

Lecture 3: Predictive models in cancer informatics

Jean-Philippe Vert

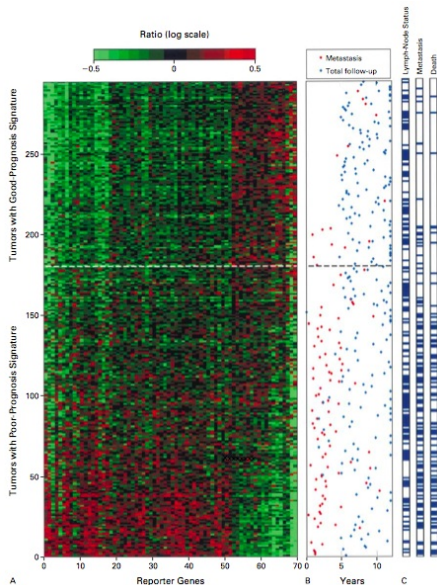
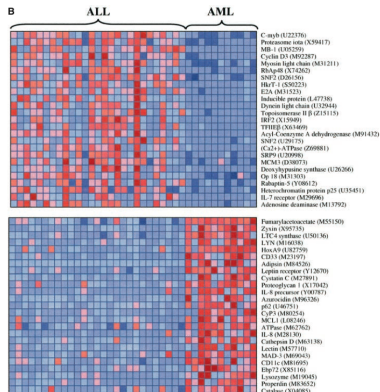
Mines ParisTech / Curie Institute / Inserm
Paris, France

"Optimization, machine learning and bioinformatics" summer
school, Erice, Sep 9-16, 2010.

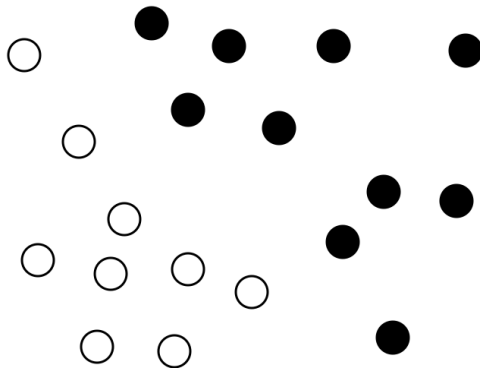
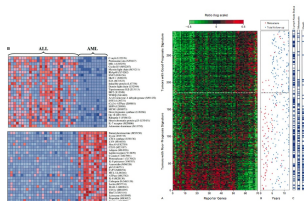
- 1 Shrinkage classifiers
- 2 Cancer prognosis from DNA copy number variations
- 3 Diagnosis and prognosis from gene expression data
- 4 Penalties for smooth classifiers
- 5 Penalties for structured feature selection
- 6 Conclusion

- 1 Shrinkage classifiers
- 2 Cancer prognosis from DNA copy number variations
- 3 Diagnosis and prognosis from gene expression data
- 4 Penalties for smooth classifiers
- 5 Penalties for structured feature selection
- 6 Conclusion

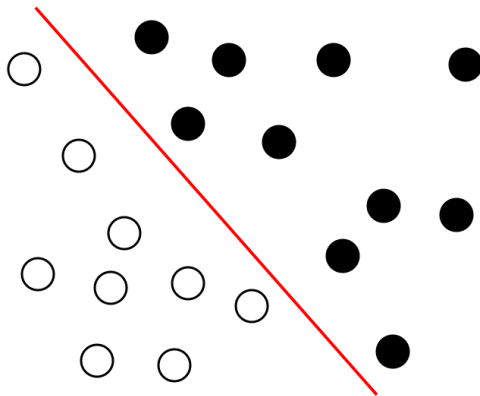
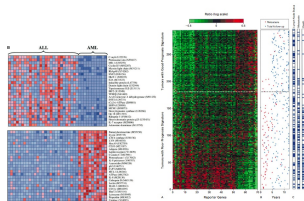
Molecular diagnosis / prognosis / theragnosis



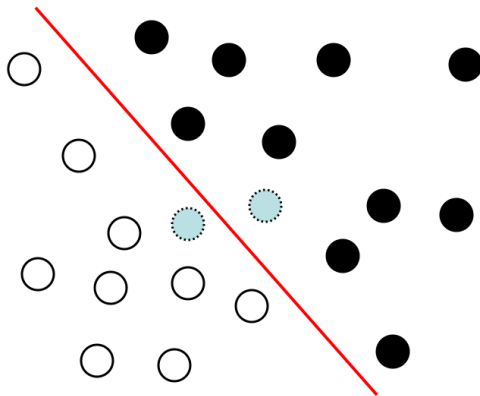
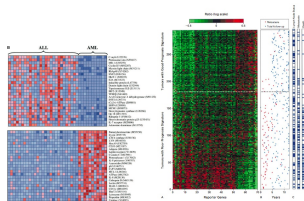
Pattern recognition, *aka* supervised classification



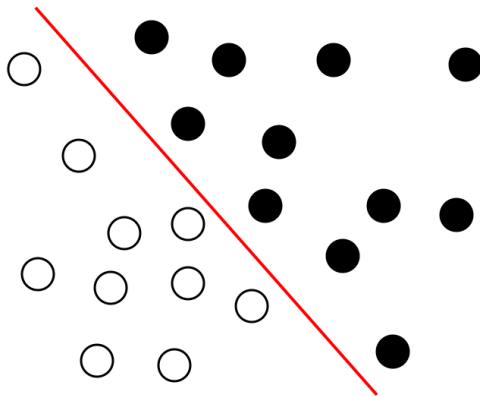
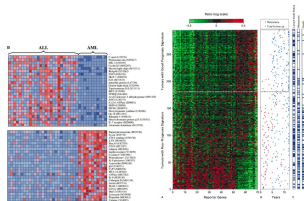
Pattern recognition, *aka* supervised classification

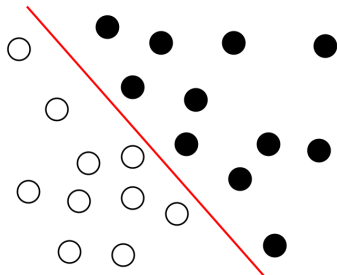


Pattern recognition, *aka* supervised classification



Pattern recognition, *aka* supervised classification





Challenges

- Few samples
- High dimension
- Structured data
- Heterogeneous data
- Prior knowledge
- Fast and scalable implementations
- Interpretable models

