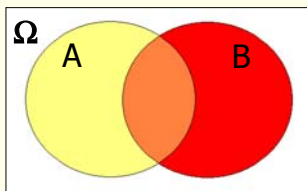


MIXTURE MODELS, BAYES AND EMPIRICAL BAYES APPROACHES

- ➔ A Sketch of Bayesian Inference
- ➔ Applications (Microarray Experiments)
 - Mixture model approach
 - Empirical Bayes methods
- ➔ Discussion

BAYESIAN INFERENCE AND MCMC

Conditional Probability (Bayes' Rule)



$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A)P(B|A)}{P(B)}$$

Bayesian Inference

y : observed data
 θ : parameters
(all unobserved quantities)

$$p(\theta|y) = \frac{p(\theta, y)}{p(y)} = \frac{p(\theta)p(y|\theta)}{p(y)}$$

posterior distribution \leftarrow $p(\theta|y) \propto p(\theta)p(y|\theta)$ \rightarrow sampling distribution
prior distribution

Marginal Posterior Distributions

$$p(\theta_1 | y) \propto \int_{\theta \neq \theta_1} p(\theta) p(y | \theta) d\theta_{\theta \neq \theta_1}$$

Gibbs Sampling

$$\theta = (\theta_1', \theta_2', \dots, \theta_r')' \rightarrow p(\theta_i | \theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_r)$$

$$\theta^{(0)} = (\theta_1^{(0)}, \theta_2^{(0)}, \dots, \theta_r^{(0)})'$$

$$\theta_1^{(1)} | \theta_2^{(0)}, \theta_3^{(0)}, \dots, \theta_r^{(0)}$$

$$\theta_2^{(1)} | \theta_1^{(1)}, \theta_3^{(0)}, \dots, \theta_r^{(0)}$$

⋮

$$\theta_r^{(1)} | \theta_2^{(1)}, \theta_3^{(1)}, \dots, \theta_{r-1}^{(1)}$$

Burn-in & Convergence

Tinning interval & Lag correlations

Sample size & Monte Carlo error

NORMAL MIXTURE MODEL

(Pan et al., 2002)

$$\begin{cases} X_{1i}, \dots, X_{mi} & \text{(gene } i, m \text{ arrays under condition 1)} \\ Y_{1i}, \dots, Y_{mi} & \text{(gene } i, m \text{ arrays under condition 2)} \end{cases}$$

$$\begin{cases} X_{ji} = \mu_{(1)i} + \varepsilon_{ji} \\ Y_{li} = \mu_{(2)i} + e_{li} \end{cases} \rightarrow \begin{cases} \mu_{(1)i} \text{ and } \mu_{(2)i} : & \text{mean expression levels for gene } \\ & \text{i under the two conditions} \\ \varepsilon_{ji} \sim (0, \sigma_{(1)i}^2) \\ e_{li} \sim (0, \sigma_{(2)i}^2) \end{cases} \text{ usually symmetric distributions}$$

$$H_0: \mu_{(1)i} = \mu_{(2)i}$$

$$Z_i = \frac{\sum_{j=1}^m X_{ji} / m}{\sigma_{(1)i}} - \frac{\sum_{l=1}^m Y_{li} / m}{\sigma_{(2)i}} = \frac{\mu_{(1)i}}{\sigma_{(1)i}} - \frac{\mu_{(2)i}}{\sigma_{(2)i}} + \frac{\sum_{j=1}^m \varepsilon_{ji}}{m\sigma_{(1)i}} - \frac{\sum_{l=1}^m e_{li}}{m\sigma_{(2)i}}$$

$$E[Z_i] = \frac{\mu_{(1)i}}{\sigma_{(1)i}} - \frac{\mu_{(2)i}}{\sigma_{(2)i}} \quad \text{and} \quad \text{Var}[Z_i] = \frac{2}{m}$$

Estimating the Null Distribution (f_0)

