

Structured feature selection for genomic data

Jean-Philippe Vert

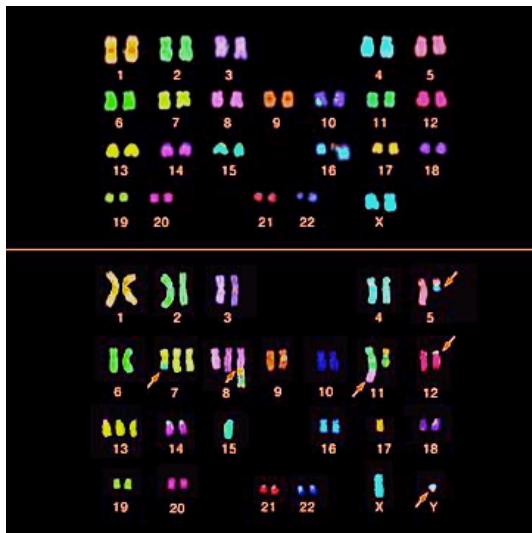
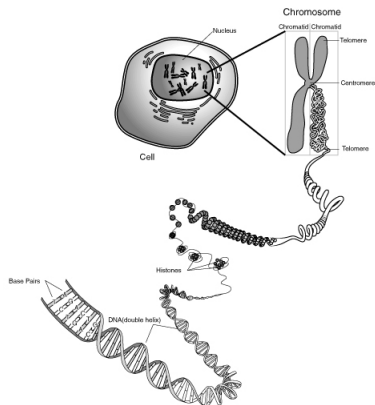
Mines ParisTech / Curie Institute / Inserm

8th World Congress in Probability and Statistics,
Istanbul, July 9-14, 2012

- 1 On chromosome abnormalities in cancer
- 2 Gene selection with prior information
- 3 Conclusion

- 1 On chromosome abnormalities in cancer
- 2 Gene selection with prior information
- 3 Conclusion

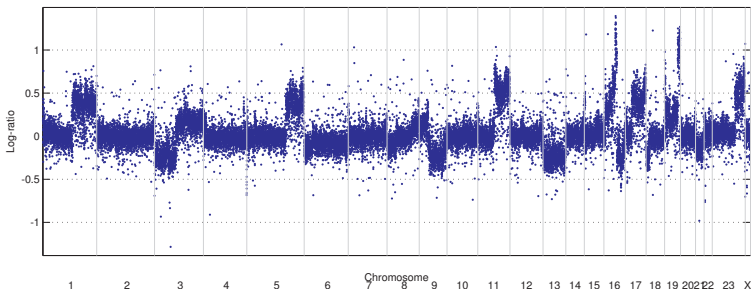
Chromosomal aberrations in cancer



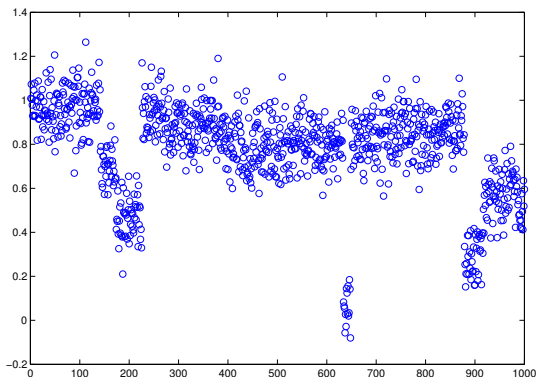
Comparative Genomic Hybridization (CGH)

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

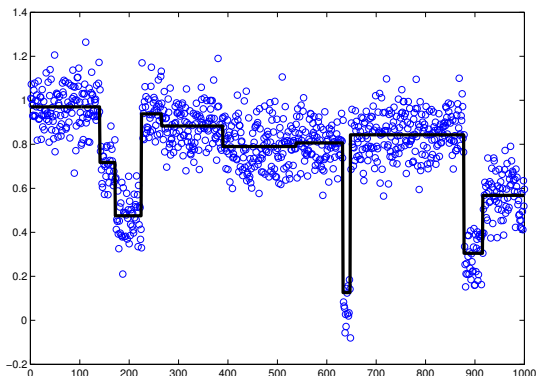


Can we identify breakpoints and "smooth" each profile?



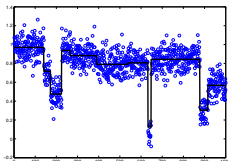
- A classical multiple change-point detection problem
- Should scale to lengths of order $10^6 \sim 10^9$

Can we identify breakpoints and "smooth" each profile?



- A classical multiple change-point detection problem
- Should scale to lengths of order $10^6 \sim 10^9$

An optimal solution?

