

Discover Common and Differential Enrichment: A Multivariate Bayesian Variable Selection Approach

Sierra M. Li and Giovanni Parmigiani

Division of Oncology Biostatistics
The Sidney Kimmel Cancer Center at Johns Hopkins

November 7, 2007

Model and notation

- Model

- Hierarchical model with cross-profile shrinkage
- Enrichment types: latent variable γ
 - no enrichment ($\gamma = 0$)
 - common enrichment: similar magnitude across profiles ($\gamma = 1$)
 - differential enrichment: different magnitude across profiles ($\gamma = 2$)

- Notation

- Gene expression: $\mathbf{Y}_{n \times q}$ response matrix
- gene-set: $\mathbf{X}_{n \times p}$ design matrix
- Enrichment measure: $\mathbf{B}_{p \times q} = (\beta_{ij})$ regression coefficients
 - β_{ij} : enrichment for gene-set i in profile j
 - β_{ij} measures the change of average expression in a gene-set
- $\boldsymbol{\mu}$: p -vector for average enrichment across profiles

Conjugate model and parameters

- Gene level:

$$Y - \mathbf{1}\alpha' - \mathbf{X}B = \Omega \sim \mathcal{N}(\mathbf{I}_n, \sigma^2 \mathbf{I}_p)$$

- Gene-set level:

$$\alpha' - \alpha'_0 \sim \mathcal{N}(h, \sigma^2 \mathbf{I}_p)$$

$$B - \mu_{p \times 1} \mathbf{1}'_{q \times 1} \sim \mathcal{N}(\mathbf{H}\gamma, \sigma^2 \mathbf{I}_p)$$

- Gene-set across profiles:

$$\mu' - \mathbf{0}' \sim \mathcal{N}(\sigma^2, \mathbf{G}\gamma)$$

- Matrices related to variable selection:

$\mathbf{H}\gamma = \mathbf{D}\gamma \mathbf{R} \mathbf{D}\gamma$ and $\mathbf{G}\gamma = \mathbf{F}\gamma \mathbf{R} \mathbf{F}\gamma$, \mathbf{R} : correlation \mathbf{D}, \mathbf{F} : diagonal matrices

$$d_i^2 = \begin{cases} \tau_{i0}^2 & \text{if } \gamma_i = 0 \text{ or } 1 \\ \tau_{i1}^2 & \text{if } \gamma_i = 2 \end{cases} \quad f_i^2 = \begin{cases} \nu_{i0}^2 & \text{if } \gamma_i = 0 \\ \nu_{i1}^2 & \text{if } \gamma_i = 1 \text{ or } 2 \end{cases}, \text{ for } i = 1, 2, \dots, p.$$

$\tau_{i0}^2 \ll \tau_{i1}^2$ and $\nu_{i0}^2 \ll \nu_{i1}^2$, trans-dimensional setting with $\tau_{i0}^2 = \nu_{i0}^2 = 0$

Prior and Posterior

- Prior:

- $\pi(\boldsymbol{\alpha}, \mathbf{B}, \boldsymbol{\gamma}, \sigma^2) = \pi(\boldsymbol{\alpha}|\boldsymbol{\gamma}, \sigma^2)\pi(\mathbf{B}|\boldsymbol{\gamma}, \sigma^2)\pi(\sigma^2)\pi(\boldsymbol{\gamma})$
- $\pi(\mathbf{B}|\boldsymbol{\gamma}, \sigma^2) = \pi(\mathbf{B}|\boldsymbol{\mu}, \boldsymbol{\gamma}, \sigma^2)\pi(\boldsymbol{\mu}|\boldsymbol{\gamma}, \sigma^2)$
- $\sigma^2 \sim IG(a, b)$
- $\boldsymbol{\gamma}$ product of independent multinomial