Nonparametric method:

- Null score for each gene

$$Z_i = \frac{X_{1i} - X_{2i} + \dots + X_{m-1,i} - X_{mi} + Y_{1i} - Y_{2i} + \dots + Y_{m-1,i} - Y_{mi}}{m\sigma_{(1)i} + m\sigma_{(2)i}}$$
$$= \frac{\epsilon_{1i} - \epsilon_{2i} + \dots + \epsilon_{m-1,i} - \epsilon_{mi} + e_{1i} - e_{2i} + \dots + e_{m-1,i} - e_{mi}}{m\sigma_{(1)i} + m\sigma_{(2)i}}$$

$$\left\{ \begin{array}{l} \epsilon_{ji} \text{ and } -\epsilon_{ji} : \text{ same (symmetric) distribution} \\ e_{ji} \text{ and } -e_{ji} : \text{ same (symmetric) distribution} \end{array} \right.$$

⇒ z_i values across all genes are used to estimate f_0

EMPIRICAL BAYES (EB) APPROACH (Newton and Kendziorski, 2003)

- Well-suited to high-dimensional inference problems ($n \times p$)
- Borrowing of information across genes ($p \gg n$)
 particularly information related to variability in the system
- Probability distributions are specified in several layers (accounting for multiple sources of variation).
- Methods provide estimates of differential expression, hypothesis tests and assessment of patterns of differential expression among multiple conditions, and rank ordering genes (cDNA or Affymetrix).
- Comparison of multiple groups: posterior probability of differential expression using a parametric mixture model.

→ Canonical EB example

- Single array (J genes)
- X_j observed (measured) expression level $\sim (\mu_j, \sigma^2)$
($j = 1, 2, \dots, J$)

(Scaled) biological +
other sources of variation

True expression level of gene j

$$\hat{\mu}_j = \text{MLE}(\mu_j) = x_j \quad \text{Least variance (among unbiased estimators)}$$

- These obvious estimates has some problems when J is relatively large (Stein, 1956)

$$\text{Squared Euclidean length of the vector } \mathbf{x} > \text{Squared length of the target profile } \boldsymbol{\mu}$$

$$\sum_j x_j^2$$

- Better estimates of $\boldsymbol{\mu}$ can be found...

Empirical Bayes argument

- Treat μ_j as random;
- Work out a summary feature of $p(\boldsymbol{\mu}|\mathbf{x})$;
- Plug in estimates of remaining parameters using the whole dataset.

e.g. $\hat{\mu}_j = wx_j + (1-w)\bar{x}$

Data-dependent weight
 $0 \leq w \leq 1$

$$\bar{x} = \sum x_j / J$$

Sharing information
among genes

→ General Model Structure: Two Conditions

- I samples
- $\mathbf{x}_j = (x_{j1}, x_{j2}, \dots, x_{ji})'$
- Probability distribution of \mathbf{x}_j ?
- Baseline hypothesis: I samples are exchangeable, i.e., x_{ji} are independent random deviations from a gene specific mean value μ_j as $f_{\text{obs}}(\cdot | \mu_j)$.

Two groups (conditions) with s_k ($k = 1, 2$) samples each; $I = s_1 + s_2$

EE_j: equivalent expression for gene j (baseline hypothesis above holds)

DE_j: differential expression for gene j (μ_{j1} and μ_{j2})

Fold change: $\rho_j = \mu_{j1} / \mu_{j2}$

Instead of computing average of fold-change estimates from pairs of arrays or by taking the ratio of average expression from each group, EB estimates of ρ_j are obtained by specifying a probability distribution on μ_{jk} and then computing some measure $\hat{\rho}_j$ (central tendency) of the posterior distribution of ρ_j (given the data).

"Significance" of Differential Expression

- Classical one-gene-at-a-time tests treats gene effects as fixed and separate.
- Here, some unknown fraction p of the genes are DE and the remainder are EE, and the state for each gene is considered to be a matter of chance.