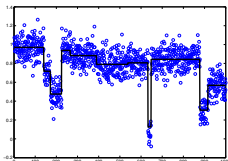


- For a signal $Y \in \mathbb{R}^p$, define an optimal approximation $\beta \in \mathbb{R}^p$ with k breakpoints as the solution of

$$\min_{\beta \in \mathbb{R}^p} \|Y - \beta\|^2 \quad \text{such that} \quad \sum_{i=1}^{p-1} \mathbf{1}(U_{i+1} \neq U_i) \leq k$$

- This is an optimization problem over the $\binom{p}{k}$ partitions...
- Dynamic programming finds the solution in $O(p^2 k)$ in time and $O(p^2)$ in memory
- But: does not scale to $p = 10^6 \sim 10^9$...

An optimal solution?

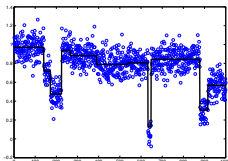


- For a signal $Y \in \mathbb{R}^p$, define an optimal approximation $\beta \in \mathbb{R}^p$ with k breakpoints as the solution of

$$\min_{\beta \in \mathbb{R}^p} \|Y - \beta\|^2 \quad \text{such that} \quad \sum_{i=1}^{p-1} \mathbf{1}(U_{i+1} \neq U_i) \leq k$$

- This is an optimization problem over the $\binom{p}{k}$ partitions...
 - Dynamic programming finds the solution in $O(p^2 k)$ in time and $O(p^2)$ in memory
 - But: does not scale to $p = 10^6 \sim 10^9$...

An optimal solution?

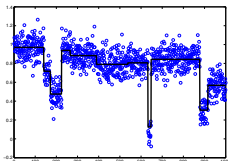


- For a signal $Y \in \mathbb{R}^p$, define an optimal approximation $\beta \in \mathbb{R}^p$ with k breakpoints as the solution of

$$\min_{\beta \in \mathbb{R}^p} \|Y - \beta\|^2 \quad \text{such that} \quad \sum_{i=1}^{p-1} \mathbf{1}(U_{i+1} \neq U_i) \leq k$$

- This is an optimization problem over the $\binom{p}{k}$ partitions...
- **Dynamic programming** finds the solution in $O(p^2k)$ in time and $O(p^2)$ in memory
- **But:** does not scale to $p = 10^6 \sim 10^9$...

An optimal solution?



- For a signal $Y \in \mathbb{R}^p$, define an optimal approximation $\beta \in \mathbb{R}^p$ with k breakpoints as the solution of

$$\min_{\beta \in \mathbb{R}^p} \|Y - \beta\|^2 \quad \text{such that} \quad \sum_{i=1}^{p-1} \mathbf{1}(U_{i+1} \neq U_i) \leq k$$

- This is an optimization problem over the $\binom{p}{k}$ partitions...
- **Dynamic programming** finds the solution in $O(p^2k)$ in time and $O(p^2)$ in memory
- **But:** does not scale to $p = 10^6 \sim 10^9$...

Promoting sparsity with the ℓ_1 penalty

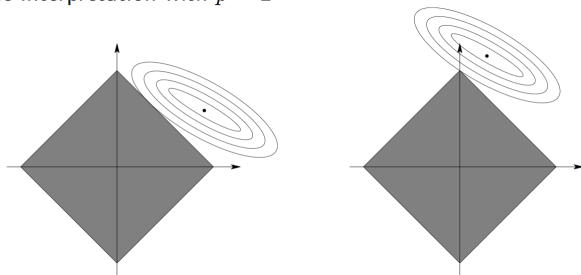
The ℓ_1 penalty (Tibshirani, 1996; Chen et al., 1998)

If $R(\beta)$ is convex and "smooth", the solution of

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

is usually **sparse**.

Geometric interpretation with $p = 2$



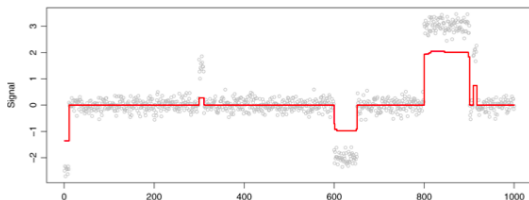
Promoting piecewise constant profiles with

- Total variation (Rudin et al., 1992; Land and Friedman, 1996):

$$\|\beta\|_{TV} = \|\nabla\beta\|_1 = \sum_{i=1}^{p-1} |\beta_{i+1} - \beta_i|$$

- Fused lasso (Tibshirani et al., 2005; Tibshirani and Wang, 2008)

$$\min_{\beta \in \mathbb{R}^p} \|Y - \beta\|^2 + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_{TV}$$



TV signal approximator as dichotomic segmentation

Algorithm 1 Greedy dichotomic segmentation

Require: k number of intervals, $\gamma(I)$ gain function to split an interval I into $I_L(I), I_R(I)$

