

# Supervised classification for structured data: Applications in bio- and chemoinformatics

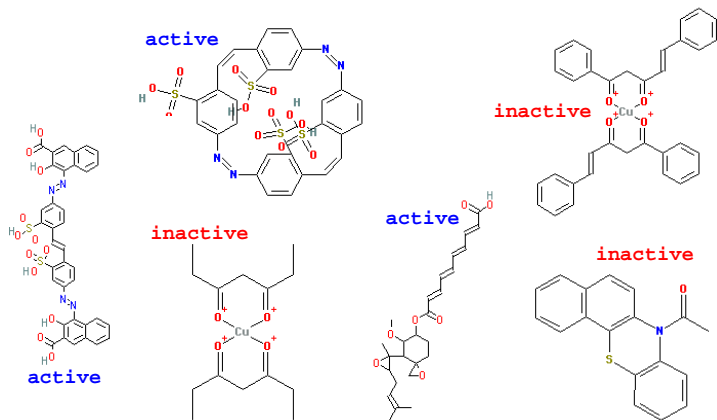
Jean-Philippe Vert

Jean-Philippe.Vert@mines-paristech.fr

Mines ParisTech / Curie Institute / Inserm

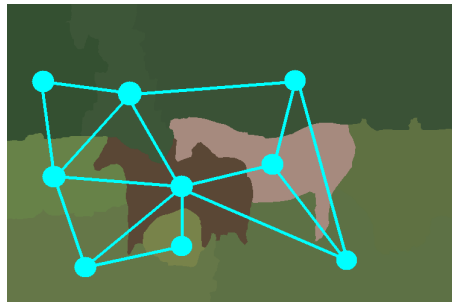
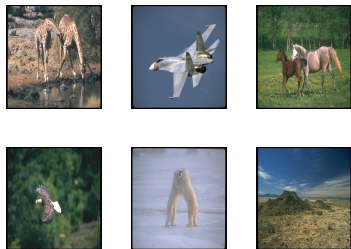
The third school on The Analysis of Patterns, Pula Science Park,  
Italy, June 1, 2009

# Virtual screening for drug discovery



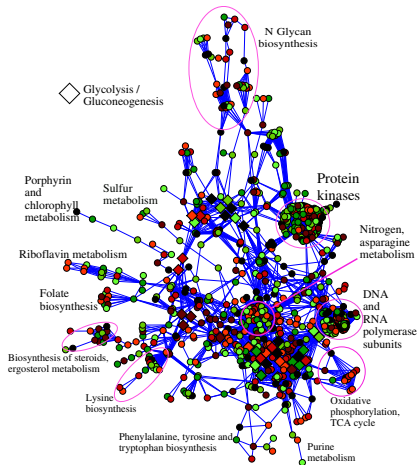
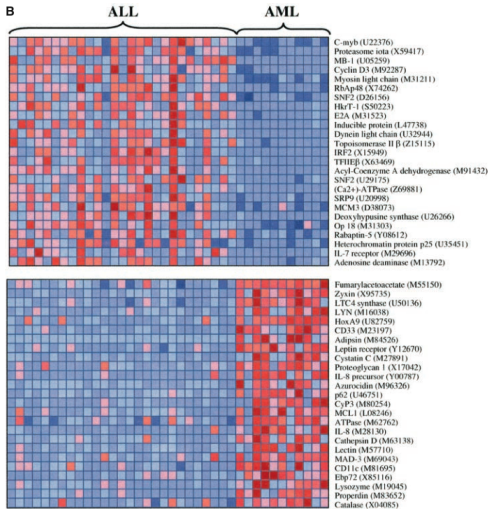
NCI AIDS screen results (from <http://cactus.nci.nih.gov>).

# Image retrieval and classification

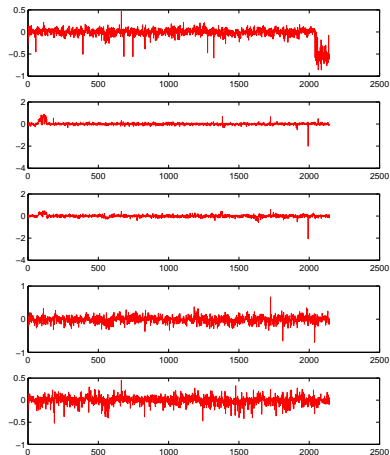
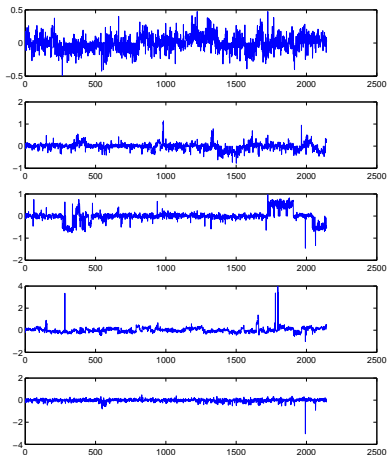


*From Harchaoui and Bach (2007).*

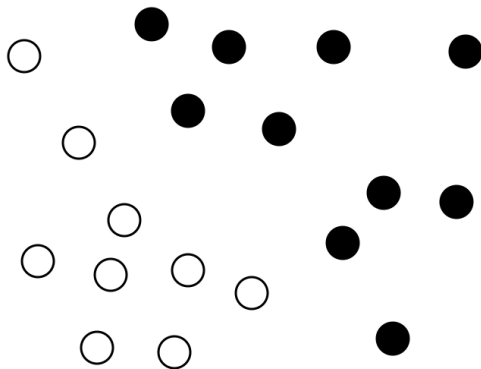
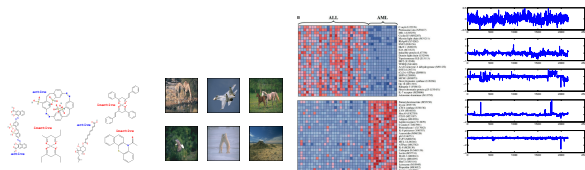
# Cancer diagnosis



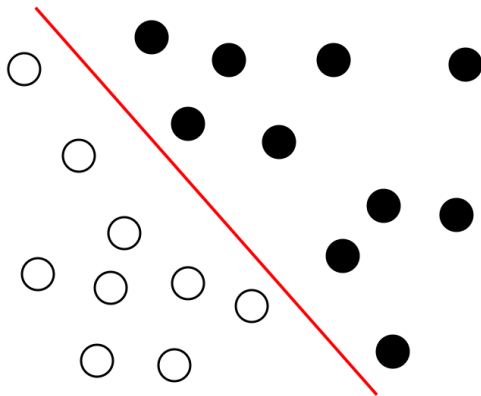
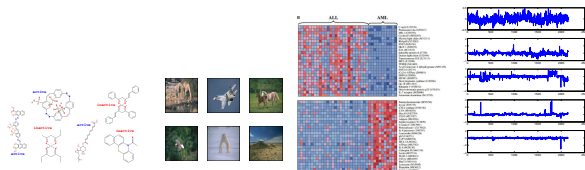
# Cancer prognosis



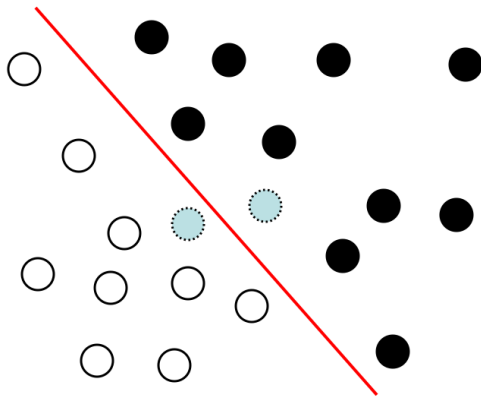
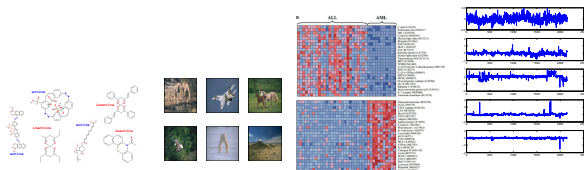
# Pattern recognition, *aka* supervised classification



# Pattern recognition, *aka* supervised classification

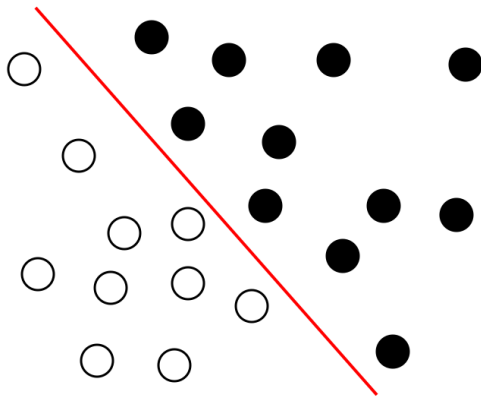
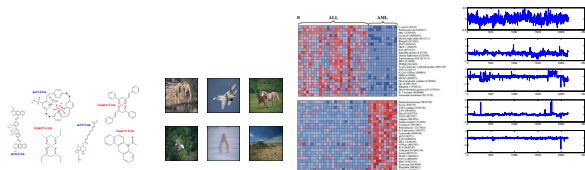


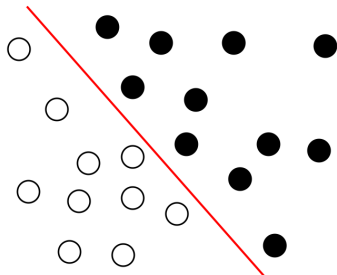
# Pattern recognition, *aka* supervised classification





# Pattern recognition, *aka* supervised classification





## Challenges

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

## The problem

- Given a set of **training instances**  $(x_1, y_1), \dots, (x_n, y_n)$ , where  $x_i \in \mathcal{X}$  are data and  $y_i \in \mathcal{Y}$  are continuous or discrete variables of interest,

- Estimate a function

$$y = f(x)$$

where  $x$  is any new data to be labeled.

- $f$  should be **accurate** and **interpretable**.

## The model

- Each sample  $x \in \mathcal{X}$  is represented by a vector of **features** (or **descriptors**, or **patterns**):

$$\Phi(x) = (\Phi_1(x), \dots, \Phi_p(x))$$

- Based on the training set we estimate a linear function:

$$f_{\beta}(x) = \sum_{i=1}^p \beta_i \Phi_i(x) = \beta^{\top} \Phi(x) .$$

## Two (related) questions

- How to **design the features**  $\Phi(x)$ ?
- How to **estimate the model**  $\beta$ ?

## The model

- Each sample  $x \in \mathcal{X}$  is represented by a vector of **features** (or **descriptors**, or **patterns**):

$$\Phi(x) = (\Phi_1(x), \dots, \Phi_p(x))$$

- Based on the training set we estimate a linear function:

$$f_{\beta}(x) = \sum_{i=1}^p \beta_i \Phi_i(x) = \beta^{\top} \Phi(x) .$$

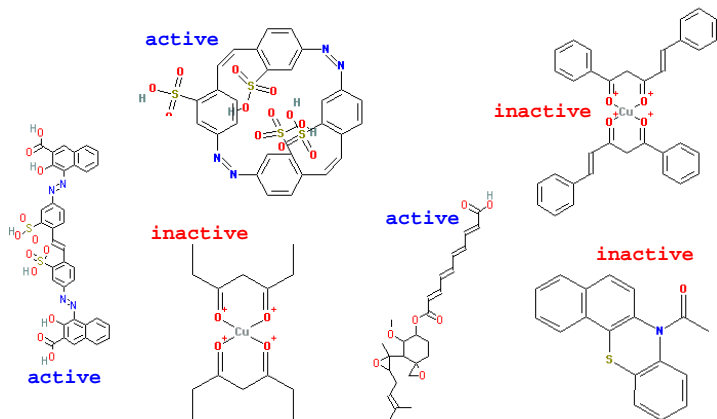
## Two (related) questions

- How to **design the features**  $\Phi(x)$ ?
- How to **estimate the model**  $\beta$ ?