EE gene j $\rightarrow \mathbf{x}_j = (x_{j1}, x_{j2}, \dots, x_{jI})$

$$f_0(\mathbf{x}_j) = \int \left(\prod_{i=1}^I f_{\text{obs}}(x_{ji} | \mu) \right) \pi(\mu) d\mu \quad (\text{predictive distribution})$$

(average over possible gene effects μ_j)

DE gene j $\rightarrow \mathbf{x}_j = (\mathbf{x}_{j1}, \mathbf{x}_{j2})$

$$f_1(\mathbf{x}_j) = f_0(\mathbf{x}_{j1}) f_0(\mathbf{x}_{j2})$$

(both continuous mixing over the unknown values of μ_j and discrete over two patterns (DE_j and EE_j) for each gene)

Marginal distribution of the data: $pf_1(\mathbf{x}_j) + (1-p)f_0(\mathbf{x}_j)$

Posterior probability of differential expression (Bayes' rule):

$$\frac{pf_1(\mathbf{x}_j)}{pf_1(\mathbf{x}_j) + (1-p)f_0(\mathbf{x}_j)}$$

→ Multiple Conditions

- 2 conditions → 2 patterns (EE_j and DE_j)
- 3 conditions → 5 possible patterns
 - (1) equivalent expression across 3 conditions
 - (3) altered expression in just one condition
 - (1) distinct expression in each condition
- 4 conditions → 15 different patterns
 - Patterns: grouping or clustering of conditions
 - Extra information might reduce the number of patterns to be considered

- (m + 1) distinct patterns for $\mathbf{x}_j = (x_{j1}, x_{j2}, \dots, x_{jm})$
- Here, the marginal distribution of the data becomes: $\sum_{k=0}^m p_k f_k(\mathbf{x}_j)$

mixing proportions
↓
- Posterior probability of expression pattern k: $\Pr(k | \mathbf{x}_j) \propto p_k f_k(\mathbf{x}_j)$

↑

product of $f_{0j}(\cdot)$ contributions of the different groups or conditions

- $\Pr(k | \mathbf{x}_j)$, $k = 0, 1, \dots, m$, can be used to identify genes with altered expression in at least one condition, to order genes within conditions, or to classify genes into distinct expression patterns.
- More generic concepts: pattern of increasing (decreasing) means over different conditions (e.g. time course experiment)

→ The Gamma-Gamma Model

- Choice of observation component and mean component, each depending on a few additional parameters θ (to be estimated from the data)

Observation component: Gamma distribution

(shape parameter $\alpha > 0$ and scale parameter $\lambda = \alpha/\mu_j \rightarrow$ mean value μ_j)

$$f_{\text{obs}}(x | \mu_j) = \frac{\lambda^\alpha x^{\alpha-1} \exp\{-\lambda x\}}{\Gamma(\alpha)}, \quad x > 0$$

Note: $CV = 1/\sqrt{\alpha}$ constant across genes j.

Mean component: $\pi(\mu_j) \sim$ Inverse Gamma distribution

$$\alpha \text{ fixed} \rightarrow \lambda = \alpha / \mu_j \sim \text{Gamma}(\alpha_0, \nu)$$

shape
scale

Parameters: $\theta = (\alpha, \alpha_0, \nu)$

$$f_0(x_1, x_2, \dots, x_I) = C \frac{\left(\prod_{i=1}^I x_i\right)^{\alpha-1}}{\left(\gamma + \sum_{i=1}^I x_i\right)^{I\alpha + \alpha_0}}, \text{ where } C = \frac{\gamma^{\alpha_0} \Gamma(I\alpha + \alpha_0)}{\Gamma^I(\alpha) \Gamma(\alpha_0)}$$

