# STATISTICAL CHALLENGES IN GENE DISCOVERY THROUGH MICROARRAY DATA ANALYSIS

#### Sreekumar. J<sup>1</sup> and Jose. K. K.<sup>2</sup>

<sup>1</sup>Central Tuber Crops Research Institute,Kerala, India <sup>2</sup>Dept. of Statistics, St. Thomas College, Pala, Kerala, India

email:sreejyothi\_in@yahoo.com

Sreekumar. J<sup>1</sup> and Jose. K. K.<sup>2</sup> Statistical Challenges in Gene Discovery Through Microarray

4 同 ト - 4 三 ト - 4 三

### Introduction

- In 1953, Francis Crick and James Watson: double-helical structure of DNA Figure
- A gene is a fragment of DNA: (start codon, e.g. AUG or GUG), (stop codon, e.g. UAG etc)
- A genome can be comprised of thousands of genes.
- A protein is a sequence of 20 different types of amino acids
- The process of protein synthesis: Central dogma of molecular biology Figure
  - transcription
  - translation

・何・ ・ヨ・ ・ヨ・

- DNA microarrays are devices that measure the gene expression of many thousands of genes simultaneously
- The new data promise to enhance fundamental understanding of life on a molecular level and may prove useful in medical diagnosis, treatment and drug design
- Novel computational tools and reliable data processing are essential for the meaningful and accurate interpretation of microarray data
- In the context of expression-based classification based on a large number of genes, the primary interest is to select a small set of genes, which have a good prediction performance

▲ 同 ▶ ▲ 目 ▶ ▲ 目

- DNA microarrays are devices that measure the gene expression of many thousands of genes simultaneously
- The new data promise to enhance fundamental understanding of life on a molecular level and may prove useful in medical diagnosis, treatment and drug design
- Novel computational tools and reliable data processing are essential for the meaningful and accurate interpretation of microarray data
- In the context of expression-based classification based on a large number of genes, the primary interest is to select a small set of genes, which have a good prediction performance

▲御▶ ▲ 国▶ ▲ 国

- DNA microarrays are devices that measure the gene expression of many thousands of genes simultaneously
- The new data promise to enhance fundamental understanding of life on a molecular level and may prove useful in medical diagnosis, treatment and drug design
- Novel computational tools and reliable data processing are essential for the meaningful and accurate interpretation of microarray data
- In the context of expression-based classification based on a large number of genes, the primary interest is to select a small set of genes, which have a good prediction performance

▲御▶ ▲ 国▶ ▲ 国▶

- DNA microarrays are devices that measure the gene expression of many thousands of genes simultaneously
- The new data promise to enhance fundamental understanding of life on a molecular level and may prove useful in medical diagnosis, treatment and drug design
- Novel computational tools and reliable data processing are essential for the meaningful and accurate interpretation of microarray data
- In the context of expression-based classification based on a large number of genes, the primary interest is to select a small set of genes, which have a good prediction performance

・何・ ・ヨ・ ・ヨ・



#### Figure: Steps in the analysis of Microarray data

Sreekumar. J<sup>1</sup> and Jose. K. K.<sup>2</sup> Statistical Challenges in Gene Discovery Through Microarray

(ロ) (四) (三) (三)

э

#### cDNA Microarray Experimental Procedure

- The process consists of four steps:
  - printing
  - preparation of the biological material
  - hybridization
  - image processing and data analysis
- Selection of a set of genes of interest
- DNA clones of interest (ESTs) are then amplified by PCR to generate a sufficient amount to allow printing onto a glass microslide(Probes)

<回> < E> < E>

- The printing is made by an arrayer
- An arrayer is a robot with a certain number of pins programmed to deposit EST aliquots in an array configuration Figure
- sample preparation, labelling: isolating RNA (e.g., samples from plants under normal condition and plants under drought stress), (Targets)
- The samples are then mixed and hybridized to the arrayed DNAs on the glass slide and any unhybridized cDNA is washed out

(ロ) (四) (三) (三)

- Slides are imaged using a laser scanner
- Typically, the two measurements (one for red and the other for green) for each spot indicate the relative abundance of the corresponding mRNA in the two tissue samples Figure
- The scanned images are analyzed using image analysis software, which evaluates the expression of a gene by quantifying the ratio of the fluorescence intensities of a spot

