes called significant in real data. The numerator is the expected number of falsely significant genes and given by the mean number across B permuted data. Median may be used instead of mean. For

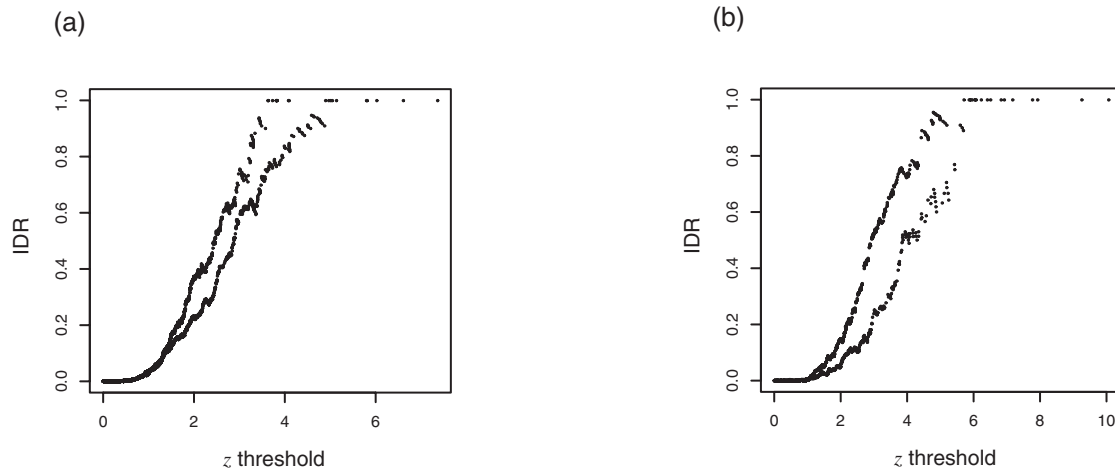


Fig. 2. Integration-discovery rate (IDR) for the (a) LC datasets (b) PC datasets. IDR was computed for $z > 0$ and $z < 0$ separately. z_{th} values were set from the ordered z values. Out of the genes where $|z| \geq z_{\text{th}}$, the genes were counted if $|z_i| < z_{\text{th}}$ for all $i = 1, \dots, k$.

the sake of illustration, using median and $B = 300$, the LC datasets identified 167 genes with a FDR of 2.39% for $z_{\text{th}} = 3.29$ ($P = 0.001$) and the PC datasets identified 285 genes with a FDR of 0.18% for $z_{\text{th}} = 3.89$ ($P = 0.0001$).

To validate our method in terms of expression data, the expression patterns of the identified genes were obtained. They can be visualized at <http://centi.kribb.re.kr/MMA>. The gene list for the PC datasets was compared with the result of Rhodes *et al.* (2002) and the difference was discussed.

INTEGRATION-DRIVEN DISCOVERY

We define integration-driven discovery (IDD) as

$$z \geq z_{\text{th}} \text{ and } \sum_{i=1}^k I(z_i \geq z_{\text{th}}) = 0, \text{ for } z > 0$$

or

$$z \leq -z_{\text{th}} \text{ and } \sum_{i=1}^k I(z_i \leq -z_{\text{th}}) = 0, \text{ for } z < 0.$$

In other words, while a gene is identified as being significantly changed in expression level according to the meta-analysis result, it might be significant in none of the individual studies (with the same statistical significance). Integration-driven discovery rates (IDRs), the ratios to total discoveries, were computed and plotted for given z_{th} values (Fig. 2). The more significant the average effect was, the higher proportion of total discoveries were accounted for by IDs. Interestingly, the most significant discoveries yielded IDRs up to 1.0. For the more realistic threshold, for example $z_{\text{th}} = 3.0$, the PC datasets yielded an IDR of 0.44 and the LC datasets yielded an IDR of 0.63. In other words, about 44–63% of the significant genes were identified purely by the meta-analysis.

It should be noted that IDs occur when combining the ‘small but consistent’ effect sizes. The weak but certainly present effects are brought together to make a result with high statistical confidence. As far as individual studies are concerned, they are potential false negatives. This decrease of false negatives for a fixed α (Type I error), that is to say, an increase of statistical power, is attributable to an increase in overall sample size. In other words, by gathering separate datasets we gain the same effect of increasing sample size. For the selection of a particular number of significant genes, meta-analysis allows us to use a smaller α than individual analyses, decreasing the chance of false positives.