Shrinkage estimators

- Define a large family of "candidate classifiers", e.g., linear predictors $f_{\beta}(x) = \beta^{\top} x$
- For any candidate classifier f_{β} , quantify how "good" it is on the training set with some **empirical risk**, e.g.:

$$R(\beta) = \frac{1}{n} \sum_{i=1}^n l(f_{\beta}(x_i), y_i).$$

- Choose β that achieves the minimum empirical risk, subject to some **constraint**:

$$\min_{\beta} R(\beta) \quad \text{subject to} \quad \Omega(\beta) \leq C.$$

Shrinkage estimators

- Define a large family of "candidate classifiers", e.g., linear predictors $f_{\beta}(x) = \beta^{\top} x$
- For any candidate classifier f_{β} , quantify how "good" it is on the training set with some **empirical risk**, e.g.:

$$R(\beta) = \frac{1}{n} \sum_{i=1}^n l(f_{\beta}(x_i), y_i).$$

- Choose β that achieves the minimum empirical risk, subject to some **constraint**:

$$\min_{\beta} R(\beta) \quad \text{subject to} \quad \Omega(\beta) \leq C.$$

Shrinkage estimators

- Define a large family of "candidate classifiers", e.g., linear predictors $f_{\beta}(x) = \beta^{\top} x$
- For any candidate classifier f_{β} , quantify how "good" it is on the training set with some **empirical risk**, e.g.:

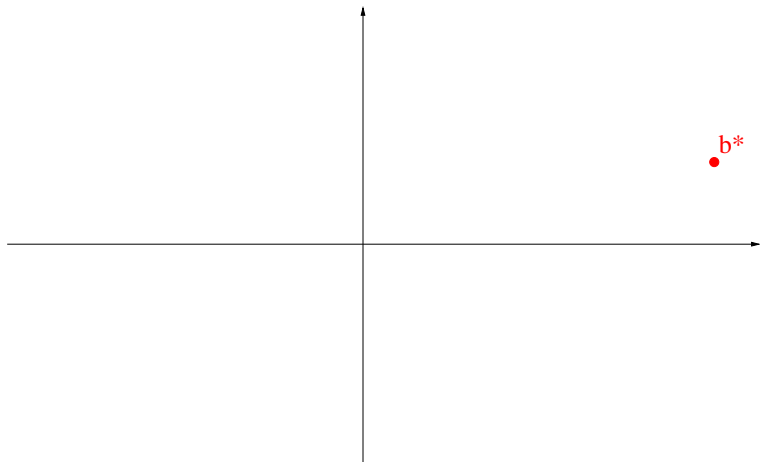
$$R(\beta) = \frac{1}{n} \sum_{i=1}^n l(f_{\beta}(x_i), y_i).$$

- Choose β that achieves the minimum empirical risk, subject to some **constraint**:

$$\min_{\beta} R(\beta) \quad \text{subject to} \quad \Omega(\beta) \leq C.$$

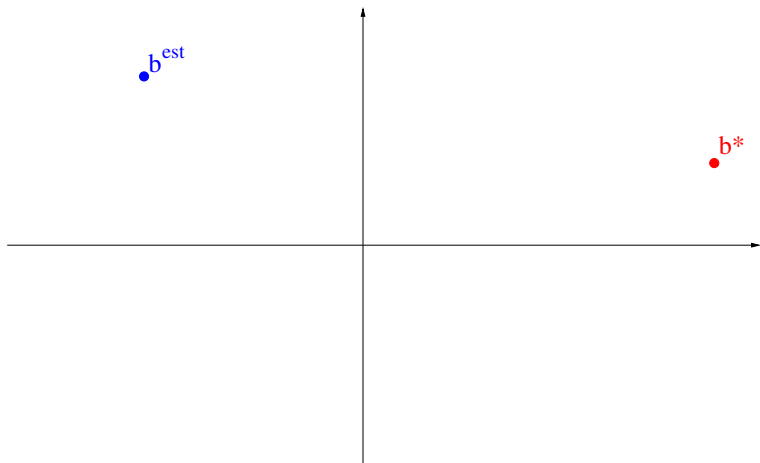
Why shrinkage classifiers?

$$\min_{\beta} R(\beta) \quad \text{subject to} \quad \Omega(\beta) \leq C.$$



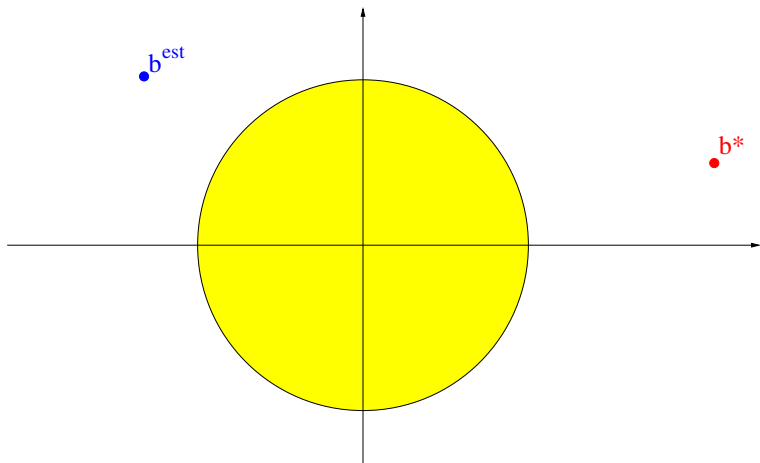
Why shrinkage classifiers?

$$\min_{\beta} R(\beta) \quad \text{subject to} \quad \Omega(\beta) \leq C.$$



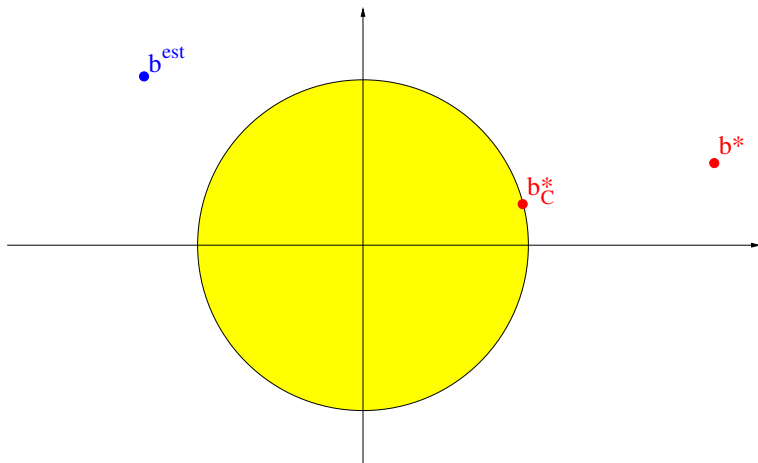
Why shrinkage classifiers?

