



Biomarkers for personalized medicine

Lars Arvastson
H. Lundbeck A/S



INTRODUCING LUNDBECK



LUNDBECK IN BRIEF

We are an international pharmaceutical company specializing in central nervous system disorders

- Founded by Hans Lundbeck in 1915
- An integrated company with core competencies in research, development, production, marketing and sales
- International presence with pharmaceuticals in more than 100 markets
- Marketed pharmaceuticals include treatments for Alzheimer's disease, depression and anxiety, epilepsy, Huntington's disease, insomnia, Parkinson's disease, and schizophrenia/bipolar disorder
- Headquarters in Copenhagen, Denmark
- Approximately 6,000 employees in 57 countries
- 2011 revenue: DKK 16 billion (approx. EUR 2.1 billion/USD 3 billion)

Molecular biomarkers – what to measure?



DNA – Gene expression – Proteins

State – Trait

Genotype – Phenotype

Cost – Quality

CNS – Periphery

Explorative – Hypothesis

Danish pharma – still doing well!



Biomarkers and personalized medicine



- The future will be more focused on personalized treatment
- Biomarkers have a central role
 - Diagnostic biomarkers
 - Prognostic biomarkers
 - Predictive biomarkers
 - Biological understanding

Gene expression analysis



- Genome wide scan
 - Micro array technology
 - ~100 000 genes
 - Low quality data
 - No prior assumptions
- Selected candidate genes
 - qPCR technology
 - ~100 genes
 - High quality data
 - Selected based on prior knowledge

Scientific questions



- Genes associated with a disease
 - Biological understanding
- Prediction of treatment response
 - Companion diagnostic
- Classification of disease state
 - Diagnosis

Associating genes with a diagnosis



$$x_{ij} = \mu_j + \beta_j^{\text{Age}} \text{Age}_i + \beta_j^{\text{Gender}} \mathbb{1}_{\{\text{Subj}=\text{Female}\}} + \beta_j^{\text{MDD}} \mathbb{1}_{\{\text{Subj has MDD}\}} + \epsilon_{ij}$$

x_{ij} = Gene expression level for gene j , subject i

- Simple t-test
- Multiple testing
 - Bonferroni
 - False Discovery Rate
- Other confounding factors
 - BMI
 - Inclusion criteria
 - Smoking
 - Alcohol

Associating genes with treatment response



$$y_i = \mu + \sum_g \beta_g x_{gi} + \sum_g \gamma_g x_{gi} \mathbb{1}_{\{\text{Active treatment of subject}\}} + \epsilon_i$$

where

y_i is the outcome depression score (adjusted for baseline score and treatment),

x_{gi} is the gene expression for gene g , patient i ,

$$\epsilon_i \sim \text{i.i.d } N(0, \sigma^2)$$

$$\text{Predictive index for patient } i: P_i(\mathbf{x}) = \sum_g \hat{\gamma}_g \cdot x_{gi}$$

Classification of disease status based on gene expression



Logistic model

$$P(\text{Subject } i \text{ has MDD} | x_{i1}, \dots, x_{in}) = \frac{\exp\left(\beta_0 + \sum_{g=1}^n \beta_g x_{ig}\right)}{1 + \exp\left(\beta_0 + \sum_{g=1}^n \beta_g x_{ig}\right)}$$

x_{ij} = Gene expression level for gene j , subject i

Mathematical toolbox



- Low dimensional data an a lot of data
 - Standard statistical toolbox
- High dimensional data and few data
 - Regularization
 - $\min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|_2 + \lambda \|\beta\|_1$
 - $\min_{\beta} (-\mathcal{L}(\beta; \mathbf{X}) + \lambda \|\beta\|_1)$
 - Selection of λ and p
 - Cross validation
 - Should be repeated
 - Permutation test
 - Should include selection of regularization parameter

Non-trivial when more genes than subjects



- Classification by logistic regression
- Matlab R2010b
- Pre-processing of data
 - Concentrations are log-transformed
 - Continuous variables are centralized to zero mean and scaled to one standard deviation
 - Binary variables defined to $\{-1, 1\}$.
 - Missing data imputed with mean or ML estimate.
- LASSO regularization
 - Regularization parameter based on cross validation
- Significance based on permutation test
- Predictive performance calculated as area under the ROC curve.
 - ROC curves calculated based on double cross validation, regularization in a inner CV loop.