4 同 ト - 4 三 ト - 4 三

#### Representation of gene expression data

Gene expression data can be represented as a real matrix

<b>X</b> 11	<b>x</b> <sub>12</sub>	 <b>x</b> 1n
<b>x</b> <sub>21</sub>	<b>x</b> <sub>22</sub>	 <b>x</b> <sub>21n</sub>
<b>x</b> <sub>m1</sub>	<b>x</b> <sub>m2</sub>	 <b>X</b> mn

- Each row in the matrix contains the expression data regarding a specified gene gene expression profile
- The vector in the j<sup>th</sup> column of the matrix X, lists the genome wide fluorescence ratios measured by the j<sup>th</sup> array
- Each array represents a condition, different time periods or tissue profile

A (1) A (1) A (1)

### Statistical issues...

- The expression levels are noisy
- Noise creaps into the data in all stages
- accidental / controllable ( Preparation of cells, environmental)
- systematic source of variation
  - array
  - spotter
  - o dye
  - imaging

▲御▶ ▲ 臣▶ ★ 臣▶

### Preprocessing

#### normalization and filtering

- Discard genes where experiments went visibly wrong
- Compute ratio between red and green to account for spot effect
- Take logarithm
- normalize to zero mean
- Variance stabilization(Transformation)
- whether the data need transformation and which?

4 同 ト - 4 三 ト - 4 三

# Visualization tools

Helps in interpreting the results of microarray experiments. The most commonly used of these are illustrated

- Heatmaps : small cells, each consisting of a colour, which represent relative expression values
- Box plots: present various statistics for a given data set. The plots consist of boxes with a central line and two tails
- MA plots: They are often used as an aid when normalizing two-colour cDNA microarrays, where  $M = \log_2 \frac{R}{G}$  and  $A = \log_2 \sqrt{RG}$
- Volcano plots: are used to look at fold change and statistical significance simultaneously
- p-value histograms: the p-values for a test of differential expression for each gene Figure

#### Error Distribution of Gene Expression data

- Gene expression data contains a large number of low expressed genes whose expression level tends to follow a skewed and definitely non normal distribution
- The lognormal distribution is frequently used to model positively skewed gene expression distributions.
- we applied the asymmetric Laplace distribution Figure
- The distribution of error helps in further interpretation and normalization techniques

A (1) > A (1) > A

#### Gene selection

- One of the important problems to be addressed in analysis of microarray data is the identification of differentially expressed gene for further investigation.
- Fold change is the simplest method for identifying differentially expressed genes
- It is known to be unreliable because statistical variability was not taken into account
- Fold change method is subject to bias if the data have not been properly normalized

・何・ ・ヨ・ ・ヨ



 Classic statistical approaches used for detecting differences between two groups include the parametric t-test

$$t_{n_1+n_2-2} = \frac{\bar{X}_1 - \bar{X}_2}{s\sqrt{\frac{1}{n_1} + \frac{1}{n_1}}}$$
(1)

- t-test utilize the variance
- However the small sample sizes : can affect the variance estimates.

4 同 ト - 4 三 ト - 4 三

- If the two groups have different variances, then the two-sample unequal variance t-statistic will be more appropriate.
- Welch method proposes an elaborate correction for degrees of freedom under unequal variances of samples

$$\nu = \left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2 / \frac{\left(\frac{s_1^2}{n_1}\right)^2}{(n_1 - 1)} + \frac{\left(\frac{s_2^2}{n_2}\right)^2}{(n_2 - 1)}$$
(2)

where  $s_1^2$ ,  $s_2^2$ ,  $n_1$  and  $n_2$  are the sample variances and the number of observations in each category.

To address the shortcomings of the *t*-test

$$t^* = rac{ar{M}}{(s+a)/\sqrt{n}}$$

(3)

- Efron et al., percentile of the distribution
- The SAM t-test (S-test) : gene variances
- The regularized *t*-test : gene-specific and global average variance: weighted average of the two
- Broberg : Minimizing a combination of estimated false positive and false negative rates over a grid of significance levels and varinace

A (1) > A (2) > A (2)

 Baldi and Long proposed : t-statistic with a Bayesian adjusted denominator cyber T program, moderated t.

$$t^* = \frac{\bar{M}}{\tilde{s}/\sqrt{n}} \tag{4}$$

A (1) × (1) × (2) × (3)

where

$$\tilde{s}^2 = \frac{\nu_0 \sigma_0^2 + (n-1)s^2}{\nu_0 + n - 2}$$

is the regularized standard deviation.  $\nu_0$  is the strength of the prior and  $\sigma_0^2$  is the background variance estimated from all genes or a set of subset genes.