- 1 Explicit computation of features : the case of graph features
- 2 Using kernels
  - Introduction to kernels
  - Graph kernels
  - Kernels for gene expression data using gene networks
- 3 Using sparsity-inducing shrinkage estimators
  - Feature selection for all subgraph indexation
  - Classification of array CGH data with piecewise-linear models
  - Structured gene selection for microarray classification
- 4 Conclusion

- 1 Explicit computation of features : the case of graph features
- 2 Using kernels
  - Introduction to kernels
  - Graph kernels
  - Kernels for gene expression data using gene networks
- 3 Using sparsity-inducing shrinkage estimators
  - Feature selection for all subgraph indexation
  - Classification of array CGH data with piecewise-linear models
  - Structured gene selection for microarray classification
- 4 Conclusion

# Motivation

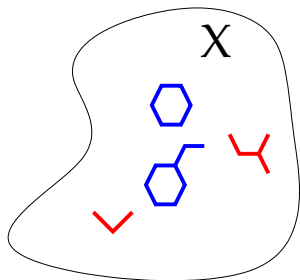


NCI AIDS screen results (from <http://cactus.nci.nih.gov>).



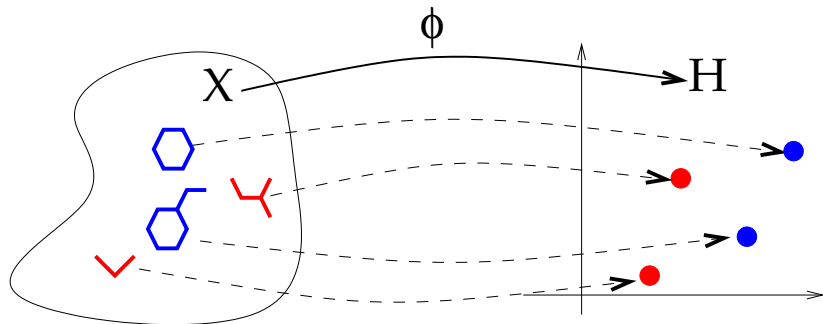
# The approach

- 1 Represent explicitly each graph  $x$  by a **vector of fixed dimension**  $\Phi(x) \in \mathbb{R}^p$ .
- 2 Use an algorithm for **regression or pattern recognition** in  $\mathbb{R}^p$ .



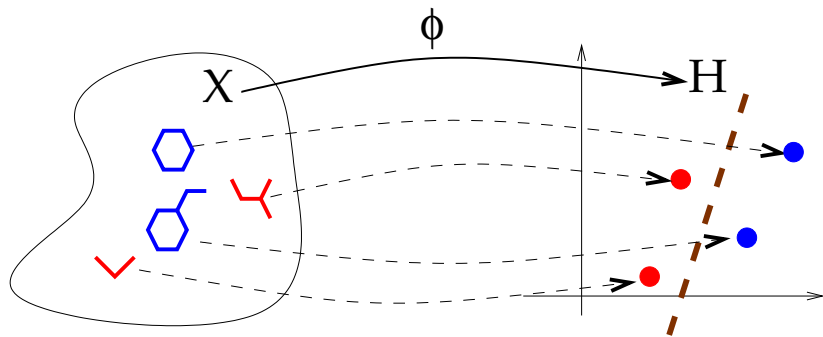
# The approach

- 1 Represent explicitly each graph  $x$  by a **vector of fixed dimension**  $\Phi(x) \in \mathbb{R}^p$ .
- 2 Use an algorithm for **regression or pattern recognition** in  $\mathbb{R}^p$ .



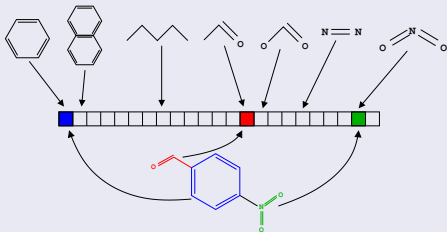
# The approach

- 1 Represent explicitly each graph  $x$  by a **vector of fixed dimension**  $\Phi(x) \in \mathbb{R}^p$ .
- 2 Use an algorithm for **regression or pattern recognition** in  $\mathbb{R}^p$ .



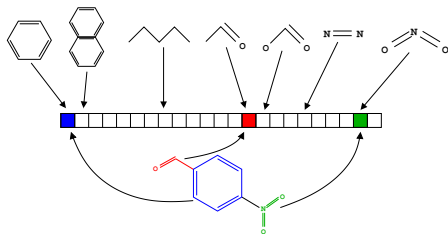
## 2D structural keys in chemoinformatics

- Index a molecule by a binary fingerprint defined by a limited set of **pre-defined** structures



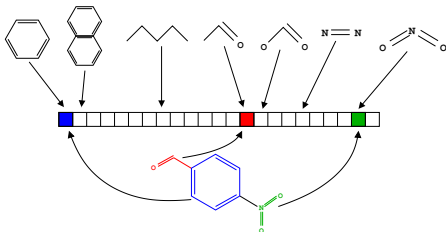
- Use a machine learning algorithms such as SVM, NN, PLS, decision tree, ...

# Challenge: which descriptors (patterns)?



- **Expressiveness**: they should retain as much information as possible from the graph
- **Computation** : they should be fast to compute
- **Large dimension** of the vector representation: memory storage, speed, statistical issues

# Indexing by substructures

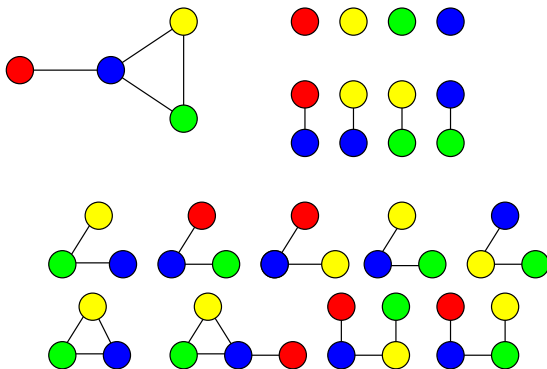


- Often we believe that **the presence substructures** are important predictive patterns
- Hence it makes sense to represent a graph by **features** that indicate the presence (or the number of occurrences) of particular substructures
- However, detecting the presence of particular substructures may be **computationally challenging**...

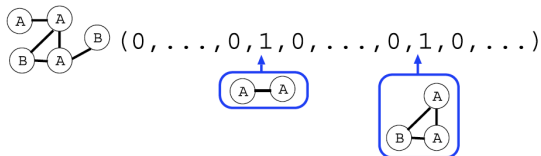
# Subgraphs

## Definition

A **subgraph** of a graph  $(V, E)$  is a connected graph  $(V', E')$  with  $V' \subset V$  and  $E' \subset E$ .



# Indexing by all subgraphs?



## Theorem

*Computing all subgraph occurrences is NP-hard.*

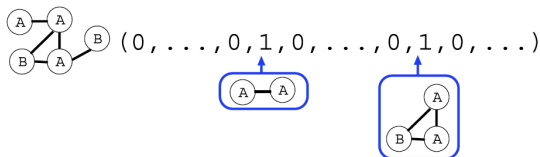
## Proof.

- The linear graph of size  $n$  is a subgraph of a graph  $X$  with  $n$  vertices iff  $X$  has an Hamiltonian path
- The decision problem whether a graph has a Hamiltonian path is NP-complete.





# Indexing by all subgraphs?



## Theorem

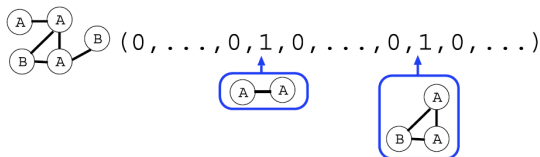
Computing all subgraph occurrences is **NP-hard**.

## Proof.

- The linear graph of size  $n$  is a subgraph of a graph  $X$  with  $n$  vertices iff  $X$  has an Hamiltonian path
- The decision problem whether a graph has a Hamiltonian path is NP-complete.



# Indexing by all subgraphs?



## Theorem

Computing all subgraph occurrences is **NP-hard**.

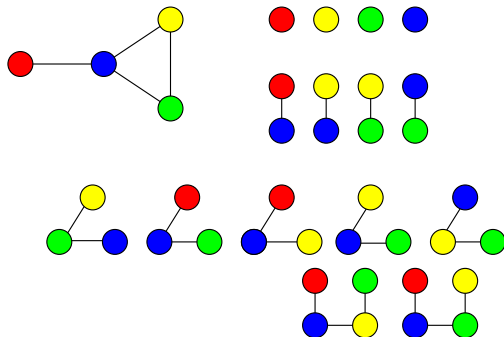
## Proof.

- The linear graph of size  $n$  is a subgraph of a graph  $X$  with  $n$  vertices iff  $X$  has an Hamiltonian path
- The decision problem whether a graph has a Hamiltonian path is NP-complete.

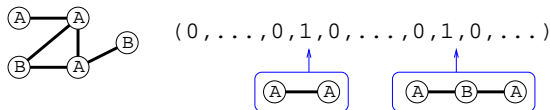


## Definition

- A **path** of a graph  $(V, E)$  is sequence of **distinct vertices**  $v_1, \dots, v_n \in V$  ( $i \neq j \implies v_i \neq v_j$ ) such that  $(v_i, v_{i+1}) \in E$  for  $i = 1, \dots, n-1$ .
- Equivalently the paths are the **linear subgraphs**.



# Indexing by all paths?



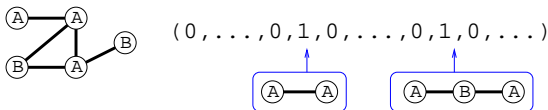
## Theorem

*Computing all path occurrences is NP-hard.*

## Proof.

Same as for subgraphs. □

# Indexing by all paths?



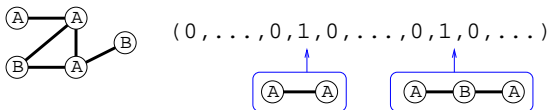
## Theorem

Computing all path occurrences is **NP-hard**.

Proof.

Same as for subgraphs. □

# Indexing by all paths?



## Theorem

Computing all path occurrences is *NP-hard*.

