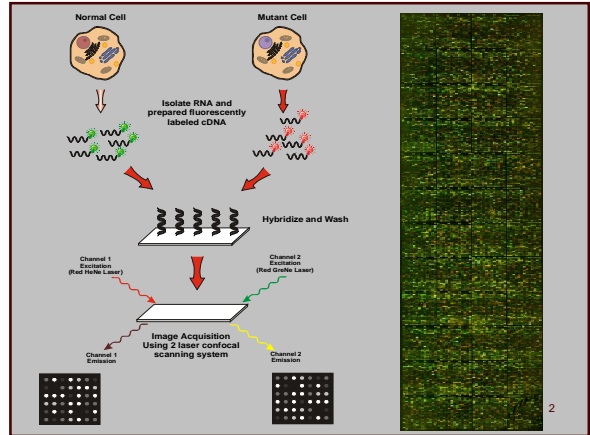


Empirical Bayes and Mixed Linear Models for Assessing Differential Expression in cDNA Microarray Experiments

Gordon Smyth
Walter and Eliza Hall Institute

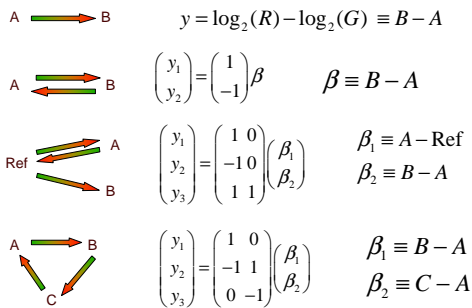


1



2

Designs → Linear Models



3

Linear Model Estimates

Obtain a linear model for each gene g
 $E(\underline{y}_g) = X \underline{\beta}_g \quad \text{var}(\underline{y}_g) = W_g^{-1} \sigma_g^2$

Estimate model by **robust regression**, **least squares** or **generalized least squares** to get

coefficients $\hat{\beta}_{gj}$
 standard deviations s_g
 standard errors $\text{se}(\hat{\beta}_{gj})^2 = c_{gj} s_g^2$

Parallel Inference for Genes

- 10,000-40,000 linear models
- Curse of dimensionality:** Need to adjust for multiple testing, e.g., control family-wise error rate (FWE) or false discovery rate (FDR)
- Boon of parallelism:** Can borrow information from one gene to another

5

Hierarchical Model

Normal Model

$$\hat{\beta}_{gj} \sim N(\beta_{gj}, c_{gj} \sigma_g^2)$$

$$s_g^2 \sim \sigma_g^2 \chi_{d_g}^2$$

Prior

$$P(\beta_{gj} \neq 0) = p$$

$$\beta_{gj} | \beta_{gj} \neq 0 \sim N(0, c_{0j} \sigma_g^2)$$

$$\sigma_g^2 \sim s_0^2 (\chi_{d_0}^2 / d_0)^{-1}$$

Reparametrization of Lönnstedt and Speed 2002

Normality, independence assumptions are wrong but convenient, resulting methods are useful

6

Posterior Statistics

Posterior variance estimators

$$\tilde{s}_g^2 = \frac{s_g^2 d_g + s_0^2 d_0}{d_g + d_0}$$

Moderated t-statistics

$$\tilde{t}_{gj} = \frac{\hat{\beta}_{gj}}{\tilde{s}_g \sqrt{c_{gj}}}$$

Eliminates large t-statistics merely from very small s

7

Marginal Distributions

The marginal distributions of the sample variances and moderated t-statistics are mutually independent

$$s_g^2 \sim s_0^2 F_{d, d_0}$$

$$\tilde{t}_g \sim \begin{cases} t_{d_0+d} & \text{with prob } 1-p \\ \sqrt{1+c_0/c} t_{d_0+d} & \text{with prob } p \end{cases}$$

Degrees of freedom add!

Known result?

Estimating Prior Parameters

Marginal moments of $\log s^2$ lead to estimators of s_0 and d_0 :

Estimate d_0 by solving

$$\psi'(d_0/2) = \text{mean} \{ n s_e^2 - \psi'(d_g/2) \}$$

where

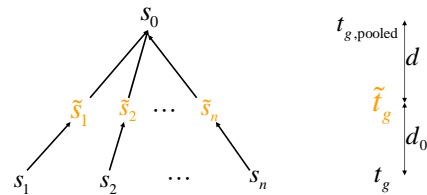
$$e_g = \log s_g^2 - \psi(d_g/2) + \log(d_g/2)$$

