# **Expression Data Exploration:**

Association, Patterns, Factors & Regression Modelling



#### Exploring gene expression data

Scale factors, *median chip* correlation on gene subsets for crude 'data quality' investigation

Selecting genes correlated with an outcome, or a gene

Image displays: Expression intensity images

Common structure in gene subsets - Underlying Factors

Principal components/Singular value decomposition

Data sets: exploration in Matlab Some data, papers linked to <u>MW info@AB52004 web site</u> More data in papers on <u>Duke genomics web site</u>

#### Some data sets linked to <u>MW info@AB52004 web site</u>

- 2001 PNAS paper: Duke breast cancer genomics (49 tumours)
- 2003 Lancet, 2004 PNAS papers: Duke-KFSYS breast cancer clinico-genomics (158 tumours, clinical data)
- 2003 Nature Genetics paper: Mice cell lines and tumours in controlled experiments on several key oncogenes
- MIT/Whitehead 1999 Leukemia data

#### **Regression models**

#### Gene expression as response

Designed experiments: e.g. Mice models: age, sex, diet, genotype

- Finding genes, effects of covariates
- Grouping genes by effects
- Patterns of coordinately expressed genes
- Signature of effect multiple genes

Predictive interest: Human disease

p genes: Multiple models in parallel



## **Regression models**

#### Gene expression as covariates (predictors)

**Molecular phenotyping:** e.g. Predict aggressive vs. benign tumour Disease susceptible vs. resistance

- Finding genes linked to response
- Grouping genes
- Patterns of coordinately expressed genes
- Signature of effect multiple genes

Predictive interest: Human disease

Multiple genes as predictors



# **Empirical Factors:**

# Principal Components ~ Singular Value Decompositions

#### SVD of Expression Data Matrices



"Tall & Skinny" - p>>n

Gene-gene associations: sample variances/covariances - centered data (row means 0)

 $(n-1)^{-1}$ **XX**'

Collinearities Co-dependencies Co-regulation

Transcription factors Cascading expression Pathways, interactions

#### SVD of Data Expression Matrices



A pxn orthonormal columns:  $A'A=I_n$ D nxn diagonal, non-negative F nxn orthogonal: F'F=FF'=I\_n

Principal Components (PCA):  $XX' = AD^2A'$ 

## SVD of Expression Matrices



F - factors: rows are *n* factors columns are *n* samples

A – loadings: A<sub>ij</sub> loads gene i on factor j

D - singular values relative importance of factors

Patterns of covariation observed among genes is "driven" by underlying factors - mediated by loadings -

## **Exploring Expression Data**

Data sets: exploration in Matlab

Gene subset selection SVD of smaller subsets "of interest" Loadings to order genes

Links to clustering - group/cluster by loadings

Scatter plots, Image displays

Signatures: Dominant pattern - Metagene pattern

# **Elements of Regression**

#### Regression models : ideas & theory (see notes)

Linear model

 $\mathbf{y} = \mathbf{H}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \qquad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{I} / \boldsymbol{\phi})$ 

Least Squares

 $\mathbf{V}^{-1} = \mathbf{H'H}$  $\hat{\boldsymbol{\beta}} = \mathbf{VH'y}$ 

(minimal) Bayes: Shrinkage priors

Decision theory Regularisation - Ridge regression Key with many predictors Relevance of zero-mean location

Prior  $\boldsymbol{\beta} \mid \boldsymbol{\phi} \sim N(\mathbf{0}, \mathbf{C}^{-1} / \boldsymbol{\phi})$ 

Posterior  $\boldsymbol{\beta} \mid \mathbf{y}, \boldsymbol{\phi} \sim N(\mathbf{b}, \mathbf{B}^{-1} / \boldsymbol{\phi})$ 

Intercept term: 'Large'  $C_{1,1}^{-1}$ other  $C_{i,i}^{-1} = \tau$  (covariates scaled) Common shrinkage towards 0 - 'controlled' by  $\tau$   $\mathbf{B} = \mathbf{C} + \mathbf{H'H}$  $\mathbf{b} = \mathbf{B}^{-1}\mathbf{H'y}$ 

#### **Regression models : full theory**

Gamma prior/posterior on precision Normal posterior for  $\beta$  becomes T

Prior  $\phi \sim Ga(a/2, b/2)$ 

Posterior

 $\phi \mid \mathbf{y} \sim Ga(a'/2, b'/2)$ a'=a+n  $b'=b+\mathbf{y}'(\mathbf{y}-\mathbf{Hb})$   $\boldsymbol{\beta} \mid \mathbf{y}, \boldsymbol{\phi} \sim N(\mathbf{b}, \mathbf{B}^{-1} / \boldsymbol{\phi})$  $\boldsymbol{\beta} \mid \mathbf{y} \sim T_{a+n}(\mathbf{b}, \mathbf{B}^{-1} s)$ 

 $p(\mathbf{y})$ 

Residual variance estimate s = (b+q)/(a+n)

Marginal likelihood (to assess shrinkage) Prediction of new samples:  $p(y_{new}|y)$ 

#### Bayesian/Shrinkage elements:

- LSE as limiting case (no shrinkage)  $\tau \rightarrow \infty$
- Shrinks when it matters weak/no association  $\tau \rightarrow 0$
- Marginal likelihood to assess shrinkage degree
- Theoretical dominance of shrinkage estimation
- Regularisation: acts against over-fitting, improves stability and robustness in prediction

#### Examples (see Matlab code explorations)

- Explore genes predictive of others collinear expression patterns, gene co-regulation (e.g., ER in breast data)
- Design effects: finding genes related to genetic interventions, environmental factors (e.g., mice cell lines and tumours, time course experiments)

## Simple Designed Experiment Example

Gene pathway studies -

Myc & Ras oncogenes E2F transcription factors

(Nat Gene 2003, Huang et al)

Genes and gene subsets (signatures) related to up-regulation of oncogenes.