## Proof.

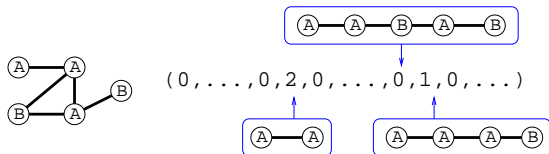
Same as for subgraphs. □

## Substructure selection

We can imagine more limited sets of substructures that lead to more computationnally efficient indexing (non-exhaustive list)

- substructures selected by **domain knowledge** (MDL fingerprint)
- all path **up to length  $k$**  (Openeye fingerprint, Nicholls 2005)
- all **shortest paths** (Borgwardt and Kriegel, 2005)
- all subgraphs **up to  $k$  vertices** (graphlet kernel, Sherashidze et al., 2009)
- all **frequent** subgraphs in the database (Helma et al., 2004)

# Example : Indexing by all shortest paths

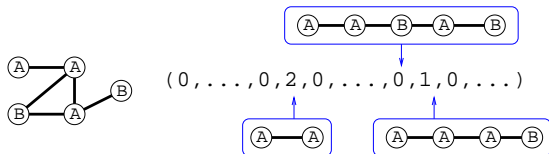


## Properties (Borgwardt and Kriegel, 2005)

- There are  $O(n^2)$  shortest paths.
- The vector of counts can be computed in  $O(n^4)$  with the Floyd-Warshall algorithm.



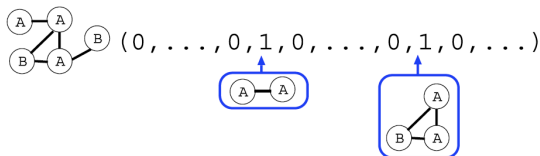
# Example : Indexing by all shortest paths



## Properties (Borgwardt and Kriegel, 2005)

- There are  $O(n^2)$  shortest paths.
- The vector of counts can be computed in  $O(n^4)$  with the Floyd-Warshall algorithm.

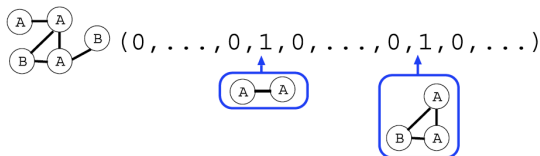
# Example : Indexing by all subgraphs up to $k$ vertices



## Properties (Shervashidze et al., 2009)

- Naive enumeration scales as  $O(n^k)$ .
- Enumeration of connected graphlets in  $O(nd^{k-1})$  for graphs with degree  $\leq d$  and  $k \leq 5$ .
- Randomly sample subgraphs if enumeration is infeasible.

# Example : Indexing by all subgraphs up to $k$ vertices



## Properties (Shervashidze et al., 2009)

- Naive enumeration scales as  $O(n^k)$ .
- Enumeration of connected graphlets in  $O(nd^{k-1})$  for graphs with degree  $\leq d$  and  $k \leq 5$ .
- Randomly sample subgraphs if enumeration is infeasible.

- Explicit computation of substructure occurrences can be **computationally prohibitive** (subgraph, paths)
- Several ideas to **reduce** the set of substructures considered
- In practice, NP-hardness may not be so prohibitive (e.g., graphs with small degrees), the strategy followed should depend on the data considered.

- 1 Explicit computation of features : the case of graph features
- 2 **Using kernels**
  - Introduction to kernels
  - Graph kernels
  - Kernels for gene expression data using gene networks
- 3 Using sparsity-inducing shrinkage estimators
  - Feature selection for all subgraph indexation
  - Classification of array CGH data with piecewise-linear models
  - Structured gene selection for microarray classification
- 4 Conclusion

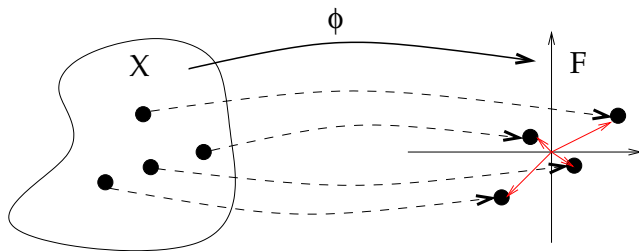
- 1 Explicit computation of features : the case of graph features
- 2 **Using kernels**
  - **Introduction to kernels**
  - Graph kernels
  - Kernels for gene expression data using gene networks
- 3 Using sparsity-inducing shrinkage estimators
  - Feature selection for all subgraph indexation
  - Classification of array CGH data with piecewise-linear models
  - Structured gene selection for microarray classification
- 4 Conclusion

# Positive definite kernels

## Definition

- Let  $\Phi(x)$  be a vector representation of the data  $x$
- The **kernel** between two graphs is defined by:

$$K(x, x') = \Phi(x)^\top \Phi(x').$$



## The trick

- Many linear algorithms for regression or pattern recognition can be **expressed only in terms of inner products** between vectors
- Computing the kernel is often **more efficient** than computing  $\Phi(x)$ , especially in high or infinite dimensions!
- Perhaps we can consider more features with kernels than with explicit feature computation?



# Learning linear classifiers with kernels

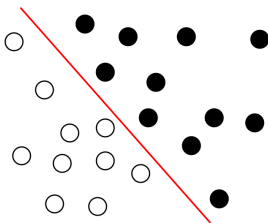
## Training the model

- Minimize an **empirical risk** on the training samples:

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

- ... subject to a **constraint** on  $\beta$ :

$$\|\beta\| \leq C.$$



## Two important strategies (not the only ones!)

- **Feature design** :

$$K(x, x') = \Phi(x)^\top \Phi(x').$$

We illustrate this idea with graph kernels.

- **Regularization design** :

$$\|\beta\| \leq C.$$

We illustrate this idea with kernels for microarray data.

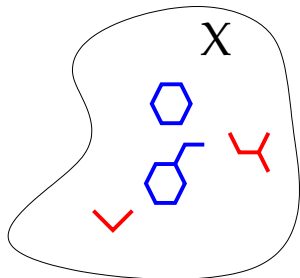
- 1 Explicit computation of features : the case of graph features
- 2 **Using kernels**
  - Introduction to kernels
  - **Graph kernels**
  - Kernels for gene expression data using gene networks
- 3 Using sparsity-inducing shrinkage estimators
  - Feature selection for all subgraph indexation
  - Classification of array CGH data with piecewise-linear models
  - Structured gene selection for microarray classification
- 4 Conclusion

# The idea

- 1 Represent **implicitly** each graph  $x$  by a vector  $\Phi(x) \in \mathcal{H}$  through the kernel

$$K(x, x') = \Phi(x)^\top \Phi(x').$$

- 2 Use a kernel method for classification in  $\mathcal{H}$ .

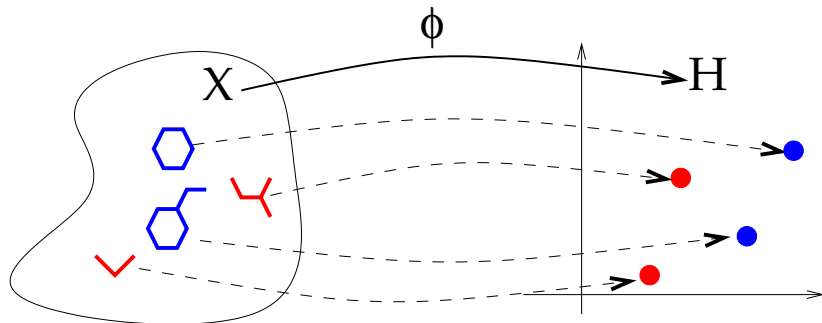


# The idea

- 1 Represent **implicitly** each graph  $x$  by a vector  $\Phi(x) \in \mathcal{H}$  through the kernel

$$K(x, x') = \Phi(x)^\top \Phi(x').$$

- 2 Use a kernel method for classification in  $\mathcal{H}$ .

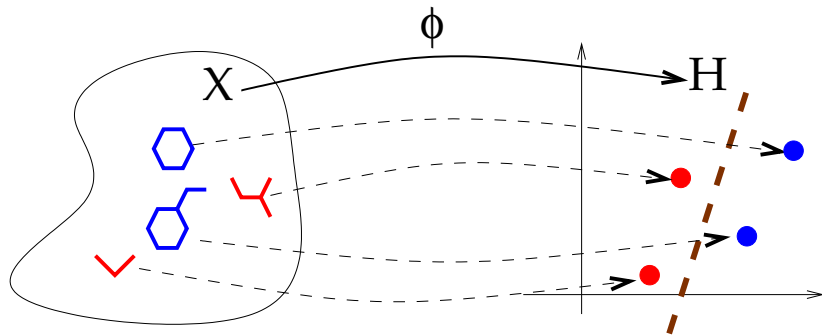


# The idea

- 1 Represent **implicitly** each graph  $x$  by a vector  $\Phi(x) \in \mathcal{H}$  through the kernel

$$K(x, x') = \Phi(x)^\top \Phi(x').$$

- 2 Use a kernel method for classification in  $\mathcal{H}$ .



# Expressiveness vs Complexity

## Definition: Complete graph kernels

A graph kernel is **complete** if it separates non-isomorphic graphs, i.e.:

$$\forall G_1, G_2 \in \mathcal{X}, \quad d_K(G_1, G_2) = 0 \implies G_1 \simeq G_2.$$

Equivalently,  $\Phi(G_1) \neq \Phi(G_2)$  if  $G_1$  and  $G_2$  are not isomorphic.

## Expressiveness vs Complexity trade-off

- If a graph kernel is not complete, then there is **no hope** to learn all possible functions over  $\mathcal{X}$ : the kernel is not **expressive** enough.
- On the other hand, kernel **computation** must be **tractable**, i.e., no more than polynomial (with small degree) for practical applications.
- Can we define **tractable** and **expressive** graph kernels?

# Expressiveness vs Complexity

## Definition: Complete graph kernels

A graph kernel is **complete** if it separates non-isomorphic graphs, i.e.:

$$\forall G_1, G_2 \in \mathcal{X}, \quad d_K(G_1, G_2) = 0 \implies G_1 \simeq G_2.$$

Equivalently,  $\Phi(G_1) \neq \Phi(G_2)$  if  $G_1$  and  $G_2$  are not isomorphic.

## Expressiveness vs Complexity trade-off

- If a graph kernel is not complete, then there is **no hope** to learn all possible functions over  $\mathcal{X}$ : the kernel is not **expressive** enough.
- On the other hand, kernel **computation** must be **tractable**, i.e., no more than polynomial (with small degree) for practical applications.
- Can we define **tractable** and **expressive** graph kernels?



## Proposition (Gärtner et al., 2003)

Computing **any complete graph kernel** is **at least as hard** as the graph isomorphism problem.

## Proof

- For any kernel  $K$  the complexity of computing  $d_K$  is the same as the complexity of computing  $K$ , because:

$$d_K(G_1, G_2)^2 = K(G_1, G_1) + K(G_2, G_2) - 2K(G_1, G_2).$$

- If  $K$  is a complete graph kernel, then computing  $d_K$  solves the graph isomorphism problem ( $d_K(G_1, G_2) = 0$  iff  $G_1 \simeq G_2$ ).  $\square$

# Complexity of complete kernels

## Proposition (Gärtner et al., 2003)

Computing **any complete graph kernel** is **at least as hard** as the graph isomorphism problem.