Example: Classifying gender based on mRNA



- Classify subjects as male/female based on gene expression profile solely.
- For each subject there is 29 gene expression levels,

$$X_{ADA}, \dots, X_{VMAT2}$$

- Predictive probability of gender based on logistic model,

$$\begin{aligned} &P(\text{Subject } i \text{ is male} | X_{ADA}, \dots, X_{VMAT2}) \\ &= 1 - P(\text{Subject } i \text{ is female} | X_{ADA}, \dots, X_{VMAT2}) \\ &= \frac{\exp(\beta_{Const} + \beta_{ADA}X_{ADA} + \beta_{VMAT2}X_{VMAT2})}{1 + \exp(\beta_{Const} + \beta_{ADA}X_{ADA} + \beta_{VMAT2}X_{VMAT2})} \end{aligned}$$

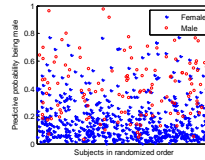
- Model parameters

$$\beta_{Const}, \beta_{ADA}, \dots, \beta_{VMAT2}$$

Example: Estimated model



- Model parameters
 - ML estimate
 - LASSO regularization
 - Cross validated regularization parameter

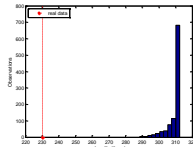


β_{Const}	= -1.6654	▼
β_{PLB}	= -0.78074	▼
β_{CPRE}	= -0.6367	▼
β_{PWRB2}	= 0.88211	▲
β_{CDBP}	= -0.57192	▼
β_{CZBA}	= 0.49142	▲
β_{NOLA}	= -0.42307	▼
β_{LA}	= -0.35866	▼
β_{JAREB}	= 0.3556	▲
β_{ADMA4}	= 0.302	▲
β_{PCR}	= -0.29322	▼
β_{PKR7}	= 0.25232	▲
β_{MIS}	= -0.25176	▼
β_{SDN}	= -0.24763	▼
β_{M31}	= 0.23745	▲
β_{VMA12}	= 0.20369	▲
β_{M32}	= -0.18548	▼
β_{M33}	= 0.1806	▲
β_{TAAA}	= -0.15273	▼
β_{M31}	= 0.12221	▲
β_{TF2}	= -0.091958	▼
β_{PGE}	= -0.081653	▼
β_{PMB}	= -0.0036507	▼
β_{SPCA}	= 0	•
β_{PCR}	= 0	•
β_{M14}	= 0	•
β_{PCR}	= 0	•
β_{MWA}	= 0	•
β_{ADA}	= 0	•
β_{LS}	= 0	•

Example: Significance



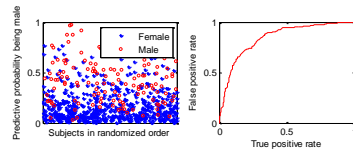
- Probability that a model would describe data equally well by chance.
- Permutation test, repeated 1000 times.
- Estimated p-value = 0/1000
- Classifier include 22 genes



Example: Performance



- Trade-off between sensitivity and specificity
- ROC curve AUC classic measure of predictive power
- AUC = 0.84 for final model on training data
- Cross validated AUC = 0.79 (10-fold repeated 10 times)



Why so complicated?



- Many genes, few subjects
- No clear signal

Method	Comment
Full sample lasso estimation	Too optimistic (performance bias)
CV	To choose smoothing parameter (α)
Repeated CV	To reduce variability in estimation due to random split
Double (outer) CV	To remove bias in performance evaluation
Repeated double CV	To reduce variance in performance evaluation
Permutation test	To give p-value for effect

Summary and Conclusions



- Data with many samples and few subjects
- Different computer intensive techniques in use
 - Simple models
 - Linear regression
 - Logistic regression
 - Regularization
 - LASSO
 - Ridge
 - L0
 - Cross validation
 - Repeated
 - Double cross validation
 - Permutation test

References



- MAQC Consortium (2010). The MicroArray Quality Control (MAQC)-II study of common practices for the development and validation of microarray-based predictive models. *Nature Biotechnology*, Vol. 28, No. 8, pages 827-841.
- Hastie, T., Tibshirani, R. and Friedman, J. (2009). *The elements of statistical learning*, New York, Springer, 2.ed.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *J. Royal. Statist. Soc B.*, Vol. 58, No. 1, pages 267-288.