

Bayesian modelling of gene expression data

Alex Lewin

Sylvia Richardson (IC Epidemiology)
Tim Aitman (IC Microarray Centre)
Philippe Broët (INSERM, Paris)

In collaboration with
Anne-Mette Hein, Natalia Bochkina (IC Epidemiology)
Helen Causton (IC Microarray Centre)
Peter Green and Graeme Ambler (Bristol)

Contents

- Introduction to microarrays
- Differential expression
- Bayesian mixture estimation of false discovery rate

Introduction to microarrays

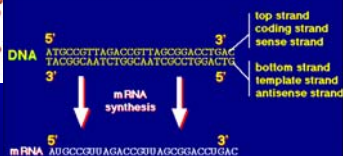
Post-genome Genetics Research

- Challenge: Identify function of all genes in genome
- DNA microarrays allow study of thousands of genes simultaneously

Gene Expression

DNA → mRNA → protein

TRANSCRIPTION



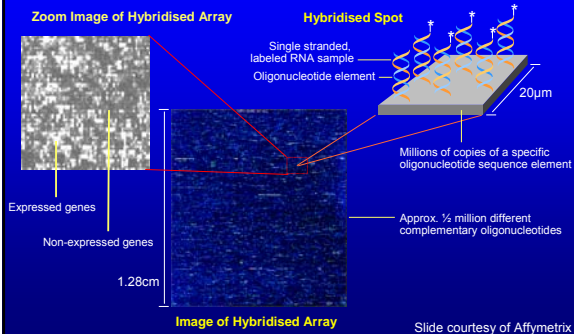
Pictures from <http://www.emc.maricopa.edu/faculty/farabee/BIOBK/BioBookTOC.html>

Hybridisation

- Known sequences of single-stranded DNA immobilised on microarray
- Tissue sample (with unknown concentration of RNA) fluorescently labelled
- Sample hybridised to array
- Excess sample washed off array
- Array scanned to measure amount of RNA present for each sequence on array



The Principle of Hybridisation



Output of Microarray

- Each gene is represented by several different DNA sequences (probes)
- Obtain intensity for each probe
- Different tissue samples on different arrays so compare gene expression for different experimental conditions

Differential Expression

AL, Sylvia Richardson,
Clare Marshall, Anne
Glazier, Tim Aitman

Microarray analysis is a multi-step process

Low-level Model
(how gene expression is estimated from signal)

Normalisation
(to make arrays comparable)

Differential Expression

Clustering Partition Model

We aim to integrate all the steps in a common statistical framework

Differential Expression 1 of 18

Bayesian hierarchical model framework

- Model different sources of variability simultaneously, within array, between array ...
- Share information in appropriate ways to get better estimates, e.g. estimation of gene specific variability.
- Uncertainty propagated from data to parameter estimates.
- Incorporate prior information into the model.

Differential Expression 2 of 18

Data Set and Biological question

Previous Work (Tim Aitman, Anne Marie Glazier)

The spontaneously hypertensive rat (SHR): A model of human insulin resistance syndromes.

Deficiency in gene Cd36 found to be associated with insulin resistance in SHR (spontaneously hypertensive rat)

Differential Expression 3 of 18

Data Set and Biological question

Microarray Data

- 3 SHR compared with 3 transgenic rats
- 3 wildtype (normal) mice compared with 3 mice with Cd36 knocked out
- ≅ 12000 genes on each array

Biological Question

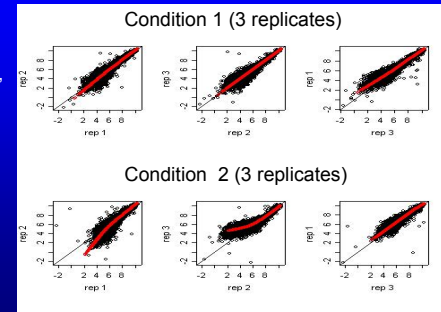
Find genes which are expressed differently in wildtype and knockout mice.