## Proof

- For any kernel  $K$  the complexity of computing  $d_K$  is the same as the complexity of computing  $K$ , because:

$$d_K(G_1, G_2)^2 = K(G_1, G_1) + K(G_2, G_2) - 2K(G_1, G_2).$$

- If  $K$  is a complete graph kernel, then computing  $d_K$  solves the graph isomorphism problem ( $d_K(G_1, G_2) = 0$  iff  $G_1 \simeq G_2$ ).  $\square$

# Subgraph kernel

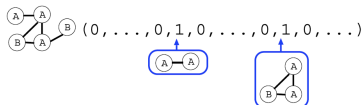
## Definition

- Let  $(\lambda_G)_{G \in \mathcal{X}}$  a set of **nonnegative** real-valued weights
- For any graph  $G \in \mathcal{X}$ , let

$$\forall H \in \mathcal{X}, \quad \Phi_H(G) = |\{G' \text{ is a subgraph of } G : G' \simeq H\}|.$$

- The **subgraph kernel** between any two graphs  $G_1$  and  $G_2 \in \mathcal{X}$  is defined by:

$$K_{\text{subgraph}}(G_1, G_2) = \sum_{H \in \mathcal{X}} \lambda_H \Phi_H(G_1) \Phi_H(G_2).$$



# Subgraph kernel complexity

Proposition (Gärtner et al., 2003)

Computing the subgraph kernel is **NP-hard**.

Proof (1/2)

- Let  $P_n$  be the path graph with  $n$  vertices.
- Subgraphs of  $P_n$  are path graphs:

$$\Phi(P_n) = ne_{P_1} + (n-1)e_{P_2} + \dots + e_{P_n}.$$

- The vectors  $\Phi(P_1), \dots, \Phi(P_n)$  are linearly independent, therefore:

$$e_{P_n} = \sum_{i=1}^n \alpha_i \Phi(P_i),$$

where the coefficients  $\alpha_i$  can be found in polynomial time (solving a  $n \times n$  triangular system).

# Subgraph kernel complexity

## Proposition (Gärtner et al., 2003)

Computing the subgraph kernel is **NP-hard**.

## Proof (1/2)

- Let  $P_n$  be the path graph with  $n$  vertices.
- Subgraphs of  $P_n$  are path graphs:

$$\Phi(P_n) = ne_{P_1} + (n-1)e_{P_2} + \dots + e_{P_n}.$$

- The vectors  $\Phi(P_1), \dots, \Phi(P_n)$  are linearly independent, therefore:

$$e_{P_n} = \sum_{i=1}^n \alpha_i \Phi(P_i),$$

where the coefficients  $\alpha_i$  can be found in polynomial time (solving a  $n \times n$  triangular system).

# Subgraph kernel complexity

## Proposition (Gärtner et al., 2003)

Computing the subgraph kernel is **NP-hard**.

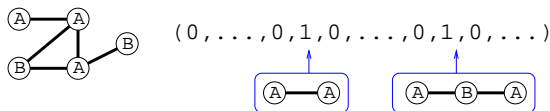
## Proof (2/2)

- If  $G$  is a graph with  $n$  vertices, then it has a path that visits each node exactly once (Hamiltonian path) if and only if  $\Phi(G)^\top e_n > 0$ , i.e.,

$$\Phi(G)^\top \left( \sum_{i=1}^n \alpha_i \Phi(P_i) \right) = \sum_{i=1}^n \alpha_i K_{\text{subgraph}}(G, P_i) > 0.$$

- The decision problem whether a graph has a Hamiltonian path is NP-complete.  $\square$

# Path kernel



## Definition

The **path kernel** is the subgraph kernel restricted to paths, i.e.,

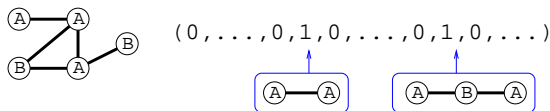
$$K_{path}(G_1, G_2) = \sum_{H \in \mathcal{P}} \lambda_H \Phi_H(G_1) \Phi_H(G_2),$$

where  $\mathcal{P} \subset \mathcal{X}$  is the set of path graphs.

Proposition (Gärtner et al., 2003)

Computing the path kernel is **NP-hard**.

# Path kernel



## Definition

The **path kernel** is the subgraph kernel restricted to paths, i.e.,

$$K_{path}(G_1, G_2) = \sum_{H \in \mathcal{P}} \lambda_H \Phi_H(G_1) \Phi_H(G_2),$$

where  $\mathcal{P} \subset \mathcal{X}$  is the set of path graphs.

## Proposition (Gärtner et al., 2003)

Computing the path kernel is **NP-hard**.

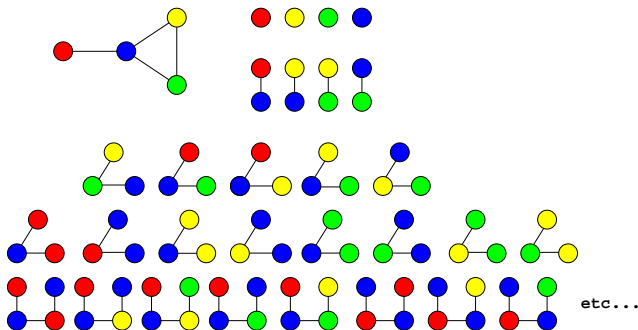


## Expressiveness vs Complexity trade-off

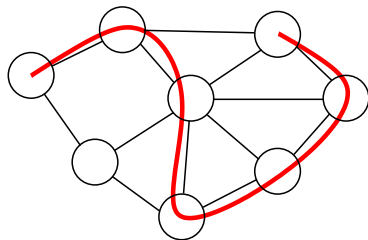
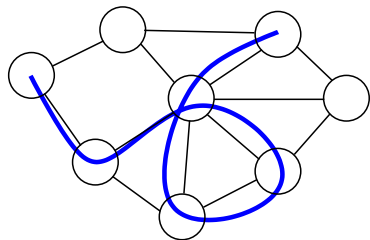
- It is **intractable** to compute **complete** graph kernels.
- It is **intractable** to compute the **subgraph kernels**.
- Restricting subgraphs to be linear does not help: it is also **intractable** to compute the **path kernel**.
- One approach to define polynomial time computable graph kernels is to have the feature space be made up of graphs **homomorphic** to subgraphs, e.g., to consider **walks** instead of paths.

## Definition

- A **walk** of a graph  $(V, E)$  is sequence of  $v_1, \dots, v_n \in V$  such that  $(v_i, v_{i+1}) \in E$  for  $i = 1, \dots, n - 1$ .
- We note  $\mathcal{W}_n(G)$  the set of walks with  $n$  vertices of the graph  $G$ , and  $\mathcal{W}(G)$  the set of all walks.



# Walks $\neq$ paths



## Definition

- Let  $\mathcal{S}_n$  denote the set of all possible **label sequences** of walks of length  $n$  (including vertices and edges labels), and  $\mathcal{S} = \cup_{n \geq 1} \mathcal{S}_n$ .
- For any graph  $\mathcal{X}$  let a **weight**  $\lambda_G(w)$  be associated to each walk  $w \in \mathcal{W}(G)$ .
- Let the feature vector  $\Phi(G) = (\Phi_s(G))_{s \in \mathcal{S}}$  be defined by:

$$\Phi_s(G) = \sum_{w \in \mathcal{W}(G)} \lambda_G(w) \mathbf{1}(s \text{ is the label sequence of } w).$$

- A walk kernel is a graph kernel defined by:

$$K_{walk}(G_1, G_2) = \sum_{s \in \mathcal{S}} \Phi_s(G_1) \Phi_s(G_2).$$

## Definition

- Let  $\mathcal{S}_n$  denote the set of all possible **label sequences** of walks of length  $n$  (including vertices and edges labels), and  $\mathcal{S} = \cup_{n \geq 1} \mathcal{S}_n$ .
- For any graph  $\mathcal{X}$  let a **weight**  $\lambda_G(w)$  be associated to each walk  $w \in \mathcal{W}(G)$ .
- Let the feature vector  $\Phi(G) = (\Phi_s(G))_{s \in \mathcal{S}}$  be defined by:

$$\Phi_s(G) = \sum_{w \in \mathcal{W}(G)} \lambda_G(w) \mathbf{1}(s \text{ is the label sequence of } w).$$

- A walk kernel is a graph kernel defined by:

$$K_{walk}(G_1, G_2) = \sum_{s \in \mathcal{S}} \Phi_s(G_1) \Phi_s(G_2).$$

## Examples

- The  **$n$ th-order walk kernel** is the walk kernel with  $\lambda_G(w) = 1$  if the length of  $w$  is  $n$ , 0 otherwise. It compares two graphs through their common walks of length  $n$ .
- The **random walk kernel** is obtained with  $\lambda_G(w) = P_G(w)$ , where  $P_G$  is a **Markov random walk on  $G$** . In that case we have:

$$K(G_1, G_2) = P(\text{label}(W_1) = \text{label}(W_2)),$$

where  $W_1$  and  $W_2$  are two independent random walks on  $G_1$  and  $G_2$ , respectively (Kashima et al., 2003).

- The **geometric walk kernel** is obtained (when it converges) with  $\lambda_G(w) = \beta^{\text{length}(w)}$ , for  $\beta > 0$ . In that case the feature space is of **infinite dimension** (Gärtner et al., 2003).

## Examples

- The  *$n$ th-order walk kernel* is the walk kernel with  $\lambda_G(w) = 1$  if the length of  $w$  is  $n$ , 0 otherwise. It compares two graphs through their common walks of length  $n$ .
- The *random walk kernel* is obtained with  $\lambda_G(w) = P_G(w)$ , where  $P_G$  is a *Markov random walk on  $G$* . In that case we have:

$$K(G_1, G_2) = P(\text{label}(W_1) = \text{label}(W_2)),$$

where  $W_1$  and  $W_2$  are two independent random walks on  $G_1$  and  $G_2$ , respectively (Kashima et al., 2003).

- The *geometric walk kernel* is obtained (when it converges) with  $\lambda_G(w) = \beta^{\text{length}(w)}$ , for  $\beta > 0$ . In that case the feature space is of *infinite dimension* (Gärtner et al., 2003).

## Examples

- The  **$n$ th-order walk kernel** is the walk kernel with  $\lambda_G(w) = 1$  if the length of  $w$  is  $n$ , 0 otherwise. It compares two graphs through their common walks of length  $n$ .
- The **random walk kernel** is obtained with  $\lambda_G(w) = P_G(w)$ , where  $P_G$  is a **Markov random walk on  $G$** . In that case we have:

$$K(G_1, G_2) = P(\text{label}(W_1) = \text{label}(W_2)),$$

where  $W_1$  and  $W_2$  are two independent random walks on  $G_1$  and  $G_2$ , respectively (Kashima et al., 2003).

- The **geometric walk kernel** is obtained (when it converges) with  $\lambda_G(w) = \beta^{\text{length}(w)}$ , for  $\beta > 0$ . In that case the feature space is of **infinite dimension** (Gärtner et al., 2003).



## Proposition

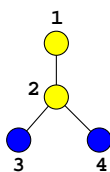
These three kernels ( $n$ th-order, random and geometric walk kernels) can be computed efficiently in **polynomial time**.

