

# Machine learning in cancer genomics

Jean-Philippe Vert

`Jean-Philippe.Vert@mines.org`

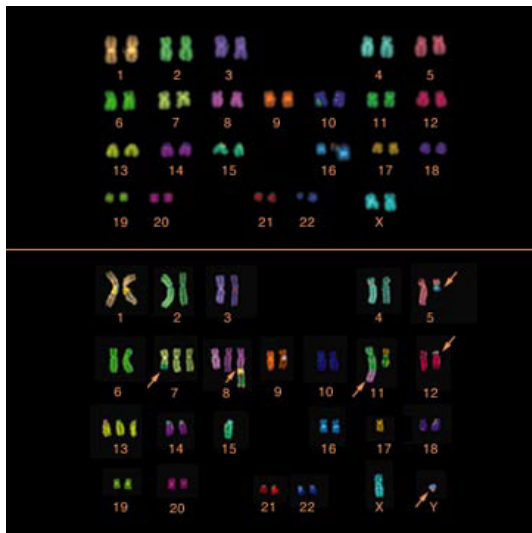
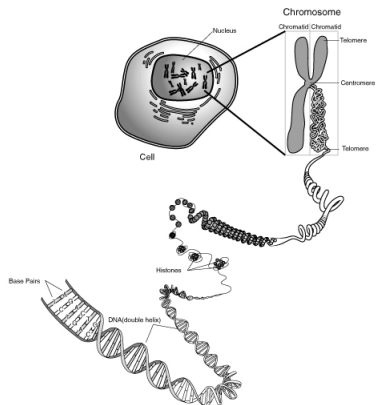
Mines ParisTech / Curie Institute / Inserm

Université catholique de Louvain, Nov 14, 2011.

- 1 Introduction
- 2 Cancer prognosis from DNA copy number variations
- 3 Diagnosis and prognosis from gene expression data
- 4 Conclusion

- 1 Introduction
- 2 Cancer prognosis from DNA copy number variations
- 3 Diagnosis and prognosis from gene expression data
- 4 Conclusion

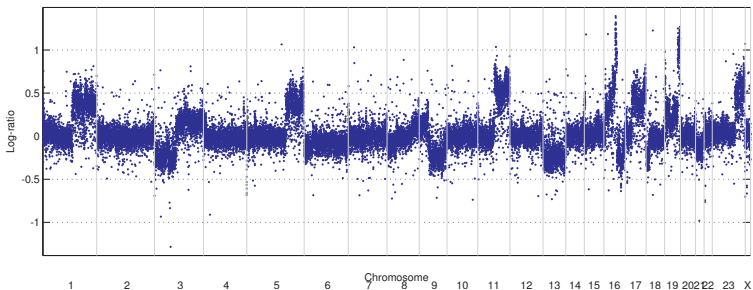
# Chromosomal aberrations in cancer



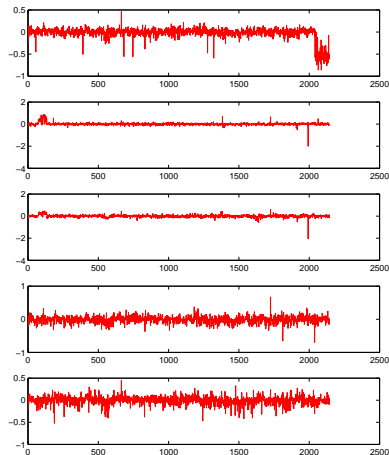
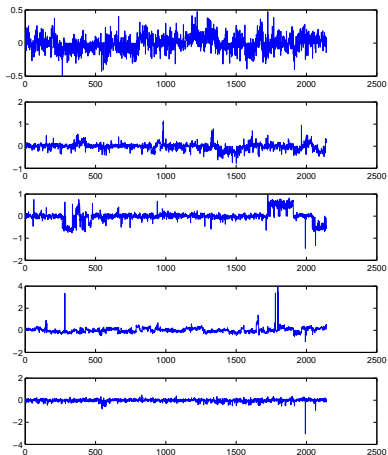
# Comparative Genomic Hybridization (CGH)

## Motivation

- Comparative genomic hybridization (CGH) data measure the **DNA copy number** along the genome
- Very useful, in particular in cancer research to observe systematically variants in DNA content

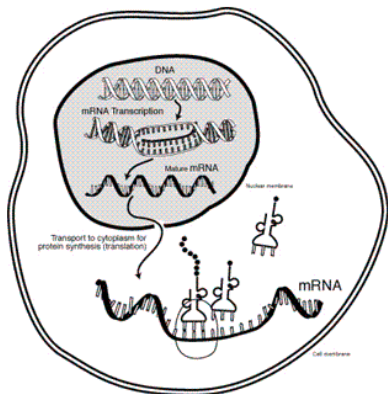


# Cancer prognosis: can we predict the future evolution?



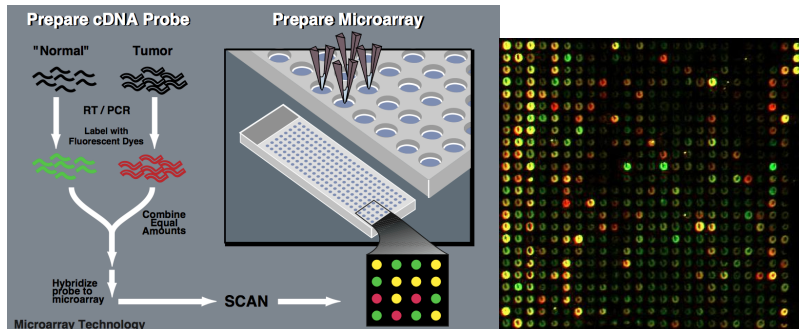
*Aggressive (left) vs non-aggressive (right) melanoma*

# DNA → RNA → protein



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

# Tissue profiling with DNA chips

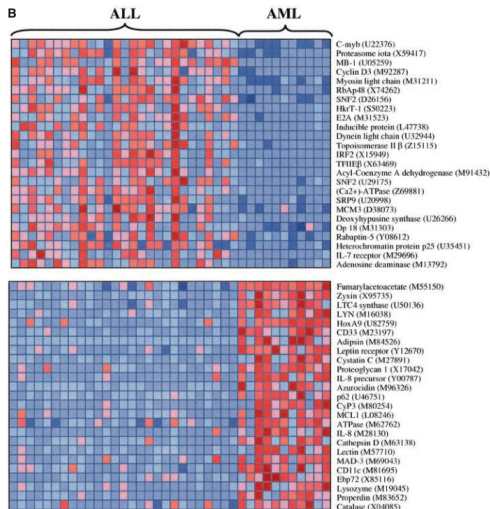


## Data

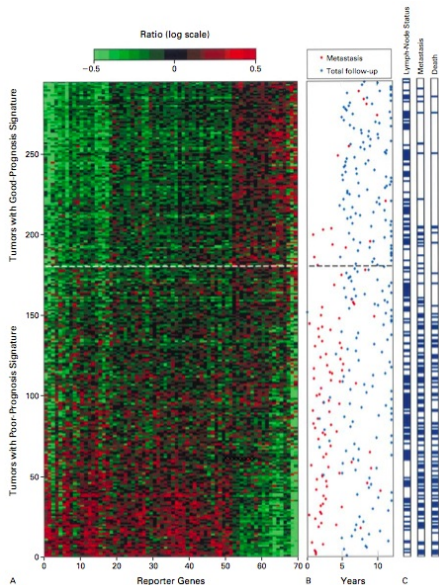
- Gene expression measures for **more than 10k genes**
- Measured typically on **less than 100 samples** of two (or more) different classes (e.g., different tumors)