$$\min_{\beta} R(\beta) \quad \text{subject to} \quad \Omega(\beta) \leq C.$$



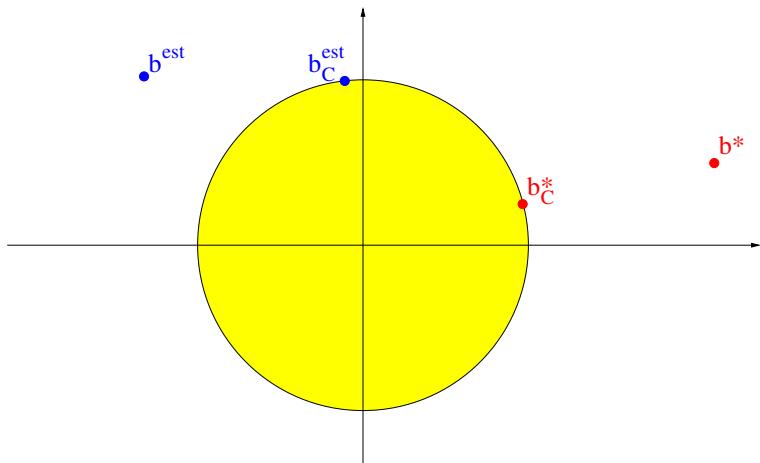
Why shrinkage classifiers?

$$\min_{\beta} R(\beta) \quad \text{subject to} \quad \Omega(\beta) \leq C.$$



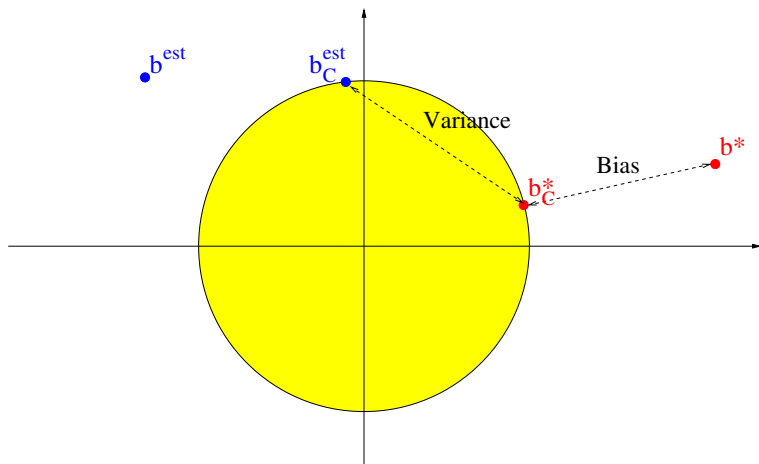
Why shrinkage classifiers?

$$\min_{\beta} R(\beta) \quad \text{subject to} \quad \Omega(\beta) \leq C.$$

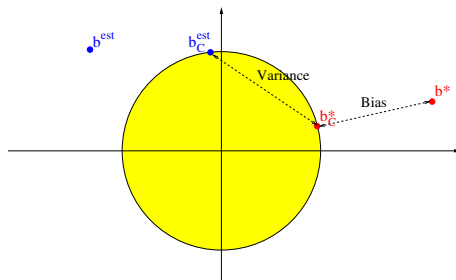


Why shrinkage classifiers?

$$\min_{\beta} R(\beta) \quad \text{subject to} \quad \Omega(\beta) \leq C.$$



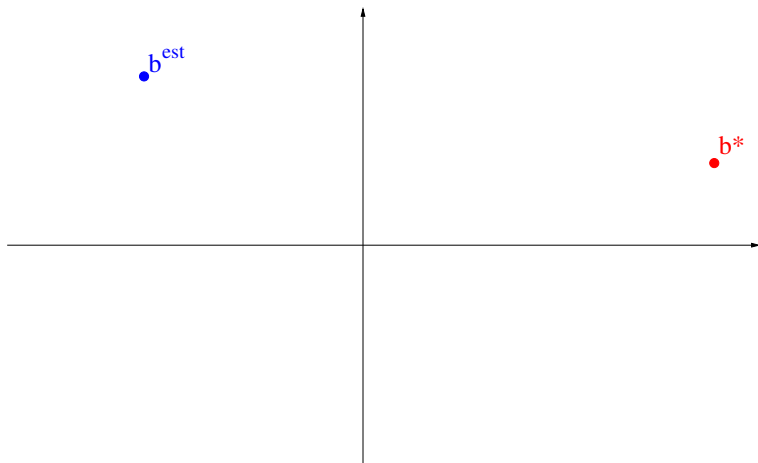
Why shrinkage classifiers?



- "Increases bias and decreases variance"
- Common choices are
 - $\Omega(\beta) = \sum_{i=1}^p \beta_i^2$ (ridge regression, SVM, ...)
 - $\Omega(\beta) = \sum_{i=1}^p |\beta_i|$ (lasso, boosting, ...)

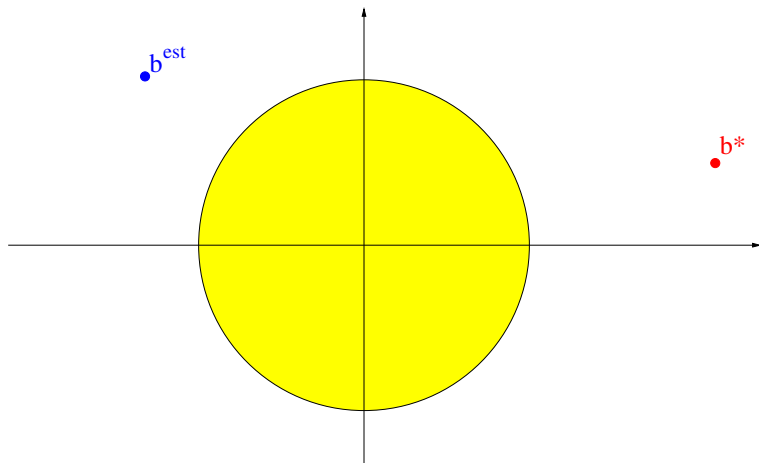
Including prior knowledge in the penalty?