1: I_0 represents the interval $[1, n]$

2: $\mathcal{P} = \{I_0\}$

3: **for** $i = 1$ to k **do**

4: $I^* \leftarrow \arg \max_{I \in \mathcal{P}} \gamma(I)$

5: $\mathcal{P} \leftarrow \mathcal{P} \setminus \{I^*\}$

6: $\mathcal{P} \leftarrow \mathcal{P} \cup \{I_L(I^*), I_R(I^*)\}$

7: **end for**

8: **return** \mathcal{P}

Theorem

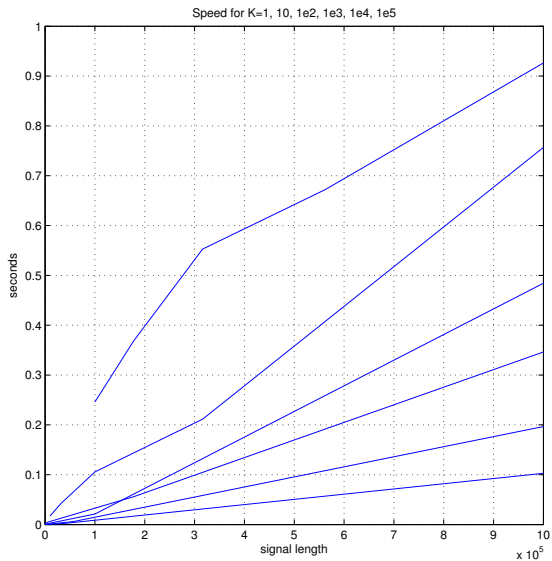
TV signal approximator performs "greedy" dichotomic segmentation.

(V. and Bleakley, 2010; see also Hoefling, 2009)

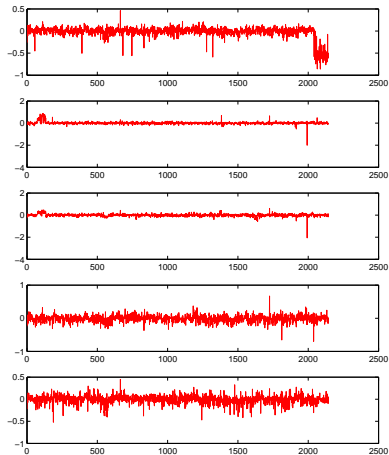
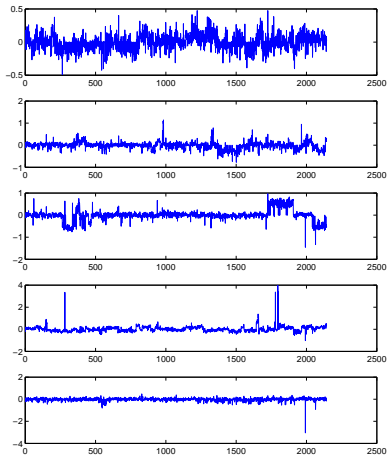
$$\min_{\beta \in \mathbb{R}^p} \|Y - \beta\|^2 + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_{TV}$$

- QP with sparse linear constraints in $O(p^2)$ -> 135 min for $p = 10^5$ (Tibshirani and Wang, 2008)
- Coordinate descent-like method $O(p)$? -> 3s for $p = 10^5$ (Friedman et al., 2007)
- For all λ with the LARS in $O(pK)$ (Harchaoui and Levy-Leduc, 2008)
- For all λ in $O(p \ln p)$ (Hoefling, 2009)
- For the first K change-points in $O(p \ln K)$ (Bleakley and V., 2010)

Speed trial : 2 s. for $K = 100$, $p = 10^7$

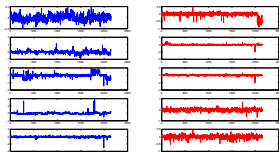


Extension: cancer prognosis



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

Fused lasso for supervised classification

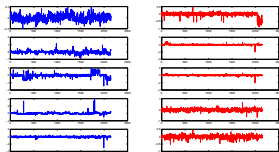


- **Idea:** find a linear predictor $f(Y) = \beta^T Y$ that best discriminates the aggressive vs non-aggressive samples, subject to the constraints that it should be sparse and piecewise constant
- **Mathematically:**

$$\min_{\beta \in \mathbb{R}^p} R(\beta) + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_{TV}$$

- **Computationally:** this is convex optimization problem that can be solved very efficiently with proximal optimization methods (V. and Bleakley, 2012)

Fused lasso for supervised classification

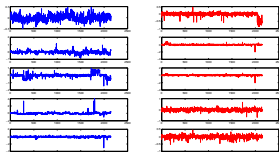


- **Idea:** find a linear predictor $f(Y) = \beta^\top Y$ that best discriminates the aggressive vs non-aggressive samples, subject to the constraints that it should be sparse and piecewise constant
- **Mathematically:**

$$\min_{\beta \in \mathbb{R}^p} R(\beta) + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_{TV}$$