➔ The Lognormal-Normal Model

common variance

$$\log(x) \sim \text{Normal}(\mu_j, \sigma^2)$$

Note: $CV = \sqrt{\exp(\sigma^2) - 1}$ (raw scale)

$$\mu_j \sim \text{Normal}(\mu_0, \tau_0^2)$$

$$f_0(\mathbf{x}) \sim \text{MN}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_n)$$

$$\boldsymbol{\Sigma}_n = \sigma^2 \mathbf{I}_n + \tau_0^2 \mathbf{M}_n$$

$$\boldsymbol{\mu}_0 = \begin{bmatrix} \mu_0 \\ \mu_0 \\ \vdots \\ \mu_0 \end{bmatrix} \quad \mathbf{M}_n = \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix}$$

➔ Model Fitting

- Maximum (marginal) likelihood of θ

$$\begin{cases} \text{GG} : \theta = (\alpha, \alpha_0, \gamma) & \text{(EM Algorithm: any gene with at least one} \\ \text{LNN} : \theta = (\mu_0, \sigma^2, \tau_0^2) & \text{negative intensity value in one condition is} \\ & \text{omitted from this step of the calculation)} \end{cases}$$

- With data x_j governed by a mixture as $\sum_{k=0}^m p_k f_k(x_j)$, indicator variables (ϕ_{jl}) are defined as 1 if the expression pattern of gene j is pattern l , and 0 otherwise.
- The complete data log likelihood is:

$$l_c(\theta) = \sum_j \left\{ \sum_{k=0}^m \phi_{jk} [\log f_k(x_j) + \log(p_k)] \right\}$$

- E-step: Using current estimate θ_0 , replace ϕ_{jl} by $\hat{\phi}_{jl}$.
($\hat{\phi}_{jl}$: posterior probability of expression pattern l for gene j)
- M-step: Estimate p_k with the arithmetic mean $\hat{\phi}_{\cdot k}$.
- Initial values: $\alpha = (1/CV)^2$ (GG) and $\sigma^2 = \log(1 + CV^2)$ (LNN)

➔ Example:

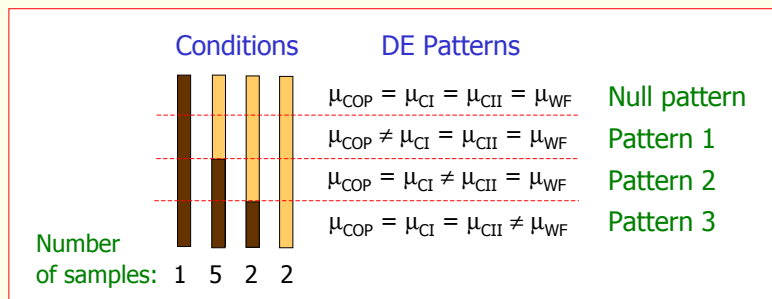
- Mammary epithelial cells in rat (model of breast cancer)
- Affymetrix U34 chip (26,379 genes)
- $I = 10$ mRNA samples
- Four inbred strains (2 parental + 2 derived congenic lines)

☞ Parental { - Copenhagen (COP): almost completely resistant to certain carcinogens
 - Wistar Furth (WF): highly susceptible

☞ Intermediate inbred lines { - CI and CII: homozygous for WF/WF or COP/COP in different (small) regions of the genome

After normalization/filtering:

$J = 25,248$ genes
 (1131 outliers)



➔ Discussion:

- Classifying genes into expression patterns by the posterior probability is an optimal procedure in the context of the mixture model: it minimizes the expected number of errors.
- Interestingly, this goal is different from the goal in classical testing, which is to bound the type I error rate and then aim to maximize the power (Newton and Kendzioriski, 2003).