$$\min_{\beta} R(\beta) \quad \text{subject to} \quad \Omega(\beta) \leq C.$$



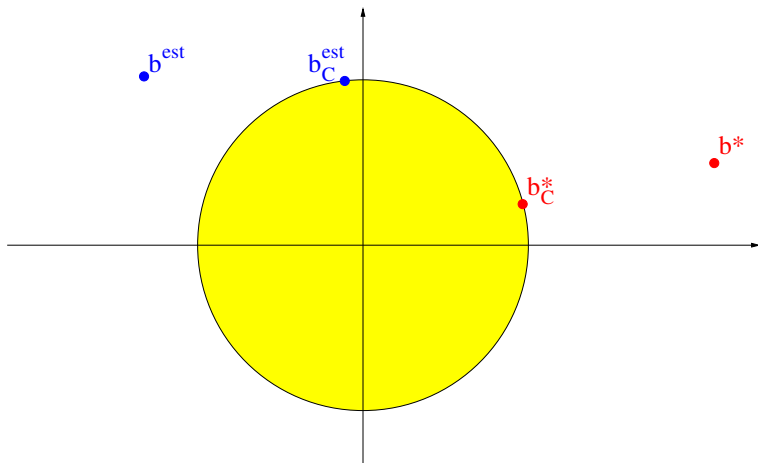
Including prior knowledge in the penalty?

$$\min_{\beta} R(\beta) \quad \text{subject to} \quad \Omega(\beta) \leq C.$$



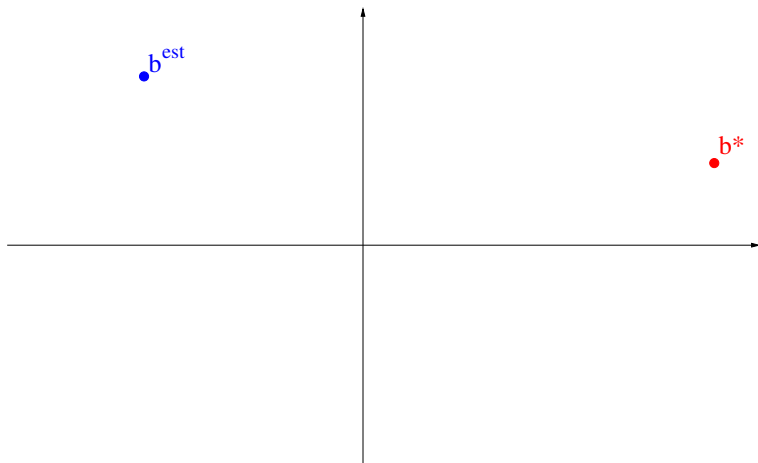
Including prior knowledge in the penalty?

$$\min_{\beta} R(\beta) \quad \text{subject to} \quad \Omega(\beta) \leq C.$$



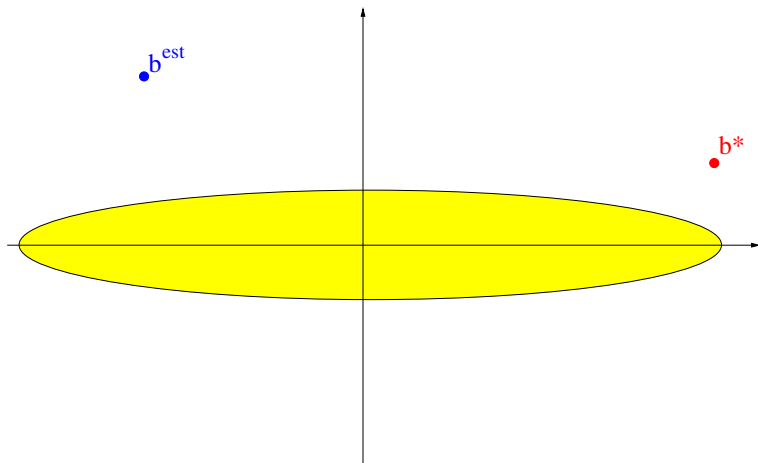
Including prior knowledge in the penalty?

$$\min_{\beta} R(\beta) \quad \text{subject to} \quad \Omega(\beta) \leq C.$$



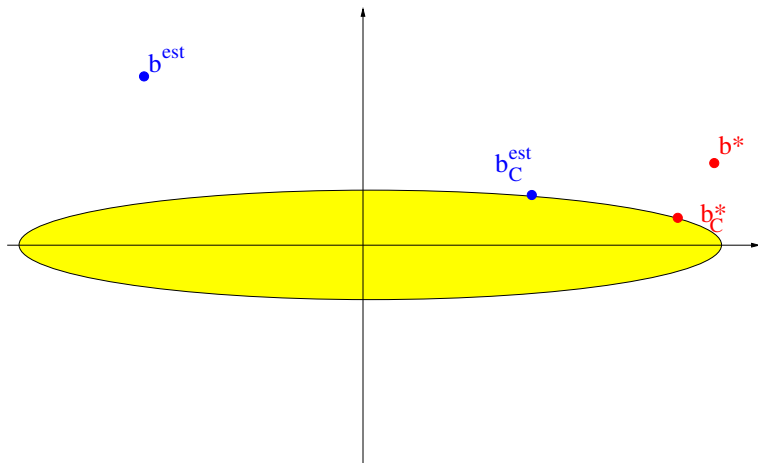
Including prior knowledge in the penalty?

$$\min_{\beta} R(\beta) \quad \text{subject to} \quad \Omega(\beta) \leq C.$$



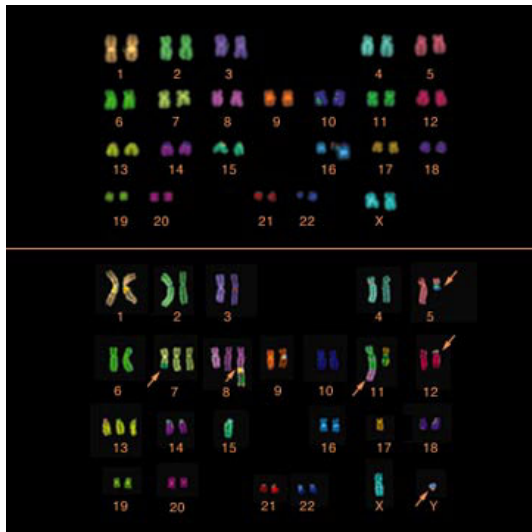
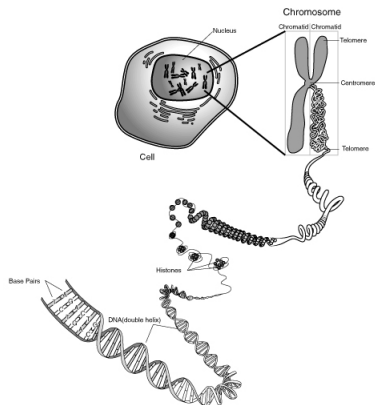
Including prior knowledge in the penalty?

$$\min_{\beta} R(\beta) \quad \text{subject to} \quad \Omega(\beta) \leq C.$$



- 1 Shrinkage classifiers
- 2 Cancer prognosis from DNA copy number variations**
- 3 Diagnosis and prognosis from gene expression data
- 4 Penalties for smooth classifiers
- 5 Penalties for structured feature selection
- 6 Conclusion