- **Computationally:** this is convex optimization problem that can be solved very efficiently with proximal optimization methods (V. and Bleakley, 2012)

Fused lasso for supervised classification

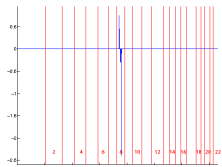
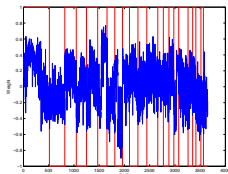
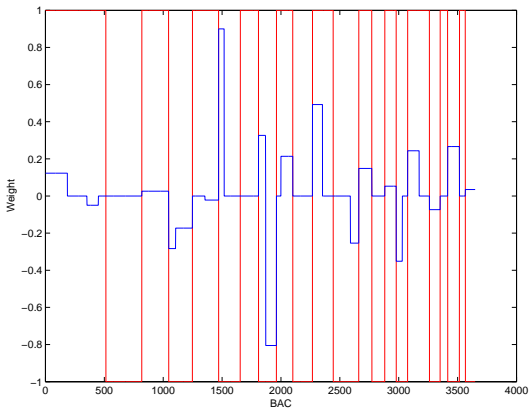


- **Idea:** find a linear predictor $f(Y) = \beta^T Y$ that best discriminates the aggressive vs non-aggressive samples, subject to the constraints that it should be sparse and piecewise constant
- **Mathematically:**

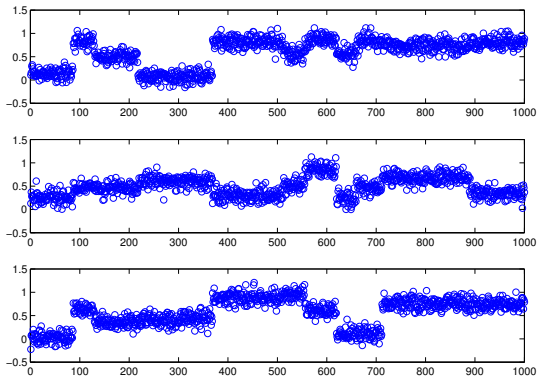
$$\min_{\beta \in \mathbb{R}^p} R(\beta) + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_{TV}$$

- **Computationally:** this is convex optimization problem that can be solved very efficiently with proximal optimization methods (V. and Bleakley, 2012)

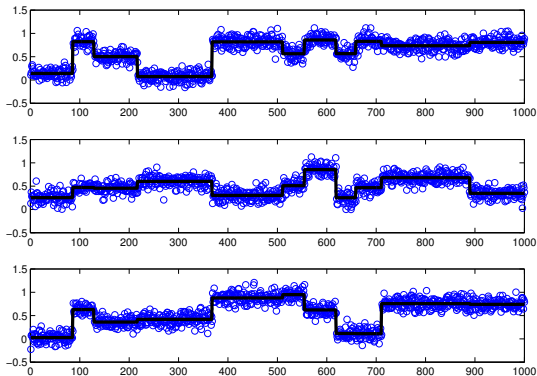
Prognostic in melanoma (Rapaport et al., 2008)



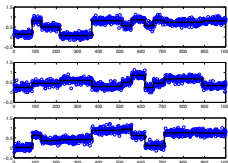
Extension: finding multiple change points shared by several profiles



Extension: finding multiple change points shared by several profiles



"Optimal" segmentation by dynamic programming



- Define the "optimal" piecewise constant approximation $\hat{U} \in \mathbb{R}^{p \times n}$ of Y as the solution of

$$\min_{U \in \mathbb{R}^{p \times n}} \|Y - U\|^2 \quad \text{such that} \quad \sum_{i=1}^{p-1} \mathbf{1}(U_{i+1, \bullet} \neq U_{i, \bullet}) \leq k$$

- DP finds the solution in $O(p^2 kn)$ in time and $O(p^2)$ in memory
- But: does not scale to $p = 10^6 \sim 10^9 \dots$

Replace

$$\min_{U \in \mathbb{R}^{p \times n}} \|Y - U\|^2 \quad \text{such that} \quad \sum_{i=1}^{p-1} \mathbf{1}(U_{i+1, \bullet} \neq U_{i, \bullet}) \leq k$$

by

$$\min_{U \in \mathbb{R}^{p \times n}} \|Y - U\|^2 \quad \text{such that} \quad \sum_{i=1}^{p-1} w_i \|U_{i+1, \bullet} - U_{i, \bullet}\| \leq \mu$$

Theorem

This is a group lasso problem!

