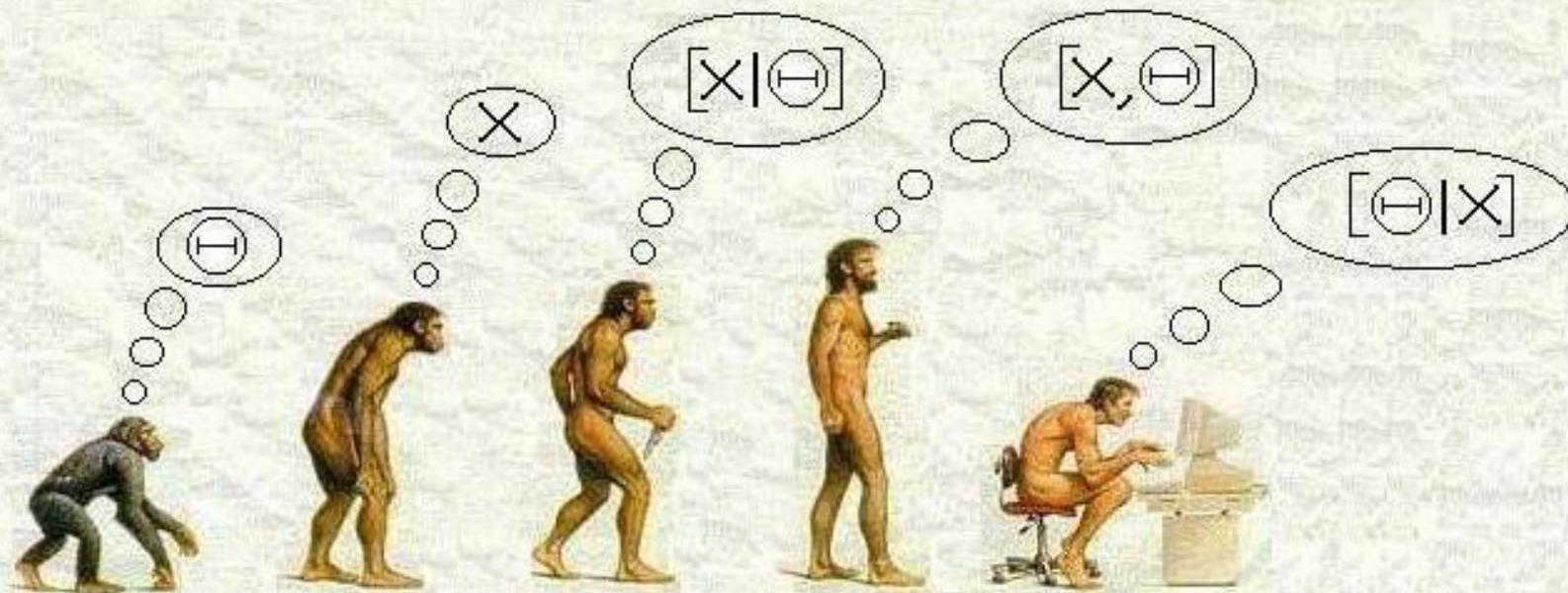




Bayesian Prediction Tree Models

(YET ANOTHER) HISTORY OF LIFE AS WE KNOW IT...

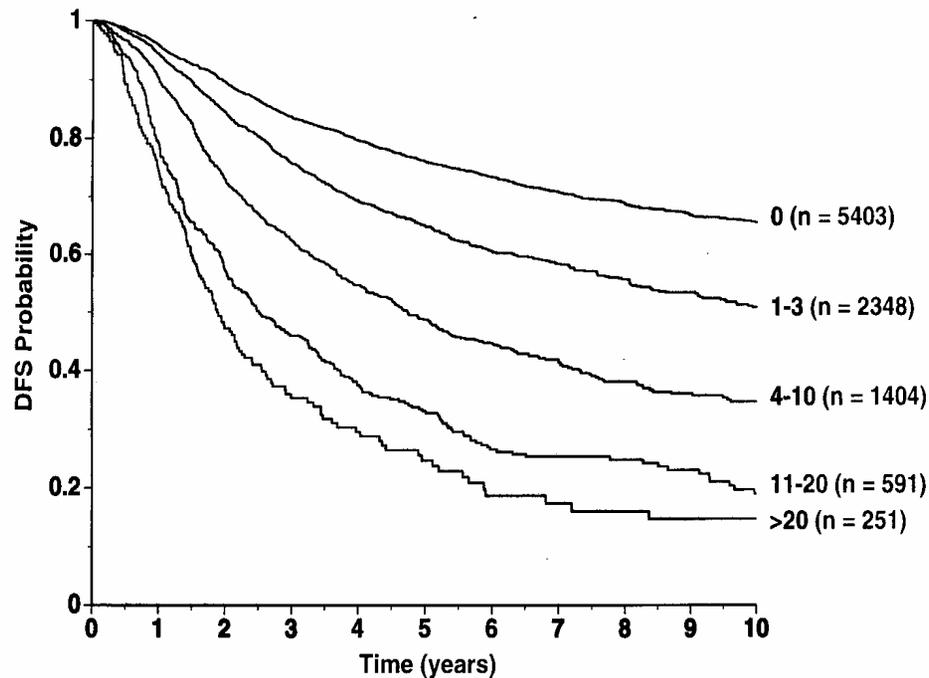


HOMO APRIORIUS **HOMO PRAGMATICUS** **HOMO FREQUENTISTUS** **HOMO SAPIENS** **HOMO BAYESIANIS**

Statistical Prediction Tree Modelling for Clinico-Genomics

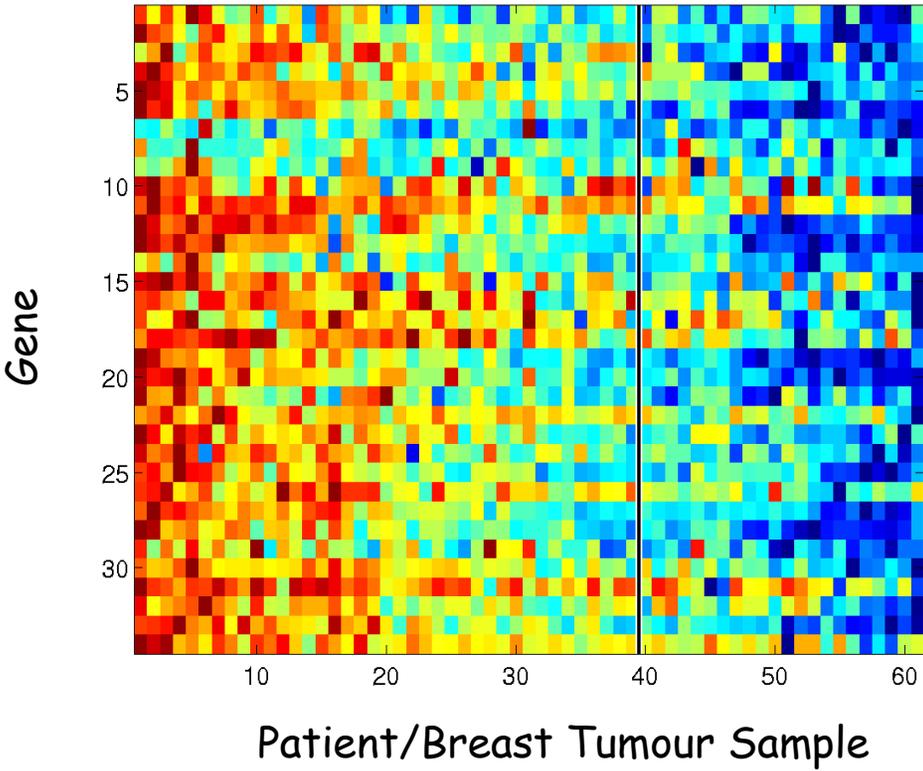
- Clinical gene expression data
 - expression signatures, profiling
- Tree models for predictive sub-typing
- Combining clinical + molecular data
- Gene discovery and prioritisation
- Breast cancer prognosis - clinical testing

Lymph Node Metastasis - Key Breast Cancer Risk Factor -



But ... lymph node dissection carries morbidity, inaccuracy

Gene Subsets Associated with Axillary Lymph Node Status



Expression Signatures - Metagenes

Multiple related genes:

- multiplicities, redundancies
- idiosyncratic noise, variability

Aggregate "Metagene" summaries:

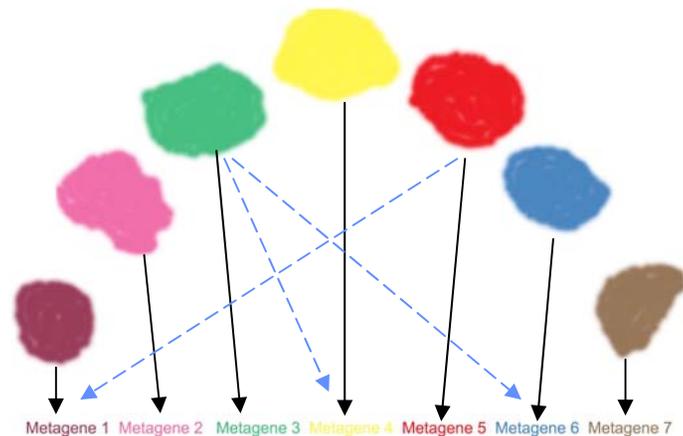
- characterize patterns
- noise reduction in estimating common "signal"
- improve statistical efficiency, robustness
- multiple "small" clusters of genes

Metagenes - Characterising Factors

$$Y=A(M)+E$$

- Dimension reduction:
Signal improvement
- Clustering, k-means
- Empirical or model-based
factor analysis

- (Latent) factors in
graphical models
- Related/intersecting
clusters



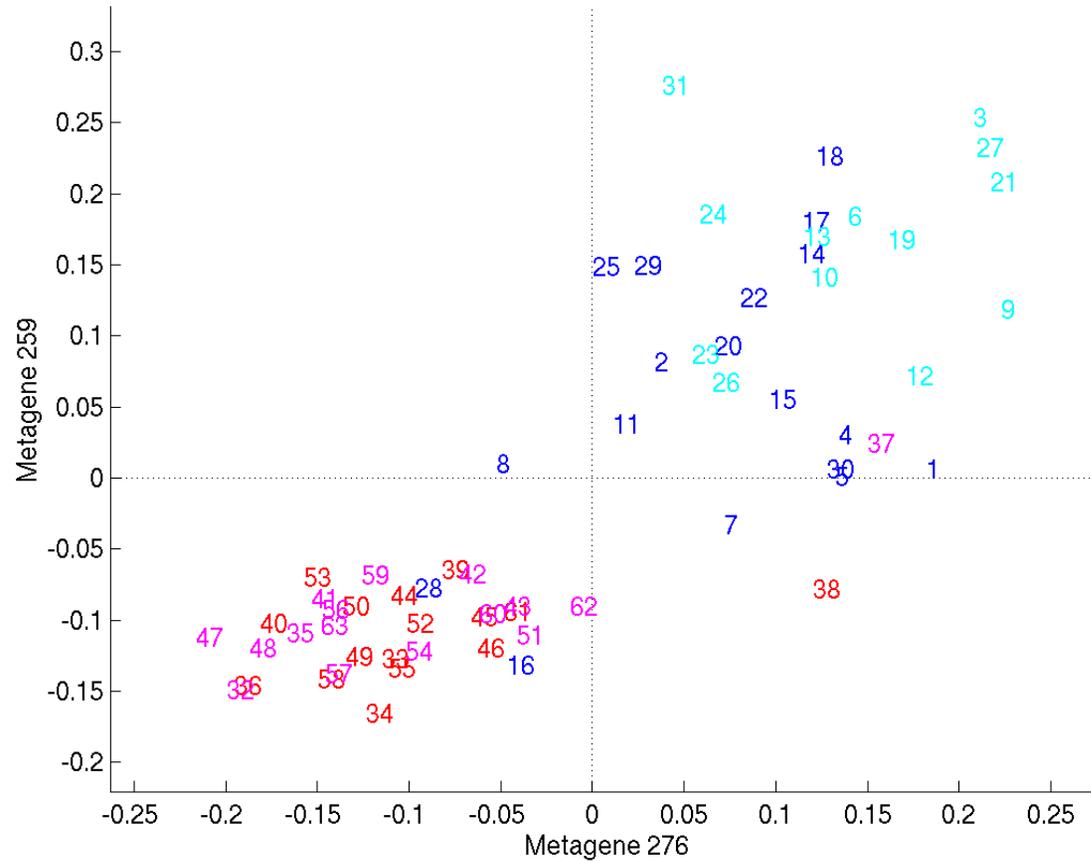
Reclustering, Covariance selection &
Graphical Models - software tools ..