Differential Expression 4 of 18

Data

Needs 'normalisation'

Spline curves shown



Differential Expression 5 of 18

Model for Differential Expression

- Expression-level-dependent normalisation
- Only 3 replicates per gene, so share information between genes to estimate gene variances
- To select interesting genes, use posterior distribution of ranks

Differential Expression 6 of 18

Bayesian hierarchical model for genes under one condition

Data: y_{gr} = log gene expression for gene g , replicate r

α_g = gene effect

$\beta_{r(g)}$ = array effect (expression-level dependent)

σ_g^2 = gene variance

• 1st level

$$y_{gr} \sim N(\alpha_g + \beta_{r(g)}, \sigma_g^2), \quad \sum_r \beta_{r(g)} = 0$$

$\beta_{r(g)}$ = function of α_g , parameters $\{\mathbf{a}\}$ and $\{\mathbf{b}\}$

Differential Expression 7 of 18

Bayesian hierarchical model for genes under one condition

• 2nd level

Priors for α_g , coefficients $\{\mathbf{a}\}$ and $\{\mathbf{b}\}$
 $\sigma_g^2 \sim \text{lognormal}(\mu, \tau)$

Hyper-parameters μ and τ can be influential.
 In a full Bayesian analysis, these are **not fixed**

• 3rd level

$$\mu \sim N(c, d) \quad \tau \sim \text{lognormal}(e, f)$$

Differential Expression 8 of 18

Details of array effects (Normalisation)

- Piecewise polynomial with unknown break points:

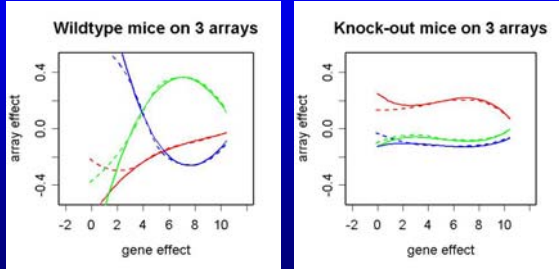
$\beta_{r(g)}$ = quadratic in α_g for $a_{rk-1} \leq \alpha_g \leq a_{rk}$
 with coeff $(b_{rk}^{(1)}, b_{rk}^{(2)})$, $k = 1, \dots, \#\text{breakpoints}$

- Locations of break points not fixed
- Must do sensitivity checks on # break points
- Cubic fits well for this data

Differential Expression 9 of 18

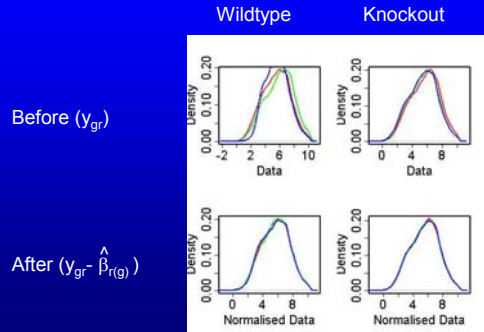
Non linear fit of array effect as a function of gene effect

— cubic
 loess



Differential Expression 10 of 18

Effect of normalisation on density

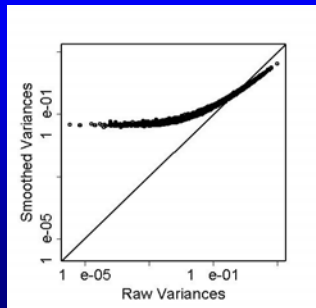


Differential Expression 11 of 18

Smoothing of the gene specific variances

• Variances are estimated using information from all G x R measurements (~12000 x 3) rather than just 3

• Variances are stabilised and shrunk towards average variance



Differential Expression 12 of 18

Bayesian Model Checking

- Check our assumption of different variance for each gene
- Predict sample variance $S_g^{2\text{new}}$ from the model for each gene
- Compare predicted $S_g^{2\text{new}}$ with observed $S_g^{2\text{obs}}$

$$\text{Bayesian p-value } \text{Prob}(S_g^{2\text{new}} > S_g^{2\text{obs}})$$

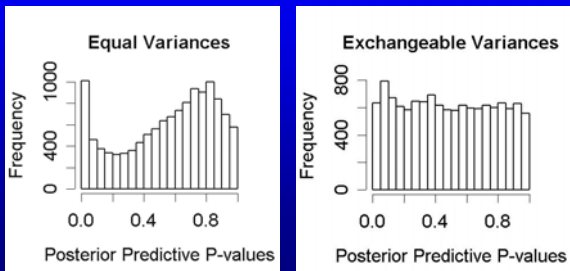
- Distribution of p-values Uniform if model is adequate
- Easily implemented in MCMC algorithm