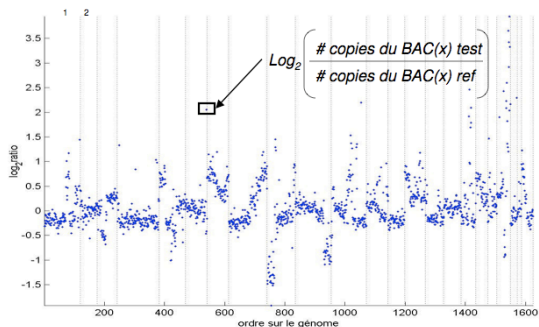
Chromosomal aberrations in cancer



Comparative Genomic Hybridization (CGH)

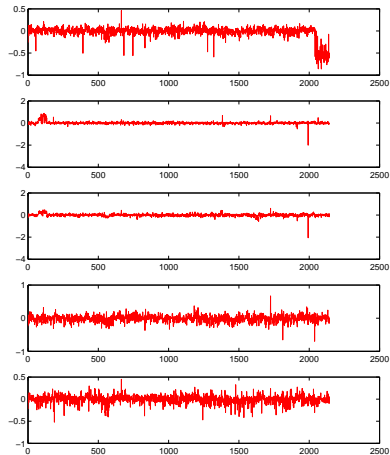
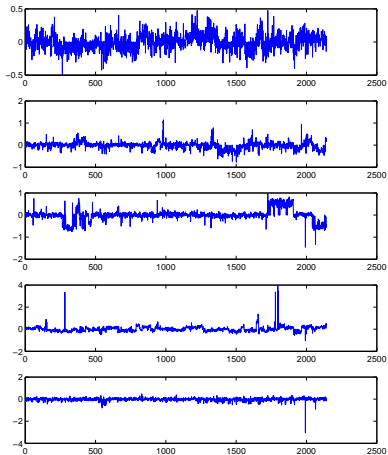
Motivation

- Comparative genomic hybridization (CGH) data measure the **DNA copy number** along the genome
- Very useful, in particular in cancer research
- Can we **classify CGH arrays** for diagnosis or prognosis purpose?



Jain et al. Genome research 2002 12:325-332

Aggressive vs non-aggressive melanoma



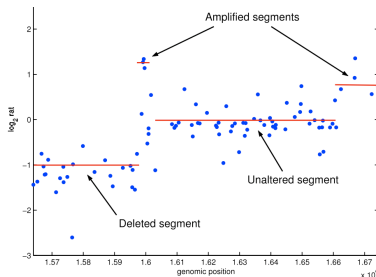
CGH array classification

Prior knowledge

- For a CGH profile $x \in \mathbb{R}^p$, we focus on linear classifiers, i.e., the sign of :

$$f_{\beta}(x) = \beta^T x .$$

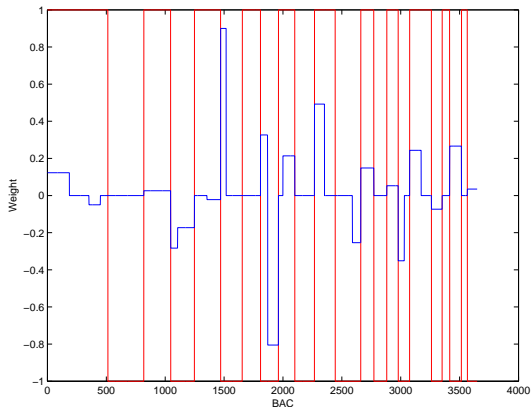
- We expect β to be
 - **sparse** : not all positions should be discriminative
 - **piecewise constant** : within a selected region, all probes should contribute equally



Fused lasso for supervised classification

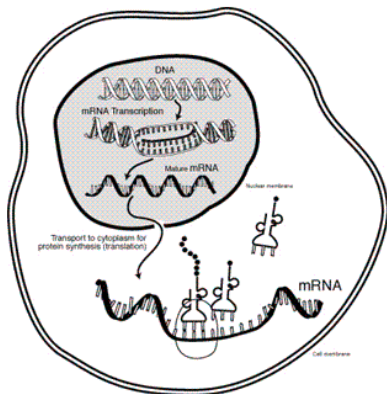
$$\min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n \ell(y_i, \beta^\top x_i) + \lambda_1 \sum_{i=1}^p |\beta_i| + \lambda_2 \sum_{i=1}^{p-1} |\beta_{i+1} - \beta_i|.$$

where ℓ is, e.g., the hinge loss $\ell(y, t) = \max(1 - yt, 0)$.



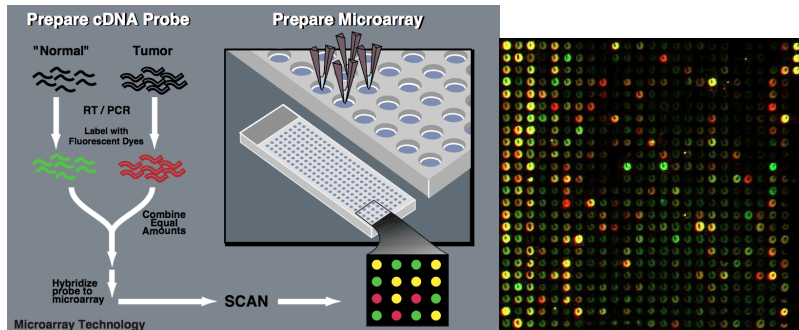
- 1 Shrinkage classifiers
- 2 Cancer prognosis from DNA copy number variations
- 3 Diagnosis and prognosis from gene expression data**
- 4 Penalties for smooth classifiers
- 5 Penalties for structured feature selection
- 6 Conclusion

DNA → RNA → protein



- CGH shows the (static) DNA
- Cancer cells have also **abnormal (dynamic) gene expression** (= transcription)

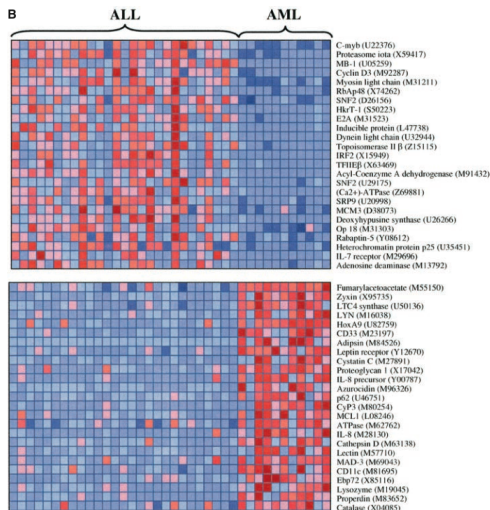
Tissue profiling with DNA chips



Data

- Gene expression measures for **more than 10k genes**
- Measured typically on **less than 100 samples** of two (or more) different classes (e.g., different tumors)