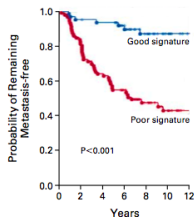
# Can we identify the cancer subtype? (diagnosis)



# Can we predict the future evolution? (prognosis)

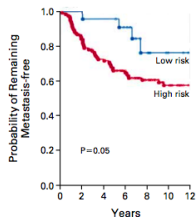


A Gene-Expression Profiling



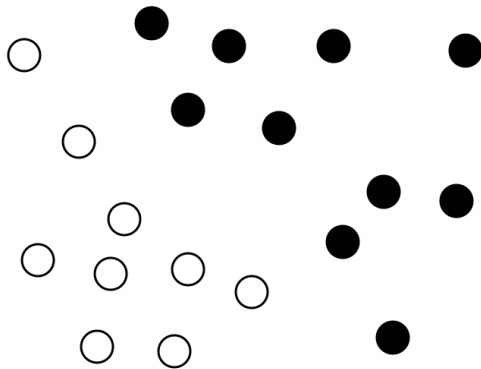
No. AT RISK	0	2	4	6	8	10	12
Good signature	60	57	54	45	31	22	12
Poor signature	91	72	55	41	26	17	9

B St. Gallen Criteria



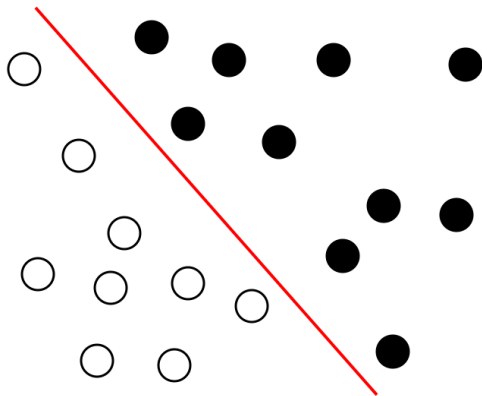
No. AT RISK	0	2	4	6	8	10	12
Low risk	22	22	21	17	9	5	2
High risk	129	107	88	69	48	34	19

# Machine learning (pattern recognition / supervised classification)



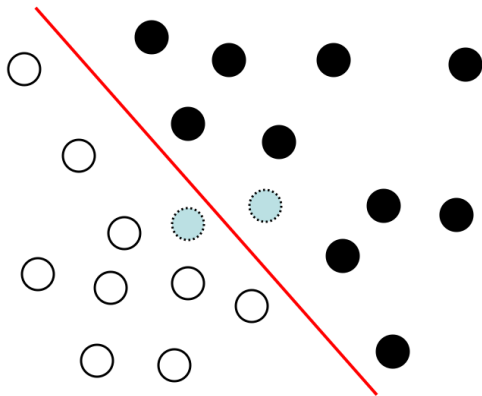
- 1 Given a **training set** of labeled data with...
- 2 **learn** a discrimination rule...
- 3 ... in order to **predict** the label of new data

# Machine learning (pattern recognition / supervised classification)



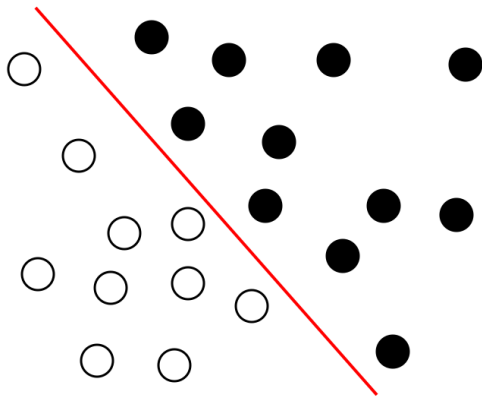
- 1 Given a **training set** of labeled data with...
- 2 **learn** a discrimination rule...
- 3 ... in order to **predict** the label of new data

# Machine learning (pattern recognition / supervised classification)

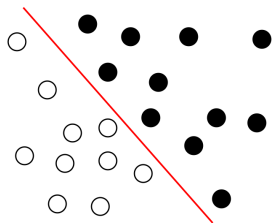


- 1 Given a **training set** of labeled data with...
- 2 **learn** a discrimination rule...
- 3 ... in order to **predict** the label of new data

# Machine learning (pattern recognition / supervised classification)



- 1 Given a **training set** of labeled data with...
- 2 **learn** a discrimination rule...
- 3 ... in order to **predict** the label of new data

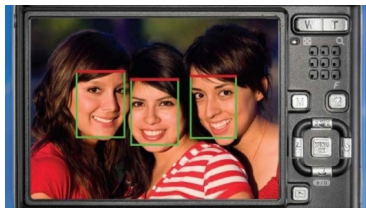


Genome annotation, systems biology, personalized medicine...

## Challenges

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

# Machine learning : tools and applications



## Many applications

Multimedia, image, video, speech recognition, web, social network, online advertising, finance, **biology, chemistry**

## Many tools

Linear discriminant analysis, logistic regression, decision trees, neural networks, support vector machines...



# ML with shrinkage estimators

- 1 Define a large family of "candidate classifiers", e.g., **linear predictors**:

$$f_{\beta}(x) = \beta^{\top} x \quad \text{for } x \in \mathbb{R}^p$$

- 2 For any candidate classifier  $f_{\beta}$ , quantify how "good" it is on the training set with some **empirical risk**, e.g.:

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

- 3 Choose  $\beta$  that achieves the minimum empirical risk, subject to some **constraint**:

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

# ML with shrinkage estimators

- 1 Define a large family of "candidate classifiers", e.g., **linear predictors**:

$$f_{\beta}(x) = \beta^{\top} x \quad \text{for } x \in \mathbb{R}^p$$

- 2 For any candidate classifier  $f_{\beta}$ , quantify how "good" it is on the training set with some **empirical risk**, e.g.:

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

- 3 Choose  $\beta$  that achieves the minimum empirical risk, subject to some **constraint**:

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

# ML with shrinkage estimators

- 1 Define a large family of "candidate classifiers", e.g., **linear predictors**:

$$f_{\beta}(x) = \beta^{\top} x \quad \text{for } x \in \mathbb{R}^p$$

- 2 For any candidate classifier  $f_{\beta}$ , quantify how "good" it is on the training set with some **empirical risk**, e.g.:

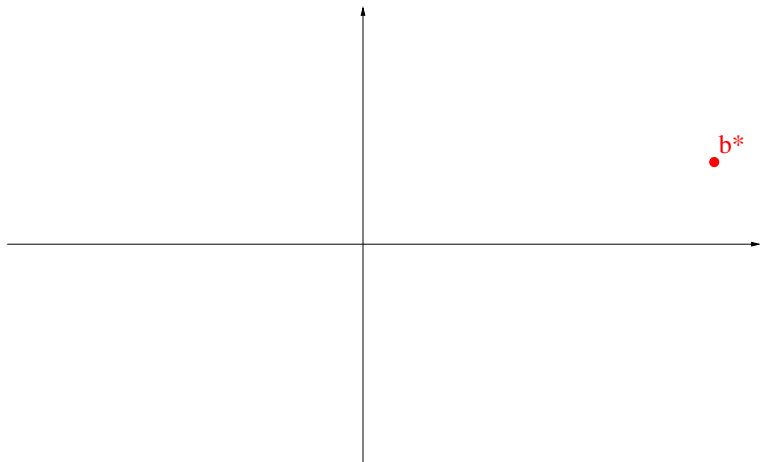
$$R(\beta) = \frac{1}{n} \sum_{i=1}^n l(f_{\beta}(x_i), y_i).$$