Smyth extended this to linear model set up

#### Generalized p value technique

Let  $X_{gij}$ ,  $i = 1, 2; j = 1, 2..., n_i$  and  $g = 1, 2..., n_g$  denote the random samples of gene expression data assumed to follow lognormal distribution. Let  $Y_{gij} = \ln(X_{gij})$ , Define,

$$\bar{Y}_{gi} = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij} \text{ and } S_{gi}^2 = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (Y_{gij} - \bar{Y}_{gi})^2, i = 1, 2.$$
for each gene  $g, g = 1, 2, ..., n_g$ . Further let  $\bar{y}_{g1}, \bar{y}_{g2}, s_{g1}^2$  and  $s_{g2}^2$ 
lenote the observed values of  $\bar{Y}_{g1}, \bar{Y}_{g2}, S_{g1}^2$  and  $S_{g2}^2$ 
espectively.

Sreekumar. J<sup>1</sup> and Jose. K. K.<sup>2</sup> Statistical Challenges in Gene Discovery Through Microarray

(ロ) (四) (三) (三)

Let

$$T_{gi} = \bar{y_{gi}} - \frac{\bar{Y_{gi}} - \mu_{gi}}{S_{gi}/\sqrt{n_{gi}}} s_{gi}/\sqrt{n_{gi}} + \frac{1}{2} \frac{\sigma_{gi}^2}{S_{gi}^2} s_{gi}^2$$
(6)

$$= \overline{y_{gi}} - \frac{g_i}{U_{gi}/\sqrt{n_i-1}} \frac{g_i}{\sqrt{n_i}} + \frac{g_i}{2} \frac{g_i}{U_{gi}^2/(n_i-1)}, \quad (7)$$

for 
$$i = 1, 2$$
 and  $g = 1, 2, ..., n_g$ ,  
where  $Z_{gi} = \sqrt{n_i} (\bar{Y}_{gi} - \mu_{gi}) / \sigma_{gi} \sim N(0, 1)$   
and  $U_{gi}^2 = (n_i - 1) S_{gi}^2 / \sigma_{gi}^2 \sim \chi_{n_i-1}^2$ 

- Define generalized pivotal quantity  $T_g = T_g^* (\eta_{g1} \eta_{g2})$ , where  $T_g^* = T_{g1} - T_{g2}$
- Hence The generalized p-value for the two sided test can be obtained as:

$$2 \times \min\left\{P(T_g \le 0), P(T_g \ge 0)\right\}$$
(8)

(ロ) (四) (三) (三)

3

# Other apparoaches

- Thomas et al., : Regression modelling
- Non parametric ( eg:Wilcoxon rank sum test)
- B statistic (log odds ratio) by Lonnstedt and Speed
- Bootstrap t–test

A (1) > A (1) > A

A typical ANOVA model under multiple conditions is

 $y_{ijkg} = \mu + A_i + D_j + AD_{ij} + G_g + AG_{ig} + VG_{kg} + DG_{jg} + \epsilon_{ijkg}$ (9)

where  $y_{ijkg}$  is the measured intensity from array *i*, dye *j*, variety *k* and gene *g* 

$$F_{1} = \frac{(rss_{0} - rss_{1})/(df_{0} - df_{1})}{rss1/df_{1}}$$
(10)

Other *F*-like statistics ( $F_2$  and  $F_3$ )

$$F_{2} = \frac{(rss_{0} - rss_{1})/(df_{0} - df_{1})}{(rss1/df1 + \sigma_{pool}^{2})/2}$$
(11)  
$$F_{3} = \frac{(rss_{0} - rss_{1})/(df_{0} - df_{1})}{2}$$
(12)

(ロ) (四) (三) (三)

 $\sigma_{pool}^2$ 

- The mixed model ANOVA treats some of the factors in the experimental design as random samples from a population and there are multiple levels of variance (biological, array, spot and residual)
- Constructing an appropriate *F*-statistics using the mixed model is tricky
- We can also apply random effects models which use BLUP (best linear unbiased prediction) for the estimation of the gene expression effects

・何・ ・ヨ・ ・ヨ・

# Multiple hypothesis testing

Table: Type I and Type II errors in multiple hypothesis testing.



 p-value: The p-value or observed significance level p is the chance of getting a test statistic as or more extreme than the observed one, under the null hypothesis H<sub>0</sub> of no differential expression.