# Product graph

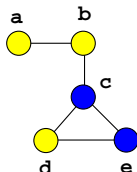
## Definition

Let  $G_1 = (V_1, E_1)$  and  $G_2 = (V_2, E_2)$  be two graphs with labeled vertices. The **product graph**  $G = G_1 \times G_2$  is the graph  $G = (V, E)$  with:

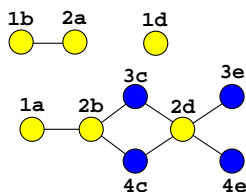
- 1  $V = \{(v_1, v_2) \in V_1 \times V_2 : v_1 \text{ and } v_2 \text{ have the same label}\}$ ,
- 2  $E = \{((v_1, v_2), (v'_1, v'_2)) \in V \times V : (v_1, v'_1) \in E_1 \text{ and } (v_2, v'_2) \in E_2\}$ .



G1



G2



G1 x G2

# Walk kernel and product graph

## Lemma

There is a **bijection** between:

- 1 The **pairs of walks**  $w_1 \in \mathcal{W}_n(G_1)$  and  $w_2 \in \mathcal{W}_n(G_2)$  with the **same label sequences**,
- 2 The **walks on the product graph**  $w \in \mathcal{W}_n(G_1 \times G_2)$ .

## Corollary

$$\begin{aligned} K_{\text{walk}}(G_1, G_2) &= \sum_{s \in \mathcal{S}} \Phi_s(G_1) \Phi_s(G_2) \\ &= \sum_{(w_1, w_2) \in \mathcal{W}(G_1) \times \mathcal{W}(G_1)} \lambda_{G_1}(w_1) \lambda_{G_2}(w_2) \mathbf{1}(l(w_1) = l(w_2)) \\ &= \sum_{w \in \mathcal{W}(G_1 \times G_2)} \lambda_{G_1 \times G_2}(w). \end{aligned}$$

# Walk kernel and product graph

## Lemma

There is a **bijection** between:

- 1 The **pairs of walks**  $w_1 \in \mathcal{W}_n(G_1)$  and  $w_2 \in \mathcal{W}_n(G_2)$  with the **same label sequences**,
- 2 The **walks on the product graph**  $w \in \mathcal{W}_n(G_1 \times G_2)$ .

## Corollary

$$\begin{aligned} K_{\text{walk}}(G_1, G_2) &= \sum_{s \in \mathcal{S}} \Phi_s(G_1) \Phi_s(G_2) \\ &= \sum_{(w_1, w_2) \in \mathcal{W}(G_1) \times \mathcal{W}(G_1)} \lambda_{G_1}(w_1) \lambda_{G_2}(w_2) \mathbf{1}(l(w_1) = l(w_2)) \\ &= \sum_{w \in \mathcal{W}(G_1 \times G_2)} \lambda_{G_1 \times G_2}(w). \end{aligned}$$

# Computation of the $n$ th-order walk kernel

- For the  $n$ th-order walk kernel we have  $\lambda_{G_1 \times G_2}(w) = 1$  if the length of  $w$  is  $n$ , 0 otherwise.
- Therefore:

$$K_{nth-order}(G_1, G_2) = \sum_{w \in \mathcal{W}_n(G_1 \times G_2)} 1.$$

- Let  $A$  be the adjacency matrix of  $G_1 \times G_2$ . Then we get:

$$K_{nth-order}(G_1, G_2) = \sum_{i,j} [A^n]_{i,j} = \mathbf{1}^\top A^n \mathbf{1}.$$

- Computation in  $O(n|G_1||G_2|d_1d_2)$ , where  $d_i$  is the maximum degree of  $G_i$ .

# Computation of random and geometric walk kernels

- In both cases  $\lambda_G(w)$  for a walk  $w = v_1 \dots v_n$  can be decomposed as:

$$\lambda_G(v_1 \dots v_n) = \lambda^i(v_1) \prod_{i=2}^n \lambda^t(v_{i-1}, v_i).$$

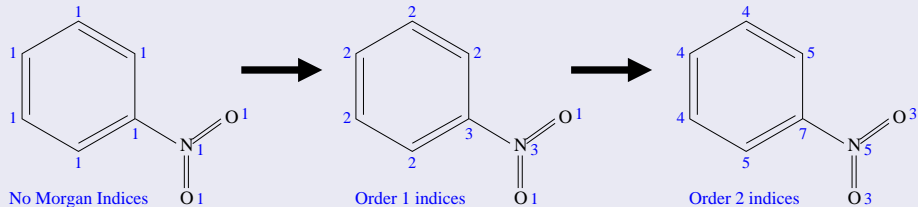
- Let  $\Lambda_i$  be the vector of  $\lambda^i(v)$  and  $\Lambda_t$  be the matrix of  $\lambda^t(v, v')$ :

$$\begin{aligned} K_{walk}(G_1, G_2) &= \sum_{n=1}^{\infty} \sum_{w \in \mathcal{W}_n(G_1 \times G_2)} \lambda^i(v_1) \prod_{i=2}^n \lambda^t(v_{i-1}, v_i) \\ &= \sum_{n=0}^{\infty} \Lambda_i \Lambda_t^n \mathbf{1} \\ &= \Lambda_i (I - \Lambda_t)^{-1} \mathbf{1} \end{aligned}$$

- Computation in  $O(|G_1|^3 |G_2|^3)$

# Extensions 1: label enrichment

## Atom relabeling with the Morgan index



- **Compromise** between **fingerprints** and **structural keys features**.
- Other **relabeling** schemes are possible (graph coloring).
- **Faster computation with more labels** (less matches implies a smaller product graph).

## Extension 2: Non-tottering walk kernel

### Tottering walks

A **tottering walk** is a walk  $w = v_1 \dots v_n$  with  $v_i = v_{i+2}$  for some  $i$ .



**Non-tottering**



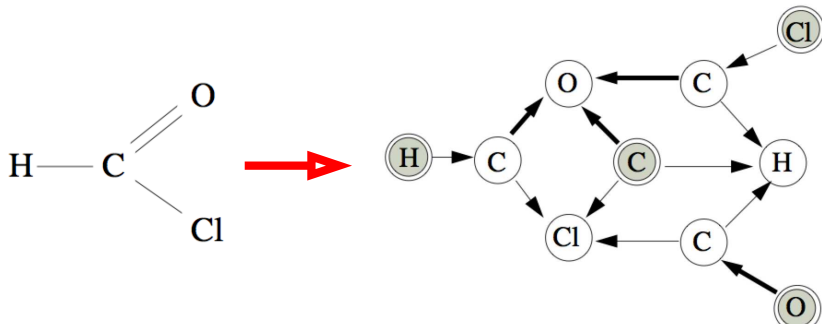
**Tottering**

- Tottering walks seem **irrelevant** for many applications
- Focusing on non-tottering walks is a way to get closer to the **path kernel** (e.g., equivalent on trees).

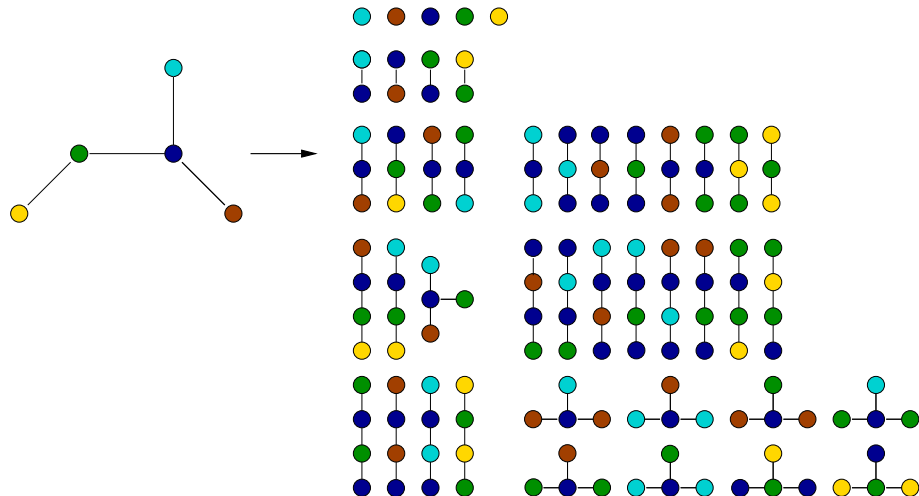


# Computation of the non-tottering walk kernel (Mahé et al., 2005)

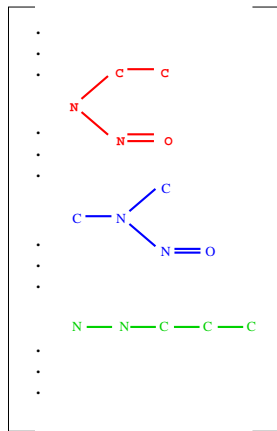
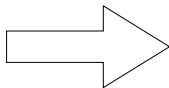
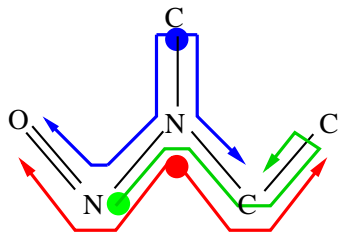
- **Second-order** Markov random walk to prevent tottering walks
- Written as a **first-order** Markov random walk on an **augmented graph**
- **Normal** walk kernel on the augmented graph (which is always a **directed** graph).



# Extension 3: Subtree kernels



# Example: Tree-like fragments of molecules



# Computation of the subtree kernel

- Like the walk kernel, amounts to compute the (weighted) number of subtrees in the **product graph**.
- Recursion: if  $\mathcal{T}(v, n)$  denotes the weighted number of subtrees of depth  $n$  rooted at the vertex  $v$ , then:

$$\mathcal{T}(v, n+1) = \sum_{R \subset \mathcal{N}(v)} \prod_{v' \in R} \lambda_t(v, v') \mathcal{T}(v', n),$$

where  $\mathcal{N}(v)$  is the set of neighbors of  $v$ .

- Can be combined with the non-tottering graph transformation as preprocessing to obtain the **non-tottering subtree kernel**.

## MUTAG dataset

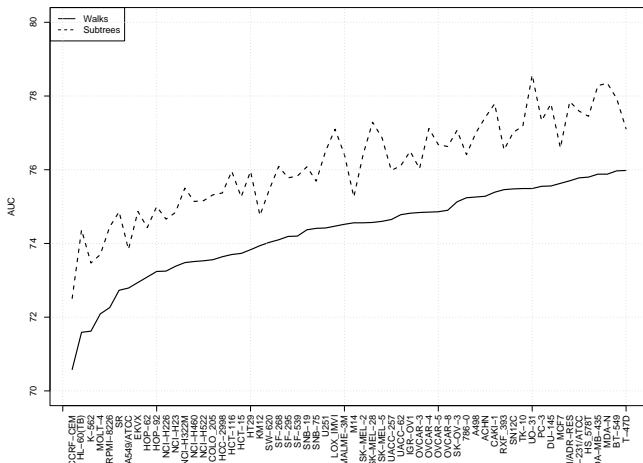
- aromatic/hetero-aromatic compounds
- high mutagenic activity /no mutagenic activity, assayed in *Salmonella typhimurium*.
- 188 compounds: 125 + / 63 -

## Results

10-fold cross-validation accuracy

Method	Accuracy
Progol1	81.4%
2D kernel	91.2%

# 2D Subtree vs walk kernels

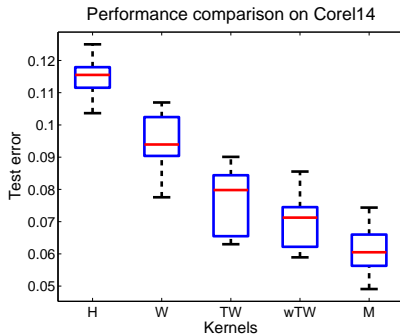
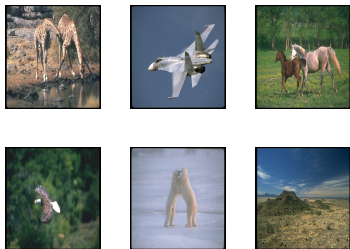


Screening of inhibitors for 60 cancer cell lines.