[MetageneCreator](#)

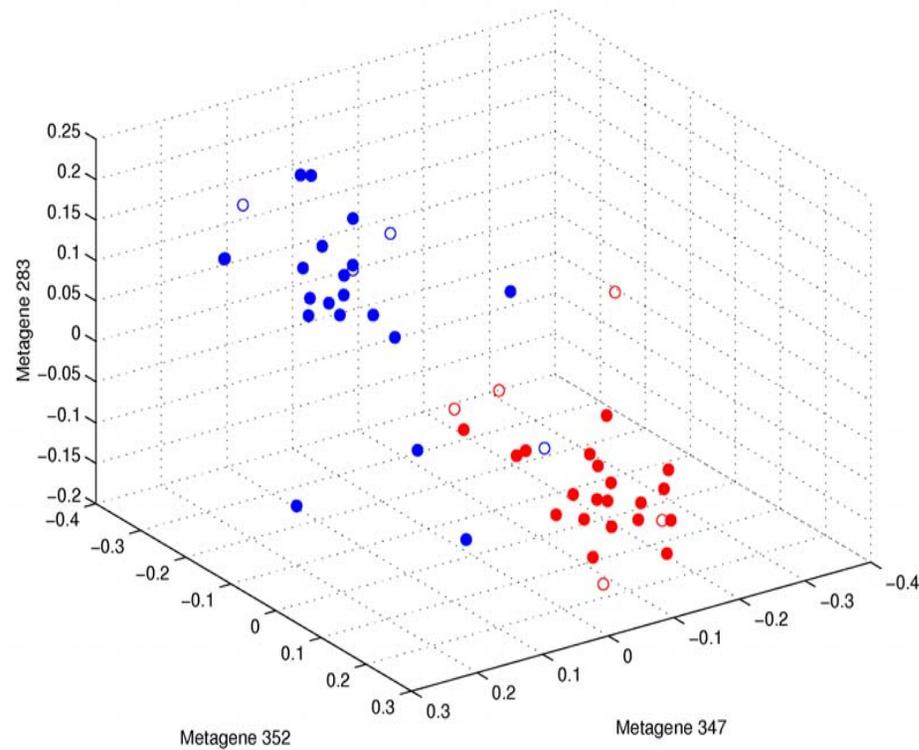
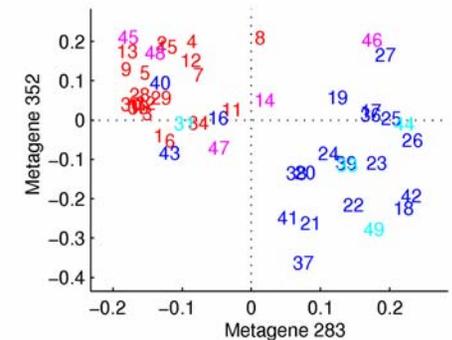
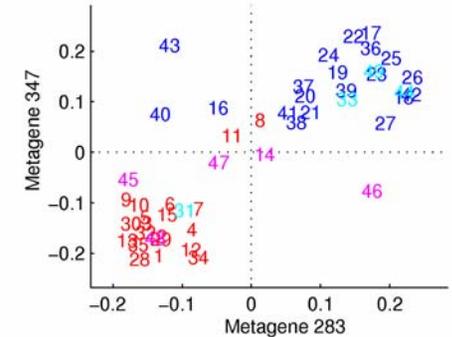
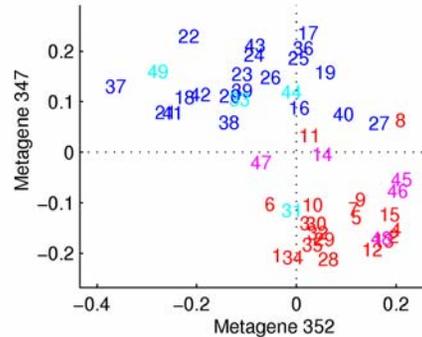
West, 2003, Bayesian Statistics 7

(Dobra, West, 2004, forthcoming)

Atherosclerosis Metagenes - Disease Susceptibility -

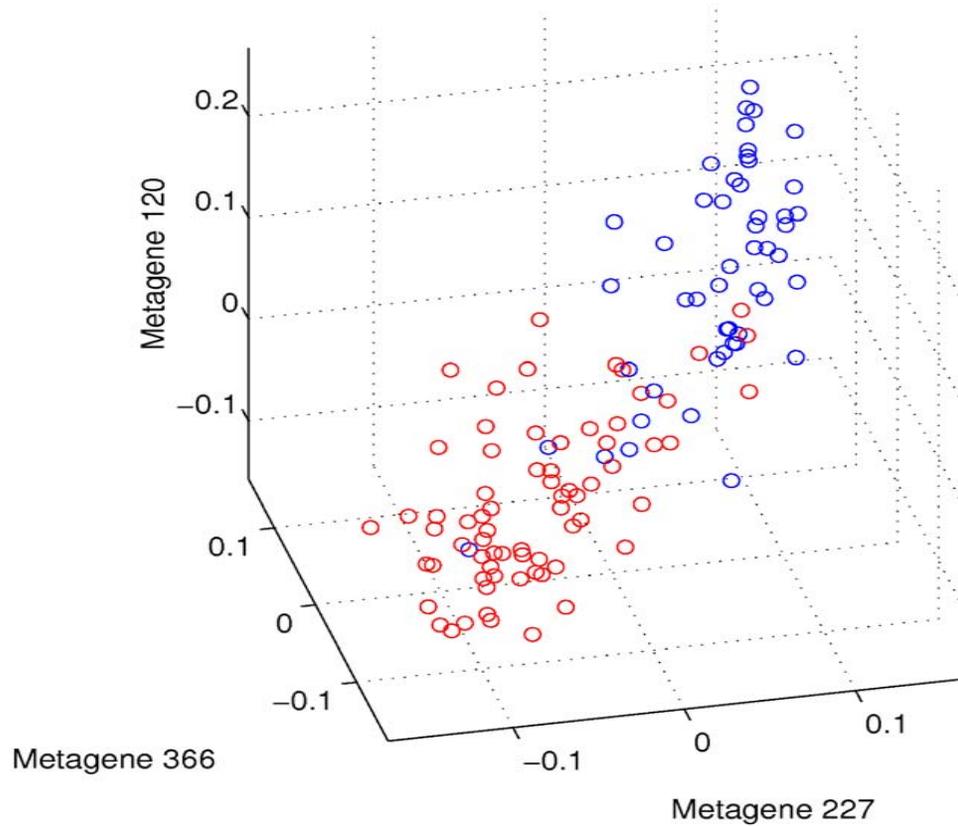


ER Metagenes in Breast Cancers



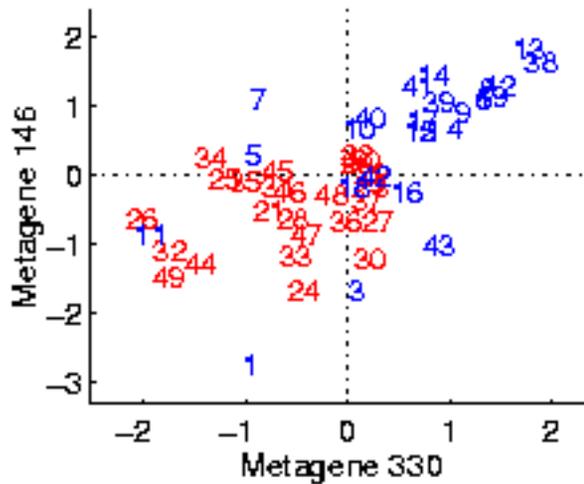
79 genes: ER regulated: Tff1, bcl-1,2, ..
 27 genes - ESR1
 39 genes - other co-regulated, androgens, ..

ER Metagenes in Breast Cancers



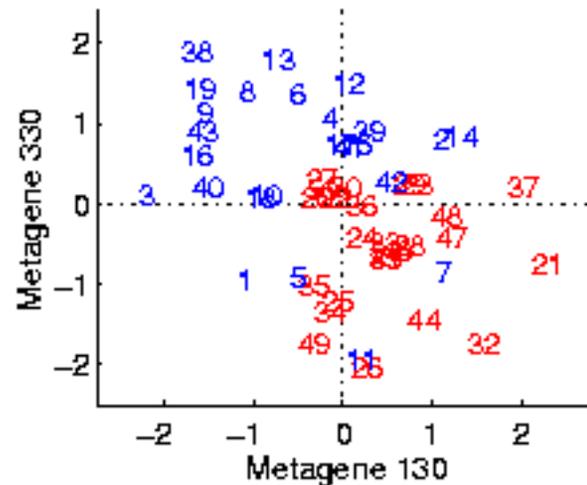
158 samples - Pittman et al, Duke/KFSYS PNAS 2004

Lymph Node Metastasis Metagenes



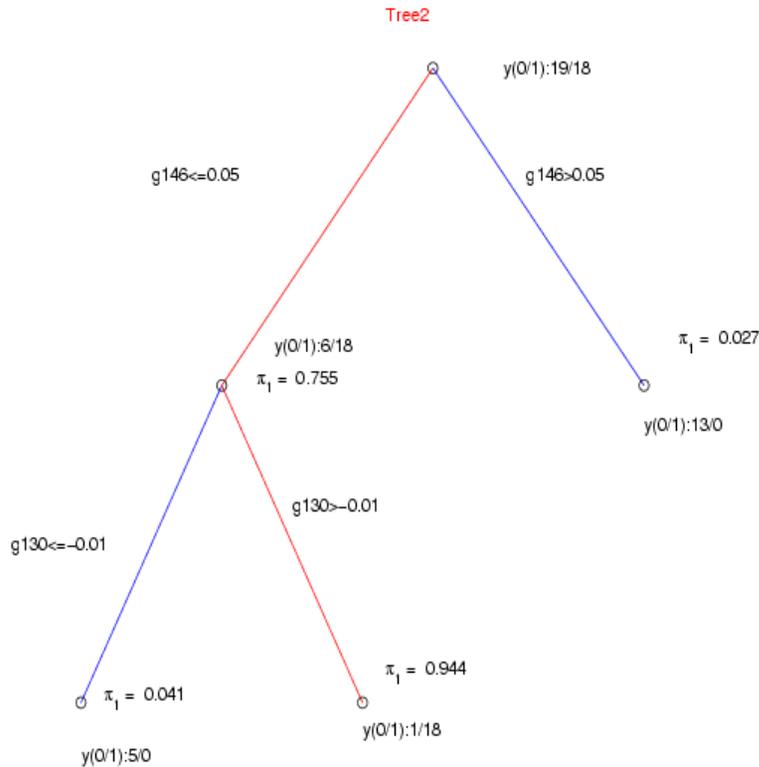
Classification Trees Models:

- *Lymph node outcomes*
- *Recurrence outcomes*
- Cross-validation assessment
- Predictive classification