- 3 Choose  $\beta$  that achieves the minimum empirical risk, subject to some **constraint**:

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

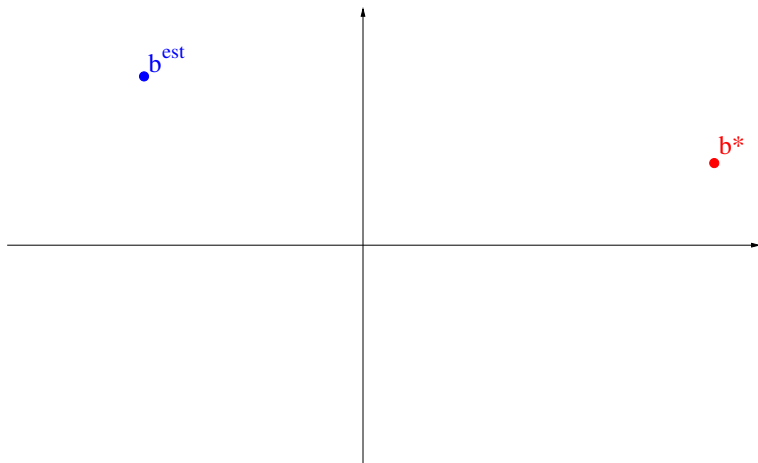
# Why shrinkage classifiers?

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



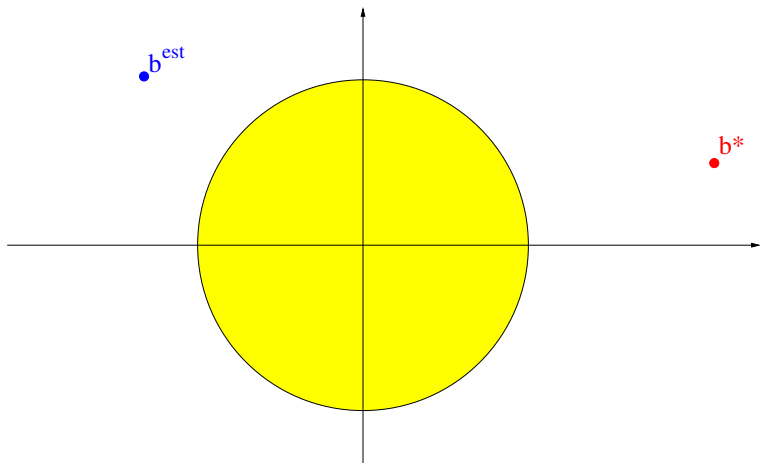
# Why shrinkage classifiers?

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



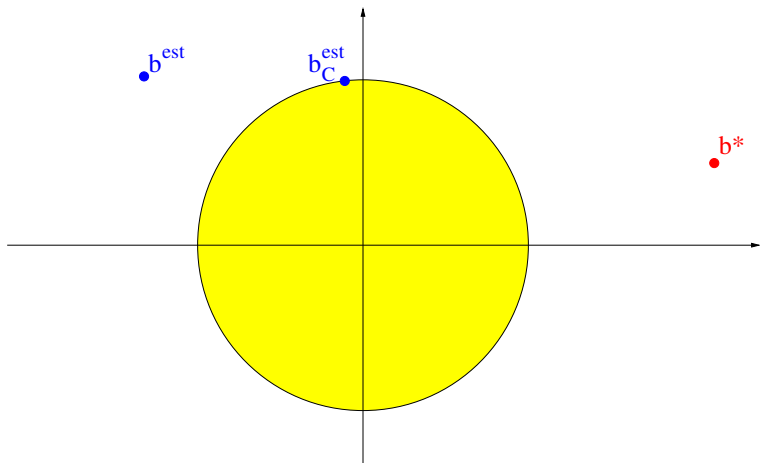
# Why shrinkage classifiers?

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



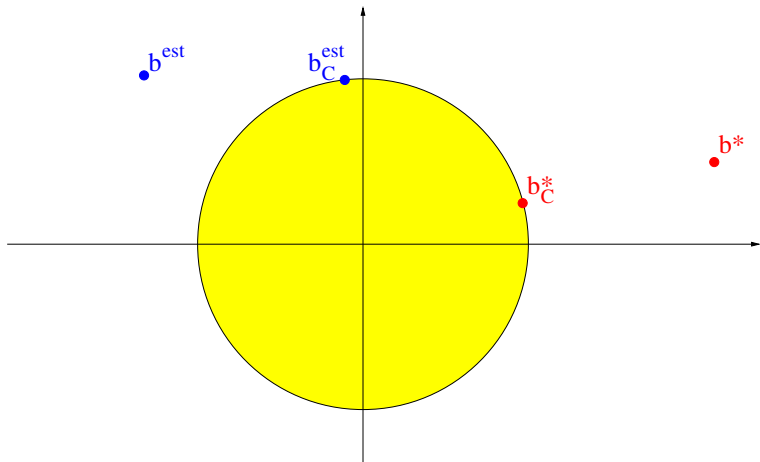
# Why shrinkage classifiers?

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



# Why shrinkage classifiers?

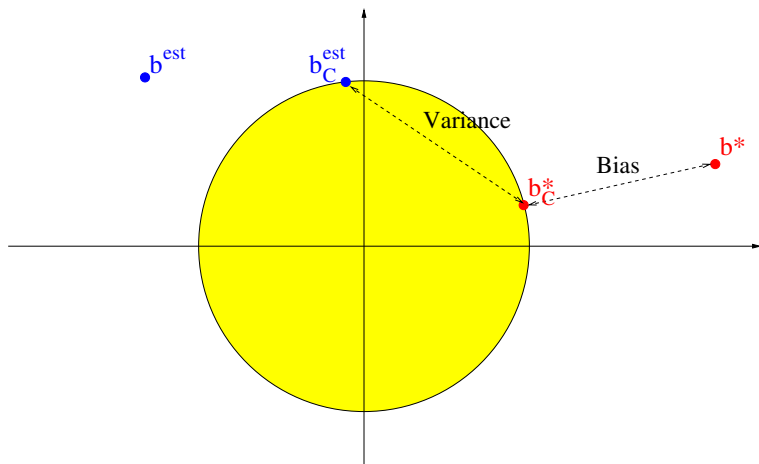
$$\min_{\beta} R(\beta) \quad \text{subject to} \quad \Omega(\beta) \leq C.$$



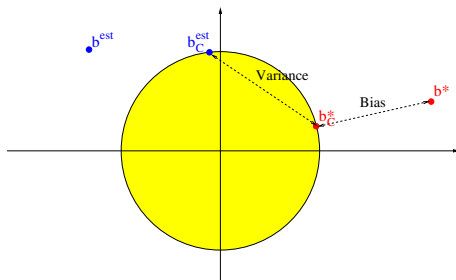


# Why shrinkage classifiers?

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



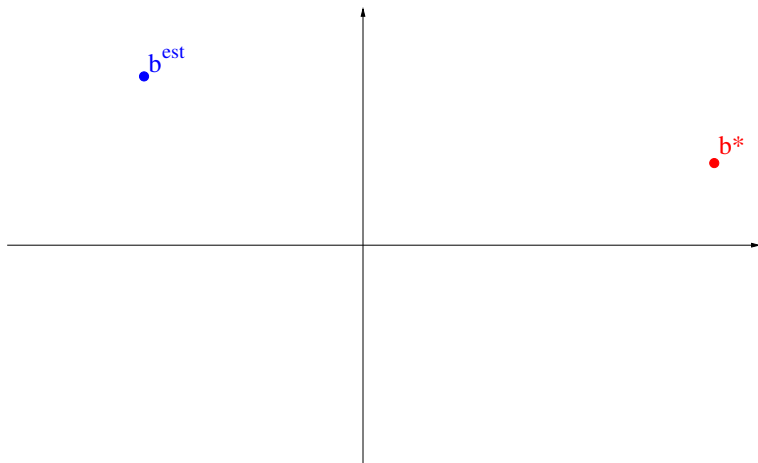
# Why shrinkage classifiers?