Differential Expression 13 of 18

Bayesian predictive p-values

Control for method: equal variance model has too little variability for the data

Exchangeable variance model is supported by the data



Differential Expression 14 of 18

Differential expression model

d_g = differential effect for gene g between 2 conditions

Joint model for the 2 conditions :

$$y_{g1r} \sim N(\alpha_g - \frac{1}{2} d_g + \beta_{r(g)1}, \sigma_{g1}^2), \quad (\text{condition 1})$$

$$y_{g2r} \sim N(\alpha_g + \frac{1}{2} d_g + \beta_{r(g)2}, \sigma_{g2}^2), \quad (\text{condition 2})$$

Prior can be put on d_g directly

Differential Expression 15 of 18

Possible Statistics for Differential Expression

$d_g \approx \log \text{fold change}$

$d_g^* = d_g / (\sigma_{g_1}^2 / 3 + \sigma_{g_2}^2 / 3)^{1/2}$ (standardised difference)

- We obtain the **joint distribution** of all $\{d_g\}$ and/or $\{d_g^*\}$
- Distributions of ranks

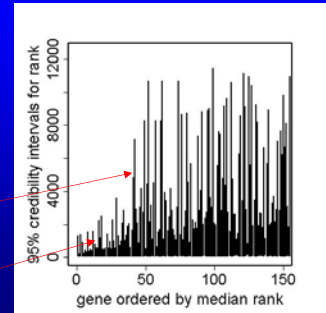
Differential Expression 16 of 18

Credibility intervals for ranks

150 genes with lowest rank (most under-expressed)

Low rank, high uncertainty

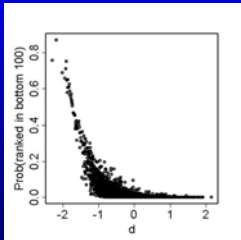
Low rank, low uncertainty



Differential Expression 17 of 18

Probability statements about ranks

Under-expression:
probability that gene is ranked in bottom 100 genes



Have to choose rank cutoff (here 100)

Have to choose how confident we want to be in saying the rank is less than the cutoff (eg prob=80%)

Differential Expression 18 of 18

Summary: Differential Expression

- Expression-level-dependent normalisation
- Only 3 replicates per gene, so share information between genes to estimate gene variances
- To select interesting genes, use posterior distribution of ranks

Bayesian estimation of False Discovery Rate

Philippe Broët, AL, Sylvia Richardson

Multiple Testing

- Testing thousands of hypotheses simultaneously
- Traditional methods (Bonferroni) too conservative
- Challenge: select interesting genes without including too many false positives.

FDR 1 of 16

False Discovery Rate

	Declare negative	Declare positive	
True negative	U	V	m_0
True positive	T	S	m_1
	N-R	R	N

Bonferroni:
FWER= $P(V>0)$

FDR= $E(V/R)$

FDR 2 of 16

FDR as a Bayesian Quantity

Storey showed that

$$E(V/R | R>0) = P(\text{truly negative} | \text{declare positive})$$

Storey starts from p-values.

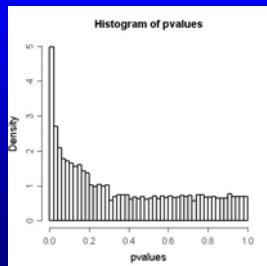
We directly estimate posterior probabilities.

FDR 3 of 16

Storey Estimate of P(null)

P-values are Uniform if all genes obey the null hypothesis

Estimate P(null) where density of p-values is approximately flat



FDR 4 of 16

Storey Estimate of FDR

Lists based on **ranking** genes (ordered rejection regions)

List i is all genes with p-value $p_g \leq p_i^{\text{cut}}$

For list i , $P(\text{declare positive} | \text{truly negative}) = p_i^{\text{cut}}$

$$\begin{aligned} \text{FDR}_i &= P(\text{truly -}) P(\text{declare +} | \text{truly -}) / P(\text{declare +}) \\ &= P(\text{null}) p_i^{\text{cut}} N/N_i \end{aligned}$$

FDR 5 of 16

Bayesian Estimate of FDR

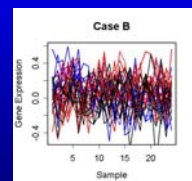
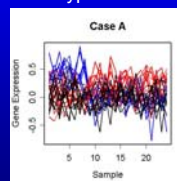
- Classify genes as under-expressed, ..., unaffected, ..., over-expressed (may be several different levels of over and under-expression)
- 'unaffected' \leftrightarrow null hypothesis
- For each gene, calculate probability of following the null distribution
- FDR = mean P(gene belonging to null) for genes declared positive