Binary Tree Models

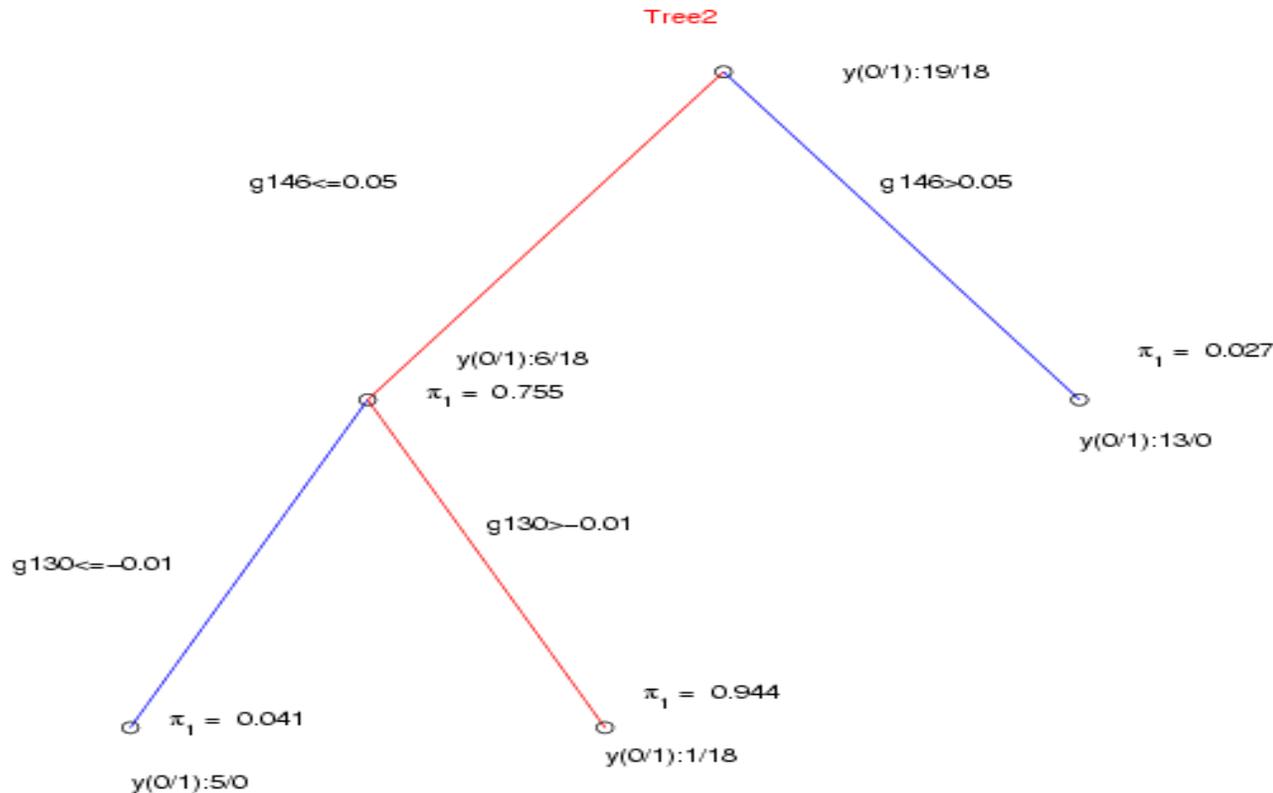
Tree Models: Classification via Recursive Partitioning



- Non-linear: interactions of covariates
- "Natural" in contexts of clinical prognosis
- Many trees:
Model & prediction uncertainty
- Prospective or
retrospective studies

Binary Outcomes: $Y=0/1$

- Retrospective Sampling (Case-Control studies) -



$$\pi = \Pr(Y = 1 | x_1 \leq \tau_1, \dots, x_k \leq \tau_k)$$

Binary Outcomes: Prospective Inference from Retrospective Model

$$\pi = \Pr(Y = 1 | x_1 \leq \tau_1, \dots, x_k \leq \tau_k)$$

$$a_{1y} = \Pr(x_1 \leq \tau_1 | Y = y), \quad y = 0, 1$$

$$a_{2y} = \Pr(x_2 \leq \tau_2 | x_1 \leq \tau_1, Y = y), \quad y = 0, 1$$

$$\frac{\pi}{1 - \pi} = \frac{\Pr(Y = 1) a_{1,1} a_{2,1}}{\Pr(Y = 0) a_{1,0} a_{2,0}}$$

Binary Outcomes: - Retrospective Model -

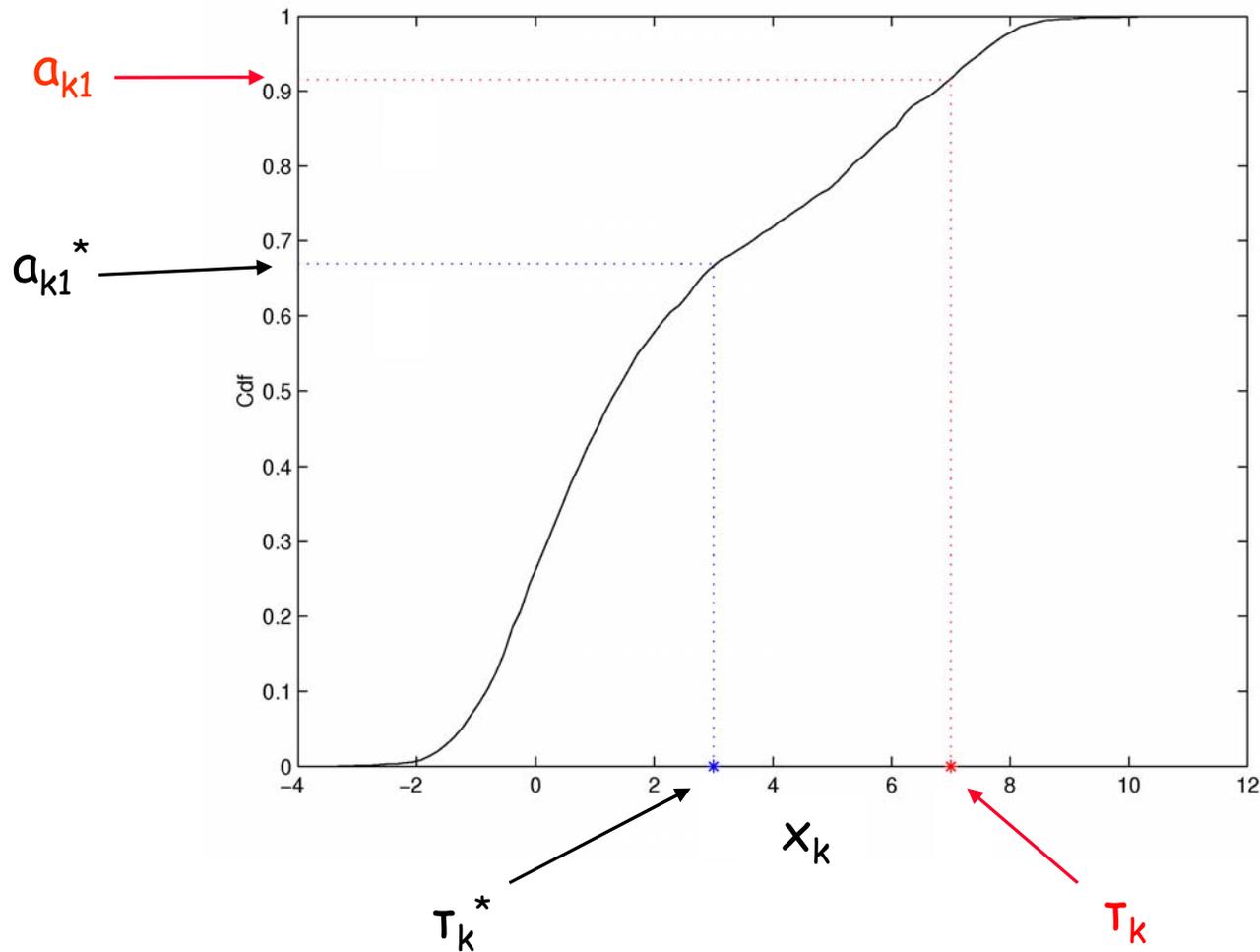
- Model conditional CDFs for predictors

$$F(x_k | x_1 \leq \tau_1, \dots, x_k \leq \tau_k, Y = y)$$

- Want consistency as thresholds vary
- Nonparametric Bayes: Dirichlet model
- Modelling in x space - joint structure
- Implies Beta priors on a_{ky}

Dirichlet Process: Beta Margins

$$F(x_k | -, y=1)$$



Assessing Predictor: Threshold Pairs for Association with Binary Outcome

Any predictor (metagene) x , threshold τ

Sample arranged in 2x2 table

Retrospective: column totals fixed

Two Bernoulli sequences: columns with "success" probabilities a_0 & a_1

Beta priors: Conjugate

| | $y=0$ | $y=1$ | |
|------------|----------|----------|-------|
| $X \leq T$ | n_{00} | n_{01} | N_0 |
| $X > T$ | n_{10} | n_{11} | N_1 |
| | M_0 | M_1 | |

Assess/test: Bayes' factor for
 $a_0 \neq a_1$ versus $a_0 = a_1$

Nonlinear association measure:
Function of threshold

Growing & Using Binary Trees

Node splits:

- assess candidate predictor:threshold pairs
- 2x2 table: conservative Bayesian tests

Inference: 1 tree

- beta posteriors for a_{iy}
- simulate: impute $\Pr(Y=1/leaf)$

Multiple trees:

- Multiple splits at any node

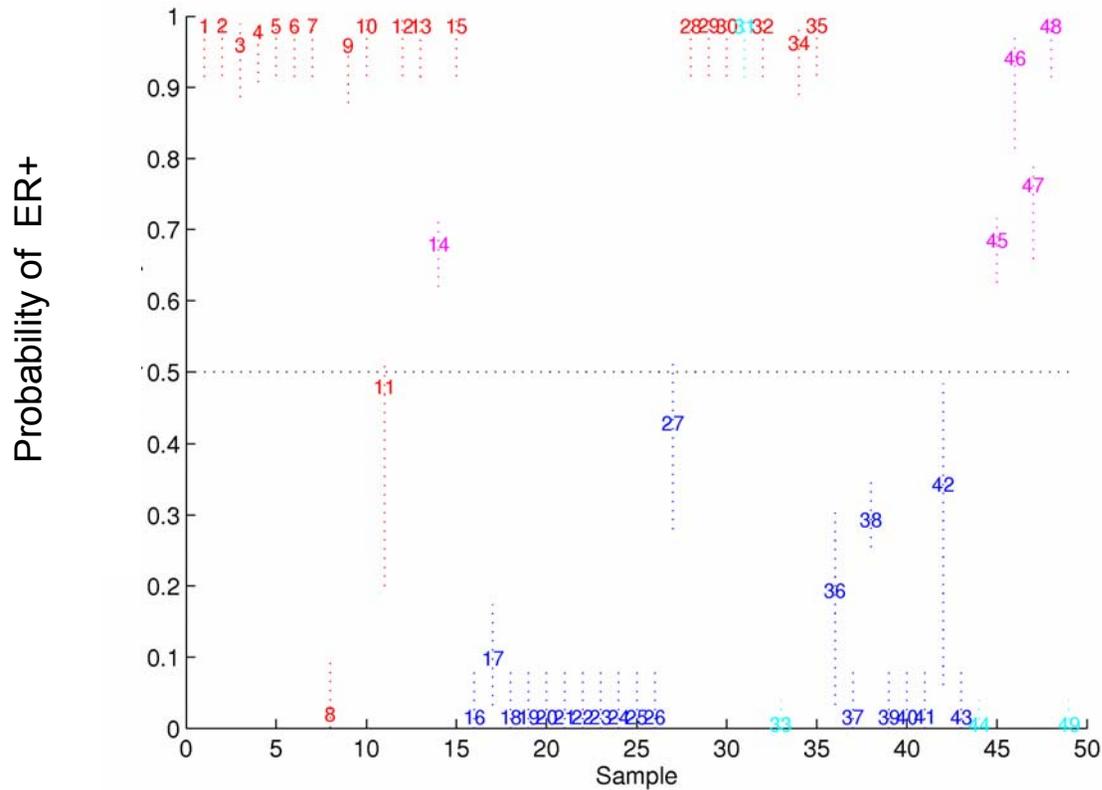
Prediction with multiple trees:

- Average across trees
- Model uncertainty
- "Smoothing" partitions

Predicting ER Status With Metagene Trees

Duke PNAS 2001 data -

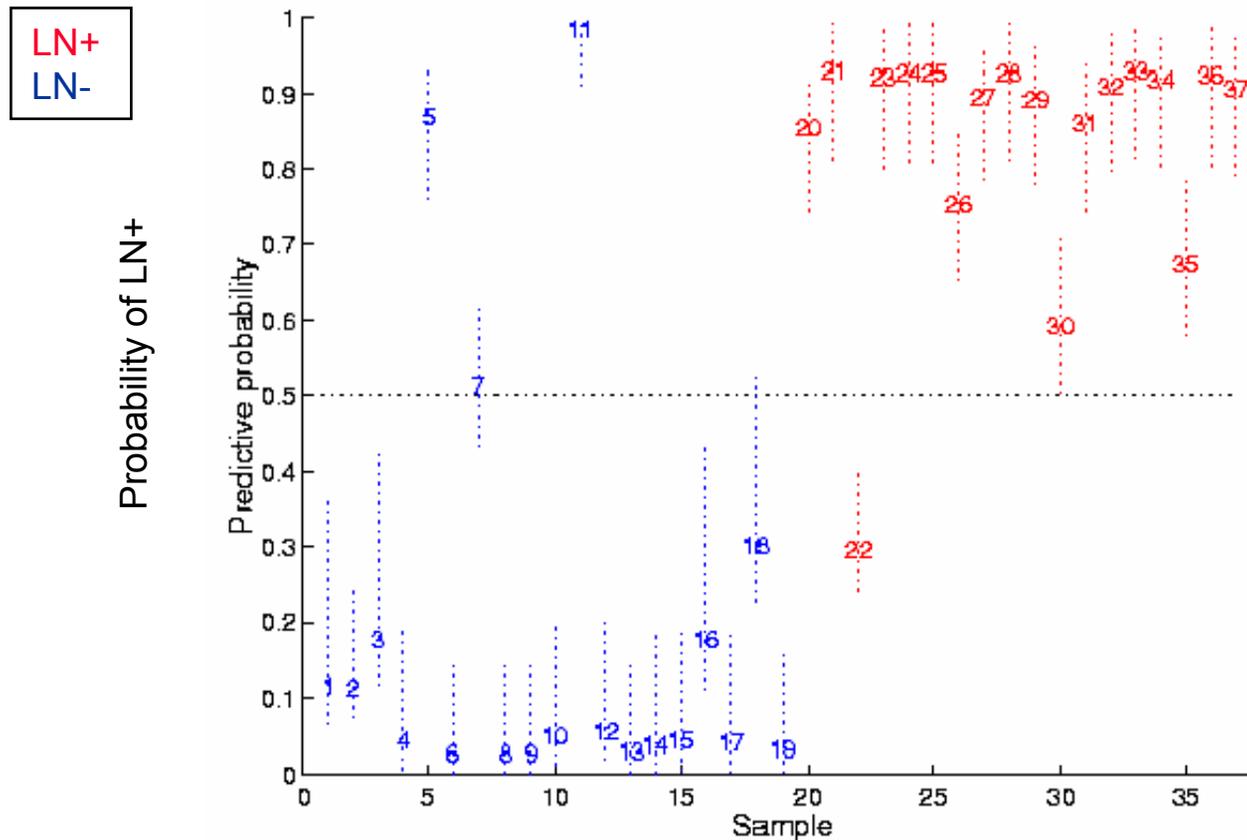
Out-of-sample cross validation



(Pittman et al, Biostatistics 2004)

Predicting Lymph Node Status With Metagenes

Out-of-sample cross validation



Gene Identification

Implicated metagenes - gene subsets

Genes correlated with key metagenes

Breast Cancer - nodal metastasis:

- Interferon pathway/inducible gene subset
- Interferons mediate anti-tumour response

Evidence of dysfunction of normal anti-tumour response?

- BRCA1, p27

Duke PNAS 2004 Breast Cancer Data

Examples (see Matlab code explorations)

Explore ER (0/1) regression using metagene predictors

Explore tree models for cancer recurrence

Look at nonlinear association with status - for each metagene across the expression (threshold) range

Fit forests of binary trees, look at 'top' trees, predictions of samples held-out

Tree Models for Survival Analysis

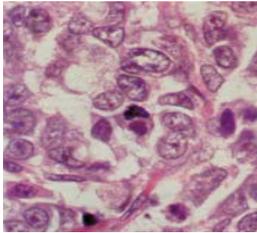
Prospective Tree Models

- Recurrence in Breast Cancer -

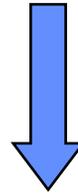
Survival tree models - Prospective

- Clinical + metagene predictors
- Survival distributions in subgroups

Challenge: Dissect Heterogeneity

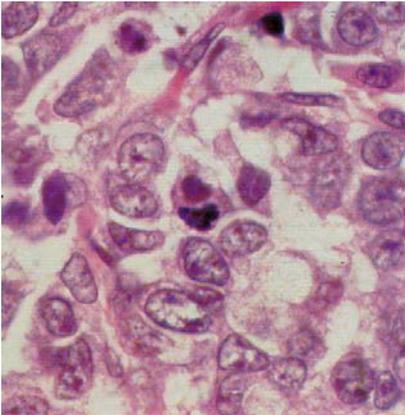


Cancer is fundamentally a **heterogeneous** disease

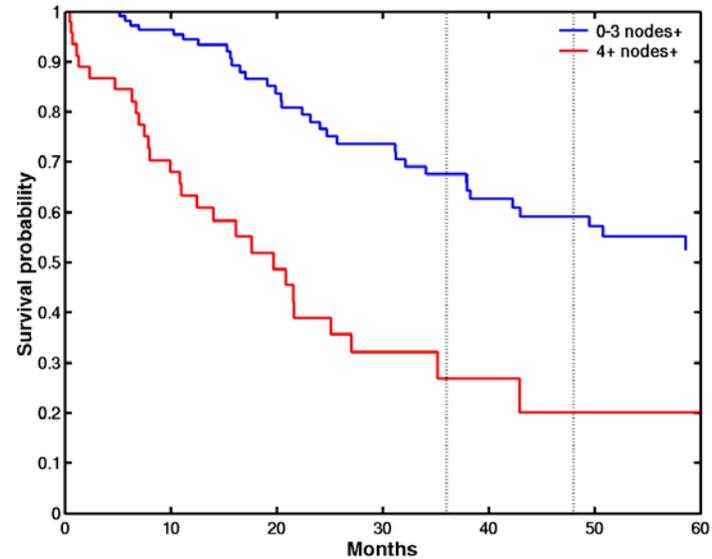


Understanding and dissecting the heterogeneity is key to prognosis and treatment

Clinical Predictors: "Low Resolution" Phenotypes



- Lymph node involvement



Refining Phenotype: Adding Molecular Information

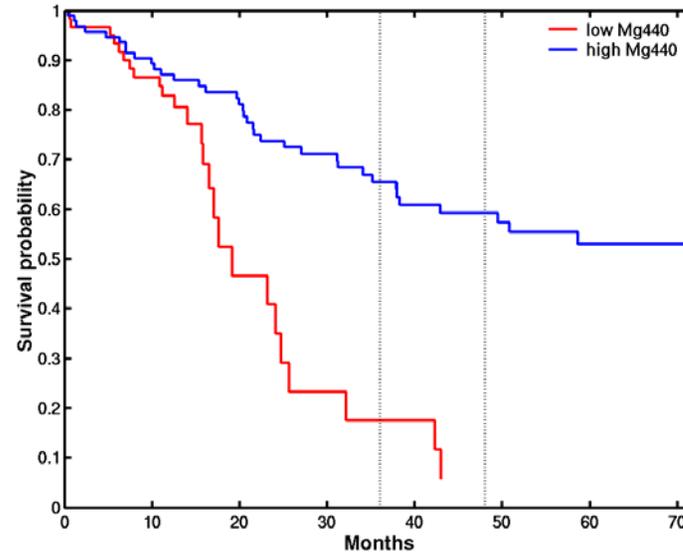
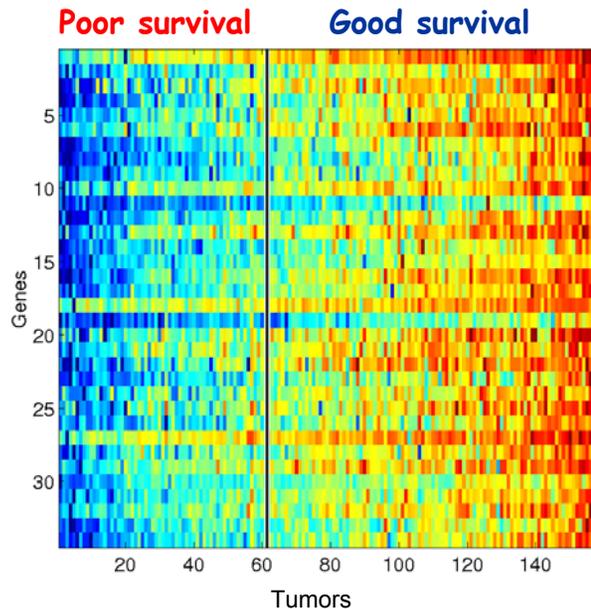
Analysis of the tumor

- *Gene expression*
- DNA copy status (CGH)
- DNA methylation
- Protein profiles
- Metabolic profiles

Analysis of the 'host'

- Genotypes
- Serum protein profiles
- Serum metabolic profiles
- Serum gene expression profiles?

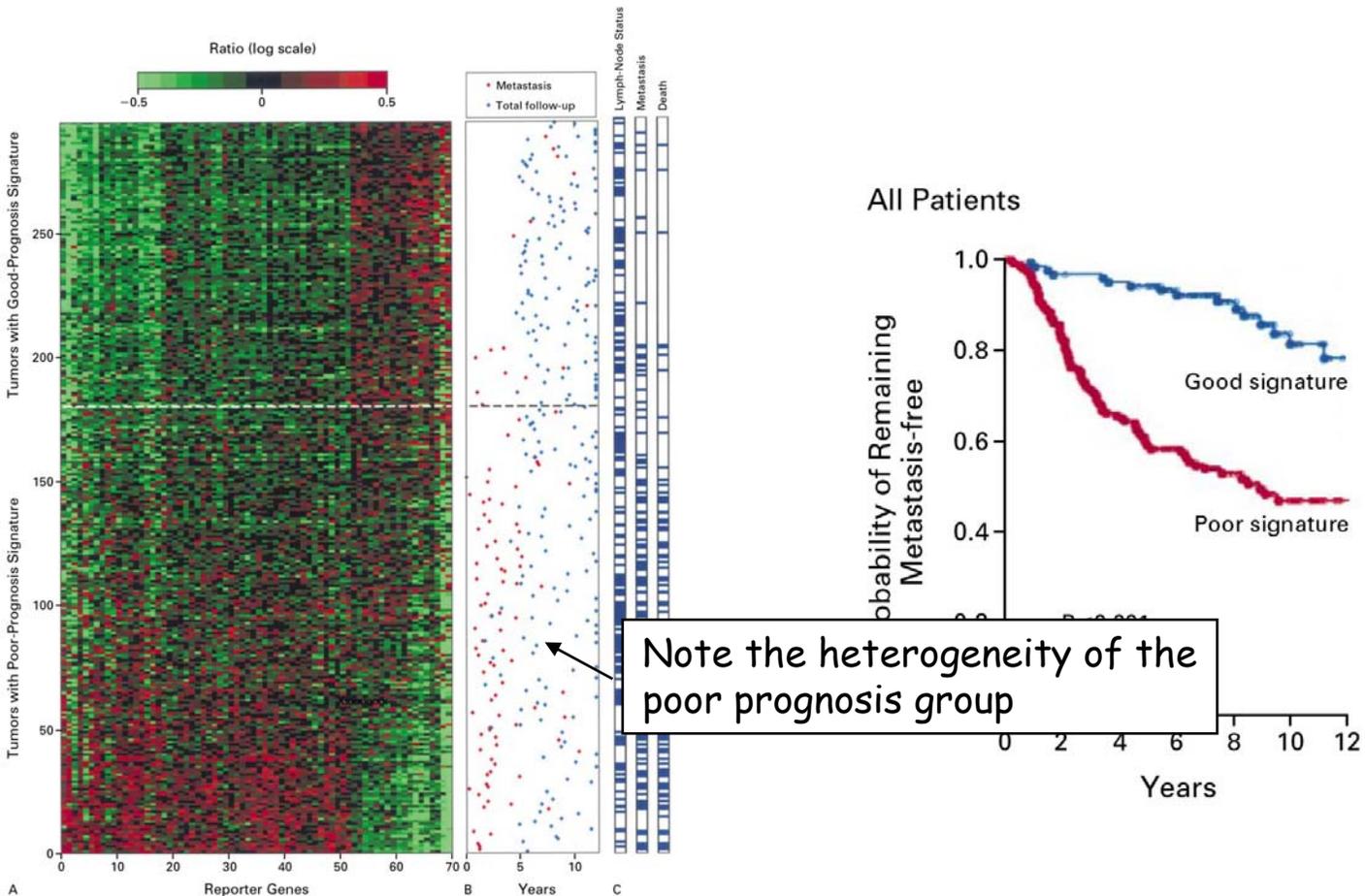
Metagenes are Predictive of Breast Cancer Recurrence



It's genomic, but 'group' prediction, yet another biomarker

A 70 Gene Predictor

van de Vijver et al., N. Engl. J. Med. 347: 1999 (2002)



Moving Gene Expression Profiling to the Clinic

genomic health | oncotype DX Breast Cancer Assay

Now we can illuminate the individual risk of breast cancer recurrence.