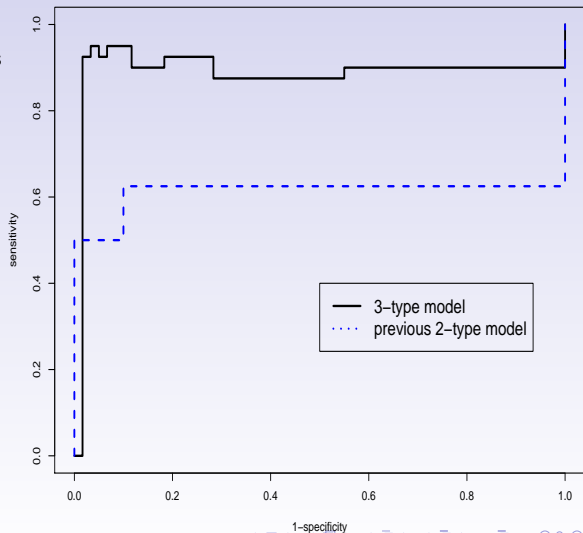
- Posterior:

- marginal posterior $\pi(\boldsymbol{\gamma}|\mathbf{Y})$
- $\pi(\boldsymbol{\gamma}|\mathbf{Y})$ involves the sum of the product of residual matrices on \mathbf{B} and $\boldsymbol{\mu}$ levels with different shrinkage

- Comment : the model framework is very general, the variable selection approach can be applied to both [gene-disease association](#) and [gene-set enrichment](#) studies

Simulation : from the model with cross-profile signal

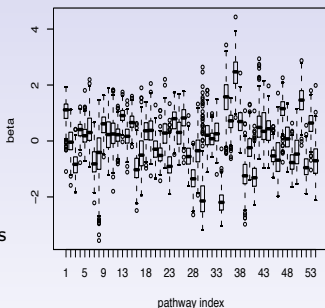
- $\mathbf{X}_{104 \times 100}$:
 - 104 yeast recombinant strains
 - randomly selection of 100 genes
 - from a real study of yeast growth under 92 different drugs
- 40 genes are in the true model:
 - 20 genes with similar association across 92 drugs
 - 20 genes with different association across 92 drugs
- generate $\mu, \mathbf{B}|\mu$ with specified gene types and large signal
- simulation from the model with no cross-profile signal: 3-type model and previous 2-type model have similar performance



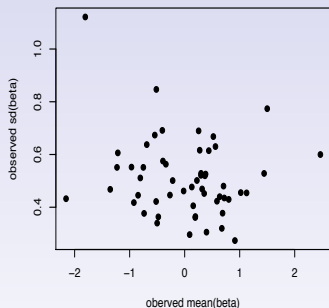
Case study

- Breast cancer by Miller et.al.
hgu133A chip
rma normalization
 - 22283 genes
251 samples
190 pathways
 - top 5% variable genes
1115 genes
- 54 sets $>$ 5 genes

ridge reg coef: real data



ridge reg coef: real data



- Non-conjugate model for signal and noise relationship in enrichment study
- A short Markov chain found no gene-set with high posterior probability close to 1
- Top differential enriched sets appear to be reasonable: pathway related to stage of the tumor and P53
- Top common enriched sets appear to be mostly in general biological processes
- The mean of a gene-set a conservative measure for enrichment?

[1] "Olfactory transduction"

[2] "Glioma"

[3] "Nitrogen metabolism"

[4] "Long-term potentiation"

[5] "Thyroid cancer"

[6] "Cell cycle"

[7] "Toll-like receptor signaling pathway"

[8] "*Neurodegenerative Disorders*"

[9] "Cell Communication"

[10] "Melanoma"

[11] "Pancreatic cancer"

[12] "Epithelial cell signaling in

Helicobacter pylori infection"

[13] "Antigen processing and

presentation"

[14] "Wnt signaling pathway"

[15] "Apoptosis"

[16] "Fatty acid metabolism"

[17] "Glutamate metabolism"

[18] "Regulation of actin cytoskeleton"

[19] "Gap junction"

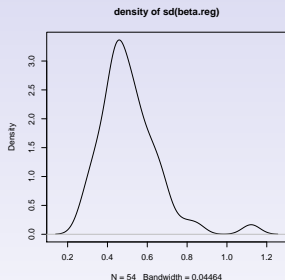
[20] "Natural killer cell mediated

cytotoxicity"

[21] "Hedgehog signaling pathway"

[22] "Adipocytokine signaling pathway"

[23] "Jak-STAT signaling pathway"



26 pathways with $sd > 0.5$
7 out of 17 p53 pathways

- Summary
 - The model framework is general
 - Marginal posterior of the latent variable type can be obtained
 - Conjugate or non-conjugate setting is flexible
 - Variable selection can be applied to both gene and gene-set levels

Thanks!