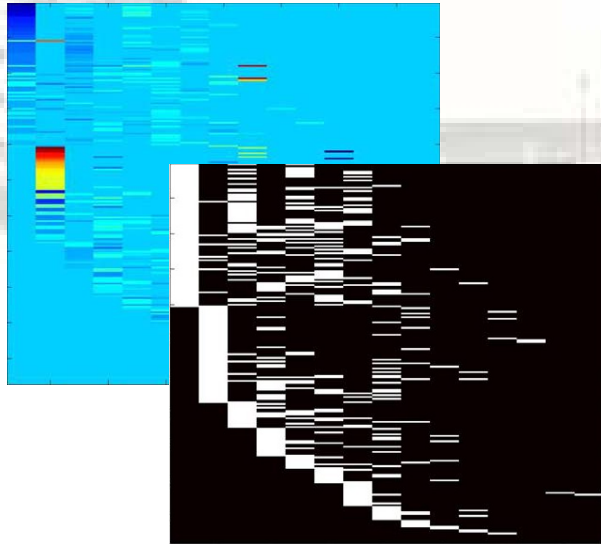# Aspects of Statistical Modelling
# & Data Analysis
# in Gene Expression Genomics

Mike West
Duke University

**#1**  Genomics, Microarrays, Data:
Big picture

**#2**  Bayesics - Regression and Shrinkage:
Gene expression as predictors

**#3**  Patterns and Factors:
Prediction via pattern profiling

**#4**  Sparse Modelling:
Regression subset-structure uncertainty

**#5**  Sparse Models and Profiling:
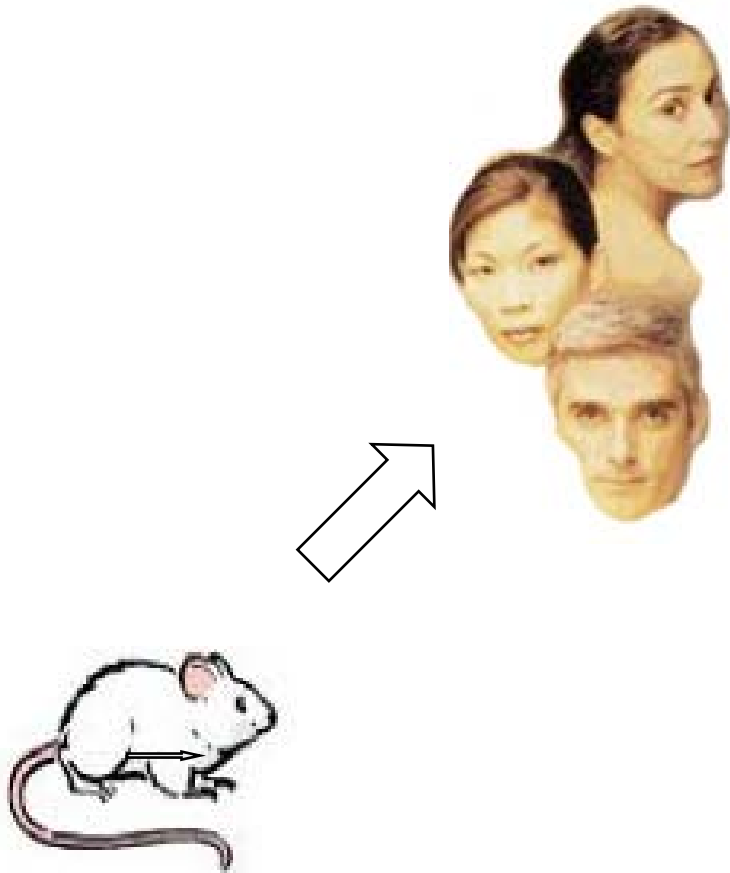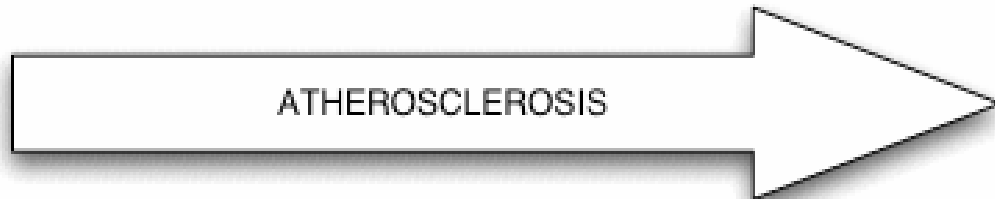Gene expression as response: Designed experiments

**#6**  Sparse Models and Profiling:
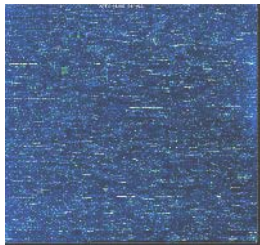Gene expression as response: Latent factor models

Gene expression profiles:

Signatures of states

ATHEROSCLEROSIS
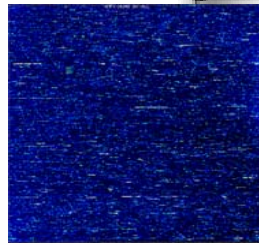
Apo E -/-, 6 wk Chow Diet

Apo E -/-, 12 wk Western Diet

PRECLINICAL DISEASE

EARLY/INTERMEDIATE DISEASE

Age, Diet, Gender, WT/ApoE

$2^4$x5+ balanced factorial design

## Human aorta gene expression:

Gene/pathway identification
Disease state prediction

## Mice models:

Disease state characterisation
Mouse ➡ Human mapping

Large-scale multifactoral design
Gene expression (aorta) response
Action is interactions

12,500 genes in parallel

$X = \beta$ WT, 6wk, chow, fem (baseline)

$+ \mu$ male

$+ \delta$ fat diet

$+ \alpha$ age=12wk/old

$+ \gamma$ ApoE genotype

$+ \mu\delta$ fat diet & male

$+ \mu\alpha$ 12wk/old & male

$+ \mu\gamma$ ApoE & male

$+ \delta\alpha, \delta\gamma, \alpha\gamma$

$+ \mu\delta\alpha, \mu\delta\gamma, \mu\alpha\gamma, \delta\alpha\gamma, \mu\delta\alpha\gamma$

+ noise

Multiplicities:  Full multivariate analysis – simultaneous inference: identify "real" effects

Bayesian multivariate Anova: Multiple shrinkage estimation

A precursor experiment: Beware the flood (RNAi, etc)

Analysis goals:    Gene identification, pattern profiling of states/interactions of design variables

Translation to human samples: Predict expression signatures
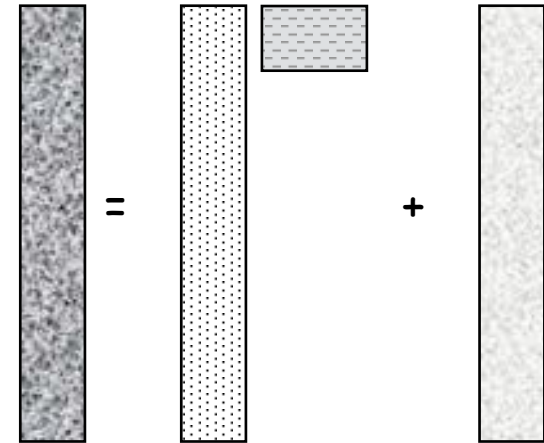
$$X = BH + N$$

Genes by samples $X_{(p \times n)}$

Gene-specific regression parameters $B_{(p \times k)} = \{\beta_{g,j}\}$

Fixed design $H_{(k \times n)}$



Design factor j:   Main effect, interaction, …   $\beta_{g,j}$

Multiple tests/comparisons – simultaneous inference

Substantial noise component

Sparsity priors:  $\#\{\beta_{g,j} \neq 0\} = $  small

$$\beta_{g,j} \sim (1 - \pi_{g,j})\delta_0(\beta_{g,j}) + \pi_{g,j}N(\beta_{g,j}|0, \tau_j)$$

Extends traditional sparsity priors:

Gene g, Design factor j:  $\pi_{g,j} \sim$  sparsity

Invites informative, structural modelling

Hierarchical models/priors within design factors

Model-based, automatic shrinkage – Simultaneous "multiple tests"

Decision theory/false discovery?
Estimation versus Decision?
How many comparisons/hypotheses?

$$\pi_{g,j}^* = Pr(\beta_{g,j} \neq 0 | X)$$

$$E(\beta_{g,j} | \beta_{g,j} \neq 0, \ X)$$

Computation and

posterior summarisation

MCMC methods

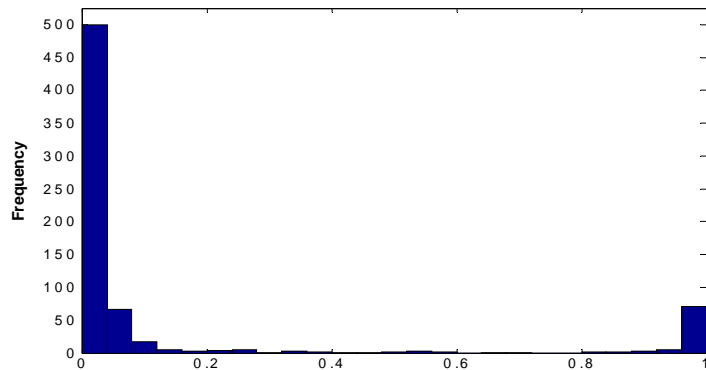New hierarchical model: $$\beta_{g,j} \sim (1 - \pi_{g,j})\delta_0(\beta_{g,j}) + \pi_{g,j}N(\beta_{g,j}|0,\tau_j)$$

$$\pi_{g,j} \sim (1-\rho_j)\delta_0(\pi_{g,j}) + \rho_j Be(\pi_{g,j}|s_j r_j, s_j(1-r_j))$$

$$p(\rho_j, \tau_j, r_j)$$



$$\pi^*_{g,j} = Pr(\beta_{g,j} \neq 0 | X)$$



Design factor j:   $\rho_j$  ~ sparsity

It matters: improved signal isolation

MCMC methods vital

(Lucas et al 05)