Another important aspect of IDD is the consistency of individual study effects. It assures us that the separate results have been validated through comparative analysis. This interstudy validation does more than simply increasing sample size, because it provides validation through the replicated experiments differing experimental variables and conditions. On the whole, the integration of microarray datasets contributed to achieving more reliable results with increased sensitivity overcoming the artifacts of single analyses.

BIOLOGICAL IMPLICATION

To demonstrate that our model translates statistics into biology, we performed a KEGG pathway database query (Ogata *et al.*, 1999) using significant genes selected in the PC datasets. Using 7342 genes present in at least two studies, $z_{\text{th}} = 3.0$ identified 697 genes containing 248 IDD genes yielding a FDR of 1.5%. Among the 405 up-regulated genes, non-IDD genes were mapped into 70 pathways while the IDD genes were mapped into

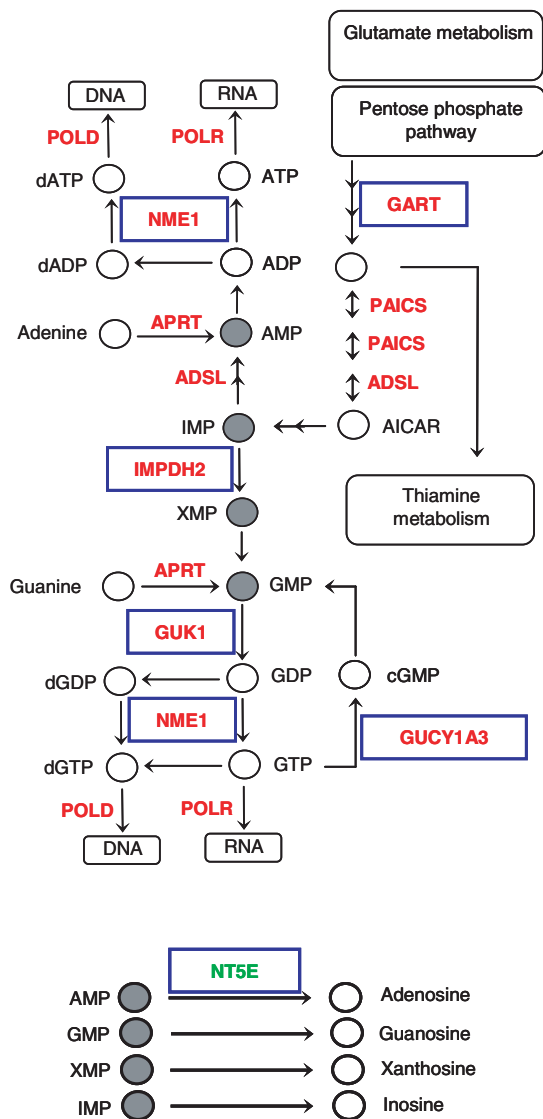


Fig. 3. Biochemical pathway of purine biosynthesis. Genes identified as over- or under-expressed by our effect size model are shown in red and green, respectively. In addition to the adenine biosynthesis pathway, our method identified the guanine biosynthesis pathway. The genes in blue boxes indicate newly discovered genes according to our method. Dephosphorylation processes of monophosphates controlled by NT5E are shown separately below. Monophosphates are shown by gray circles. All the gene names were represented by gene symbols.

51 pathways. Remarkably, in about 70% (35 pathways) among the 51 pathways, the IDD genes appeared along with at least one of the non-IDD genes. This finding implies that the IDD genes enriched the gene set involved in particular pathway-dysregulation in prostate cancer. In other words, IDD would be helpful for building a network of genes dysregulated in related biological processes, by providing sensitive discoveries.