> Home > Genomic Health > Contact Us > Register | > Search GO | Take me to:

Privacy Terms and Conditions Printer-friendly

Healthcare Professionals

- > About Oncotype DX
- > The Role of Oncotype DX
- > Oncotype DX Studies
- > Oncotype DX Logistics
- > Frequently Asked Questions

Welcome to oncotype DX[™] Breast Cancer Assay for Healthcare Professionals

Oncotype DX[™] is a clinically validated diagnostic assay that provides a quantitative assessment of the likelihood of distant breast cancer recurrence. Oncotype DX brings significant new information to enhance treatment planning and will drive further progress in the fight against breast cancer. [Click here](#) to learn more about Oncotype DX.

Oncotype DX is now available and the Genomic Health Reference Laboratory is accepting patient samples.

[Click here](#) to register for updates if you wish to be notified when new information about Oncotype DX is added to this site.

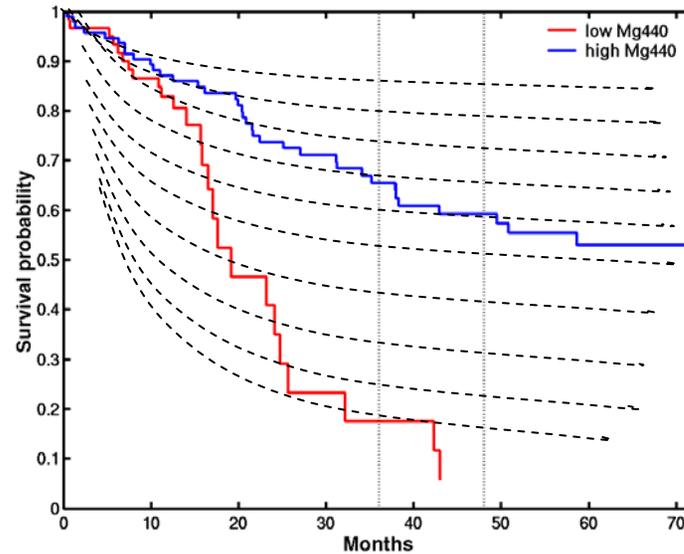
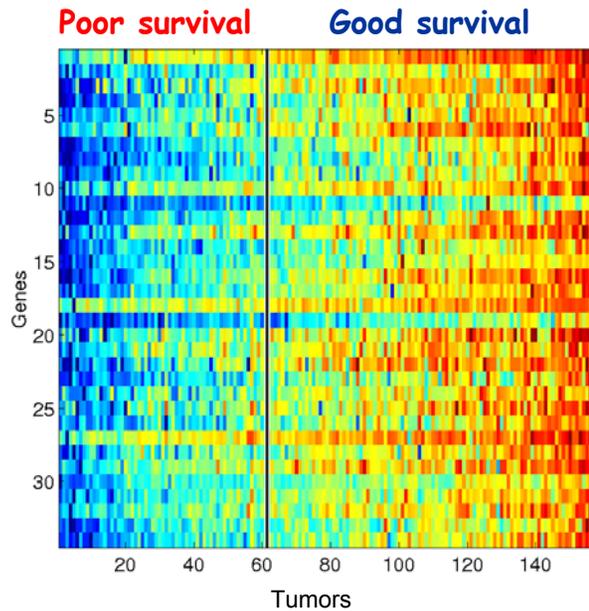
What's New?

12/04/2003
[The NSABP and Genomic Health Announce Positive Results from Large-Scale, Prospective Validation Study to Quantify Breast Cancer Recurrence in Newly Diagnosed Patients](#)

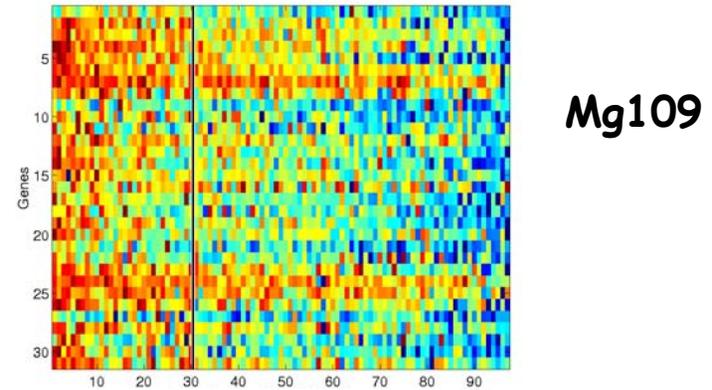
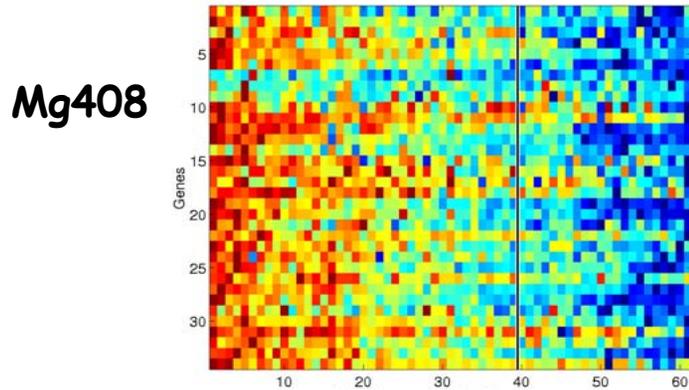
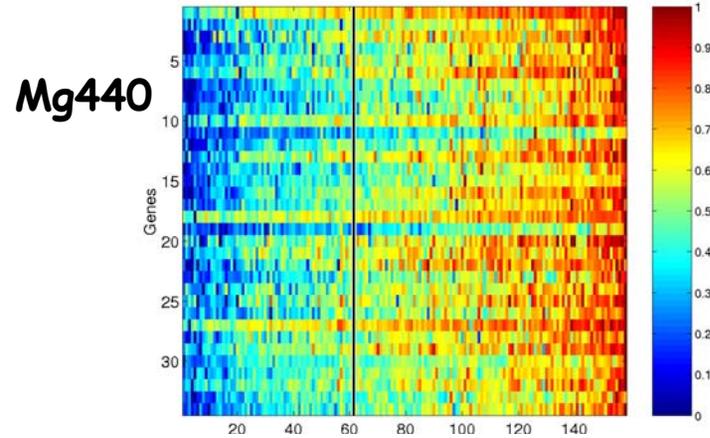
Register for updates
[Register](#)

for more Information
About Oncotype DX, call
(866) ONCOTYPE
662-6897

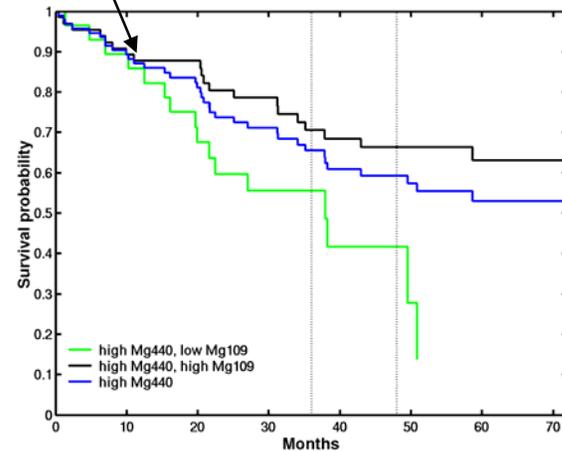
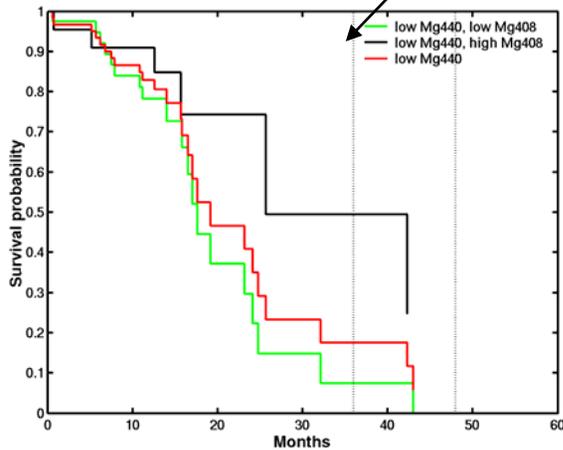
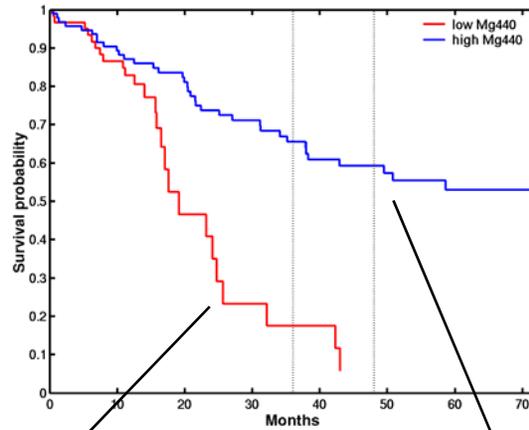
Personalised Prediction - Dissecting Heterogeneity -



Dissecting Heterogeneity Using Multiple Metagenes



Subtyping by Multiple Metagenes - Can Improve Risk Discrimination -



Predictive Survival Tree Models

Leaf j : $Y \sim \text{Weib}(\mu_j, \alpha)$

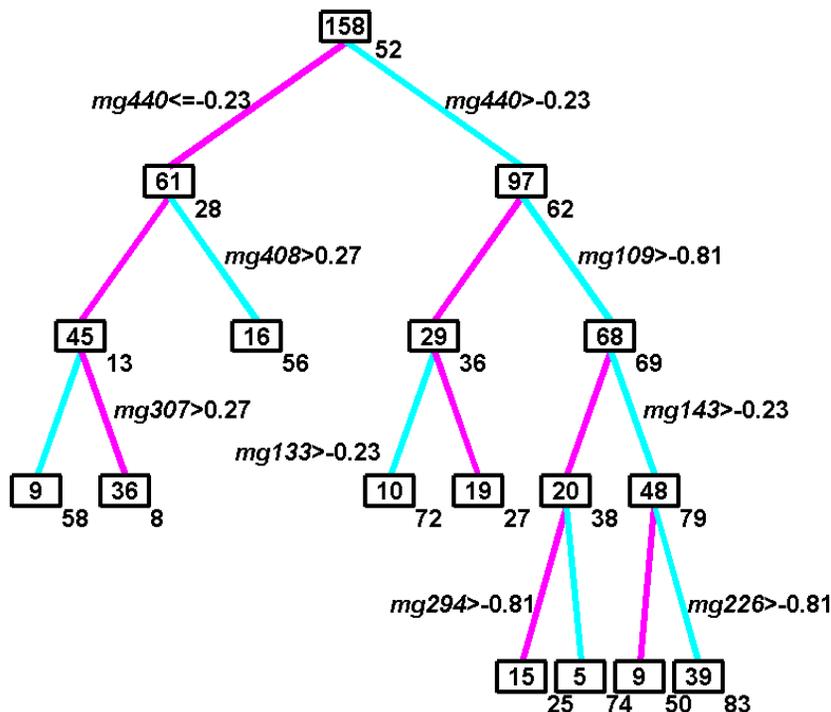
Hierarchical (gamma) priors:

$$p(\mu_j | m) = \text{Ga}(\mu_j | c, c/m_j)$$

- mean m_j has hyperprior
- m_j "similar" for sibling leaves
- shrinkage/data sharing

Marginal likelihood of tree:

$$p(Y | \text{Tree}, \alpha) = \prod_{\text{leaves } j} p(Y_j, \alpha)$$



(Pittman et al, PNAS 2004)