Regressions example:
"housekeeping gene control factors"

Gene-Sample-Study artifacts

Oncogene interventions experiments



(Lucas et al 05)

**Major but Sparse effects:**
**Selective impact across genes**



o  Gene expression data

x  Fitted component on
   control factor 1

(Lucas et al 05)

No control factors



With control factors
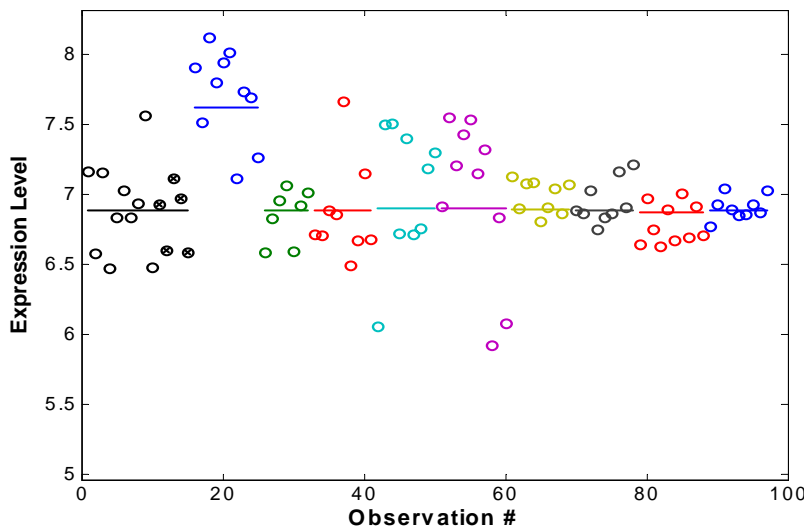
Oncogene experiments - Myc up-regulated gene:

NFE2L1

Nuclear factor (erythroid-derived 2)-like 1
Involved in the regulation of apoptosis

Directly induced by Myc

consensus binding sequence
(T/C)GCGCA(C/T)GCGC(A/G)

Myc binding site

GENES & DEVELOPMENT (2003-01-15)

(Lucas et al 05)

Sparse ANOVA for noise variances

$$X = BH + N$$

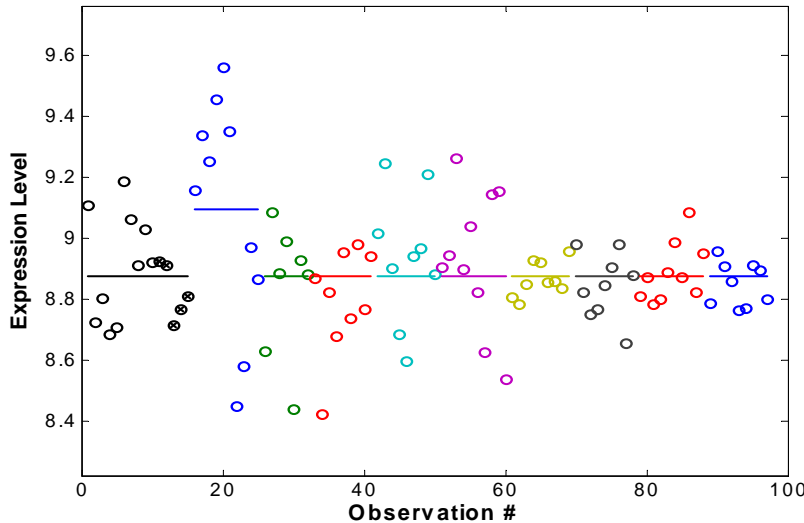$$N = \{\nu_{g,i}\}, \quad \nu_{g,i} \sim N(0, \exp(\kappa_{g,i}))$$