# Image classification (Harchaoui and Bach, 2007)

## COREL14 dataset

- 1400 natural images in 14 classes
- Compare kernel between histograms (H), walk kernel (W), subtree kernel (TW), weighted subtree kernel (wTW), and a combination (M).



## What we saw

- Kernels do **not allow** to overcome the NP-hardness of subgraph patterns
- They allow to work with approximate subgraphs (walks, subtrees), in infinite dimension, thanks to the **kernel trick**
- However: using kernels makes it difficult to **come back to patterns** after the learning stage



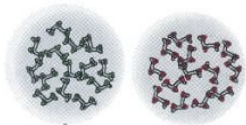
- 1 Explicit computation of features : the case of graph features
- 2 Using kernels**
  - Introduction to kernels
  - Graph kernels
  - Kernels for gene expression data using gene networks**
- 3 Using sparsity-inducing shrinkage estimators
  - Feature selection for all subgraph indexation
  - Classification of array CGH data with piecewise-linear models
  - Structured gene selection for microarray classification
- 4 Conclusion

# Microarrays measure gene expression

Make cDNA reverse transcript  
Label cDNAs with fluorescent dyes

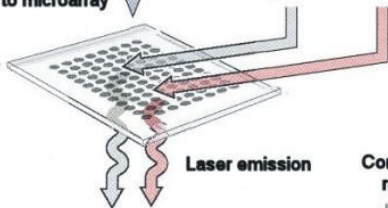
Control

Experimental

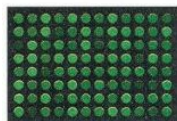


Hybridization to microarray

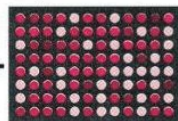
Laser excitation at dye-specific Hz



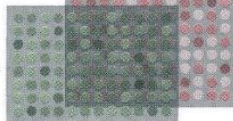
**Red = "up-regulation"**  
**Green = "down-regulation"**  
**Black = constitutive expression**



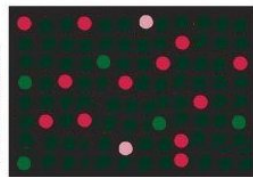
+



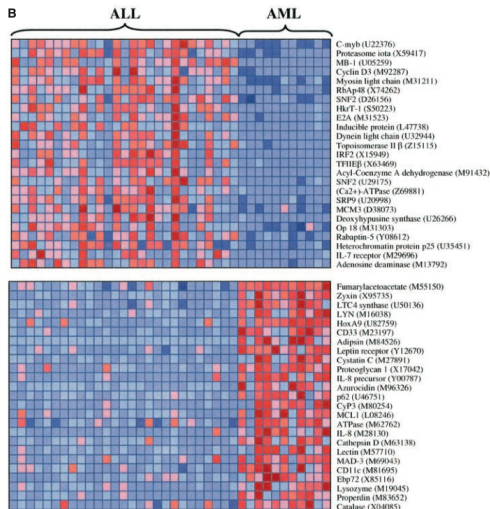
=



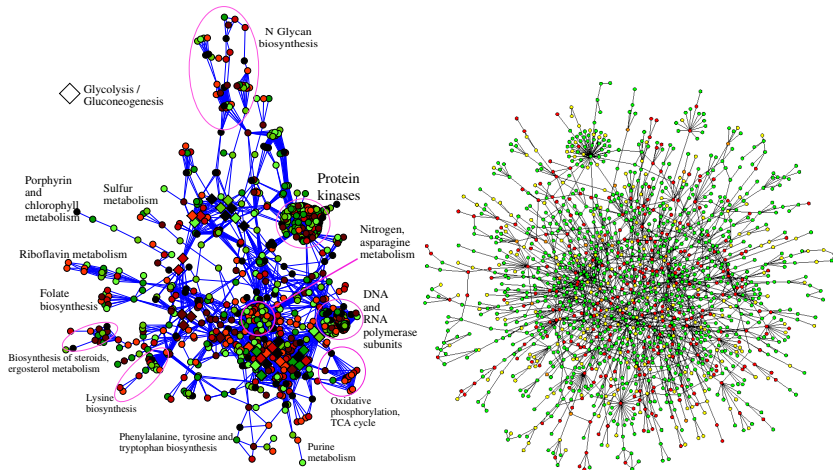
Computer calculates ratio of intensity



# Cancer classification from microarray data

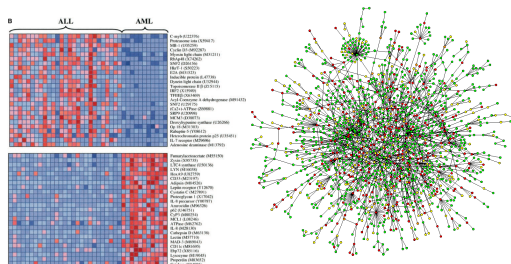


# Gene networks

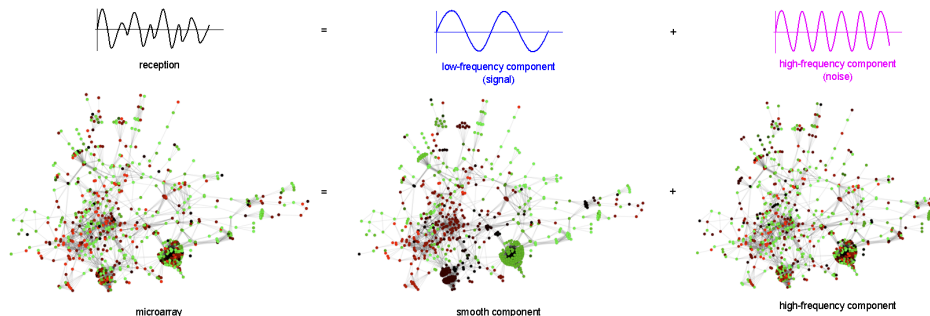


## Motivation

- 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**



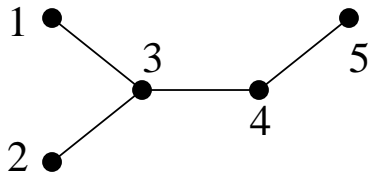
# An idea



- 1 Use the gene network to extract the “important information” in gene expression profiles by **Fourier analysis** on the graph
- 2 Learn a linear classifier on the **smooth components**

## Definition

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



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

- $L$  is positive semidefinite
- The **eigenvectors**  $\mathbf{e}_1, \dots, \mathbf{e}_n$  of  $L$  with eigenvalues  $0 = \lambda_1 \leq \dots \leq \lambda_n$  form a basis called **Fourier basis**
- For any  $f : V \rightarrow \mathbb{R}$ , the **Fourier transform** of  $f$  is the vector  $\hat{f} \in \mathbb{R}^n$  defined by:

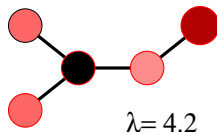
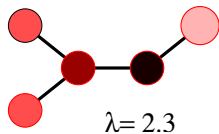
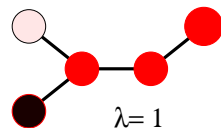
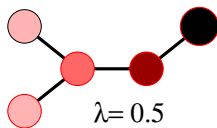
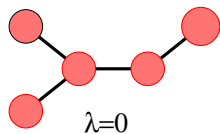
$$\hat{f}_i = f^\top \mathbf{e}_i, \quad i = 1, \dots, n.$$

- The **inverse Fourier formula** holds:

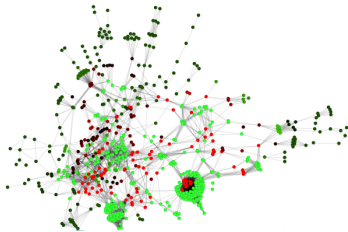
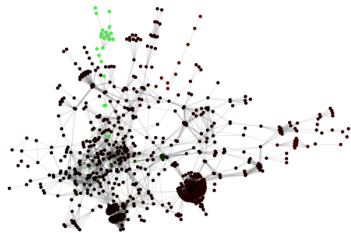
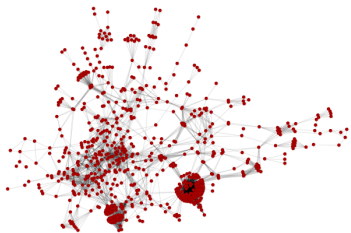
$$f = \sum_{i=1}^n \hat{f}_i \mathbf{e}_i.$$



# Fourier basis



# Fourier basis



## Definition

- Let  $\phi : \mathbb{R}^+ \rightarrow \mathbb{R}^+$  be **non-increasing**.
- A smoothing operator  $S_\phi$  transform a function  $f : V \rightarrow \mathbb{R}$  into a smoothed version:

$$S_\phi(f) = \sum_{i=1}^n \hat{f}_i \phi(\lambda_i) e_i.$$

## Examples

- Identity operator ( $S_\phi(f) = f$ ):

$$\phi(\lambda) = 1, \quad \forall \lambda$$

- Low-pass filter:

$$\phi(\lambda) = \begin{cases} 1 & \text{if } \lambda \leq \lambda^*, \\ 0 & \text{otherwise.} \end{cases}$$

- Attenuation of high frequencies:

$$\phi(\lambda) = \exp(-\beta\lambda).$$

## Examples

- Identity operator ( $S_\phi(f) = f$ ):

$$\phi(\lambda) = 1, \quad \forall \lambda$$

- Low-pass filter:

$$\phi(\lambda) = \begin{cases} 1 & \text{if } \lambda \leq \lambda^*, \\ 0 & \text{otherwise.} \end{cases}$$

- Attenuation of high frequencies:

$$\phi(\lambda) = \exp(-\beta\lambda).$$

## Examples

- Identity operator ( $S_\phi(f) = f$ ):

$$\phi(\lambda) = 1, \quad \forall \lambda$$

- Low-pass filter:

$$\phi(\lambda) = \begin{cases} 1 & \text{if } \lambda \leq \lambda^*, \\ 0 & \text{otherwise.} \end{cases}$$

- Attenuation of high frequencies:

$$\phi(\lambda) = \exp(-\beta\lambda).$$

## Working with smoothed profiles

- Classical methods for linear classification and regression with a ridge penalty solve:

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

- Applying these algorithms on the smooth profiles means solving:

$$\min_{\beta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n l(\beta^\top \mathbf{S}_\phi(f_i), y_i) + \lambda \beta^\top \beta.$$

## Lemma

This is equivalent to:

$$\min_{v \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n l(v^\top f_i, y_i) + \lambda \sum_{i=1}^p \frac{\hat{v}_i^2}{\phi(\lambda_i)},$$

hence the linear classifier  $v$  is **smooth**.