As an illustration, we examined the purine biosynthesis pathway that had been provided as biological evidence in Rhodes *et al.* (2002). It is well known that nucleotide biosynthesis is necessary for DNA synthesis in cell division. As Figure 3 shows, our model detected a considerable number of up-regulated genes and one important down-regulated gene involved in this pathway. The genes we discovered include all the genes found in the previous study without exception (Rhodes *et al.*, 2002) and six additional genes (blue boxes in Fig. 3). Furthermore, discovery of these genes extended the network to include the guanine monophosphate biosynthesis pathway as well. Two of them (GART and GUC1) were IDD genes. Especially, GART is a multifunctional enzyme working in the three consecutive steps linking the pentose phosphate pathway to the adenine phosphate pathway.

For the LC datasets, we found several interesting IDD genes, which called for further investigation. The genes will appear in a separate paper with discussion of the biological implications of up- or down-regulated genes in association with hepatocellular carcinoma.

BAYESIAN META-ANALYSIS

Bayesian methods go naturally with the concept of meta-analysis. Additionally, the hierarchical formulation of a REM would be extended for Bayesian interpretations (DuMouchel and Harris, 1983; Morris and Normand, 1992; Smith *et al.*, 1995; Normand, 1999). We can write the kernel of the joint posterior density of $V = \{\mu, \theta_1, \dots, \theta_k, \tau^2\}$ as

$$p(V|y, s^2) = p(\theta|y, s^2)p(\mu, \tau^2|\theta) \propto \prod_i p(\theta_i|y_i, s_i^2)p(\theta_i|\mu, \tau^2)\pi(\mu)\pi(\tau^2),$$

where $y = (y_1, \dots, y_k)$, $s^2 = (s_1^2, \dots, s_k^2)$ and $\theta = (\theta_1, \dots, \theta_k)$. $\pi(\mu)$ and $\pi(\tau^2)$ are non-informative priors given as $\mu \sim N(0, 10^6)$ and $1/\tau^2 \sim \text{gamma}(0.001, 0.001)$. At first, we assumed that the study effects (θ_i) of the PC datasets arise from a Student's t distribution with 3 degrees of freedom in order to permit the tails to be heavier than a normal distribution. Estimators were given by

$$\begin{aligned} \hat{\tau}_B^2 &= \int \tau^2 p(V|y, s^2) d\theta_i d\mu d\tau^2, \\ \hat{\mu}_B &= \int \mu p(V|y, s^2) d\theta_i d\tau^2 d\mu, \\ \hat{\theta}_i^B &= \int \theta_i p(V|y, s^2) d\theta_j d\mu d\tau^2 d\theta_i, \end{aligned}$$

where $j = 1, \dots, k, j \neq i$. We evaluated the integrals numerically using Markov chain Monte Carlo (MCMC)

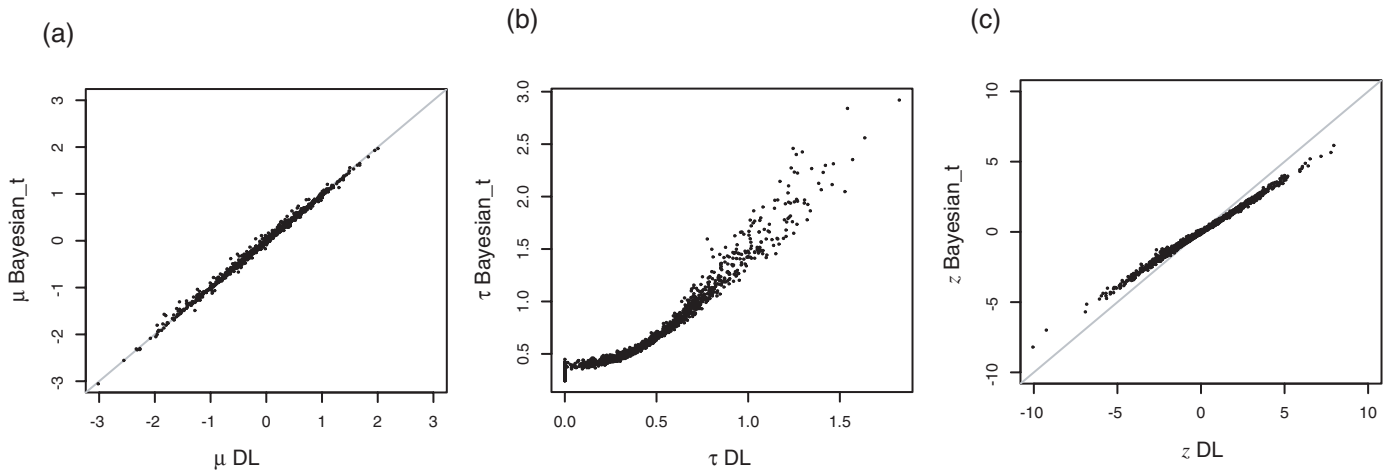


Fig. 4. Comparison of the Bayesian approach and the method of moments. Scatter plot between two estimates for the (a) average effect size μ (b) between-study variance τ^2 (c) z score of the average effect size.