- *The estimated segmentation converges to the true segmentation when the number of profiles increases (if the noise is not too large)*
- *We can solve it efficiently in $O(npk)$*

Speed trial

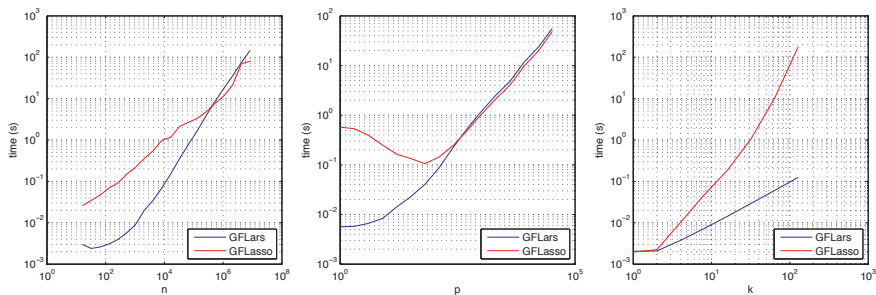


Figure 2: **Speed trials for group fused LARS (top row) and Lasso (bottom row).** *Left column:* varying n , with fixed $p = 10$ and $k = 10$; *center column:* varying p , with fixed $n = 1000$ and $k = 10$; *right column:* varying k , with fixed $n = 1000$ and $p = 10$. Figure axes are log-log. Results are averaged over 100 trials.

Consistency for a single change-point

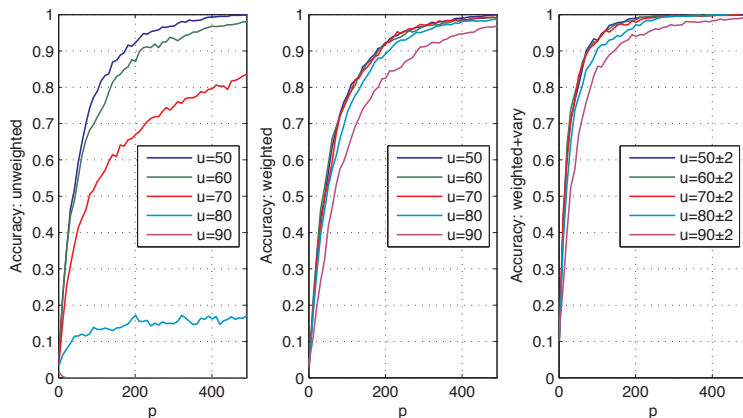


Figure 3: **Single change-point accuracy for the group fused Lasso.** Accuracy as a function of the number of profiles p when the change-point is placed in a variety of positions $u = 50$ to $u = 90$ (left and centre plots, resp. unweighted and weighted group fused Lasso), or: $u = 50 \pm 2$ to $u = 90 \pm 2$ (right plot, weighted with varying change-point location), for a signal of length 100.

Consistency for many change-points

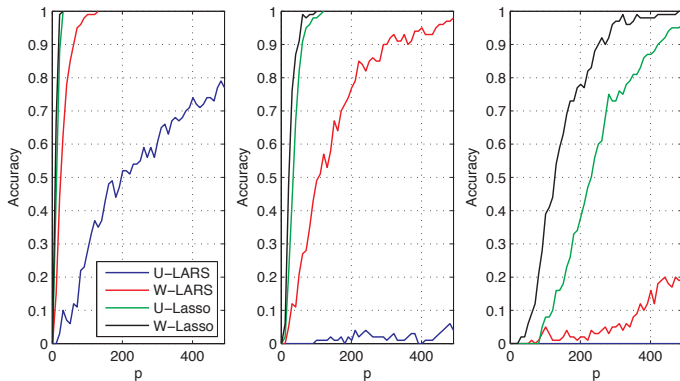
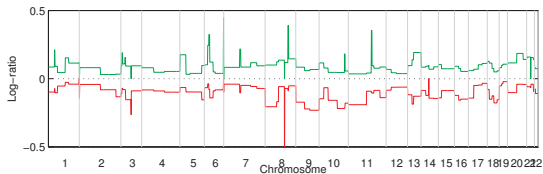
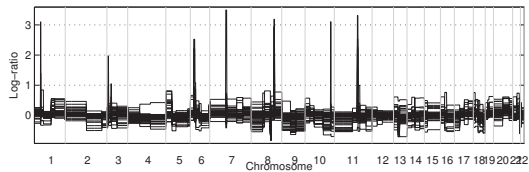
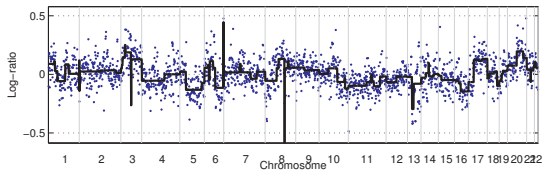