Tissue classification from microarray data



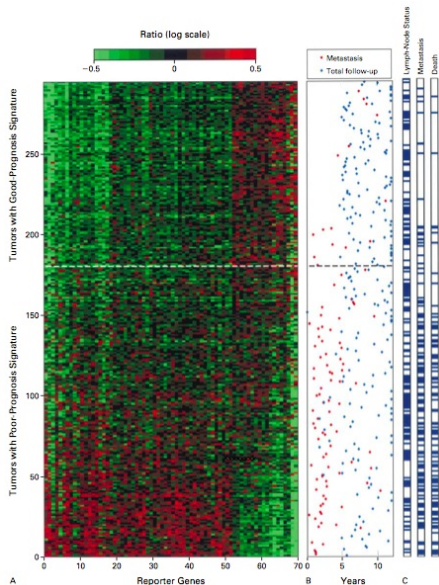
Goal

- Design a **classifier** to automatically assign a class to future samples from their expression profile
- **Interpret** biologically the differences between the classes

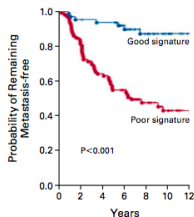
Difficulty

- Large dimension
- Few samples

Prognosis from microarray data (MAMMAPRINT)



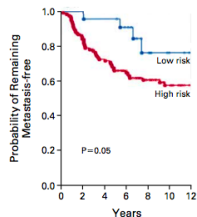
A Gene-Expression Profiling



No. AT RISK

Good signature	60	57	54	45	31	22	12
Poor signature	91	72	55	41	26	17	9

B St. Gallen Criteria



No. AT RISK

Low risk	22	22	21	17	9	5	2
High risk	129	107	88	69	48	34	19

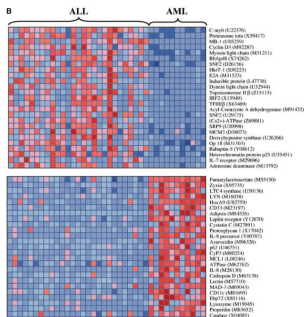
The idea

- We look for a limited set of genes that are sufficient for prediction.
- Equivalently, the linear classifier will be **sparse**

Motivations

- **Bet on sparsity**: we believe the "true" model is sparse.
- **Interpretation**: we will get a biological interpretation more easily by looking at the selected genes.
- **Statistics**: by restricting the class of classifiers, we **increase the bias but decrease the variance**. This should be helpful in large dimensions.

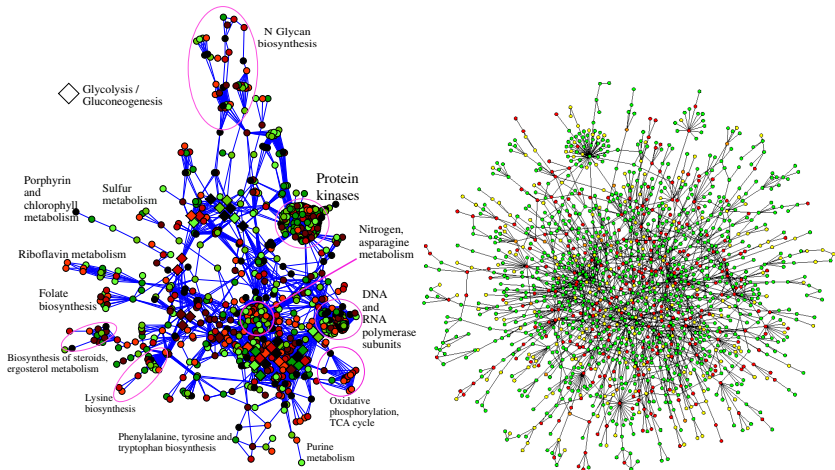
But...



Challenging the idea of gene signature

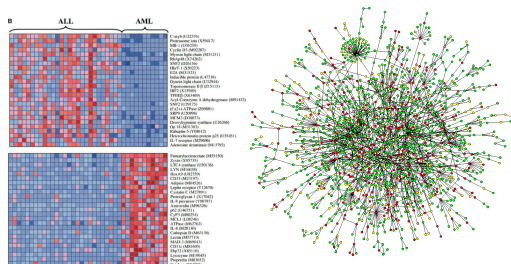
- We often observe little **stability** in the genes selected...
- Is gene selection the most **biologically relevant** hypothesis?
- What about thinking instead of "**pathways**" or "**modules**" **signatures**?

Gene networks



Motivation

- Basic biological functions usually involve the **coordinated action of several proteins**:
 - Formation of **protein complexes**
 - Activation of metabolic, signalling or regulatory **pathways**
- Many pathways and protein-protein interactions are **already known**
- Hypothesis**: the weights of the classifier should be “coherent” with respect to this **prior knowledge**



$$\min_{\beta} R(\beta) + \lambda \Omega_G(\beta)$$

Hypothesis

We would like to design penalties $\Omega_G(\beta)$ to promote one of the following hypothesis:

- **Hypothesis 1**: genes near each other on the graph should have **similar weights** (but we do not try to select only a few genes), i.e., the classifier should be **smooth** on the graph
- **Hypothesis 2**: genes selected in the signature should be **connected** to each other, or be in **a few known functional groups**, without necessarily having similar weights.

- 1 Shrinkage classifiers
- 2 Cancer prognosis from DNA copy number variations
- 3 Diagnosis and prognosis from gene expression data
- 4 Penalties for smooth classifiers**
- 5 Penalties for structured feature selection
- 6 Conclusion

Prior hypothesis

Genes near each other on the graph should have **similar weights**.

An idea (Rapaport et al., 2007)

$$\Omega_{\text{spectral}}(\beta) = \sum_{i \sim j} (\beta_i - \beta_j)^2,$$

$$\min_{\beta \in \mathbb{R}^p} R(\beta) + \lambda \sum_{i \sim j} (\beta_i - \beta_j)^2.$$

Prior hypothesis

Genes near each other on the graph should have **similar weights**.