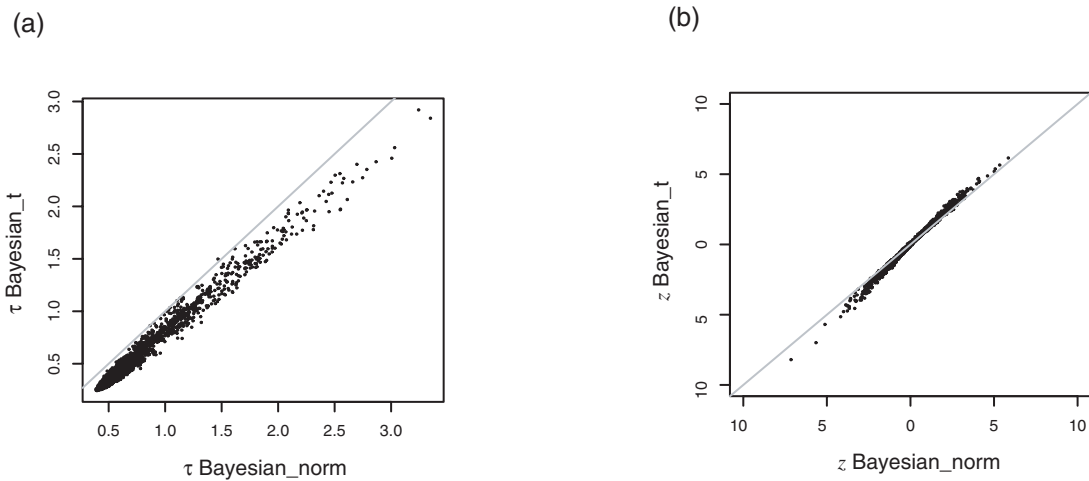


Fig. 5. Comparison of the Bayesian t model and the normal model. Scatter plot between two estimates for the (a) between-study variance τ^2 (b) z score of the average effect size.

via the BUGS (Spiegelhalter *et al.*, 1995) software package. The estimates were compared gene by gene with the estimates obtained using $\hat{\tau}_{DL}^2$ (Fig. 4). The estimates for the average effect size almost completely coincided with one another (Fig. 4a). But the method of moments estimator tended to underestimate τ^2 relative to the Bayesian (Fig. 4b) and consequently overestimate the final z score (Fig. 4c). Whereas, as expected, the normal model for θ_i overestimated the variance τ^2 relative to the t model (Fig. 5a) and slightly underestimated the z score (Fig. 5b).

Therefore, the Bayesian t model could be considered as a compromise between the Bayesian normal model and the classical method. It could correct z outliers overestimated

in the DerSimonian and Laird's method by integrating them into the prior distribution. Simultaneously, it seems that the t model could handle large variability more properly than the normal model. This capability is important in the respect that the variability is a hallmark of microarray data created from a multitude of experimental variables.

Here we used a Student's t with 3 degrees of freedom. As degrees of freedom of the t prior increases, the result will become similar to the normal model. With appropriate prior information and more detailed modeling, incorporating a multivariate normal model if desired, the Bayesian approach would offer a more flexible and robust modeling strategy.

CONCLUSIONS

An explosion of microarray studies has been calling for a method of systematic integration of multiple studies. Moreover, as presented in this paper, microarray is one of the fields that is able to make the most of meta-analysis. Thousands of tests for the homogeneity of study effects can easily reveal overall study-to-study variation. Replication of experiments, carefully designed if possible, would increase the statistical power to detect small effects that may be false negatives in single analyses. Completely independent datasets, with different measures and platforms, would offer interstudy validation to detect consistent effects across different types of experiments. A Bayesian approach, a natural and promising way to combine information from multiple studies, would appear to offer a more flexible and robust modeling strategy. To apply multivariate meta-analysis, taking account of possible interactions and correlations amongst genes, may be the next challenge of meta-analysis on microarrays.