- "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, ...)

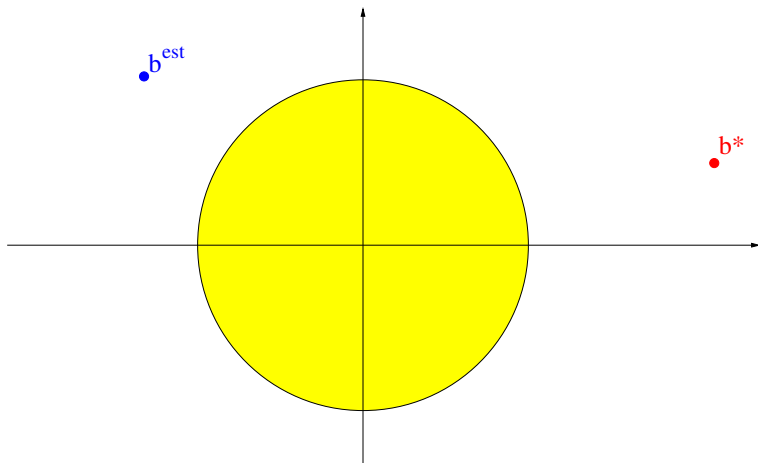
# Including prior knowledge in the penalty?

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



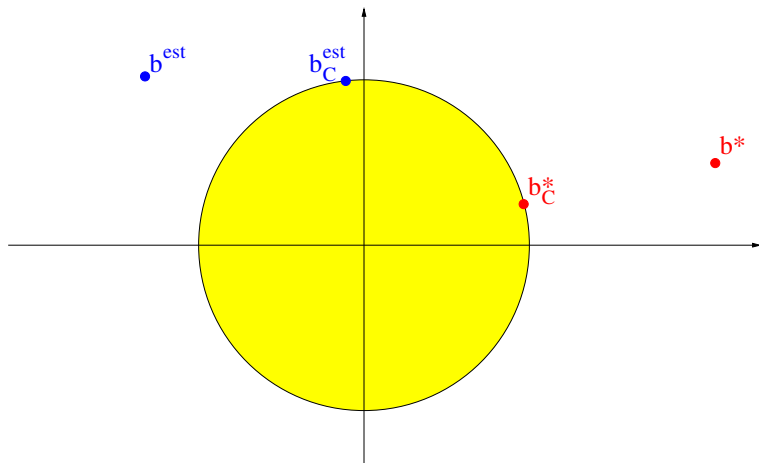
# Including prior knowledge in the penalty?

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



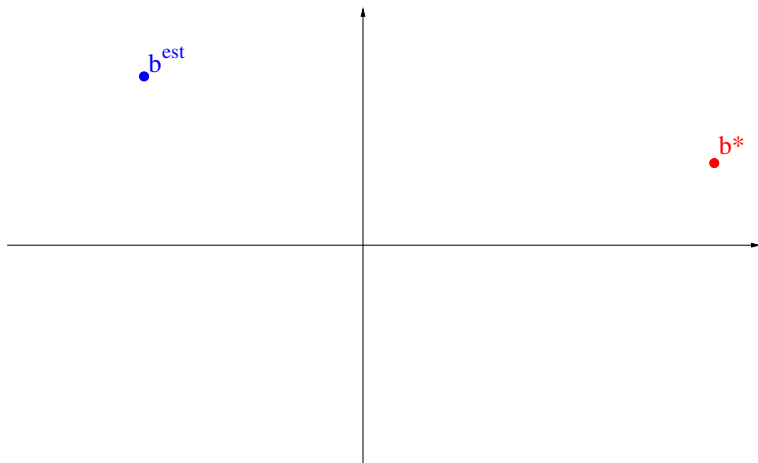
# Including prior knowledge in the penalty?

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



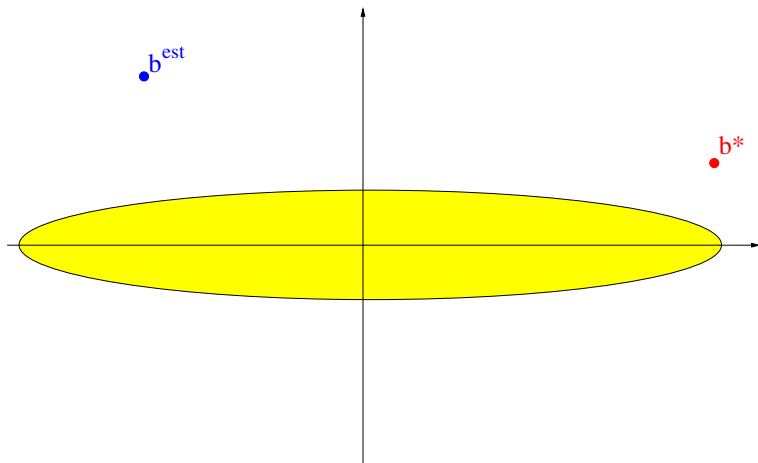
# Including prior knowledge in the penalty?

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



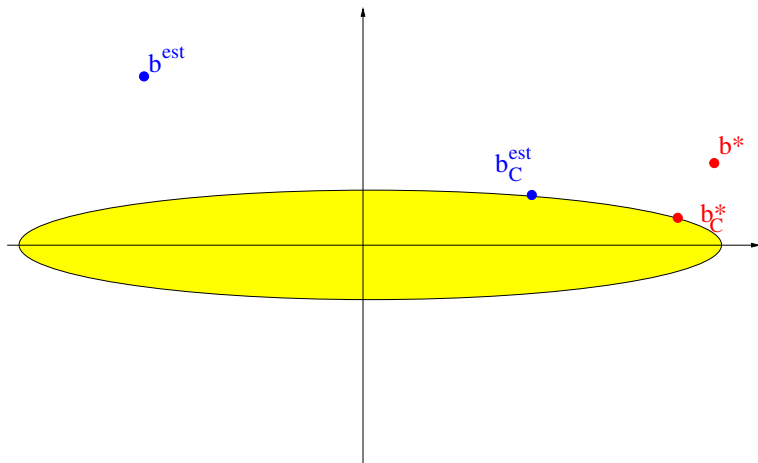
# Including prior knowledge in the penalty?