Growing & Using Survival Trees

Node splits:

- assess candidate predictor:threshold pairs
- 2 Weibulls vs 1?
- Bayes' test

Multiple trees:

- Multiple splits at any node

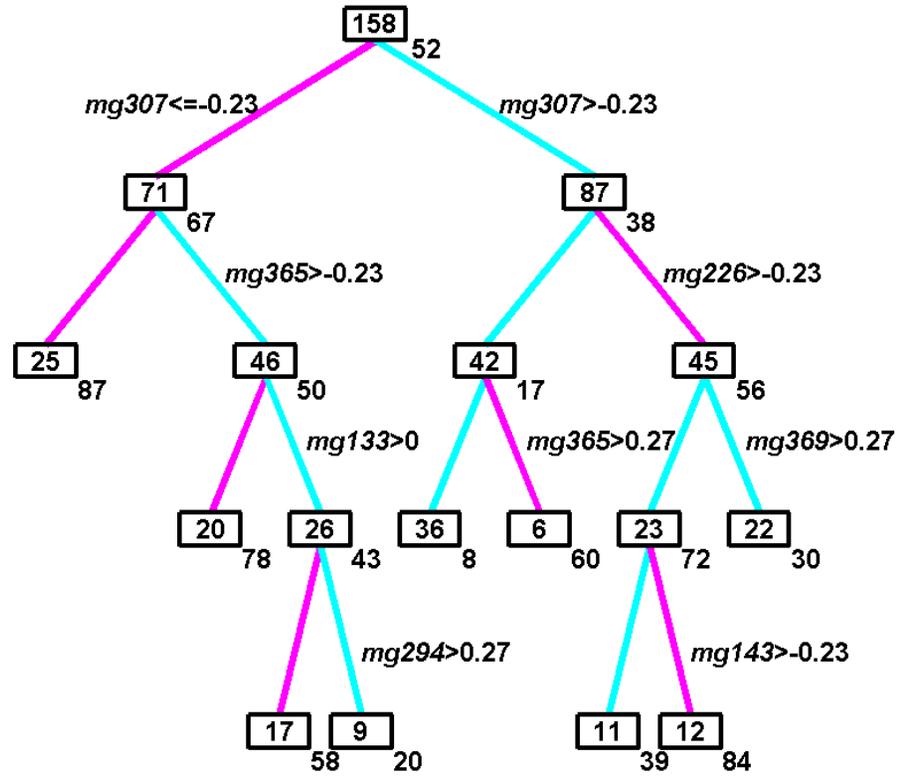
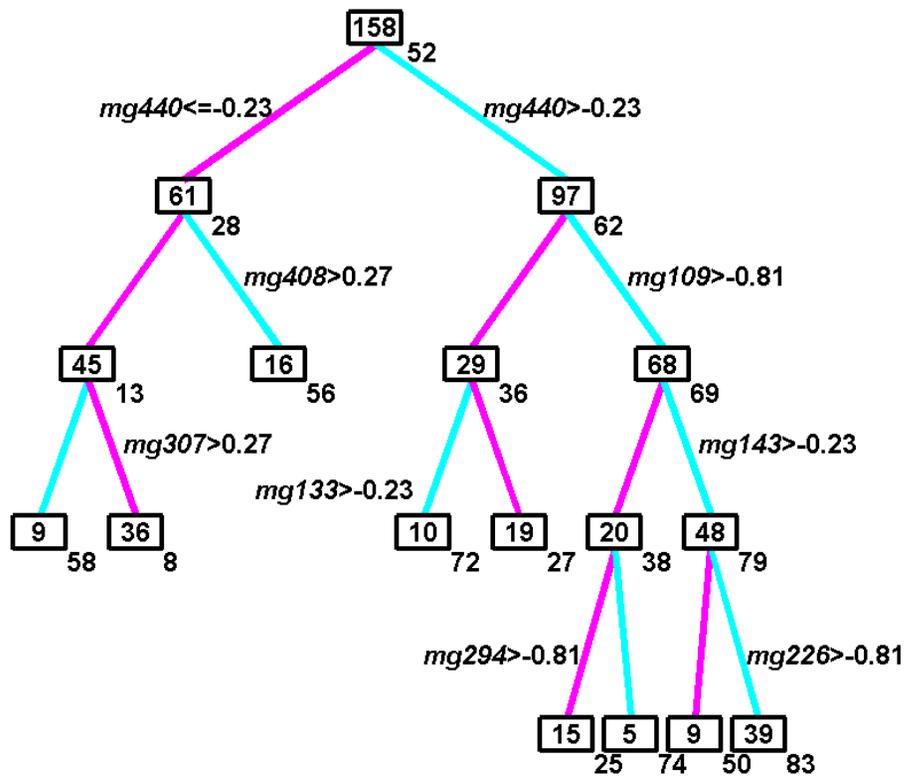
Inference: 1 tree

- Empirical Bayes for m_j, α
- Pareto predictions in leaf

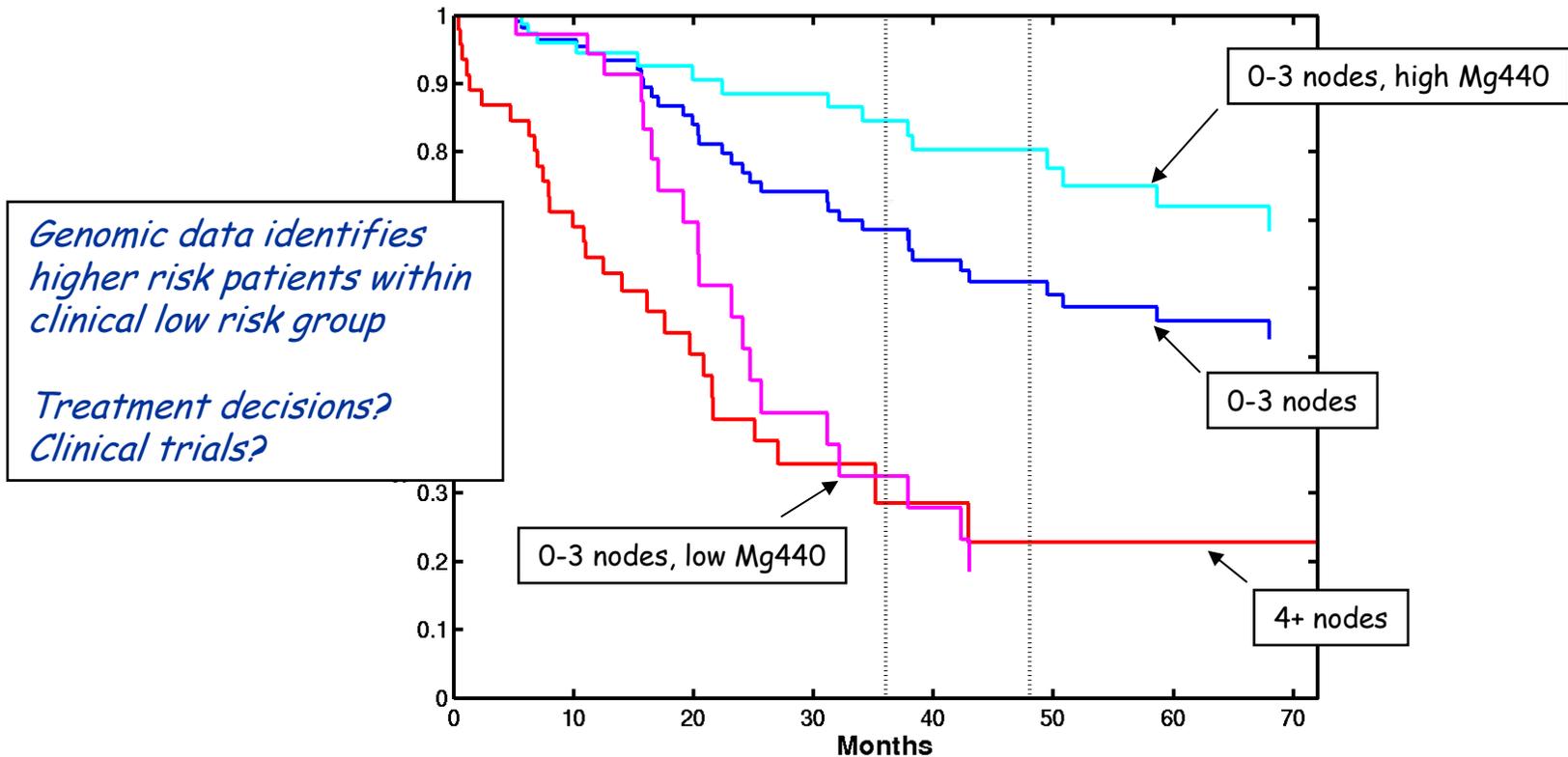
Prediction with multiple trees:

- Average across trees
- Model uncertainty
- Simulation

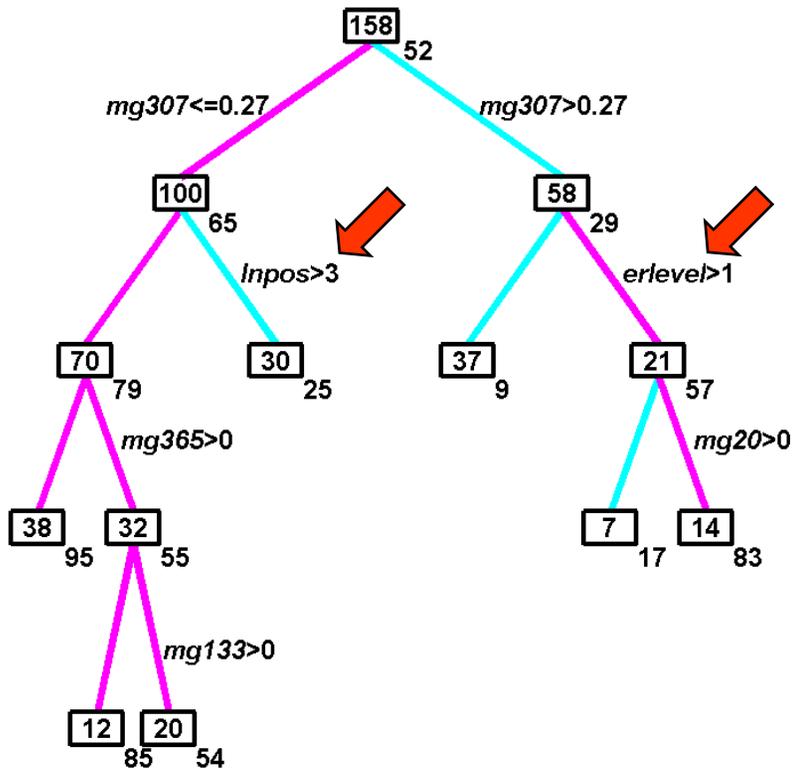
Forests of Metagene Trees



Clinico-Genomics for Modifying and Refining Risk Stratification



Forests of Clinico-Genomic Trees



*Select from full sets of
clinical and genomics
potential predictor variables*

multiple trees

*variable subset combinations
- co-occurrence*

Personalised Prognosis: Prediction of Individual Outcomes

Patient 158:

Tumor size 1.5cm

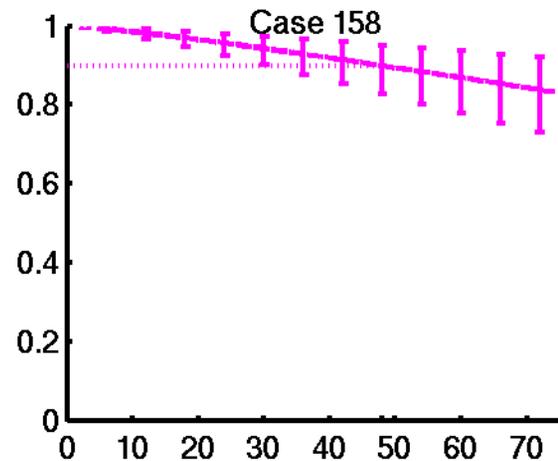
LN=1

ER-

Mg440= 0.25

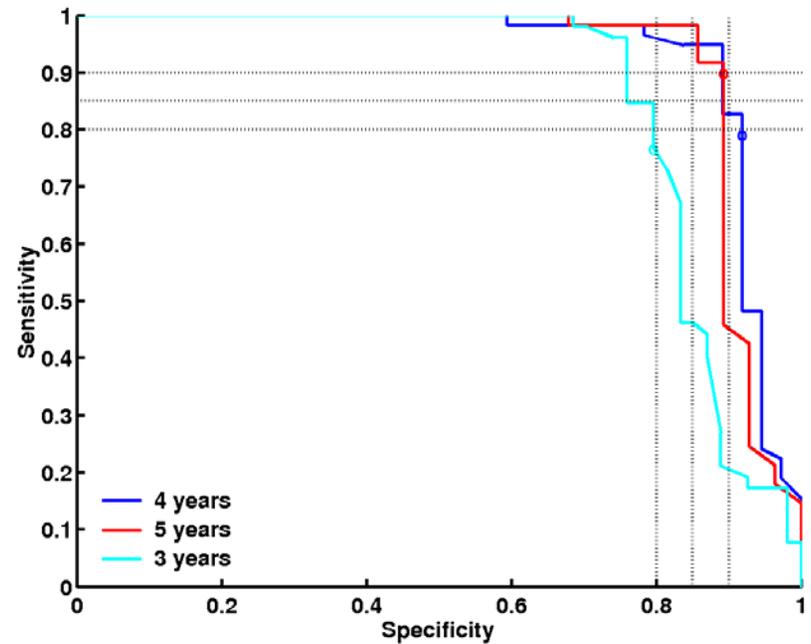
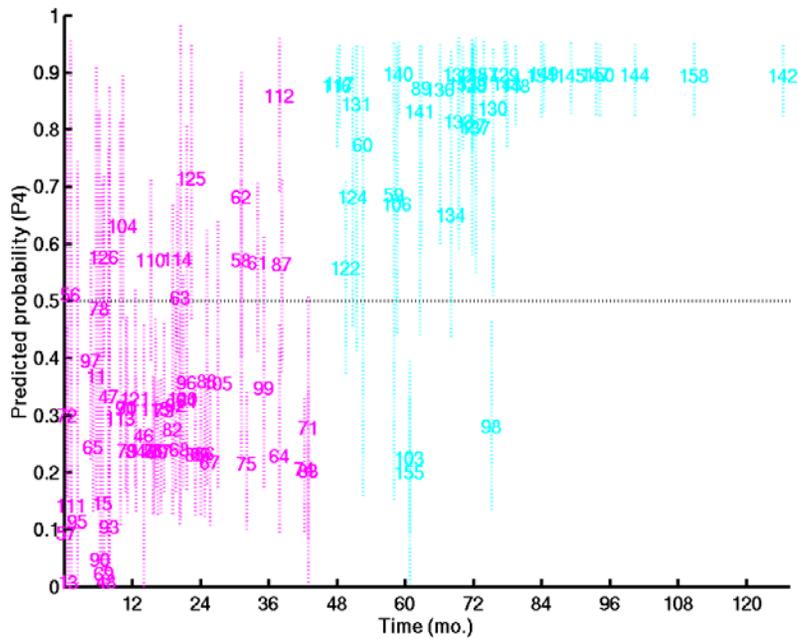
Mg20= -.04

... etc



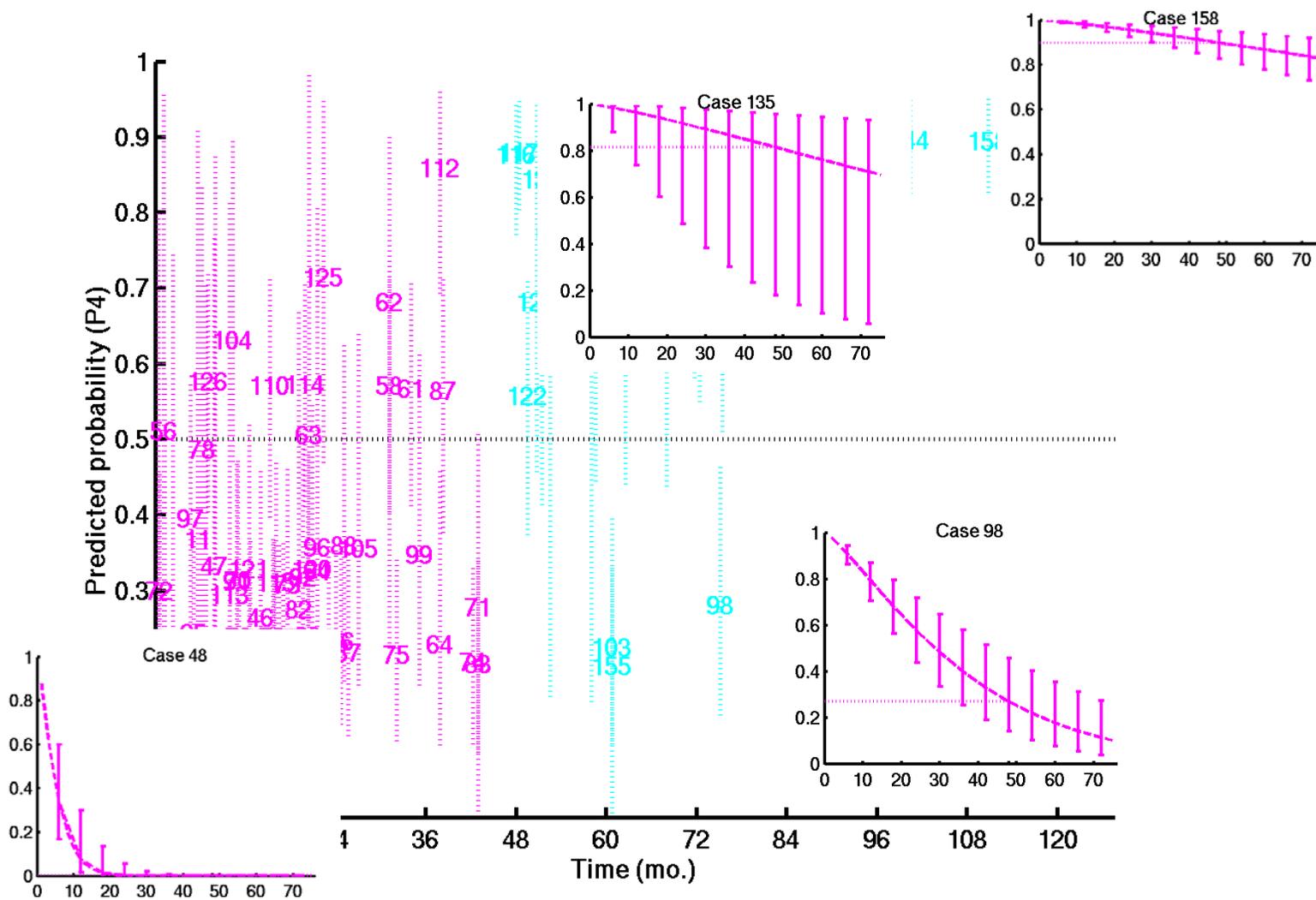
Predictions: Uncertainty & Assessment

Out-of-sample cross validation

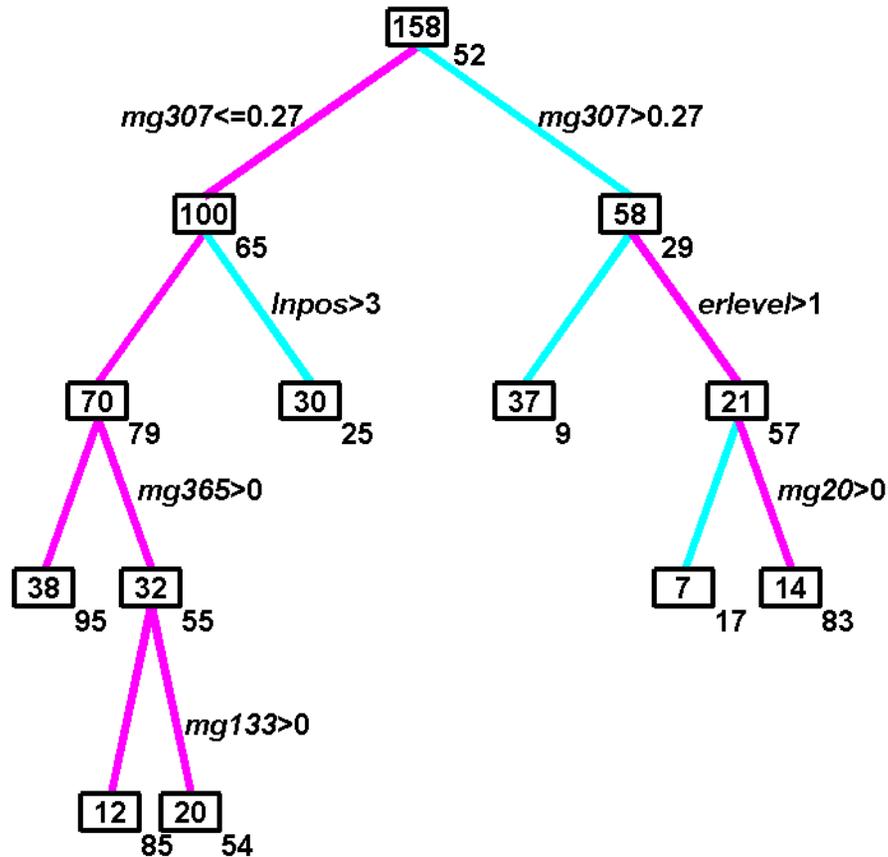


Recur+
Recur-

Personalised Prognosis



Biology behind Metagenes



*Erb-b2/Her-2
nu metagene*



Gene Identification

- Biological Signatures -

Erb-B2/Her2-nu metagene

Several ER metagenes

Recurrence metagenes: Myc oncogene, E2Fs

Lymph node & recurrence predictors

- *metagenes predictive of lymph node status*

 - ... *"surrogate" for node counts ...*

- *BRCA1, p27*

- *Interferon inducible genes*

Some Current Related Studies

- Validation studies
- Combining data - Multi-center studies
- Ethnic variation in genomic patterns
- Treatment trials - Improved risk stratification
- Treatment response predictors

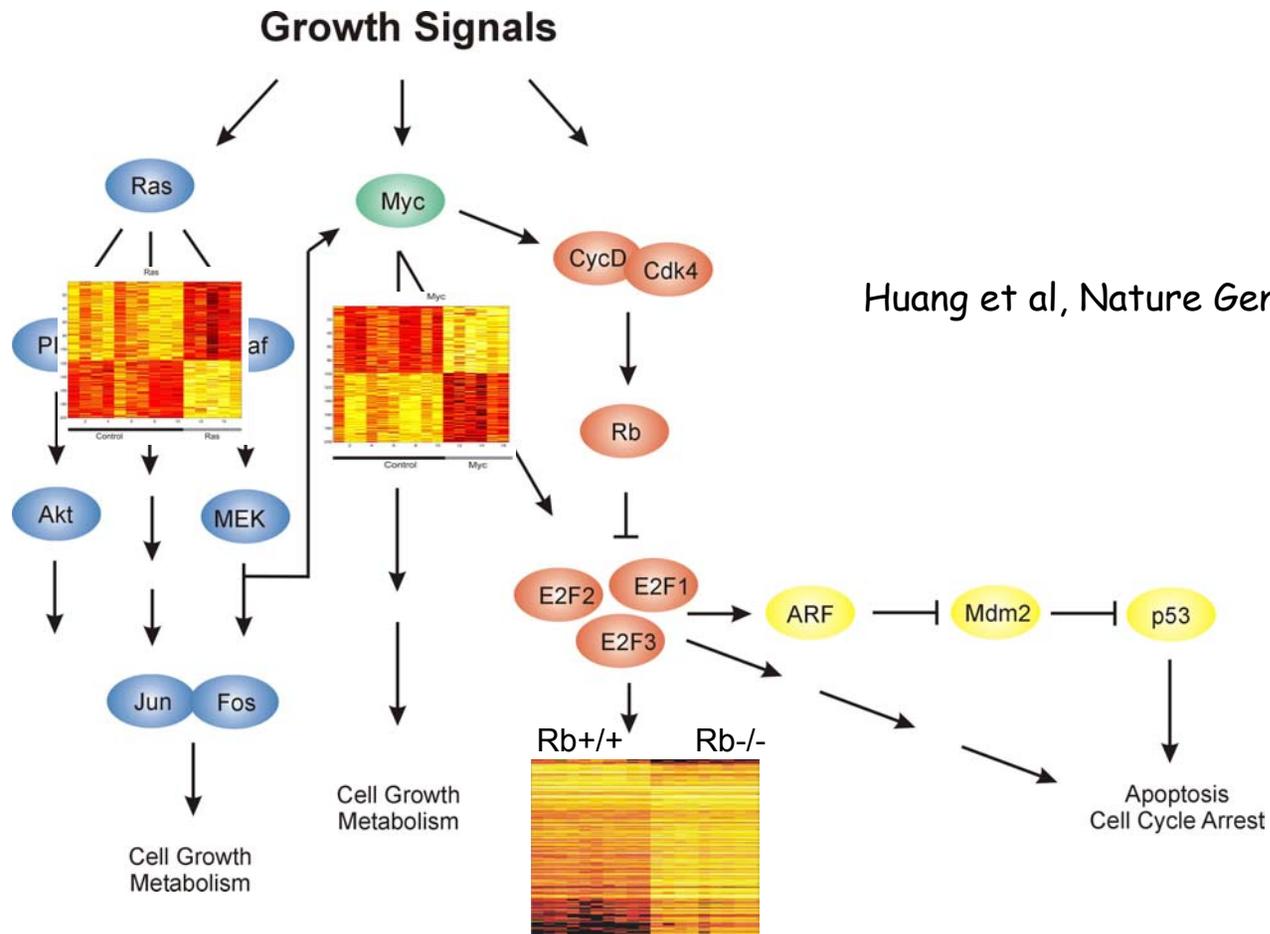
Patterns related to genes involved in cell proliferation -