## Proof

- Let  $v = \sum_{i=1}^n \phi(\lambda_i) e_i e_i^\top \beta$ , then

$$\beta^\top \mathcal{S}_\phi(f_i) = \beta^\top \sum_{i=1}^n \hat{f}_i \phi(\lambda_i) e_i = f_i^\top v.$$

- Then  $\hat{v}_i = \phi(\lambda_i) \hat{\beta}_i$  and  $\beta^\top \beta = \sum_{i=1}^n \frac{\hat{v}_i^2}{\phi(\lambda_i)^2}$ .



## Lemma

This is equivalent to:

$$\min_{v \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n l(v^\top f_i, y_i) + \lambda \sum_{i=1}^p \frac{\hat{v}_i^2}{\phi(\lambda_i)},$$

hence the linear classifier  $v$  is **smooth**.

## Proof

- Let  $v = \sum_{i=1}^n \phi(\lambda_i) e_i e_i^\top \beta$ , then

$$\beta^\top S_\phi(f_i) = \beta^\top \sum_{i=1}^n \hat{f}_i \phi(\lambda_i) e_i = f_i^\top v.$$

- Then  $\hat{v}_i = \phi(\lambda_i) \hat{\beta}_i$  and  $\beta^\top \beta = \sum_{i=1}^n \frac{\hat{v}_i^2}{\phi(\lambda_i)^2}$ .

## Smoothing kernel

Kernel methods (SVM, kernel ridge regression..) only need the **inner product** between smooth profiles:

$$\begin{aligned}K(f, g) &= S_\phi(f)^\top S_\phi(g) \\&= \sum_{i=1}^n \hat{f}_i \hat{g}_i \phi(\lambda_i)^2 \\&= f^\top \left( \sum_{i=1}^n \phi(\lambda_i)^2 \mathbf{e}_i \mathbf{e}_i^\top \right) g \\&= f^\top K_\phi g,\end{aligned}\tag{1}$$

with

$$K_\phi = \sum_{i=1}^n \phi(\lambda_i)^2 \mathbf{e}_i \mathbf{e}_i^\top.$$

# Examples

- For  $\phi(\lambda) = \exp(-t\lambda)$ , we recover the **diffusion kernel**:

$$K_\phi = \exp_M(-2tL).$$

- For  $\phi(\lambda) = 1/\sqrt{1+\lambda}$ , we obtain

$$K_\phi = (L + I)^{-1},$$

and the penalization is:

$$\sum_{i=1}^n \frac{\hat{v}_i^2}{\phi(\lambda_i)} = v^\top (L + I) v = \|v\|_2^2 + \sum_{i \sim j} (v_i - v_j)^2.$$

- For  $\phi(\lambda) = \exp(-t\lambda)$ , we recover the **diffusion kernel**:

$$K_\phi = \exp_M(-2tL).$$

- For  $\phi(\lambda) = 1/\sqrt{1+\lambda}$ , we obtain

$$K_\phi = (L + I)^{-1},$$

and the penalization is:

$$\sum_{i=1}^n \frac{\hat{v}_i^2}{\phi(\lambda_i)} = \mathbf{v}^\top (L + I) \mathbf{v} = \|\mathbf{v}\|_2^2 + \sum_{i \sim j} (v_i - v_j)^2.$$

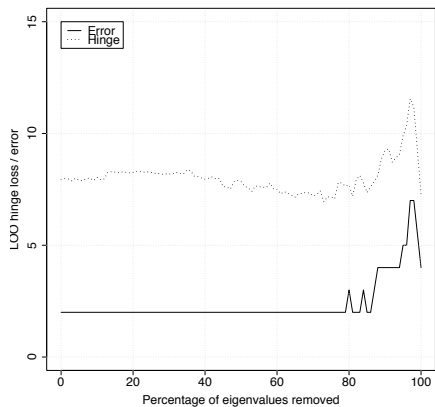
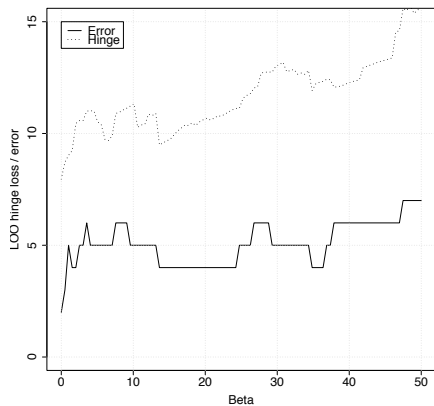
## Expression

- Study the effect of low irradiation doses on the yeast
- 12 non irradiated vs 6 irradiated
- Which pathways are involved in the response at the transcriptomic level?

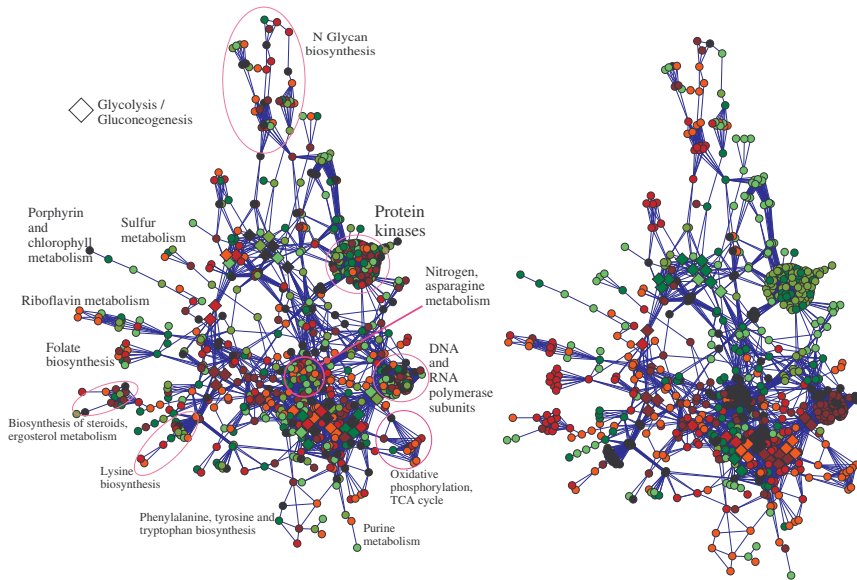
## Graph

- KEGG database of metabolic pathways
- Two genes are connected if they code for enzymes that catalyze successive reactions in a pathway (**metabolic gene network**).
- 737 genes, 4694 vertices.

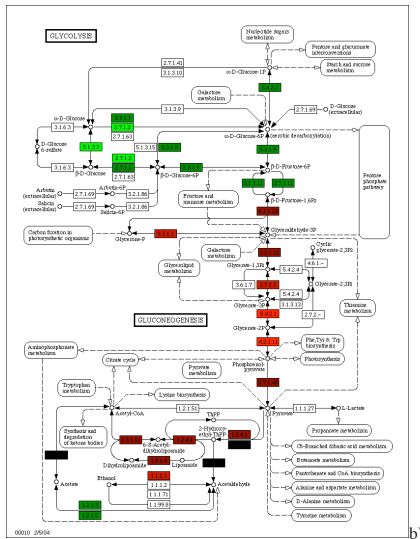
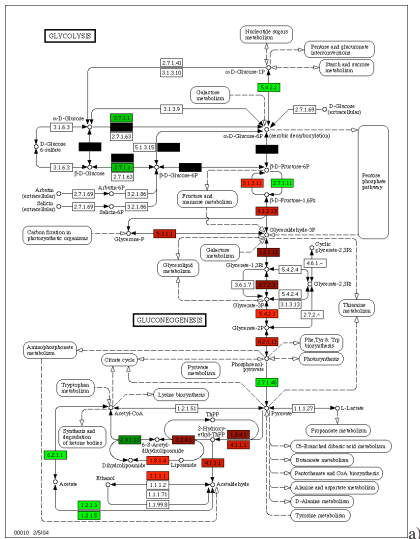
# Classification performance



# Classifier



# Classifier



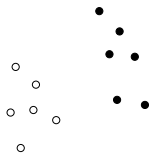


- With kernels we are able to **soft constrain** the shape of the classifier through regularization, e.g.:

$$\min_{v \in \mathbb{R}^p} R_{emp}(v) + \lambda \sum_{i=1}^p \frac{\hat{v}_i^2}{\phi(\lambda_i)},$$

- This is related to **priors** in Bayesian learning
- The resulting classifier is **interpretable**, even without selection of a specific list of features.

- 1 Explicit computation of features : the case of graph features
- 2 Using kernels
  - Introduction to kernels
  - Graph kernels
  - Kernels for gene expression data using gene networks
- 3 Using sparsity-inducing shrinkage estimators
  - Feature selection for all subgraph indexation
  - Classification of array CGH data with piecewise-linear models
  - Structured gene selection for microarray classification
- 4 Conclusion



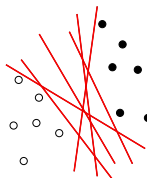
## Training the model

- Minimize an **empirical risk** on the training samples:

$$\min_{\beta \in \mathbb{R}^{p+1}} R_{emp}(\beta) = \frac{1}{n} \sum_{i=1}^n l(f_{\beta}(x_i), y_i),$$

- ... subject to some **constraint** on  $\beta$ , e.g.:

$$\Omega(\beta) \leq C.$$



## Training the model

- Minimize an **empirical risk** on the training samples:

$$\min_{\beta \in \mathbb{R}^{p+1}} R_{emp}(\beta) = \frac{1}{n} \sum_{i=1}^n l(f_{\beta}(x_i), y_i),$$

- ... subject to some **constraint** on  $\beta$ , e.g.:

$$\Omega(\beta) \leq C.$$

# Example : Norm Constraints

## The approach

A common method in statistics to learn with few samples in high dimension is to **constrain the Euclidean norm of  $\beta$**

$$\Omega_{\text{ridge}}(\beta) = \|\beta\|_2^2 = \sum_{i=1}^p \beta_i^2,$$

(ridge regression, support vector machines, kernel methods...)