BAYESIAN MODELS (Ibrahim et al., 2002)

➔ Data:

- Affymetrix Hu6800; 7,000+ probe sets (\approx 5,600 distinct genes)
- 14 individuals (samples): 4 normal + 10 endometrial cancer tissues
- Filtering/normalization: 3,214 genes (final data set)

<p>A: absent</p> <p>M: marginal</p> <p>P: present</p>	<p>A, M and negative values were set to $c_0 = 20$ (arbitrary value)</p>
<p>non expressed genes</p>	

Expression level $x = \begin{cases} c_0 & \text{with probability } p \\ c_0 + y & \text{with probability } (1 - p) \end{cases}$

➔ Model settings:

$$y_{jig} \sim \text{log normal} \begin{cases} i = 1, 2, \dots, n_j \text{ (sample)} \\ j = 1, 2 \text{ (tissue type)} \\ g = 1, 2, \dots, G \text{ (gene)} \end{cases}$$

$$p(y_{jig} | \mu_{jg}, \sigma_{jg}^2) = (2\pi)^{-1/2} y_{jig}^{-1} \sigma_{jg}^{-1} \exp\left\{-\frac{1}{2\sigma_{jg}^2} [\log(y_{jig}) - \mu_{jg}]^2\right\}$$

$$\delta_{jig} = \begin{cases} 1 & \text{if } x_{jig} = c_0 \\ 0 & \text{otherwise} \end{cases} \quad p_{jg} = \Pr(x_{jig} = c_0) \equiv \Pr(\delta_{jig} = 1)$$

$$L(\theta | \mathbf{x}, \delta) = \prod_{j=1}^2 \prod_{i=1}^{n_j} \prod_{g=1}^G p_{jg}^{\delta_{jig}} (1 - p_{jg})^{1 - \delta_{jig}} p(y_{jig} | \mu_{jg}, \sigma_{jg}^2)^{1 - \delta_{jig}}$$

• Differentially expressed genes:

$$\begin{aligned} \psi_{jg} &= E_{\delta,y}[c_0 \delta_{jig} + (1 - \delta_{jig})(c_0 + y_{jig}) | p_{jg}, \mu_{jg}, \sigma_{jg}^2] \\ &= c_0 p_{jg} + (1 - p_{jg})[c_0 + E(y_{jig} | \mu_{jg}, \sigma_{jg}^2)] \\ &= c_0 p_{jg} + (1 - p_{jg})[c_0 + \exp\{\mu_{jg} + \sigma_{jg}^2 / 2\}] \end{aligned}$$

- Posterior distribution of $\xi_g = \Psi_{2g} / \Psi_{1g}$ for comparing expression level means between normal and cancer tissues.

☞ posterior mean, standard deviation, quantiles, and probabilities such as $\Pr(\xi_g > 1 | \mathbf{x}, \delta)$.

- Priors: $\mu_{jg} | \sigma_{jg}^2, \mu_{j0} \sim N(\mu_{j0}, \tau_0 \sigma_{jg}^2 / \bar{n}_j)$
 $\mu_{j0} \sim N(m_{j0}, v_{j0}^2)$
 $\sigma_{jg}^2 \sim \text{Inv.Gamma}(a_{j0}, b_{j0})$
 - $\tau_0 > 0$
 - $\bar{n}_j = \frac{1}{G} \sum_{g=1}^G \left(n_j - \sum_{i=1}^{n_j} \delta_{jig} \right)$
 - a_{j0} fixed
 - $b_{j0} \sim \text{Gamma}(q_{j0}, t_{j0})$ (or fixed)