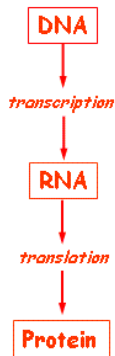
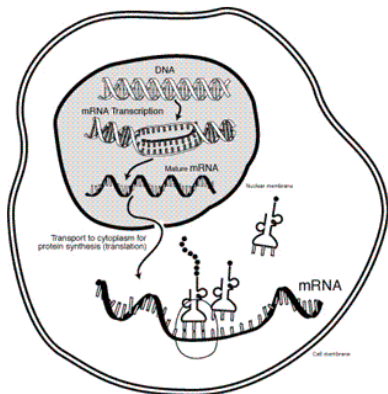
Figure 4: **Multiple change-point accuracy.** Accuracy as a function of the number of profiles p when change-points are placed at the nine positions $\{10, 20, \dots, 90\}$ and the variance σ^2 of the centered Gaussian noise is either 0.05 (left), 0.2 (center) and 1 (right). The profile length is 100.

Application: detection of frequent abnormalities



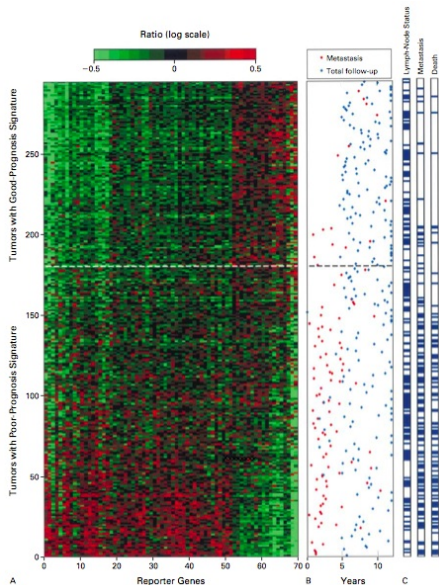
- 1 On chromosome abnormalities in cancer
- 2 Gene selection with prior information
- 3 Conclusion

DNA → RNA → protein

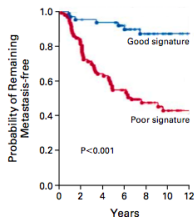


- CGH shows the (static) DNA
- Cancer cells have also **abnormal (dynamic) gene expression** (= transcription)