Compare & characterise oncogenic states

Why?



#### Simple Myc/Ras/E2F Example

y =	۵	Wild type mice (baseline
	+m	if Myc up-expressed
	+r	if Ras up-expressed
	+mr	if Myc & Ras up
	+e1, or +e2, or +e3	if E2F1, 2 or 3 up
	3+	

Design matrix H 0/1 entries

Same for each gene

Parallel processing to compute posterior, predictive summaries

Examples (see Matlab code explorations)

A more elaborate example: <u>Mice & Heart Disease</u>

# Factor Regression

#### Linear Regression on SVD Factors

Linear model with genes as predictors	$\mathbf{y} = \mathbf{H}\boldsymbol{\beta} + \boldsymbol{\epsilon}$	$\mathbf{H} = \mathbf{X'}$
SVD implies	$\mathbf{y} = \mathbf{F'}\mathbf{\theta} + \mathbf{\varepsilon}$ $\mathbf{\theta} = \mathbf{D}\mathbf{A'}\mathbf{\beta}$	Regression on <i>n&lt;<p< i=""> factors Massive dimension reduction Many-one: No unique inverse</p<></i>
	$\boldsymbol{\beta} = \mathbf{A}\mathbf{D}^{-1}\boldsymbol{\theta}$	Common "optimal" least-norm inverse

Priors: Normal priors are consistent/coherent Generalised g-priors (West, Bayesian Statistics 7, 2003)

#### Linear Regression on SVD Factors



Implied posterior  $\boldsymbol{\beta} \mid \mathbf{y}, \boldsymbol{\phi} \sim N(\mathbf{A}\mathbf{D}^{-1}\mathbf{g}, \mathbf{A}(\mathbf{D}\mathbf{G}\mathbf{D})^{-1}\mathbf{A}')$ 

Extension of shrinkage prior: Multiple shrinkage

Different degrees of shrinkage for different factor dimensions

Estimation of multiple shrinkage factors?

$$\mathbf{C}^{-1} = diag(\tau_1, \dots, \tau_n)$$

#### MCMC - Gibbs Sampling in Factor Regression

## Multiple shrinkage prior model

Simulate: Iteratively resample conditional posteriors

Sample means, histograms MC approximation of posterior quantities of interest



#### **MCMC - Gibbs Sampling in Factor Regression**

Illustrative/exploratory example: Breast cancer predicting ER gene RNA levels to identify Metagene predictor of ER -pathway related genes -

PNAS 2004 data and code

Second, non-genomic example:

Predicting fat content of cookies from infrared mass spectroscopy - covariates are spectra/finely sampled curves -(Bayesian Statistics 7 paper/M West)

39 training samples, 39 to predict p=300 wavelengths on NIR reflectance spectra: Aim to produce model to predict fat of future cookies

#### Cookie Data



#### Cookie Data



#### Data and fitted/predicted values

\* training, o predicted/validation

![](_page_23_Figure_4.jpeg)

Bayesian applot of residuals

#### Cookie Data

![](_page_24_Figure_1.jpeg)

#### **Coherent Theoretical Basis for SVD Regression?**

Model & design data

$$y = H\beta + \varepsilon$$
$$H' = X = ADF$$

Transform

$$\mathbf{y} = \mathbf{F'} \boldsymbol{\theta} + \boldsymbol{\epsilon}$$
$$\boldsymbol{\theta} = \mathbf{D} \mathbf{A'} \boldsymbol{\beta}$$

 $\boldsymbol{\theta} \sim N(\boldsymbol{0}, \boldsymbol{C}^{-1})$ 

![](_page_25_Picture_6.jpeg)

Conceptual, technical questions:

- *n* parameters θ sample size dependent?
- prior on  $\theta$  design-dependency
- $\cdot$  non-unique reverse map to  $\beta$

e.g., predict new  $y_{n+1}$  at new design point ...

... which prior to use?

## **Theoretical Foundation: Latent Factor Models**

![](_page_26_Figure_1.jpeg)

West 2003, Valencia 7

Latent Factor Regression Models - & SVD Regression Limiting Case -

$$\mathbf{x}_{i} = \mathbf{B}\boldsymbol{\lambda}_{i} + \mathbf{v}_{i}$$
$$y_{i} \sim N(\boldsymbol{\lambda}_{i}'\boldsymbol{\theta}, 1/\phi)$$

Latent/common structure in gene expression patterns ...

... is predictive of outcome, phenotype

Everything is normal, so: 
$$E(y|x) = x'\beta$$
  
 $\beta = M\Theta$ 

Coherent model:

prior on  $\boldsymbol{\Theta}$  transfers uniquely to  $\boldsymbol{\beta}$ 

Limiting case: SVD ... β≈AD<sup>-1</sup> θ

West 2003, Valencia 7

# - many open research questions -

## Model fitting: Parametrisation of **B** Identification MCMC Sparse Models: **B** sparse One gene - one or a few "pathways" One "pathway" - few or many genes, but not all Mixture priors on B<sub>ii</sub> ... $B_{ii} \sim pI(B_{ii} = 0) + (1 - p)N(0, \omega)$

West 2003, Valencia 7

#### Factors in Breast Cancer Data SVD ~ Sparse latent factor model

![](_page_29_Figure_1.jpeg)

Next up:

# Binary Regressions & Molecular Phenotyping