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



# Including prior knowledge in the penalty?

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



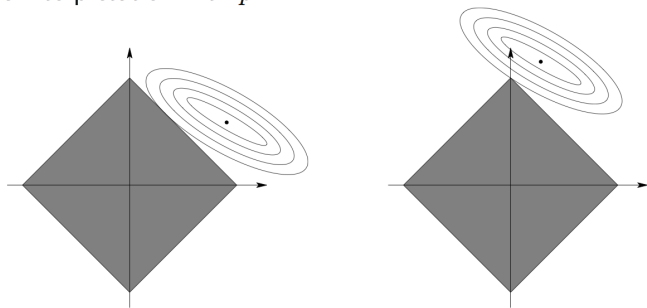


# Further benefit: sparsity-inducing penalties

(Lasso)

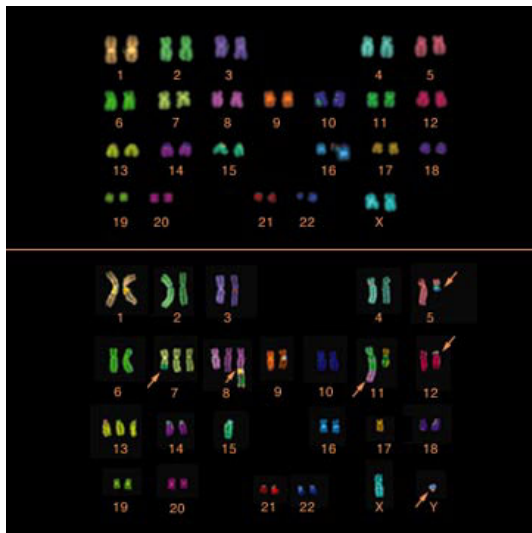
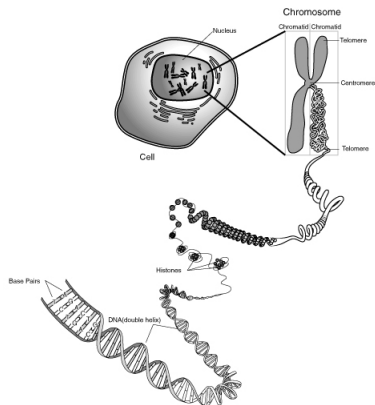
$$\min_{\beta} R(\beta) + \lambda \sum_{i=1}^p |\beta_i|$$

Geometric interpretation with  $p = 2$

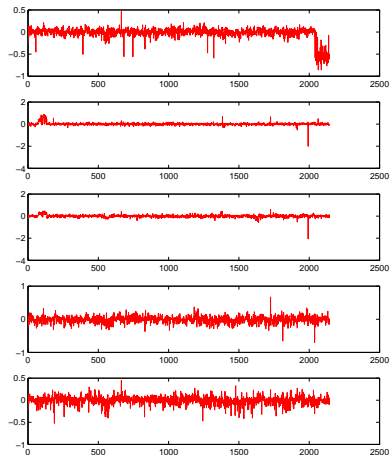
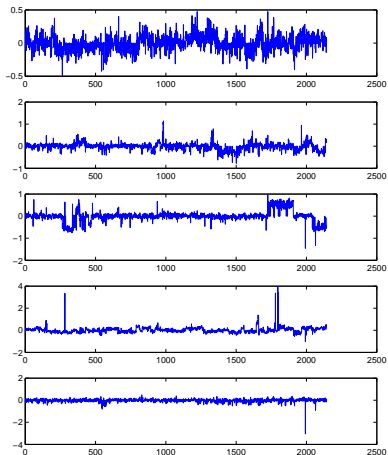


- 1 Introduction
- 2 Cancer prognosis from DNA copy number variations
- 3 Diagnosis and prognosis from gene expression data
- 4 Conclusion

# Chromosomal aberrations in cancer



# Cancer prognosis: can we predict the future evolution?



*Aggressive (left) vs non-aggressive (right) melanoma*

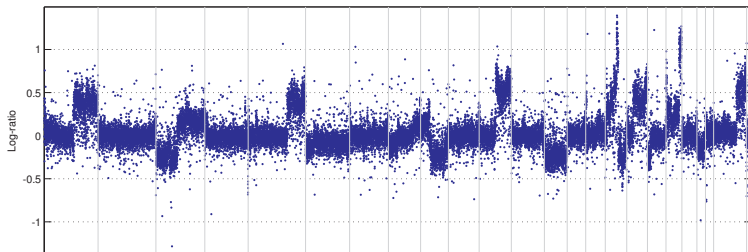
# CGH array classification

## Prior knowledge

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

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

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



# Promoting sparsity with the $\ell_1$ penalty

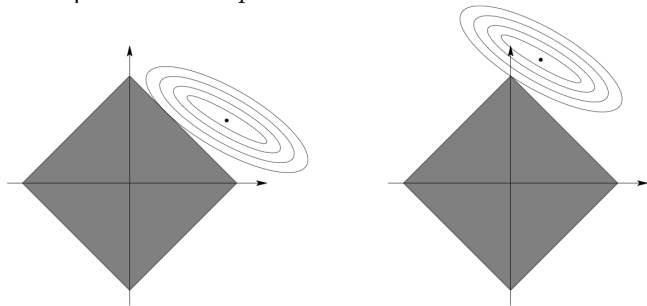
The  $\ell_1$  penalty (Tibshirani, 1996; Chen et al., 1998)

The solution of

$$\min_{\beta \in \mathbb{R}^p} R(\beta) + \lambda \sum_{i=1}^p |\beta_i|$$

is usually sparse.

Geometric interpretation with  $p = 2$



# Promoting piecewise constant profiles penalty

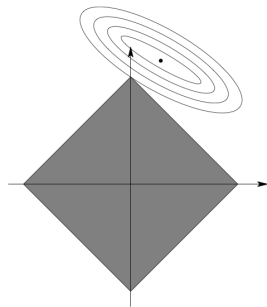
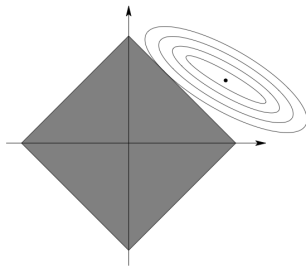
The variable fusion penalty (Land and Friedman, 1996)

The solution of

$$\min_{\beta \in \mathbb{R}^p} R(\beta) + \lambda \sum_{i=1}^{p-1} |\beta_{i+1} - \beta_i|$$

is usually piecewise constant.

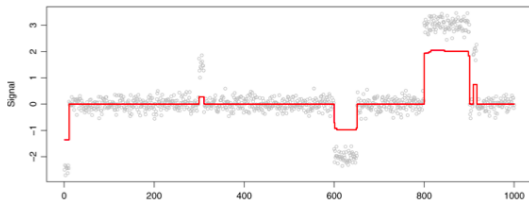
Geometric interpretation with  $p = 2$



# Fused Lasso signal approximator (Tibshirani et al., 2005)

$$\min_{\beta \in \mathbb{R}^p} \sum_{i=1}^p (y_i - \beta_i)^2 + \lambda_1 \sum_{i=1}^p |\beta_i| + \lambda_2 \sum_{i=1}^{p-1} |\beta_{i+1} - \beta_i|.$$

- First term leads to **sparse** solutions
- Second term leads to **piecewise constant** solutions





# Fused lasso for supervised classification (Rapaport et al., 2008)

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

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

## Implementation

- When  $\ell$  is the hinge loss (fused SVM), this is a **linear program** -> up to  $p = 10^3 \sim 10^4$
- When  $\ell$  is convex and smooth (logistic, quadratic), efficient implementation with **proximal methods** -> up to  $p = 10^8 \sim 10^9$