$$K_{(p \times n)} = \{\kappa_{g,i}\}$$

$$K = \Theta H + N_K$$

Sparsity prior for $\Theta_{(p \times n)} = \{\theta_{g,j}\}$

Oncogene interventions experiments: False discovery
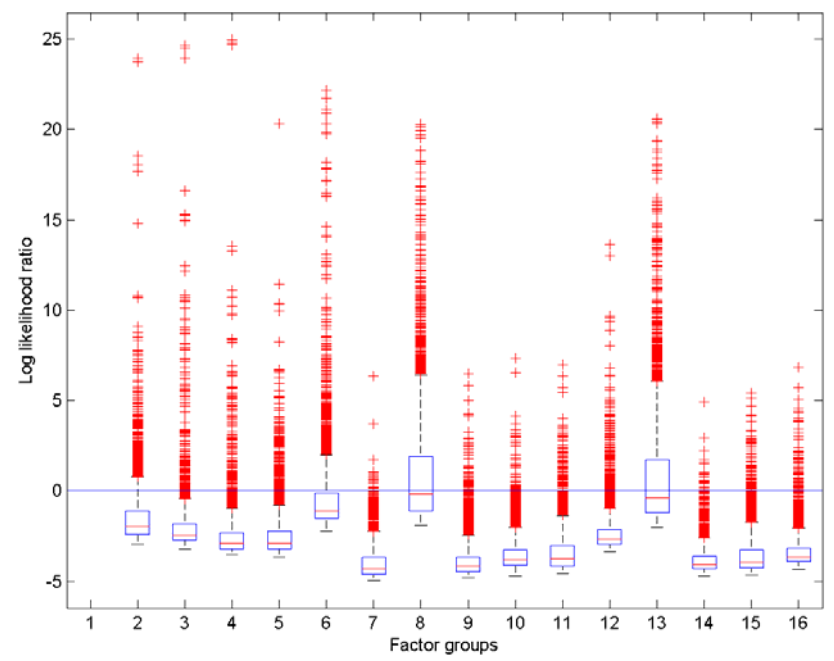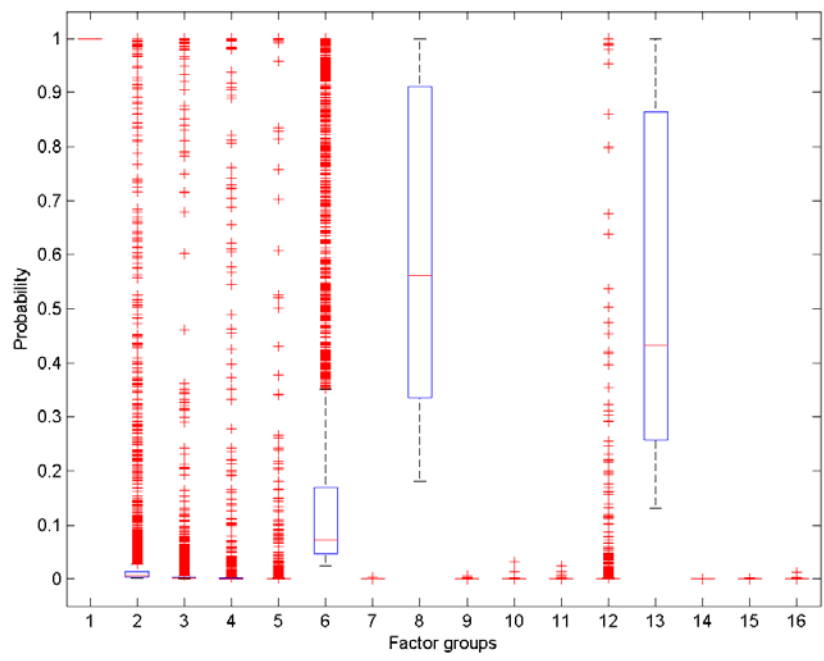
MCMC methods vital

(Lucas et al 05)

Probabilities and log-likelihood ratios

- SHRINKAGE

$$\pi_{g,j}^* = Pr(\beta_{g,j} \neq 0 | X)$$

'significant' effects define
characterising metagenes



ApoE.Age

65 genes

ApoE.Diet

321 genes

ApoE.Age.Diet

292 genes

## Mapping genes from mouse to human
### normalise/evaluate/extrapolate signatures
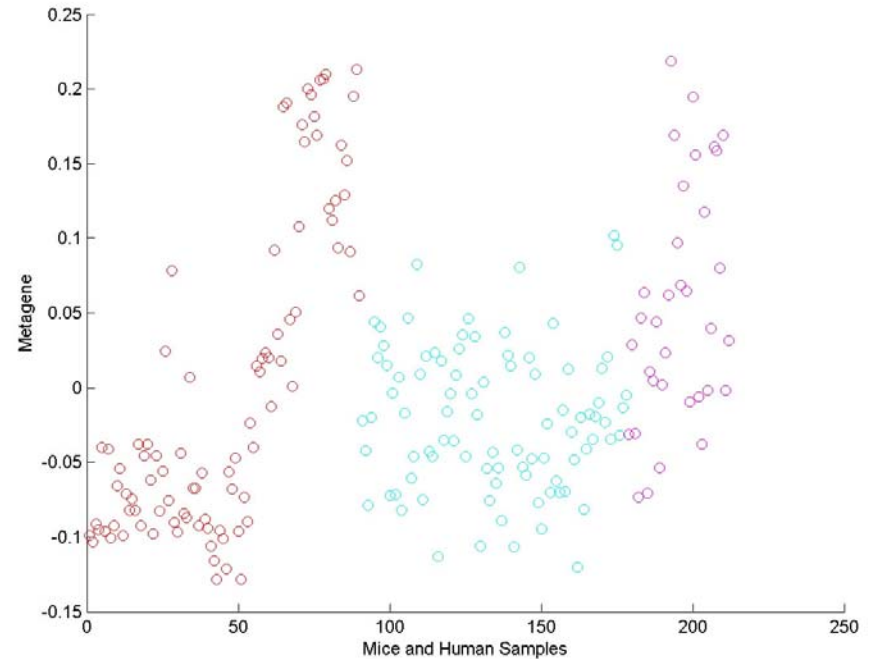