- > improved biological metagene characterisations
- > refined "subtyping":

Myc, Ras, Rb/E2F, ..., metagenes

(Nature Genetics 2003, Cancer Research 2003)

Signaling Pathways That Control Cell Proliferation



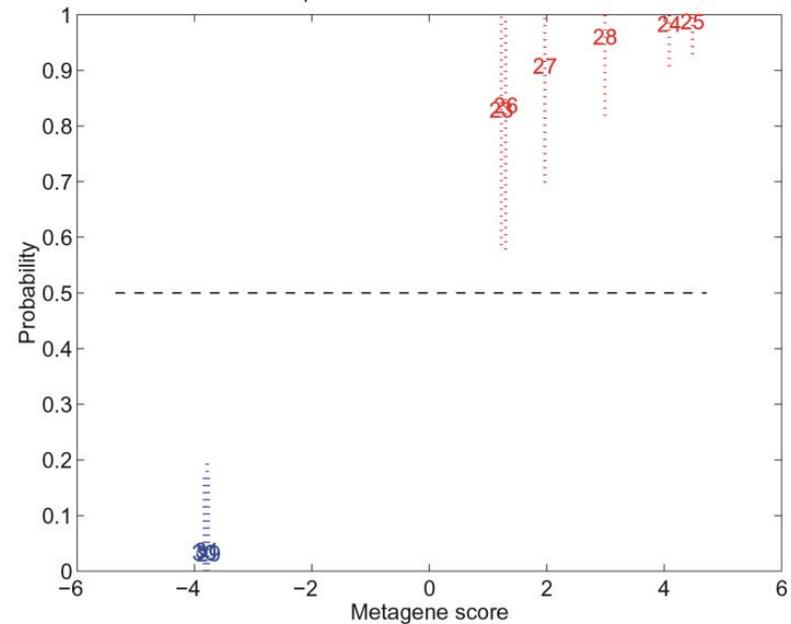
An Rb Metagene Predicts Loss of Rb Function in Tumors

Human

Rb+/- → Rb-/- Retinoblastoma

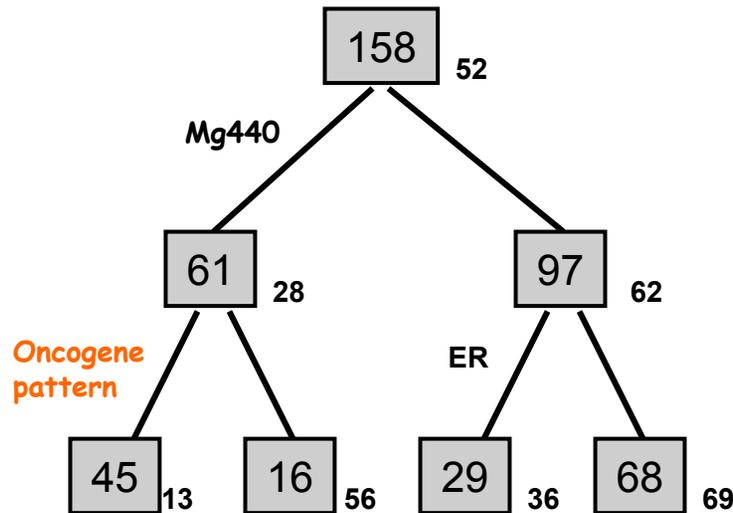
Mouse

Rb+/- → Rb-/- Pituitary tumors
Thyroid tumors



Black et al, Cancer Research 2003

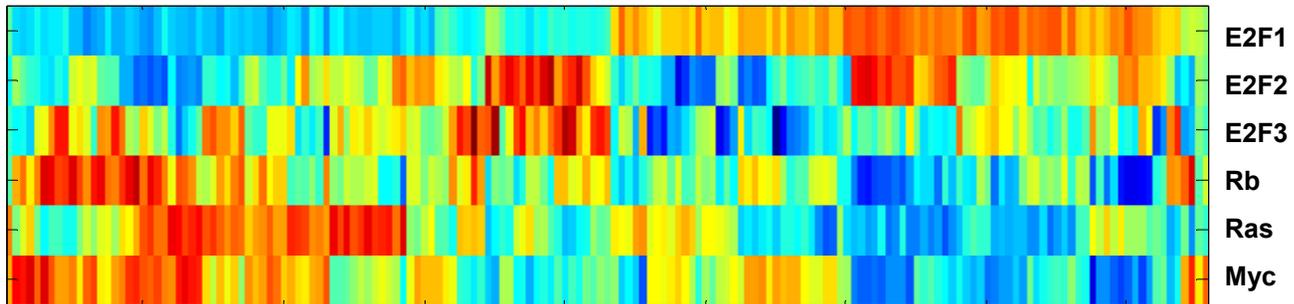
Potential for Improved Characterisation and Prediction



... reflecting multiple interacting aspects of tumour genomics

... and also potential to aid in identification of therapeutic targets

Pathway Predictions and MetaPathways



Tumor

Using Pathway Information to Guide Treatment?

Metastatic Breast Cancer Patients

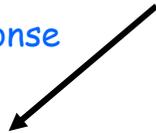


Biopsy of metastatic tumour

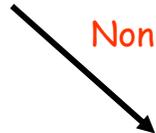


Treat with Drug A

Response



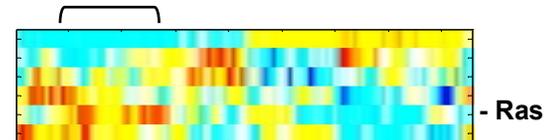
Non-response



Pathway Prediction



Pathway Specific Drug



Current Foci

- Statistics & Computation -

- Model refinements

 - smooth partition trees; aggregation of predictors

- Stochastic search

 - MCMC analysis, and search & annealing to explore space of candidate trees, and posterior over trees

- Cluster/distributed computing

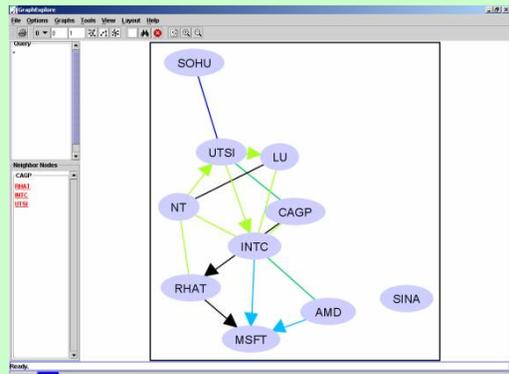
- Defining and evaluating patterns

 - improved metagene models: factor methods, large-scale graphical models, association network

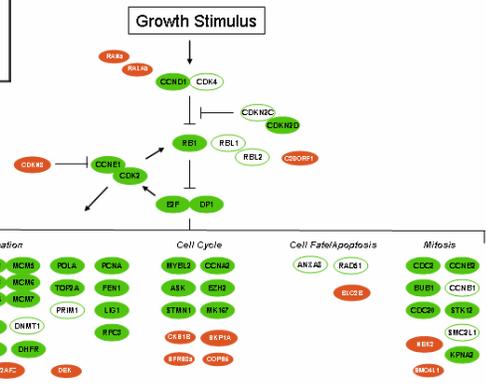
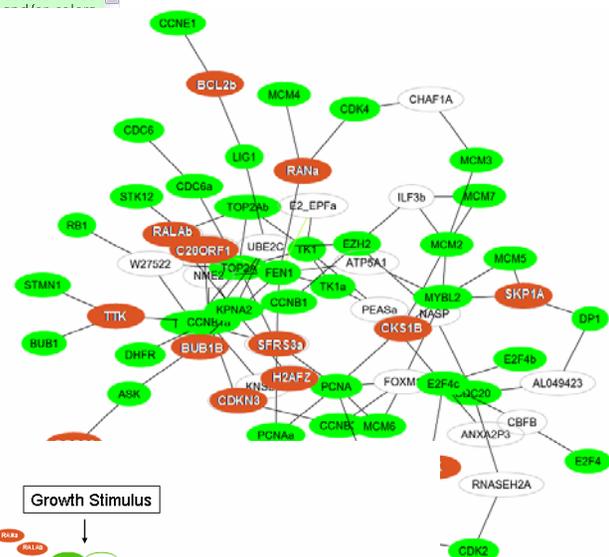
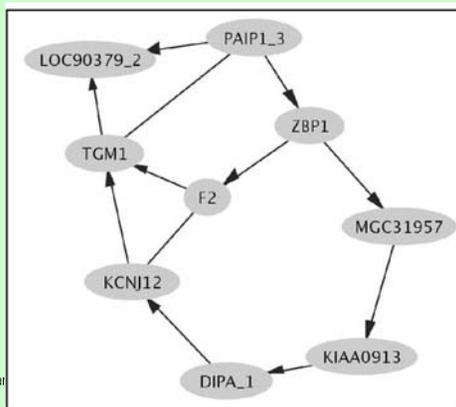
GraphExplore

A very important and attractive feature of GraphExplore is that it lets you choose specific colors and shapes for each object in the network. You can access these functions from the menu by going to **Options > Node > Color** or **Options > Node > Shape**. Moreover, GraphExplore lets you assign the same shapes to groups/clusters of objects. These groups can be different than the clusters you loaded with a new project. Therefore objects with the same shape or color, objects with the same color can identify another clustering and both of these clusterings can be different than the clusters loaded with the degree of flexibility is necessary to create meaningful displays of objects having different functions.

You can begin by typing "*" in the **Query** box and select **Graphs > Subgraph** to create a display of the entire network. Remark that all the objects all have the same default color (light blue).

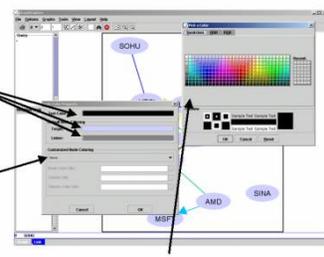


Bring up the colors dialog box by going to **Options > Node > Color**. You can format.



Check here to select new default colors for text, targets and linkers.

Click here to further customize the colors of nodes.



Large-scale graphical model search and evaluation
 Inference on large, sparse inverse covariance matrix
 (Dobra et al, JMVA 2004)

Duke University

Joseph Nevins

Erich Huang

Holly Dressman

Penni Black

Guang Yao

Molecular Genetics
Cancer Genomics

Jennifer Pitman

Adrian Dobra

Quanli Wang

Chris Hans

Carlos Carvalho

Statistics

Computational & Applied Genomics Program

Koo Foundation-Sun Yat Sen Cancer Center

Andrew Huang, Skye Cheng, Mei-Hua Tsou