# Fused lasso for supervised classification (Rapaport et al., 2008)

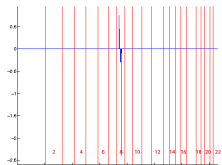
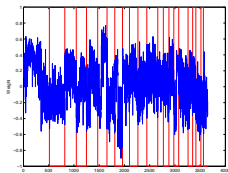
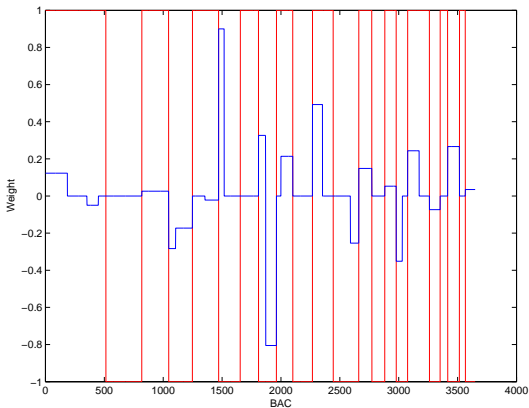
$$\min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n \ell(y_i, \beta^\top x_i) + \lambda_1 \sum_{i=1}^p |\beta_i| + \lambda_2 \sum_{i=1}^{p-1} |\beta_{i+1} - \beta_i|.$$

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

## Implementation

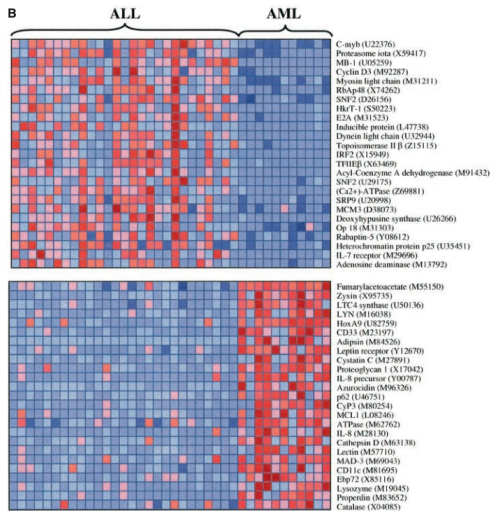
- When  $\ell$  is the hinge loss (fused SVM), this is a **linear program** -> up to  $p = 10^3 \sim 10^4$
- When  $\ell$  is convex and smooth (logistic, quadratic), efficient implementation with **proximal methods** -> up to  $p = 10^8 \sim 10^9$

# Example: predicting metastasis in melanoma

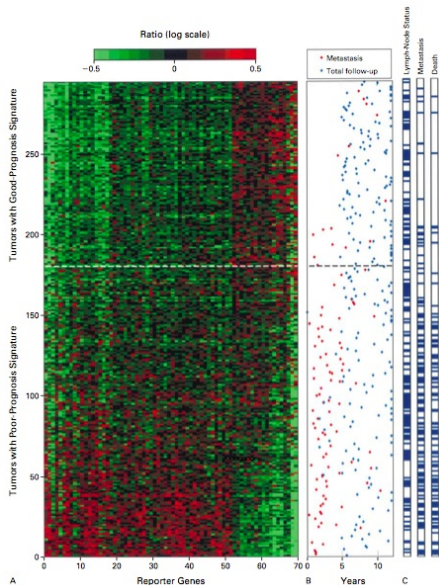


- 1 Introduction
- 2 Cancer prognosis from DNA copy number variations
- 3 Diagnosis and prognosis from gene expression data**
- 4 Conclusion

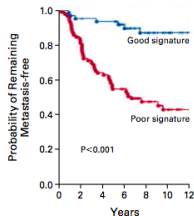
# Diagnosis



# Prognosis

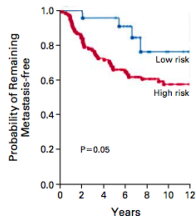


A Gene-Expression Profiling



No. AT RISK	0	2	4	6	8	10	12
Good signature	60	57	54	45	31	22	12
Poor signature	91	72	55	41	26	17	9

B St. Gallen Criteria

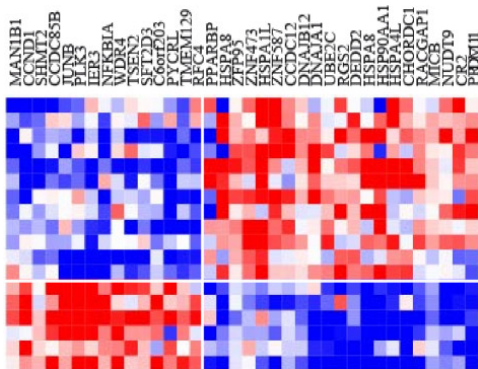


No. AT RISK	0	2	4	6	8	10	12
Low risk	22	22	21	17	9	5	2
High risk	129	107	88	69	48	34	19

# Gene selection, molecular signature

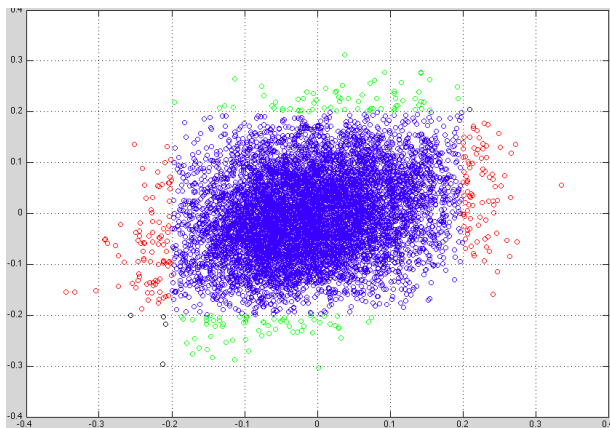
## The idea

- We look for a **limited set** of genes that are sufficient for prediction.
- Selected genes should inform us about the underlying biology



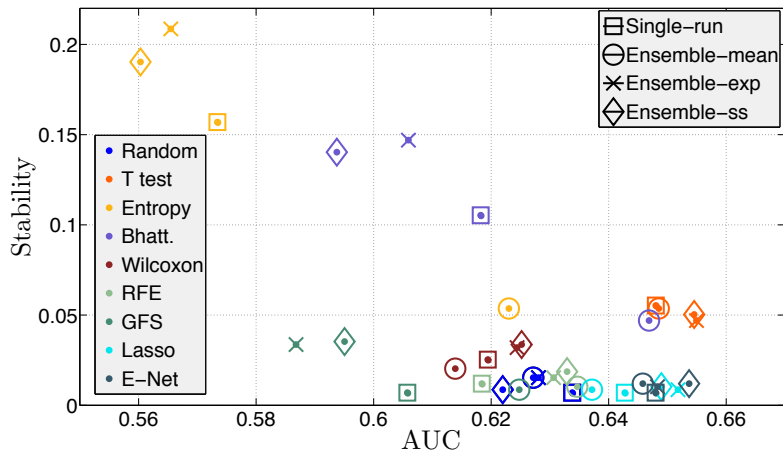
# But... instability of selected features

- Wang dataset:  $n = 286$ ,  $p = 8141$
- Pearson correlation with the output on 2 random subsamples of 143 samples:



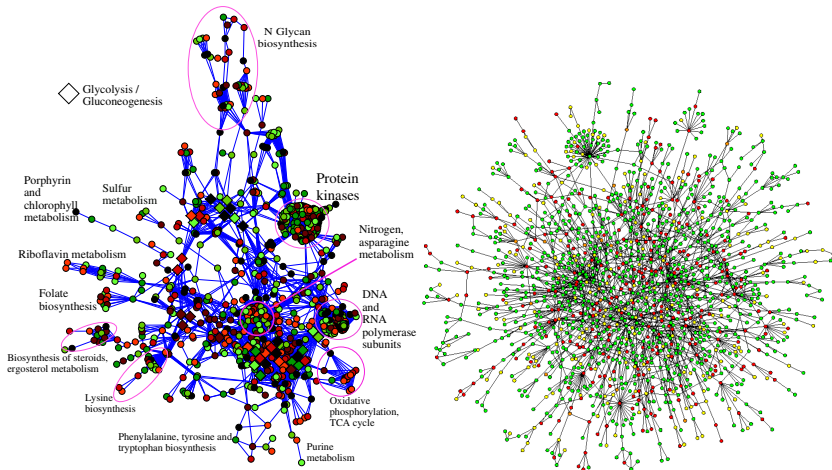


# Comparison of feature selection methods...



*Haury et al. (2011)*

# 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**
- We know these groups through functional groups and protein networks