Human D+/-

Aorta lesions, advanced atherosclerosis

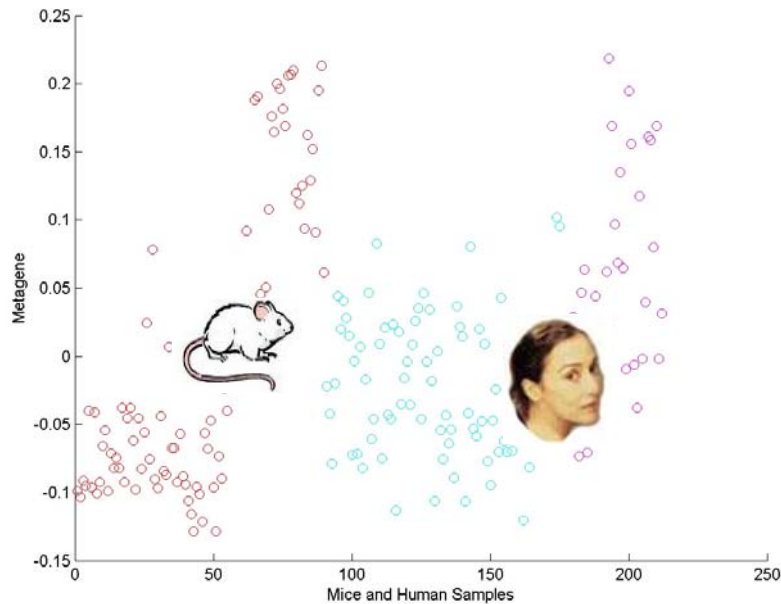Mouse ApoE-

Mouse WT/ApoE+

**ApoE.Age**



**Intersection of ApoE, ApoE+{Age,Diet}**

Improved pathway characterisation:
implicated genes

**Disease Risk Signature**



More/better phenotypes in humans

Blood/serum assays for gene expression

Metabolites in serum and blood:
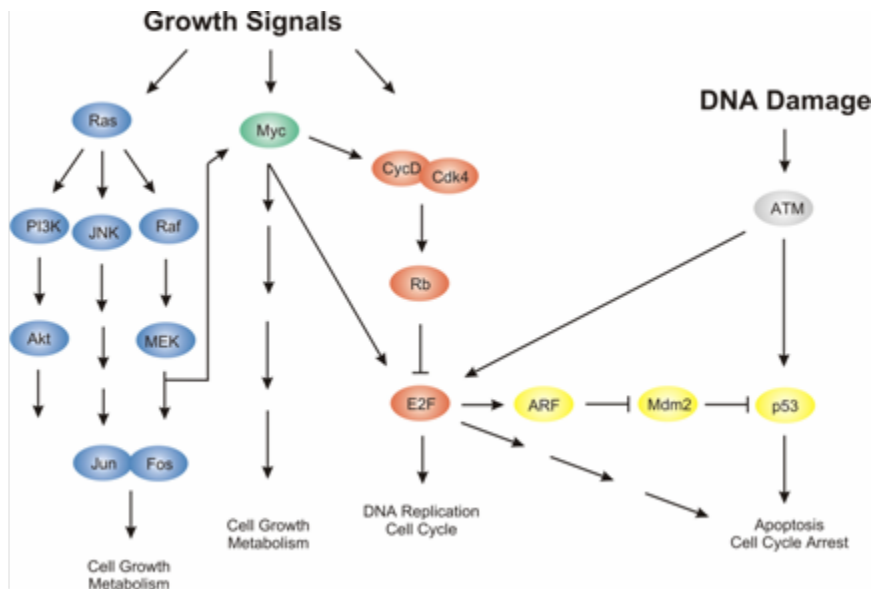metabolic/genomic fusion

(Karra et al 04; Seo et al 05)

Decompositions of p(x)

Latent structure underlying associations

Cancer Studies:

Multiple deregulated pathway components

Latent Factors:

intersecting sub-pathways

**One sample**

- column p-vector

$$x = B\lambda + \nu \qquad \nu \sim N(0, \Psi)$$

**Vector of k<<p latent**

- underlying –

factor variables

**Idiosyncratic variation**

**Latent factors:**

$$\lambda \sim N(0, T)$$

**Model of covariance matrix:**

$$V(x) = BTB' + \Psi$$

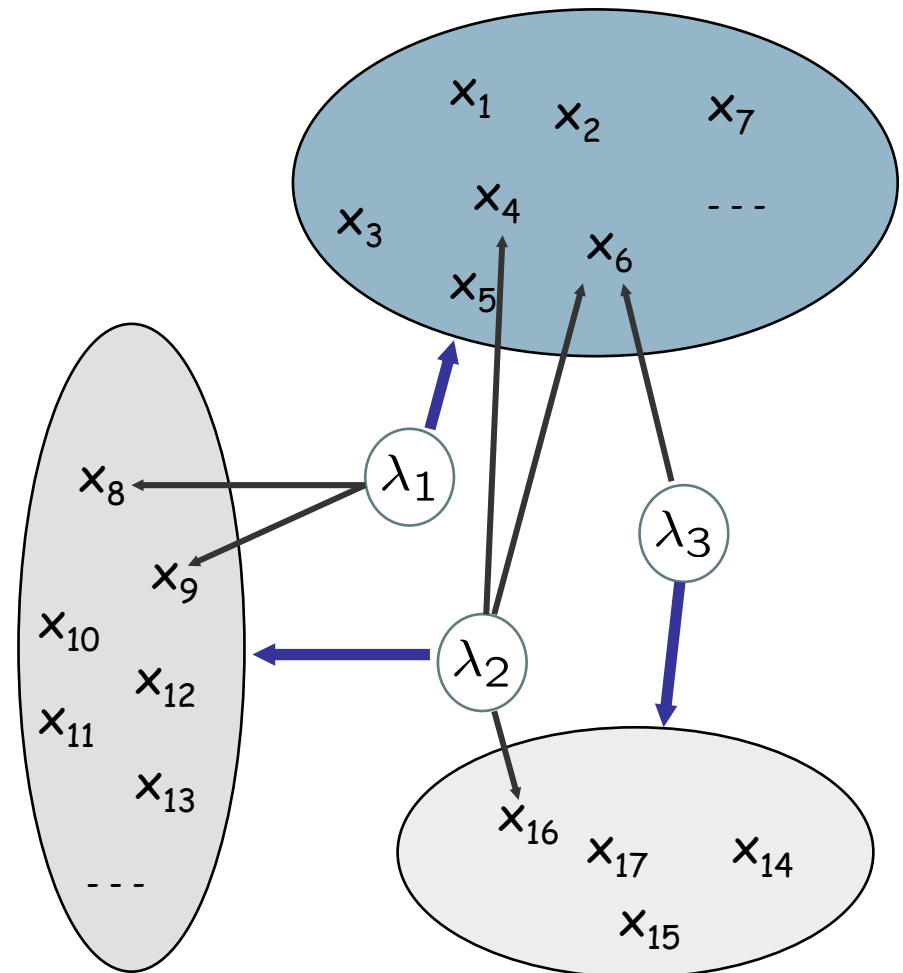(West 03, Valencia 7; +Aguilar 01 JBES; +Lopes 04 Stat Sinica)

# Sparse Models:

One factor – few or many variables

One variable – 0,1, or few factors

$$B = \{\beta_{g,j}\}$$

Row (variable) g, factor j:

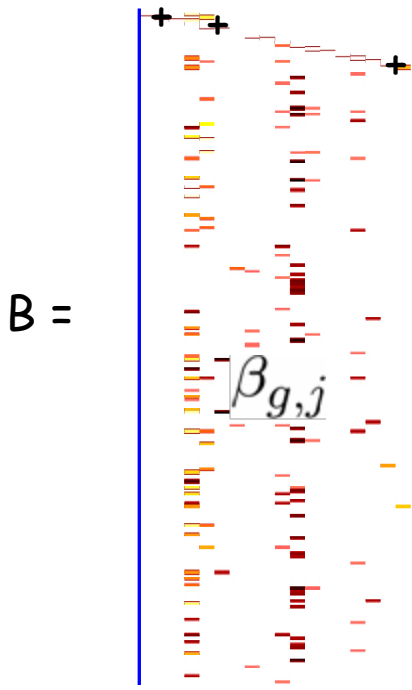$$\#\{\beta_{g,j} \neq 0\} = \quad 0,1, .., \text{small}$$

(West 2003, Valencia 7)

## Uncertain sparsity patterns:

$$\beta_{g,j} \sim (1 - \pi_{g,j})\delta_0(\beta_{g,j}) + \pi_{g,j}N(\beta_{g,j}|0, \tau_j)$$

$$\pi_{g,j} \sim (1-\rho_j)\delta_0(\pi_{g,j}) + \rho_j Be(\pi_{g,j}|s_j r_j, s_j(1-r_j))$$

$$p(\rho_j, \tau_j, r_j)$$

B =

$\beta_{g,j}$

### Structure:

Parametrisation of **B** - Identification

"founders" of factors

(West 2003, Valencia 7)

$$X = B\Lambda + N$$

$$B = [b_1, \ldots, b_k], \quad \Lambda = [\lambda_1, \ldots, \lambda_n]$$

$$\Pi = \{\pi_{g,j}\}$$

$$p(b_j | \text{else}) = \prod_{g=1}^{p} p(\beta_{g,j} | \text{else})$$

Other (hyper) parameters "easy"

$$p(\Pi | \text{else}) = \prod_{g=1,j=1}^{p,k} p(\pi_{g,j} | \text{else})$$

Not parallelizable:

Serial Gibbs/MCMC

$$p(\Lambda | \text{else}) = \prod_{i=1}^{n} p(\lambda_i | \text{else})$$

Parallel within Gibbs iterates

Aspects of p(B,Π|X)

Sample correlations

Fitted correlations
in B'TB+Ψ

Covariance decompositions:

$$BTB' = t_1 b_1 b_1' + t_2 b_2 b_2' + \cdots$$

F-decomp of H-GATA3

F-decomp of TFF1

F-decomp of HNF3a

Data decompositions:

$$x_g = \beta_{g,1}\lambda_1 + \beta_{g,2}\lambda_2 + \cdots$$

Breast cancer gene expression
Factor 2: ER, 9: HER2

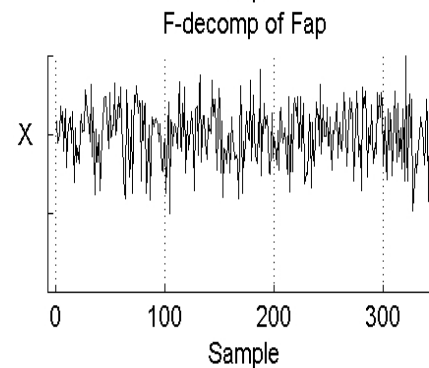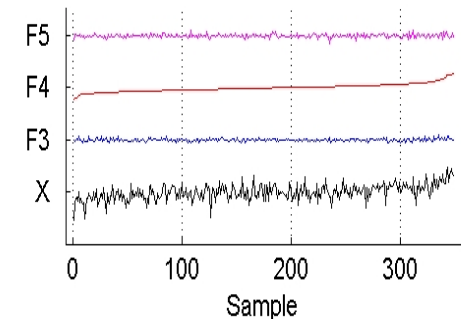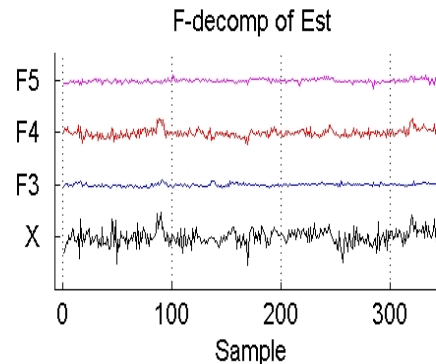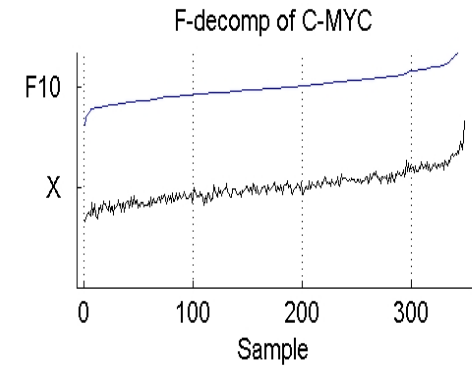Discovery of gene-factor associations

Isolation of non –associated genes

... noise/idiosyncratic

Factors
=
Patterns,
Signatures

Response variables z

Evaluate p(z,x)

Predict z

z linked to some latent factors in x space?

… and to some individual x variables?

Response factors $\lambda_z$

(West 2003, Valencia 7; Carvalho et al, 2005)

SVD

Sparse Model



"bad samples" factors

A "cleaned-up" ER factor

Factor models "clean up" SVD

SVD "noisy" factors

(West 2003, Valencia 7)

ER

HER2

Expression assays of hormonal status

Jointly predicted: Multiple factors

4+ studies

Jointly predicted:

Multiple tests/assays

E2F Pathway:

40 initial genes

· Replication
· Proliferation
· Cell cycle:check point
· Apoptosis
· "Structure"

Factor model fit to breast cancer

Thresholded gene-factor probabilities & effects

Evolutionary analysis: Add/delete genes, factors

Self-organising factor analysis

"Greedy" Stochastic Search methods – MCMC cousins

(Carvalho et al 05; Hans et al 05)

Evolutionary gene set enrichment
Experiments: gene set responses
to perturbations?
Promotor TF interactions