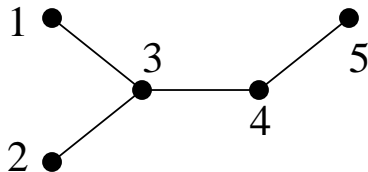
An idea (Rapaport et al., 2007)

$$\Omega_{\text{spectral}}(\beta) = \sum_{i \sim j} (\beta_i - \beta_j)^2,$$

$$\min_{\beta \in \mathbb{R}^p} R(\beta) + \lambda \sum_{i \sim j} (\beta_i - \beta_j)^2.$$

Definition

The Laplacian of the graph is the matrix $L = D - A$.



$$L = D - A = \begin{pmatrix} 1 & 0 & -1 & 0 & 0 \\ 0 & 1 & -1 & 0 & 0 \\ -1 & -1 & 3 & -1 & 0 \\ 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & 1 & 1 \end{pmatrix}$$

Theorem

The function $f(x) = \beta^\top x$ where β is solution of

$$\min_{\beta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n l(\beta^\top x_i, y_i) + \lambda \sum_{i \sim j} (\beta_i - \beta_j)^2$$

is equal to $g(x) = \gamma^\top \Phi(x)$ where γ is solution of

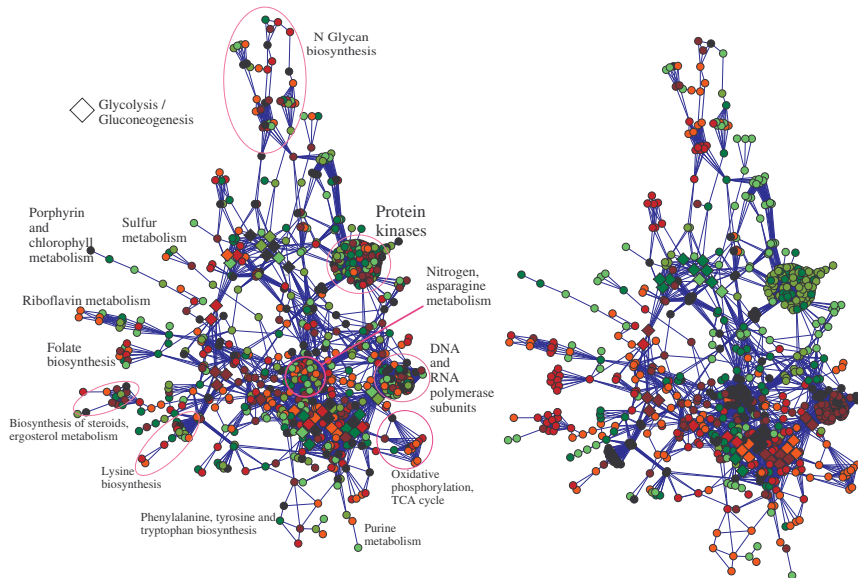
$$\min_{\gamma \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n l(\gamma^\top \Phi(x_i), y_i) + \lambda \gamma^\top \gamma,$$

and where

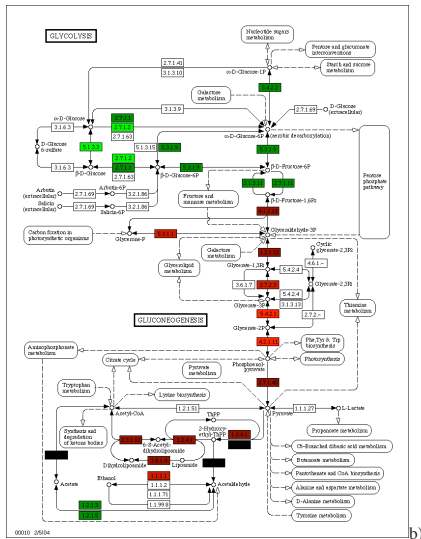
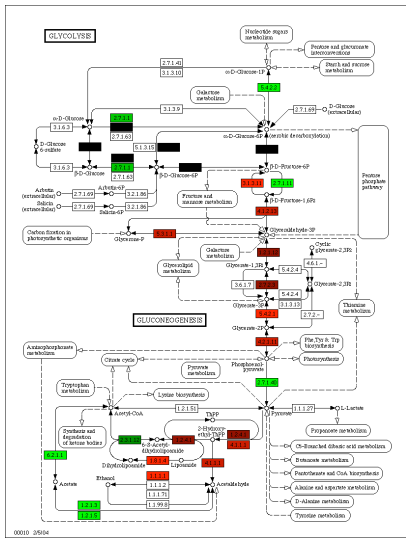
$$\Phi(x)^\top \Phi(x') = x^\top K_G x'$$

for $K_G = L^*$, the pseudo-inverse of the graph Laplacian.

Classifiers



Classifier



$$\Phi(x)^\top \Phi(x') = x^\top K_G x'$$

with:

- $K_G = (c + L)^{-1}$ leads to

$$\Omega(\beta) = c \sum_{i=1}^p \beta_i^2 + \sum_{i \sim j} (\beta_i - \beta_j)^2 .$$

- The diffusion kernel:

$$K_G = \exp_M(-2tL) .$$

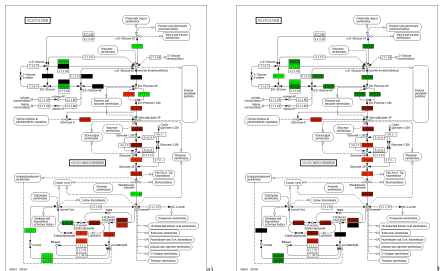
penalizes high frequencies of β in the Fourier domain.

- Gene selection + Piecewise constant on the graph

$$\Omega(\beta) = \sum_{i \sim j} |\beta_i - \beta_j| + \sum_{i=1}^p |\beta_i|$$

- Gene selection + smooth on the graph

$$\Omega(\beta) = \sum_{i \sim j} (\beta_i - \beta_j)^2 + \sum_{i=1}^p |\beta_i|$$



- We are happy to see pathways appear.
- However, in some cases, connected genes should have "opposite" weights (inhibition, pathway branching, etc...)
- **How to capture pathways without constraints on the weight similarities?**

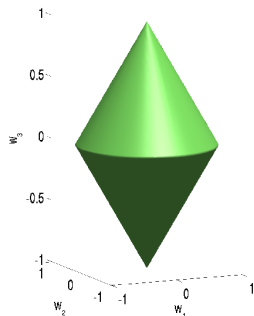
- 1 Shrinkage classifiers
- 2 Cancer prognosis from DNA copy number variations
- 3 Diagnosis and prognosis from gene expression data
- 4 Penalties for smooth classifiers
- 5 Penalties for structured feature selection**
- 6 Conclusion

Selecting pre-defined groups of variables

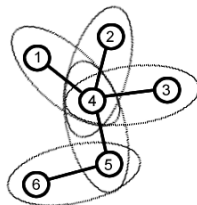
Group lasso (Yuan & Lin, 2006)