## Shrinkage estimators with prior knowledge

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

How to design penalties  $\Omega(\beta)$  to encode the following hypotheses:

- 1 Connected genes on a network should have similar weights
- 2 Select few genes that are connected or belong to same predefined functional groups

## Motivation

- Basic biological functions usually involve the **coordinated action of several proteins**:
  - Formation of **protein complexes**
  - Activation of metabolic, signalling or regulatory **pathways**
- We know these groups through functional groups and protein networks

## Shrinkage estimators with prior knowledge

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

How to design penalties  $\Omega(\beta)$  to encode the following hypotheses:

- 1 Connected genes on a network should have similar weights
- 2 Select few genes that are connected or belong to same predefined functional groups

# Hypothesis 1: connected genes on a network should have similar weights

- Smooth weights on the graph (or more generally graph kernels)

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

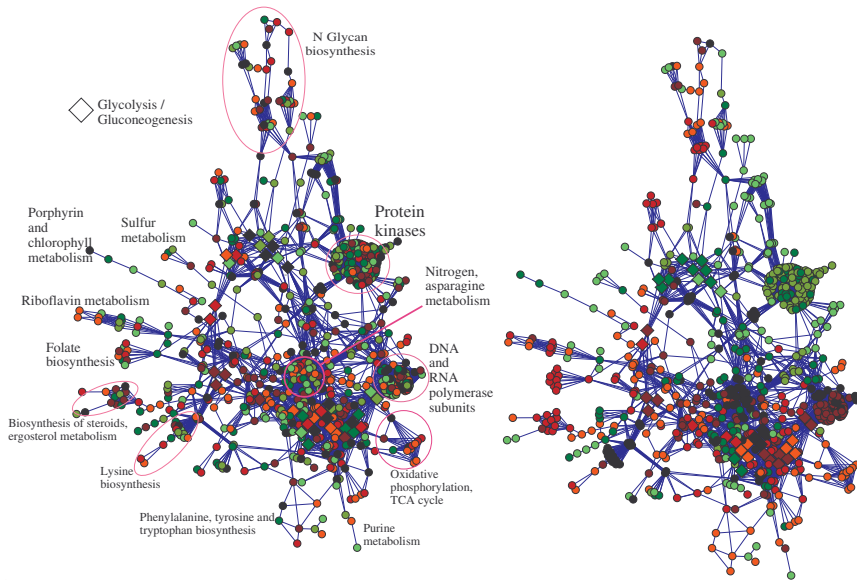
- Gene selection + smooth on the graph

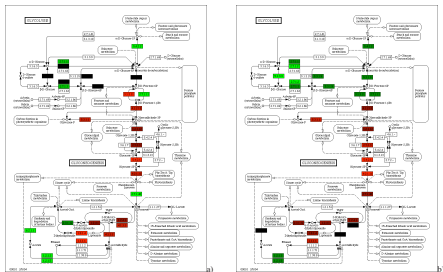
$$\Omega(\beta) = \sum_{i \sim j} (\beta_i - \beta_j)^2 + \sum_{i=1}^p |\beta_i|$$

- Gene selection + Piecewise constant on the graph (total variation)

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

# Illustration





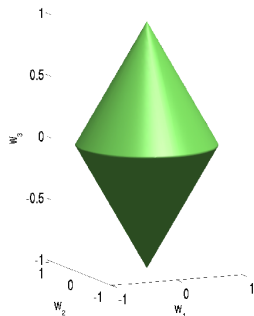
- We are happy to see pathways appear.
- However, in some cases, connected genes should have "opposite" weights (inhibition, pathway branching, etc...)
- **How to capture pathways without constraints on the weight similarities?**

# 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*:

$$\Omega_{group}(\beta) = \sum_g \|\beta_g\|_2$$



$$\begin{aligned}\Omega(\beta_1, \beta_2, \beta_3) &= \|(\beta_1, \beta_2)\|_2 + \|\beta_3\|_2 \\ &= \sqrt{\beta_1^2 + \beta_2^2} + |\beta_3|\end{aligned}$$

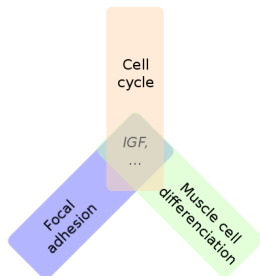


# Group Lasso when groups overlap

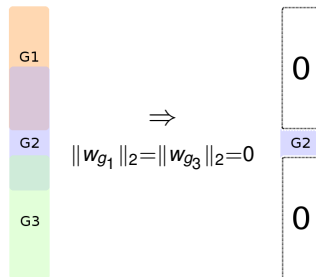
When groups overlap, the group Lasso

$$\Omega_{group}(\beta) = \sum_g \|\beta_g\|$$

sets groups to 0  $\Rightarrow$  the support of the solution is the complement of a union of groups



IGF selection  $\Rightarrow$  selection of unwanted groups



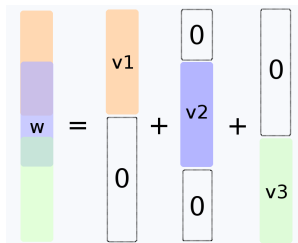
Removal of *any* group containing a gene  $\Rightarrow$  the weight of the gene is 0.

# The latent group Lasso (Jacob et al., 2009)

$$\Omega_{latent}(\beta) = \sup_{\alpha \in \mathbb{R}^p : \forall g, \|\alpha_g\| \leq 1} \alpha^\top \beta$$

or, equivalently:

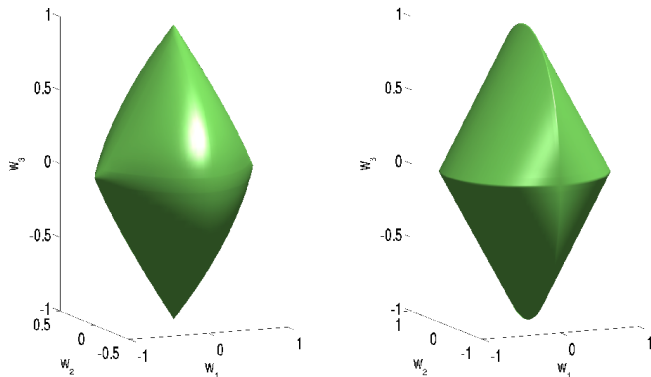
$$\Omega_{latent}(\beta) \triangleq \begin{cases} \min_v \sum_{g \in \mathcal{G}} \|v_g\|_2 \\ \beta = \sum_{g \in \mathcal{G}} v_g \\ \text{supp}(v_g) \subseteq g. \end{cases}$$



## Properties

- Resulting support is a *union* of groups in  $\mathcal{G}$ .
- Possible to select one variable without selecting all the groups containing it.
- Equivalent to group lasso when there is no overlap

# Group Lasso vs latent group Lasso



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

## Consistency in group support (Jacob et al., 2009)

- Let  $\bar{\mathbf{w}}$  be the true parameter vector.
- Assume that there exists a unique decomposition  $\bar{\mathbf{v}}_g$  such that  $\bar{\mathbf{w}} = \sum_g \bar{\mathbf{v}}_g$  and  $\Omega_{\text{latent}}(\bar{\mathbf{w}}) = \sum \|\bar{\mathbf{v}}_g\|_2$ .
- Consider the regularized empirical risk minimization problem  $L(\mathbf{w}) + \lambda \Omega_{\text{latent}}(\mathbf{w})$ .