ACKNOWLEDGEMENTS

This work was supported by grant number FG-5-01 Frontier Functional Human Genome Project from the Ministry of Science and Technology of Korea.

REFERENCES

- Benjamini, Y. and Hochberg, Y. (1995) *J. R. Stat. Soc. B*, **57**, 289–300.
- Cochran, B.G. (1954) The combination of estimates from different experiments. *Biometrics*, **10**, 101–129.
- Cooper, H. and Hedges, L.V. (1994) *The Handbook of Research Synthesis*. Russell Sage Foundation, New York.
- DerSimonian, R. and Laird, N.M. (1986) Meta-analysis in clinical trials. *Controlled Clinical Trials*, **7**, 177–188.
- Dhanasekaran, S.M., Barrette, T.R., Ghosh, D., Shah, R., Varambally, S., Kurachi, K., Pienta, K.J., Rubin, M.A. and Chinnaiyan, A.M. (2001) Delineation of prognostic biomarkers in prostate cancer. *Nature*, **412**, 822–826.
- DuMouchel, W.H. and Harris, J.E. (1983) Bayes methods for combining the results of cancer studies in humans and other species. *J. Am. Stat. Assoc.*, **78**, 293–315.
- Hedges, L.V. and Olkin, I. (1985) *Statistical Methods for Meta-analysis*. Academic Press, Orlando.
- Luo, J., Duggan, D.J., Chen, Y., Sauvageot, J., Ewing, C.M., Bittner, M., L. Trent, J.M. and Issacs, W.B. (2001) Human prostate cancer and benign prostatic hyperplasia: molecular dissection by gene expression profiling. *Cancer Res.*, **61**, 4683–4688.
- Magee, J.A., Araki, T., Patil, S., Ehrig, T., True, L., Humphrey, P.A., Catalona, W.J., Watson, M.A. and Milbrandt, J. (2001) Expression profiling reveals hepsin overexpression in prostate cancer. *Cancer Res.*, **61**, 5692–5696.
- Mood, A.M., Graybill, F.A. and Boes, D.C. (1995) *Introduction to the Theory of Statistics*. McGraw-Hill, New York.
- Morris, C.N. and Normand, S.L. (1992) Hierarchical models for combining information and for meta-analyses. In Bernardo, J.M., Berger, J.O., Dawid, A.P. and Smith, A.F.M. (eds), *Bayesian Statistics*, 4th edn, Oxford University Press, New York, pp. 321–344.
- Normand, S.L. (1999) Tutorial in biostatistics—Meta-analysis: formulating, evaluating, combining, and reporting. *Statist. Med.*, **18**, 321–359.
- Ogata, H., Goto, S., Sato, K., Fujibuchi, W., Bono, H. and Kanehisa, M. (1999) KEGG: Kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.*, **27**, 29–34.
- Rhodes, D.R., Barrette, T.R., Rubin, M.A., Ghosh, D. and Chinnaiyan, A.M. (2002) Meta-analysis of microarrays: interstudy validation of gene expression profiles reveals pathway dysregulation in prostate cancer. *Cancer Res.*, **62**, 4427–4433.
- Smith, T.C., Spiegelhalter, D.J. and Thomas, A. (1995) Bayesian approaches to random-effects meta-analysis: a comparative study. *Statist. Med.*, **14**, 2685–2699.
- Spiegelhalter, D.J., Thomas, A., Best, N.G. and Gilks, W.R. (1995) *BUGS: Bayesian Inference Using Gibbs Sampling, version 0.50*. MRC Biostatistics Unit, Cambridge.
- 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.
- Welsh, J.B., Sapinoso, L.M., Su, A.I., Kern, S.G., Wang-Rodriguez, J., Moskaluk, C.A., Frierson, Jr, H.F. and Hampton, G.M. (2001) Analysis of gene expression identifies candidate markers and pharmacological targets in prostate cancer. *Cancer Res.*, **61**, 5974–5978.