If groups of covariates are likely to be selected together, the l_1/l_2 -norm induces sparse solutions *at the group level*:

$$\Omega_{group}(w) = \sum_g \|w_g\|_2$$



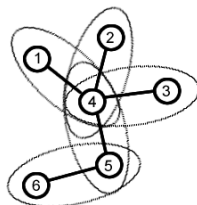
$$\Omega(w_1, w_2, w_3) = \|(w_1, w_2)\|_2 + \|w_3\|_2$$



- **Hypothesis:** selected genes should form connected components on the graph
- Two solutions (Jacob et al., 2009):

$$\Omega_{group}(\beta) = \sum_{i \sim j} \sqrt{\beta_i^2 + \beta_j^2},$$

$$\Omega_{overlap}(\beta) = \sup_{\alpha \in \mathbb{R}^p: \forall i \sim j, \|\alpha_i^2 + \alpha_j^2\| \leq 1} \alpha^T \beta.$$

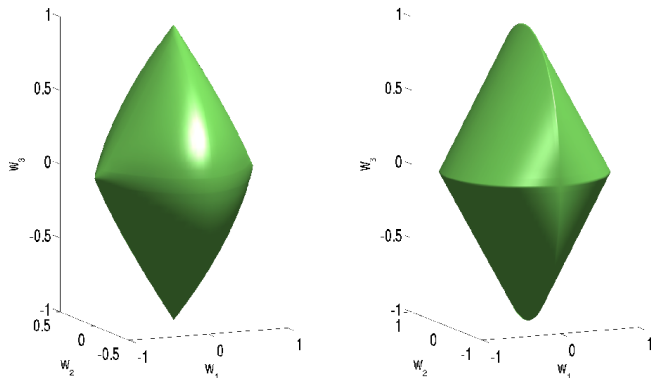


- **Hypothesis:** selected genes should form connected components on the graph
- Two solutions (Jacob et al., 2009):

$$\Omega_{group}(\beta) = \sum_{i \sim j} \sqrt{\beta_i^2 + \beta_j^2},$$

$$\Omega_{overlap}(\beta) = \sup_{\alpha \in \mathbb{R}^p: \forall i \sim j, \|\alpha_i^2 + \alpha_j^2\| \leq 1} \alpha^\top \beta.$$

Overlap and group unity balls



Balls for $\Omega_{\text{group}}^{\mathcal{G}}(\cdot)$ (middle) and $\Omega_{\text{overlap}}^{\mathcal{G}}(\cdot)$ (right) for the groups $\mathcal{G} = \{\{1, 2\}, \{2, 3\}\}$ where w_2 is represented as the vertical coordinate.

Summary: Graph lasso vs kernel

- Graph lasso:

$$\Omega_{\text{graph lasso}}(\mathbf{w}) = \sum_{i \sim j} \sqrt{w_i^2 + w_j^2}.$$

constrains the **sparsity**, not the values

- Graph kernel

$$\Omega_{\text{graph kernel}}(\mathbf{w}) = \sum_{i \sim j} (w_i - w_j)^2.$$

constrains the values (**smoothness**), not the sparsity

Breast cancer data

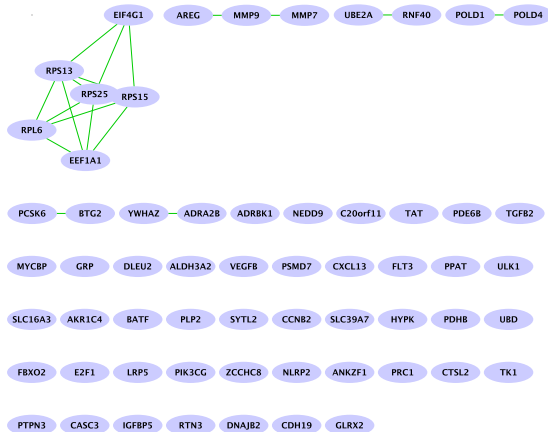
- Gene expression data for 8,141 genes in 295 breast cancer tumors.
- Canonical pathways from MSigDB containing 639 groups of genes, 637 of which involve genes from our study.

METHOD	l_1	$\Omega_{\text{OVERLAP}}^G(\cdot)$
ERROR	0.38 ± 0.04	0.36 ± 0.03
MEAN \ddagger PATH.	130	30

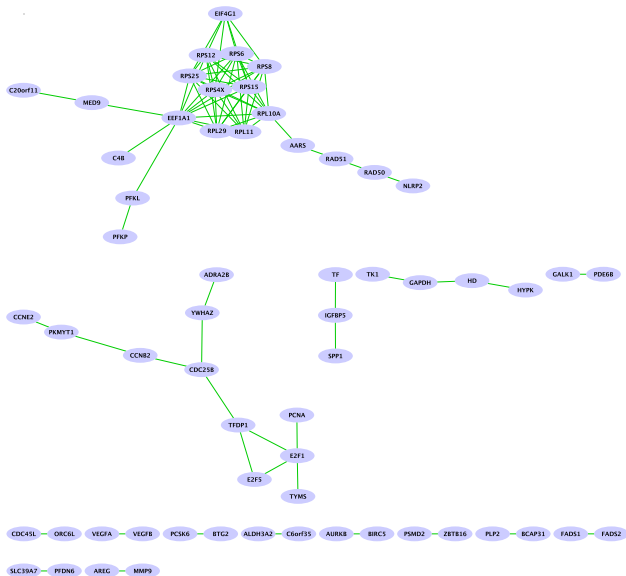
- Graph on the genes.

METHOD	l_1	$\Omega_{\text{graph}}(\cdot)$
ERROR	0.39 ± 0.04	0.36 ± 0.01
AV. SIZE C.C.	1.03	1.30

Lasso signature



Graph Lasso signature



- 1 Shrinkage classifiers
- 2 Cancer prognosis from DNA copy number variations
- 3 Diagnosis and prognosis from gene expression data
- 4 Penalties for smooth classifiers
- 5 Penalties for structured feature selection
- 6 Conclusion**

- Modern machine learning methods for regression / classification lend themselves well to the **integration of prior knowledge** in the penalization / regularization function.
- Several **computationally efficient** approaches (structured LASSO, kernels...)
- Tight collaborations with domain experts can help develop specific learning machines for specific data
- Natural extensions for **data integration**

People I need to thank



Franck Rapaport (MSKCC), Emmanuel Barillot, Andrei Zynoviev Kevin Bleakley, Anne-Claire Haury (Institut Curie / ParisTech), Laurent Jacob (UC Berkeley) Guillaume Obozinski (INRIA)