Then

- under appropriate mutual incoherence conditions on  $X$ ,
- as  $n \rightarrow \infty$ ,
- with very high probability,

the optimal solution  $\hat{\mathbf{w}}$  admits a unique decomposition  $(\hat{\mathbf{v}}_g)_{g \in \mathcal{G}}$  such that

$$\{g \in \mathcal{G} | \hat{\mathbf{v}}_g \neq 0\} = \{g \in \mathcal{G} | \bar{\mathbf{v}}_g \neq 0\}.$$

## Consistency in group support (Jacob et al., 2009)

- Let  $\bar{w}$  be the true parameter vector.
- Assume that there exists a unique decomposition  $\bar{v}_g$  such that  $\bar{w} = \sum_g \bar{v}_g$  and  $\Omega_{\text{latent}}(\bar{w}) = \sum \|\bar{v}_g\|_2$ .
- Consider the regularized empirical risk minimization problem  $L(w) + \lambda \Omega_{\text{latent}}(w)$ .

Then

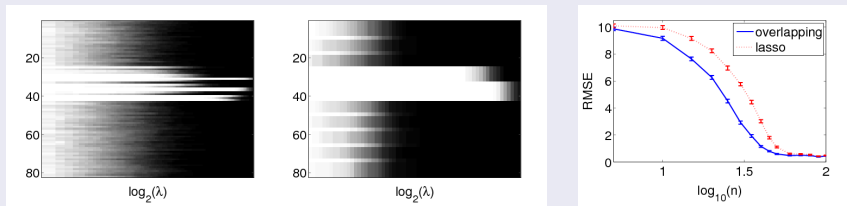
- under appropriate mutual incoherence conditions on  $X$ ,
- as  $n \rightarrow \infty$ ,
- with very high probability,

the optimal solution  $\hat{w}$  admits a unique decomposition  $(\hat{v}_g)_{g \in \mathcal{G}}$  such that

$$\{g \in \mathcal{G} | \hat{v}_g \neq 0\} = \{g \in \mathcal{G} | \bar{v}_g \neq 0\}.$$

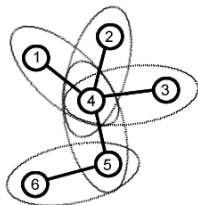
## Synthetic data: overlapping groups

- 10 groups of 10 variables with 2 variables of overlap between two successive groups :  $\{1, \dots, 10\}, \{9, \dots, 18\}, \dots, \{73, \dots, 82\}$ .
- Support: union of 4<sup>th</sup> and 5<sup>th</sup> groups.
- Learn from 100 training points.



Frequency of selection of each variable with the lasso (left) and  $\Omega_{\text{latent}}(\cdot)$  (middle), comparison of the RMSE of both methods (right).

# Graph lasso vs kernel on graph



- Graph lasso:

$$\Omega_{group}(\beta) = \sum_{i \sim j} \sqrt{\beta_i^2 + \beta_j^2} \quad \text{or} \quad \Omega_{latent}(\beta) = \sup_{\alpha \in \mathbb{R}^p : \forall i \sim j, \sqrt{\alpha_i^2 + \alpha_j^2} \leq 1} \alpha^\top \beta$$

constrains the **sparsity**, not the values

- Graph kernel

$$\Omega_{graph\ kernel}(\beta) = \sum_{i \sim j} (\beta_i - \beta_j)^2.$$

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.

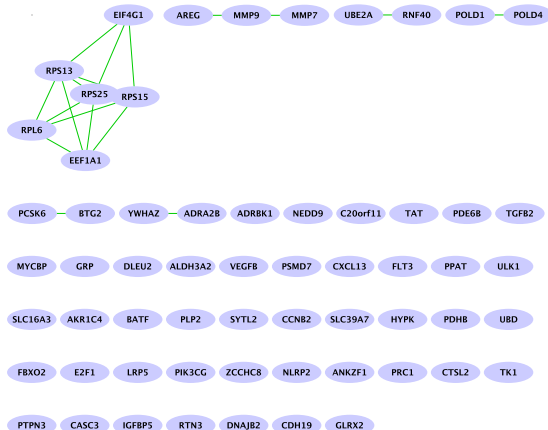
METHOD	$\ell_1$	$\Omega_{\text{LATENT}}(\cdot)$
ERROR	$0.38 \pm 0.04$	$0.36 \pm 0.03$
MEAN $\#$ PATH.	130	30

- Graph on the genes.

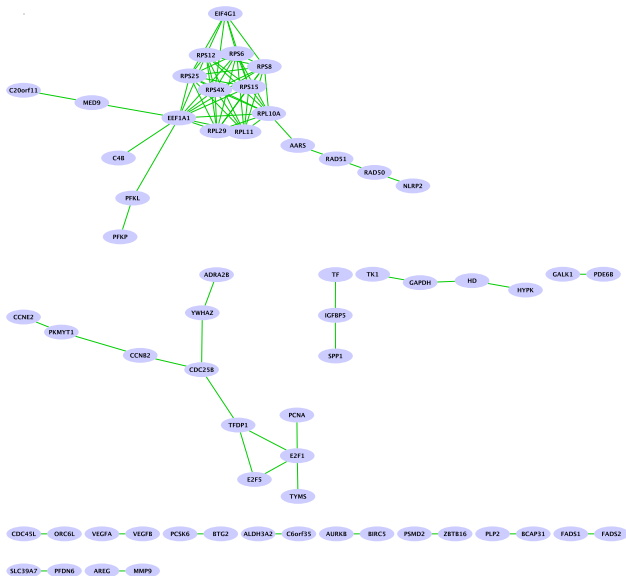
METHOD	$\ell_1$	$\Omega_{\text{graph}}(\cdot)$
ERROR	$0.39 \pm 0.04$	$0.36 \pm 0.01$
AV. SIZE C.C.	1.03	1.30



# Classical lasso signature



# Graph Lasso signature

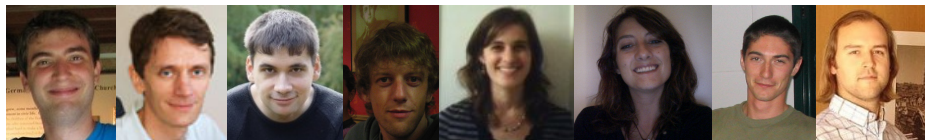


- 1 Introduction
- 2 Cancer prognosis from DNA copy number variations
- 3 Diagnosis and prognosis from gene expression data
- 4 Conclusion**

# Conclusion

- Machine learning offers **many powerful tools** to learn predictive models from large sets of complex data
- **Specific developments** are required to solve complex problems that arise in bio-informatics
- **Integration of prior knowledge** in the penalization / regularization function is an efficient approach to fight the curse of dimension
- Requires **interdisciplinary collaborations** to incorporate expert knowledge at the heart of learning algorithms
- Many other applications not covered in this presentation!

# Acknowledgements!



Franck Rapaport (MSKCC), Emmanuel Barillot, Andrei Zynoviev, Kevin Bleakley (INRIA), Fantine Mordelet (Duke), Anne-Claire Haury, Laurent Jacob (UC Berkeley) Guillaume Obozinski (INRIA)