FDR 6 of 16

Gene Expression Profiles

Each gene has repeat measurements under several conditions: gene profile

Null hypothesis: no change across conditions

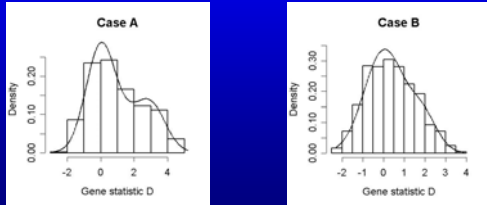


Summarize profile by F-statistic (one for each gene)

FDR 7 of 16

Transformation of F-statistics

Transform $F \rightarrow D$ approx. Normal if no change across conditions



FDR 8 of 16

Bayesian mixture model

Mixture model specification $\xrightarrow{\text{NULL}}$ $\xrightarrow{\text{ALTERNATIVE}}$

$$D_g \sim w_0 N(0, \sigma_0^2) + \sum_{j=1,k} w_j N(\mu_j, \sigma_j^2)$$

μ_j ordered, uniform on upper range

k , unknown number of components \rightarrow alternative is modelled semi-parametrically

Results integrated over different values of k

FDR 9 of 16

Bayes Estimate of FDR

Latent variable $z_g = 0, 1, \dots, k$ with prob w_0, w_1, \dots, w_k

$P(\text{gene } g \text{ in null} \mid \text{data}) = P(z_g = 0 \mid \text{data})$

For **any** given list L containing N_L genes,

$$\text{FDR}_L = 1/N_L \sum_{g \text{ on list } L} P(\text{gene } g \text{ belonging to null} \mid \text{data})$$

FDR 10 of 16

Compare Estimates of FDR

Storey $\text{FDR}_i = 1/N_i p_i^{\text{cut}} P(\text{null}) N$

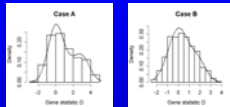
Bayes $\text{FDR}_i = 1/N_i \sum_{\text{list } i} P(\text{gene in null} \mid \text{data})$

NB Bayes estimate can be calculated for any list of genes, not just those based on ranking genes

FDR 11 of 16

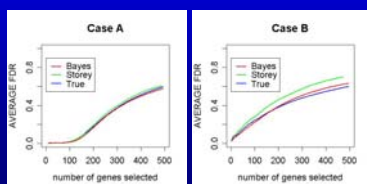
Results for Simulated Data (Ordered lists of genes)

Simulation study: 50 simulations of each set of profiles



Usual methods (Storey q-value and SAM) overestimate FDR

Bayes mixture estimate of FDR is closer to true value



FDR 12 of 16

Breast Cancer Data

- Study of gene expression changes among 3 types of tumour: BRCA1, BRCA2 and sporadic tumours (Hedenfalk et al).
- Gene profiles across tumours summarized by F-statistics, transformed to D.
- Estimate FDR using Bayesian mixture model.

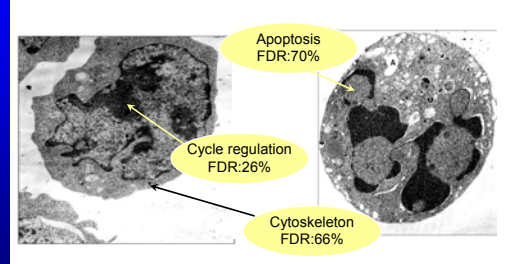
FDR 13 of 16

FDR for subsets of genes

- Fit data for **all** genes (2471) using the Bayesian mixture model
- Estimate FDR for pre-defined groups of genes with known functions:
 - apoptosis (26 genes)
 - cell cycle regulation (21 genes)
 - cytoskeleton (25 genes)

FDR 14 of 16

Results for subsets of genes



FDR 15 of 16

Slide from Philippe Broët

Summary: FDR

- Good estimate of FDR and FNR
- Semi-parametric model for differentially expressed genes.
- Obtain posterior probability for each gene.
- Can calculate FDR, FNR for any list of genes.

FDR 16 of 16

Summary

Differential Expression

Expression-level-dependent normalisation
Borrow information across genes for variances
Joint distribution of ranks

False Discovery Rate

Flexible mixture gives good estimate of FDR

Future work

Mixture prior on log fold changes, with uncertainty propagated to mixture parameters