Breast cancer prognosis



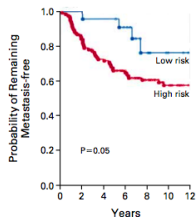
A Gene-Expression Profiling



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



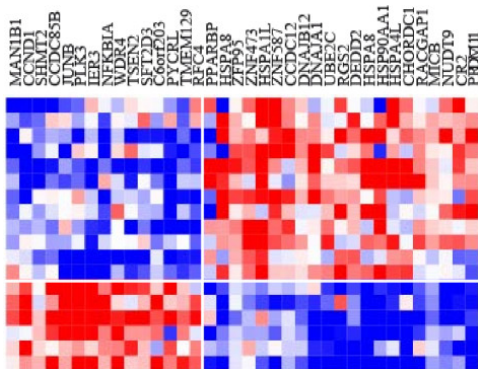
No. AT RISK

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

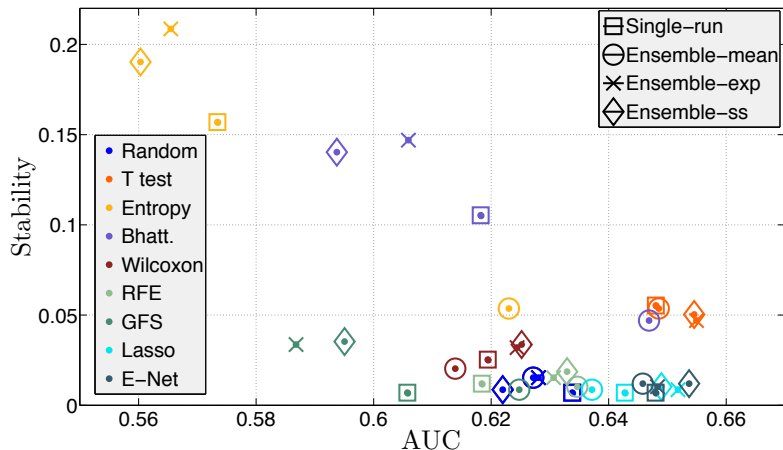
Gene selection, molecular signature

The idea

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

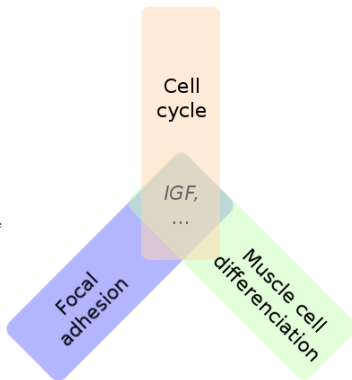
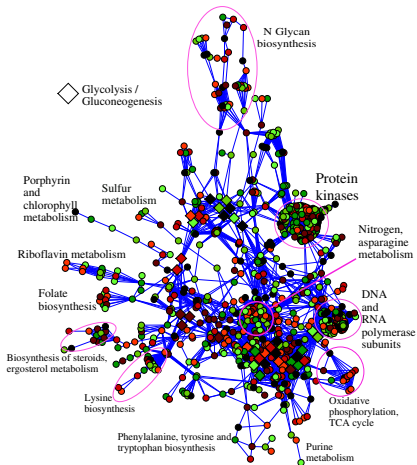


Lack of stability of signatures



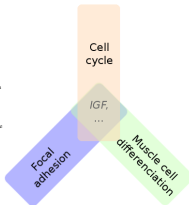
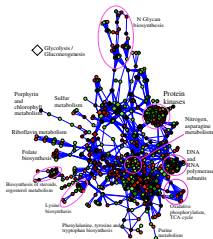
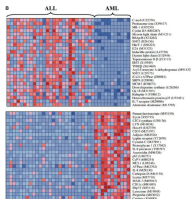
Haury et al. (2011)

Gene networks, gene groups



Structured feature selection

- Basic biological functions usually involve the **coordinated action of several proteins**:
 - Formation of **protein complexes**
 - Activation of metabolic, signalling or regulatory **pathways**
- How to perform **structured feature selection**, such that selected genes
 - belong to only a few groups?
 - form a small number of connected components on the graph?

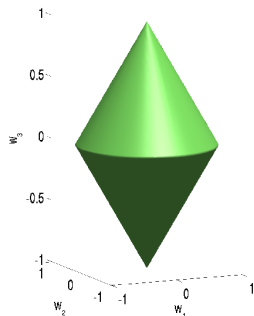


Selecting pre-defined groups of variables

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*:

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

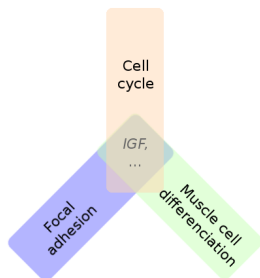


$$\Omega(w_1, w_2, w_3) = \|(w_1, w_2)\|_2 + \|w_3\|_2$$

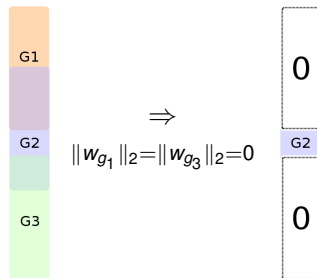
Group lasso with overlapping groups

Idea 1: shrink groups to zero (Jenatton et al., 2009)

- $\Omega_{group}(w) = \sum_g \|w_g\|_2$ sets groups to 0.
- One variable is selected \Leftrightarrow all the groups to which it belongs are selected.



IGF selection \Rightarrow selection of unwanted groups

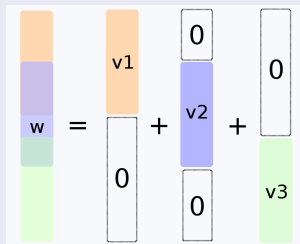


Removal of *any* group containing a gene \Rightarrow the weight of the gene is 0.

Group lasso with overlapping groups

Idea 2: latent group Lasso (Jacob et al., 2009)

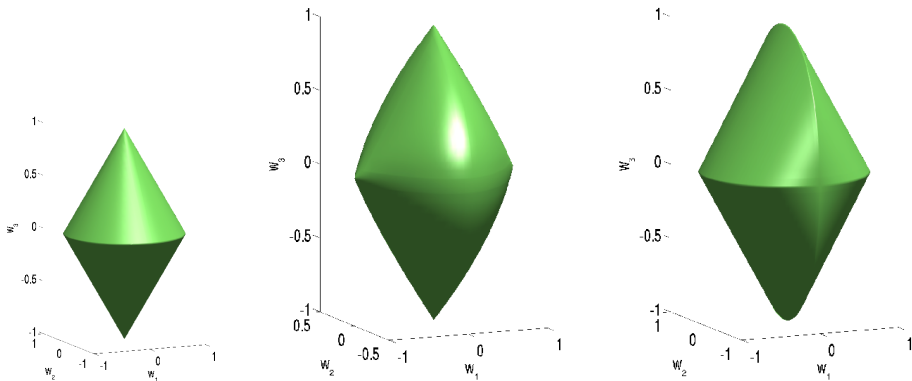
$$\Omega_{\text{latent}}^{\mathcal{G}}(w) \triangleq \begin{cases} \min_v \sum_{g \in \mathcal{G}} \|v_g\|_2 \\ w = \sum_{g \in \mathcal{G}} v_g \\ \text{supp}(v_g) \subseteq g. \end{cases}$$



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

Overlap and group unity balls



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

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 $\Omega_{\text{latent}}^{\mathcal{G}}(\bar{w}) = \sum \|\bar{v}_g\|_2$.
- Consider the regularized empirical risk minimization problem $L(w) + \lambda \Omega_{\text{latent}}^{\mathcal{G}}(w)$.