A ∰ ► A ≡ ► A

The commonly used Type I error rates in multiple hypotheses testing are:

 Per comparison error rate (PCER): the expected value of the number of Type I errors over the number of hypotheses,

$$PCER = \frac{E(V_n)}{m}$$

Per-family error rate (PFER): the expected number of Type I errors,

 $\mathsf{PFE} = \mathsf{E}(\mathsf{V}_n)$ 

Family-wise error rate: the probability of at least one Type I error

$$FEWR = Pr(V_n = 1)$$

 False discovery rate (FDR) is the expected proportion of Type I errors among the rejected hypotheses

$$\mathsf{FDR} = \mathsf{E}(\frac{\mathsf{V}_n}{\mathsf{R}_n};\mathsf{R}_n > 0) = \mathsf{E}(\frac{\mathsf{V}_n}{\mathsf{R}_n}|\mathsf{R}_n > 0) \times \mathsf{Pr}(\mathsf{R}_n > 0)$$

Positive false discovery rate (pFDR): the rate that discoveries are false

$$\mathsf{pFDR} = \mathsf{E}(\frac{\mathsf{V}_n}{\mathsf{R}_n} | \mathsf{R}_n > 0)$$

▲ □ ▷ < @ ▷ < 분 ▷ < 분 ▷ 분 의 오 ↔</p>
Statistical Challenges in Gene Discovery Through Microarray

#### There are single step procedures for controlling FWER

- Bonferroni single step adjusted p-values
- Sidak single-step adjusted p-values
- Westfall and Young single-step minP adjusted p-values

▲□▶ ▲ 三▶ ▲ 三▶

3

#### False Discovery Rate

FDR is defined to be the expected value of the ratio of the number of incorrectly rejected hypotheses and the total of number of rejected hypotheses. Controlling FDR proves more powerful than FWER and has become increasingly adopted for genomic studies.

### **Bayesian variable selection**

- The literature on microarray data is mainly based on two distributions: the log-normal and the gamma distributions
- That often appear to be effective when used in a Bayesian hierarchical framework
- In the case of empirical Bayes studies, inference is usually made on some quantities related to the posterior distribution of the parameter of interest, or of a certain type of hypothesis.
- Lee *et al.*, proposed a hierarchical bayesian model and employed latent variables to specialize the model to a regression setting and applied variable selection to select differentially expressed genes

#### Conclusion and Future recommendations

- Multiple methods can produce list of differentially expressed genes
- Beware of multiple hypothesis testing problems
- In many areas (normalization algorithms and false-discovery-rate estimation procedures), the need for thoroughly evaluating existing techniques currently seems to outweigh the need to develop new techniques.
- There may be no single best way to represent microarray data
- More research is needed on how to best examine intersections between sets of findings and evaluate complex multi-component hypotheses. Bayesian approaches might be especially advantageous here
- For all statistical procedures, the fact that transcripts are not necessarily independent should be considered. The potential impact of this on the performance of procedures should be assessed, and ways to accommodate this are needed.



#### Figure: The structure of DNA

Back to presentation

Sreekumar. J<sup>1</sup> and Jose. K. K.<sup>2</sup> Statistical Challenges in Gene Discovery Through Microarray

-12



Back to presentation

Sreekumar. J<sup>1</sup> and Jose. K. K.<sup>2</sup>

Statistical Challenges in Gene Discovery Through Microarray

-1

(日) (四) (王) (王) (王)

# Microarray spotting device and examples of commonly used print heads



Back to presentation

Sreekumar. J<sup>1</sup> and Jose. K. K.<sup>2</sup> Statistical Challenges in Gene Discovery Through Microarray

(日) (종) (종) (종)

э



- red : more expressed in test sample
- green : more expressed in reference sample
- vellow : equally expressed
- black : no expression

< 17 ×



Back to presentation

Sreekumar. J<sup>1</sup> and Jose. K. K.<sup>2</sup>

Statistical Challenges in Gene Discovery Through Microarray

≣⇒

-1

・ロト ・四ト ・ヨト ・



Figure: Histogram

Back to presentation

Sreekumar. J<sup>1</sup> and Jose. K. K.<sup>2</sup> Statistical Challenges in Gene Discovery Through Microarray

◆□▶ ◆□▶ ◆三▶ ◆三▶ -

æ –