Duke/NCI Systems Biology Center

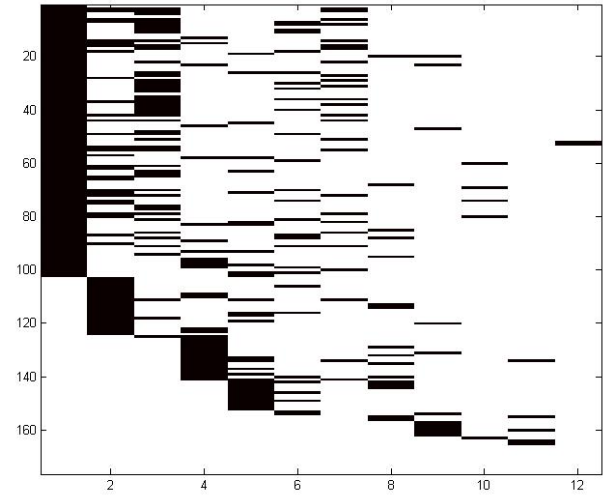| 2 (apoptosis) | 4 (DNA replication) |
|---|---|
| MDM2 | MCM6 |
| GSN | VCAM1 |
| AGC1 | CYP2A13 |
| RBM8A | PITX1 |
| MTHFS | MCM2 |
| CUGBP2 | DDX39 |
| PKN2 | MCM7 |
| SSR4 | GSN |
| FANCG | GSTM1 |
| NME3 | CCNE1 |
| POU4F1 | MCM3 |
| AGC1 | MCM4 |
| MDM2 | MFGE8 |
| CRHR1 | ABCA3 |
| H1FX | CDKN2A |
| RPS3A | MCM5 |
| ABCB8 | KIAA1026 |
| RGS12 | TOMM70A |
| GLG1 | CDC2 |
| DOC-1R | SAS |
| TNPO3 | POLR2H |
| MDM2 | CSNK1D |
| MAPT | NME3 |
| LOR | CDC6 |
| GUCA1A | CHC1 |
| GRIA1 | TNFRSF14 |
| CDC34 | BACH |
| COL11A2 | MVD |
| MYC | FANCG |
| TBL3 | LIG1 |
| BTF3 | SF3B4 |
| UCP3 | CCNE1 |
| LBA1 | ABCF3 |
| CDKN2C | |
| HTR6 | |
| CDC6 | |
| CYP2A13 | |
| KHDRBS1 | |
| KIAA0284 | |
| PEX5 | |
| CYP2A6 | |
| LTK | |
| SSTR3 | |
| MDM2 | |

Factors

Genes

HER2

ER

Myc.a

Myc.a

Myc.b

Myc.b

bcta-cat

E2F1

E2F1

E2F2

E2F2

E2F4

Biologically "Hard" vs "Soft" interactions

TF binding sites

P-P interactions

Expression expts

Priors over sparsity probabilities
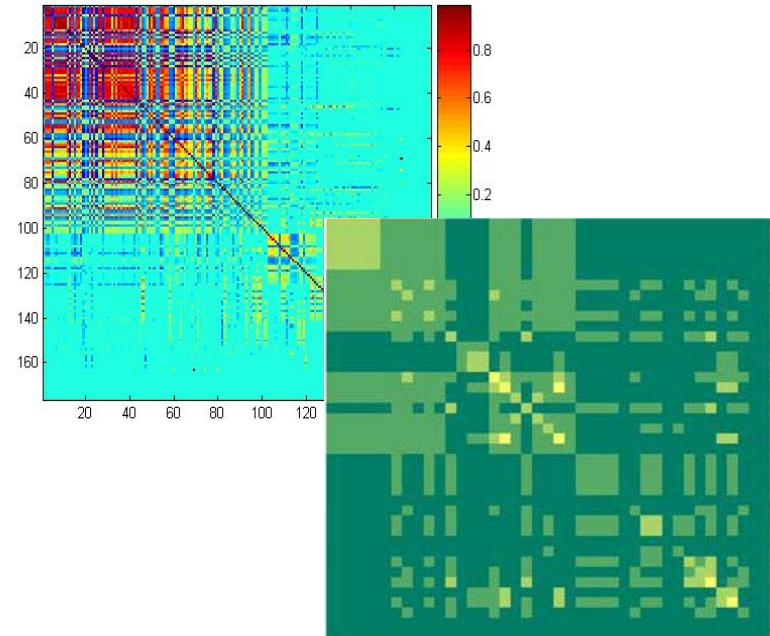
Status, relative "activation", deregulation

Priors over non-zero effects

## Graphs of sparse factor models – covariance or association graphs

### Graphical models – precision (inverse covariance) graphs





( Dobra et al 04, J Multivariate Analysis;

Jones et al 05, Statistical Science;

Rich et al 05, Cancer Research )

Multiple cancer data sets – "Meta-analysis"

Constant/consonant structure

- network "motifs"

(Zhou et al USC/PNAS 05)

## Sparsity models for:

- Design effects
- Control/correction factor effects
- Latent structure: Complex associations
- Response prediction

## Combined model
(Carvalho et al 05; Lucas et al 05)

Computation, Model Search:        MCMC and MCMC-inspired Stochastic Search

These slides:
www.isds.duke.edu/~mw/downloads/SemStat05

Papers, software, many links:
www.isds.duke.edu/~mw

ABS04 web site: Lecture slides, stats notes, papers, data, links:
www.isds.duke.edu/~mw/ABS04

Integrated Cancer Biology Program
icbp.genome.duke.edu

Genome Institute @ Duke
www.genome.duke.edu

Chris Hans

Joe Nevins

Carlos Carvalho

Beatrix Jones

Quanli Wang

Joe Lucas

**www.isds.duke.edu/~mw**