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

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

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 $\Omega_{\text{latent}}^{\mathcal{G}}(\bar{w}) = \sum \|\bar{v}_g\|_2$.
- Consider the regularized empirical risk minimization problem $L(w) + \lambda \Omega_{\text{latent}}^{\mathcal{G}}(w)$.

Then

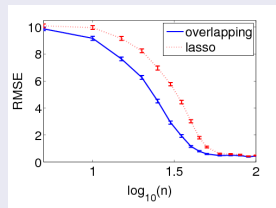
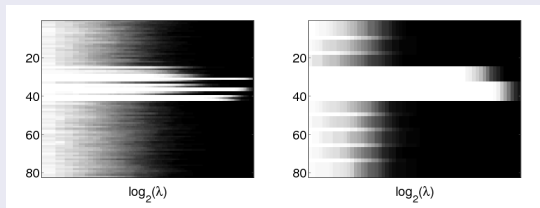
- under appropriate mutual incoherence conditions on X ,
- as $n \rightarrow \infty$,
- with very high probability,

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

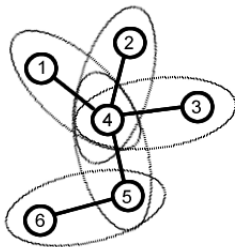
$$\{g \in \mathcal{G} | \hat{v}_g \neq 0\} = \{g \in \mathcal{G} | \bar{v}_g \neq 0\}.$$

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 4th and 5th groups.
- Learn from 100 training points.



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



Two solutions

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

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

Breast cancer data

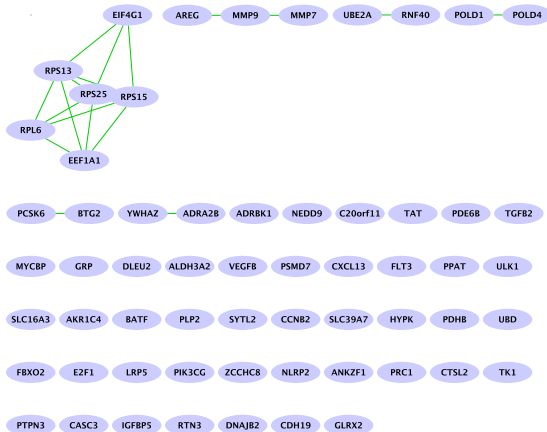
- Gene expression data for 8,141 genes in 295 breast cancer tumors.
- Canonical pathways from MSigDB containing 639 groups of genes, 637 of which involve genes from our study.

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

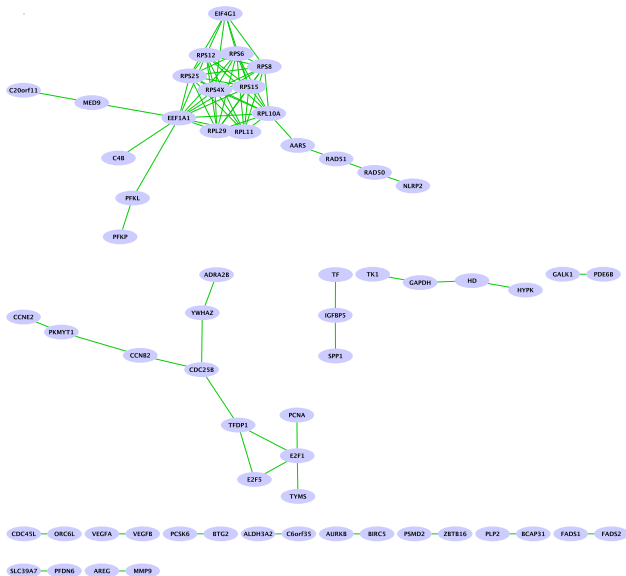
- Graph on the genes.

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 On chromosome abnormalities in cancer
- 2 Gene selection with prior information
- 3 Conclusion**

Conclusions

- Feature / pattern selection in high dimension is central for many applications
- Convex sparsity-inducing penalties are useful; efficient implementations + consistency results



Kevin Bleakley (INRIA), Laurent Jacob (UC Berkeley) Guillaume Obozinski (INRIA), Anne-Claire Haury (ParisTech)



European Research Council