- Prior correlation between μ_{jg} :

$$\begin{bmatrix} \mu_{jg} \\ \mu_{jg'} \end{bmatrix} \sim N(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*) \rightarrow \boldsymbol{\mu}^* = \begin{bmatrix} m_{j0} \\ m_{j0} \end{bmatrix} \quad \boldsymbol{\Sigma}^* = \begin{bmatrix} \frac{\tau_0 \sigma_{jg}^2}{\bar{n}_j} + v_{j0}^2 & v_{j0}^2 \\ v_{j0}^2 & \frac{\tau_0 \sigma_{jg'}^2}{\bar{n}_j} + v_{j0}^2 \end{bmatrix}$$

$$\text{corr}(\mu_{jg}, \mu_{jg'} | \sigma_{jg}^2, \sigma_{jg'}^2, v_{j0}) = \frac{v_{j0}^2}{\left\{ \left[\frac{\tau_0 \sigma_{jg}^2}{\bar{n}_j} + v_{j0}^2 \right] \left[\frac{\tau_0 \sigma_{jg'}^2}{\bar{n}_j} + v_{j0}^2 \right] \right\}^{1/2}}$$

Note: If $v_{j0}^2 \rightarrow \infty$ or $\bar{n}_j \rightarrow \infty \Rightarrow \text{corr} \rightarrow 1$

$$e_{jg} = \text{logit}(p_{jg}) = \log\left(\frac{p_{jg}}{1-p_{jg}}\right) \sim N(u_{j0}, k_{j0} w_{j0}^2)$$

$$u_{j0} \sim N(\hat{u}_{j0}, h_{j0} w_{j0}^2)$$

(k_{j0}, w_{j0}, h_{j0} : specified hyperparameters)

- Data-based guide values for hyperparameters:

$$m_{j0} = \frac{1}{N_j} \sum_{i=1}^{n_j} \sum_{g=1}^G (1 - \delta_{jig}) \log(y_{jig}) \quad \text{where} \quad N_j = \sum_{i=1}^{n_j} \sum_{g=1}^G (1 - \delta_{jig})$$

(gene sample mean on the natural logarithm scale for tissue type j)

$$v_{j0}^2 = \eta_{j0} \text{MSG}_j \quad \text{where} \quad \text{MSG}_j = \frac{1}{G-1} \sum_{g=1}^G n_{jg} (m_{jg0} - m_{j0})^2$$

$$\text{and} \quad m_{jg0} = \frac{\sum_{i=1}^{n_j} (1 - \delta_{jig}) \log(y_{jig})}{\sum_{i=1}^{n_j} (1 - \delta_{jig})} \quad n_{jg} = n_j - \sum_{i=1}^{n_j} \delta_{jig}$$

$$t_{j0}^{-1} = d_{j0} \text{MSE}_j \quad \text{where} \quad \text{MSE}_j = \frac{1}{N_j - G} \sum_{g=1}^G \sum_{i=1}^{n_j} (1 - \delta_{jig}) (\log(y_{jig}) - m_{jg0})^2$$

- Gene selection:

- ① Compute the posterior distributions of all ξ_g and $\gamma_g = \Pr(\xi_g > 1/D)$.
- ② Select a threshold value γ_0 (possible values $\gamma_0 = .7, .8, .9, \text{ and } .95$).
- ③ Once a set of genes is declared different, set the mean parameters for the tissue types to be unequal for that gene in the model.
- ④ Create several submodels using different values of γ_0 and using the step 3 and evaluate them (Bayesian criterion).
- ⑤ Select the best-fitting model.

REFERENCES

- Ibrahim, J. G., Chen, M. H., Gray, R. J. (2002) Bayesian models for gene expression with DNA microarray data. *Journal of the American Statistical Association* 97: 88-99.
- Newton, M. A. et al. (2001) On differential variability of expression ratios: Improving statistical inference about gene expression changes. *Journal of Computational Biology* 8: 37-52.
- Newton, M. A, Kendziorski, C. (2003) Parametric Empirical Bayes Methods for Microarrays. In Parmigiani, G; E. S. Garrett; R. A. Irizarry and S. L. Zeger (Eds). *The Analysis of Gene Expression Data: Methods and Software*. Springer, New York. p. 254-271.
- Pan, W., Lin, J., Lee, C. T. (2002) How many replicates of arrays are required to detect gene expression changes in microarray experiments? A mixture model approach. *Genome Biology* 3(5):research0022.1-0022.10.