Finally

$$s_0^2 = \exp \{ \bar{e} + \psi(d_0/2) - \log(d_0/2) \}$$

9

Shrinkage of Standard Deviations

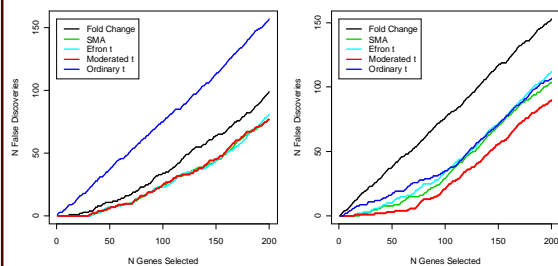


The data decides whether \tilde{t}_g should be closer to

$t_{g, \text{pooled}}$ or to t_g

10

Simulations

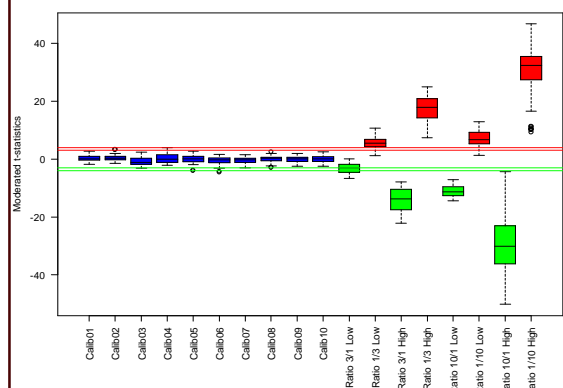


σ^2 similar

σ^2 very different

11

Scorecard Controls



Posterior Odds

Posterior probability of differential expression for any gene is

$$\frac{p(\beta \neq 0 | \hat{\beta}, s^2)}{p(\beta = 0 | \hat{\beta}, s^2)} = \frac{p}{1-p} \left(\frac{c}{c+c_0} \right)^{1/2} \left\{ \frac{\tilde{t}^2 + d + d_0}{\tilde{t}^2 \frac{c}{c+c_0} + d + d_0} \right\}^{\frac{1+d+d_0}{2}}$$

Monotonic function of \tilde{t}^2 for constant d

Reparametrization of Lönnstedt and Speed 2002

Quantile Estimation of c_0

Let r be rank of $|\tilde{t}_g|$ in descending order, and let $F(\cdot)$ be the distribution function of the t-distribution. Can estimate c_0 by equating empirical to theoretical quantiles:

$$2 \left[pF \left(-\sqrt{\frac{c_g}{c_g+c_0}} |\tilde{t}_g|; d_0 + d_g \right) + (1-p)F(-|\tilde{t}_g|; d_0 + d_g) \right] = \frac{r-0.5}{n}$$

Get overall estimator of c_0 by averaging the individual estimators from the top $p/2$ proportion of the $|\tilde{t}_g|$

Duplicate spots

- Replicate spots of each gene on same array, assume duplicates at regular spacing
- Assume spatial component of correlation between duplicates is same for each gene
- Estimate spatial correlation from **consensus** estimator across genes

15

Posterior F-tests

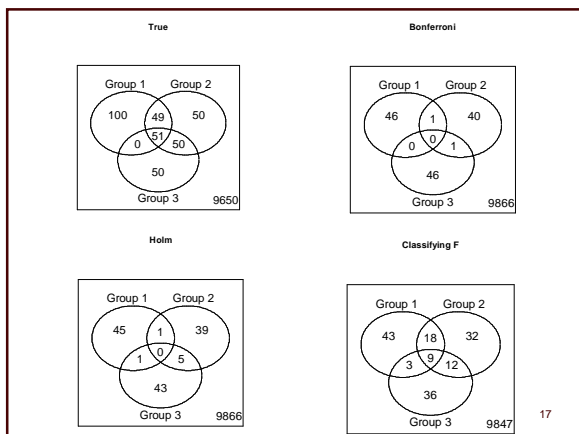
If $\beta_g = 0$

then

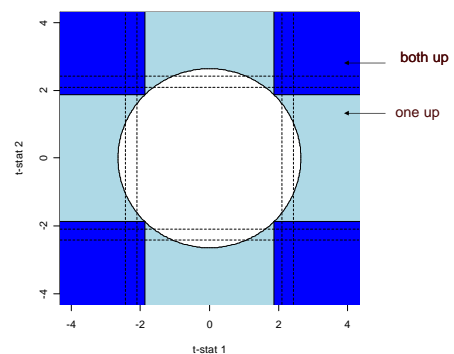
$$\frac{\hat{\beta}_g^T X^T W X \hat{\beta}_g}{\tilde{s}_g^2} \sim F_{k, d+d_0}$$

Non-null prior on β doesn't enter

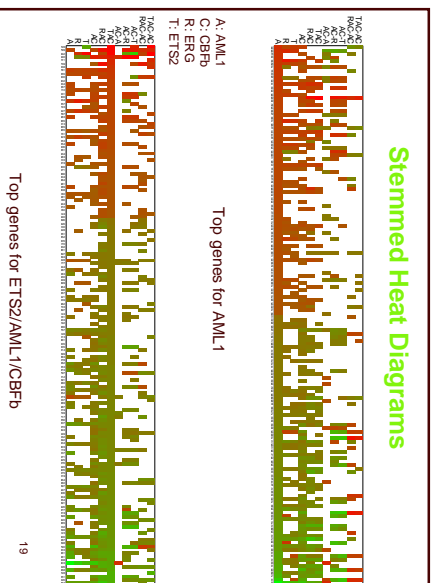
16



F-Tests as Classification Problem



Stemmed Heat Diagrams



Acknowledgements

- WEHI Bioinformatics*
 - Terry Speed
 - Ingrid Lönnstedt
 - Matt Ritchie
- WEHI Scott Lab*
 - Joelle Michaud
 - Hamish Scott
- AGRF*
 - Steve Wilcox
 - Cathy Jensen

20