### Pros

- Good performance in classification

### Cons

- Limited interpretation (small weights)
- No prior biological knowledge

# Example : Feature Selection

## The approach

Constrain most weights to be 0, i.e., **select a few genes** whose expression are sufficient for classification.

$$\Omega_{\text{Best subset selection}}(\beta) = \|\beta\|_0 = \sum_{i=1}^p \mathbf{1}(\beta_i > 0).$$

This is usually a NP-hard problem, many greedy variants have been proposed (filter methods, wrapper methods)

### Pros

- Good performance
- **Biomarker** selection
- Interpretability

### Cons

- NP-hard
- Gene selection **not robust**
- No use of prior knowledge

# Example : Sparsity inducing convex priors

## The approach

Constrain most weights to be 0 through a convex non-differentiable penalty:

$$\Omega_{\text{LASSO}}(\beta) = \|\beta\|_1 = \sum_{i=1}^p |\beta_i| .$$

- Several variants exist, e.g., **elastic net** penalty ( $\|\beta\|_1 + \|\beta\|_2$ ), ... )

## Pros

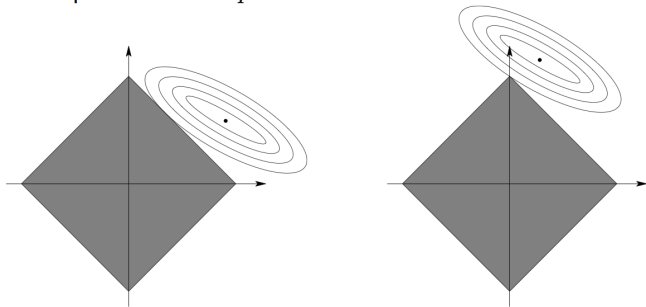
- Good performance
- **Biomarker** selection
- Interpretability

## Cons

- Gene selection **not robust**
- No use of prior knowledge

# Why LASSO leads to sparse solutions

Geometric interpretation with  $p = 2$

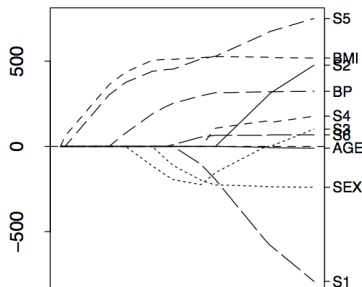




# Efficiently computation of the regularization path

$$\min_{\beta \in \mathbb{R}^{p+1}} R^n(f_\beta) = \sum_{i=1}^n (f_\beta(\mathbf{x}_i) - \mathbf{y}_i)^2 + \lambda \sum_{i=1}^p |\beta_i| \quad (2)$$

- No explicit solution, but this is just a quadratic program.
- **LARS** (Efron et al., 2004) provides a fast algorithm to compute the solution for all  $\lambda$ 's simultaneously (regularization path)



## The idea

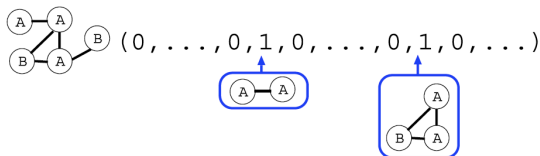
- If we have a specific **prior knowledge** about the “correct” weights, it can be included in  $\Omega$  in the constraint:

Minimize  $R_{emp}(\beta)$  subject to  $\Omega(\beta) \leq C$ .

- If we design a **convex** function  $\Omega$ , then the algorithm boils down to a convex optimization problem (usually **easy to solve**).
- Similar to priors in Bayesian statistics

- 1 Explicit computation of features : the case of graph features
- 2 Using kernels
  - Introduction to kernels
  - Graph kernels
  - Kernels for gene expression data using gene networks
- 3 Using sparsity-inducing shrinkage estimators
  - **Feature selection for all subgraph indexation**
  - Classification of array CGH data with piecewise-linear models
  - Structured gene selection for microarray classification
- 4 Conclusion

# Motivation



- Indexing by all subgraphs is appealing but **intractable** in practice (both explicitly and with the kernel trick)
- Can we work **implicitly** with this representation using **sparse** learning, e.g., LASSO regression or boosting?
- This may lead to both **accurate predictive** model and the identification of **discriminative patterns**.
- The iterations of LARS or boosting amount to an **optimization** problem over subgraphs, which may be solved efficiently using graph mining technique...

- **Weak learner** = decision stump indexed by subgraph  $H$  and  $\alpha = \pm 1$ :

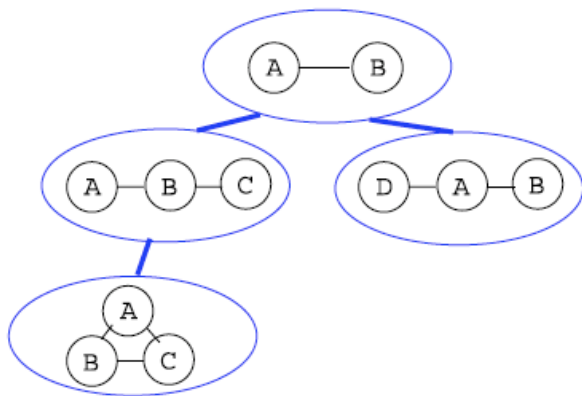
$$h_{\alpha,H}(G) = \alpha \Phi_H(G)$$

- **Boosting**: at each iteration, for a given distribution  $d_1 + \dots + d_n = 1$  over the training points  $(G_i, y_i)$ , select a weak learner (subgraph  $\tilde{H}$ ) which **maximizes the gain**

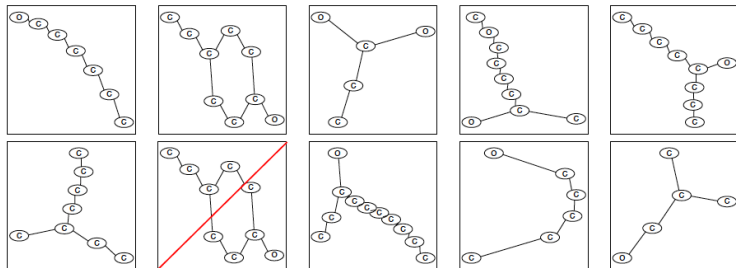
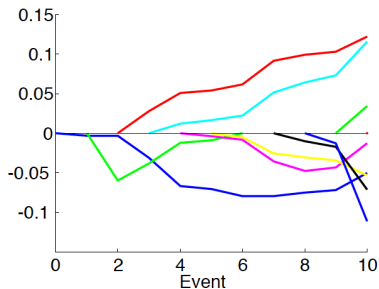
$$\text{gain}(H, \alpha) = \sum_{i=1}^n y_i h_{\alpha,H}(G_i).$$

- This can be done "efficiently" by branch-and-bound over a **DFS code tree** (Yan and Han, 2002).

# The DFS code tree



# Graph LASSO regularization path (Tsuda, 2007)

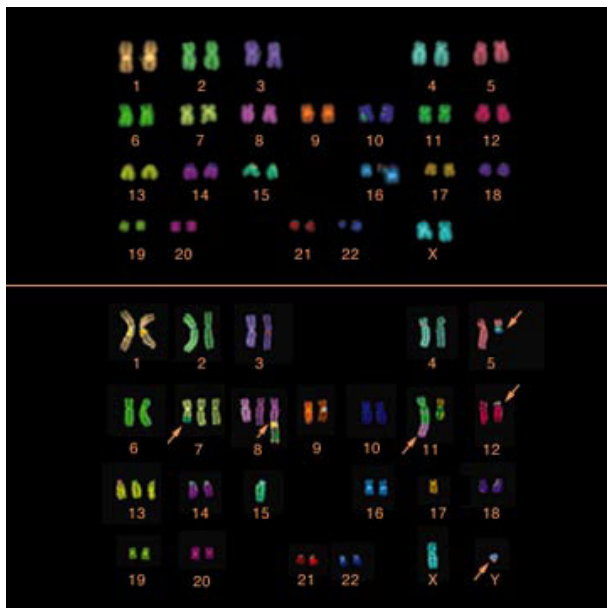


- Sparse learning is practically feasible in the space of graphs indexed by all subgraphs
- Leads to subgraph selection
- Several extensions
  - LASSO regularization path (Tsuda, 2007)
  - gboost (Saigo et al., 2009)
- A beautiful and promising marriage between machine learning and data mining



- 1 Explicit computation of features : the case of graph features
- 2 Using kernels
  - Introduction to kernels
  - Graph kernels
  - Kernels for gene expression data using gene networks
- 3 Using sparsity-inducing shrinkage estimators
  - Feature selection for all subgraph indexation
  - **Classification of array CGH data with piecewise-linear models**
  - Structured gene selection for microarray classification
- 4 Conclusion

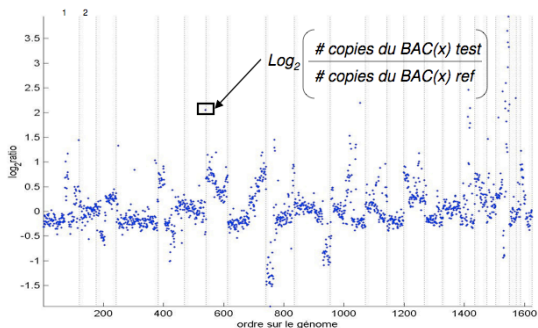
# Chromosomal aberrations in cancer



# Comparative Genomic Hybridization (CGH)

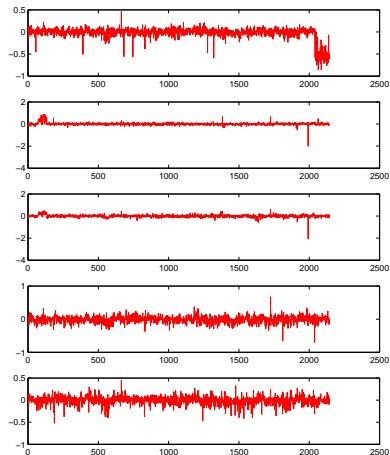
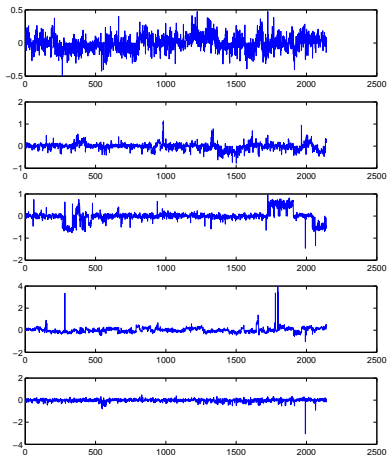
## Motivation

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



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

# Aggressive vs non-aggressive melanoma



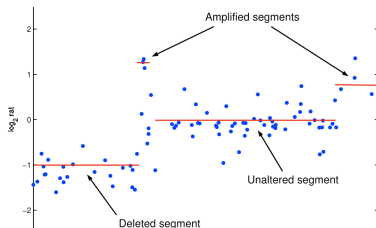
# Classification of array CGH

## Prior knowledge

- Let  $\mathbf{x}$  be a CGH profile
- We focus on linear classifiers, i.e., the sign of :

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

- We expect  $\beta$  to be
  - **sparse** : only a few positions should be discriminative
  - **piecewise constant** : within a region, all probes should contribute equally

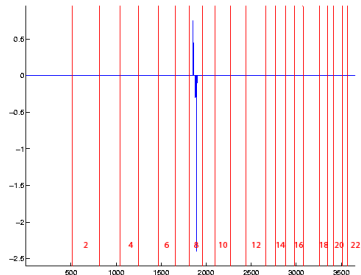
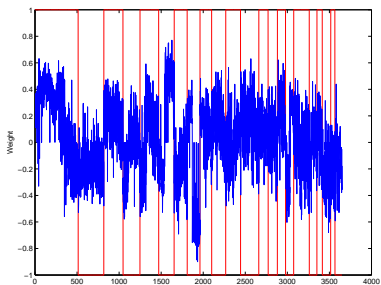
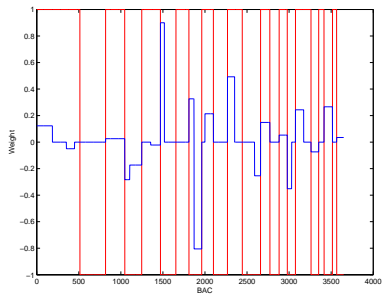


## The fused LASSO penalty (Tibshirani et al., 2005)

$$\Omega_{fusedlasso}(\beta) = \sum_i |\beta_i| + \sum_{i \sim j} |\beta_i - \beta_j|.$$

- First term leads to **sparse** solutions
- Second term leads to **piecewise constant** solutions
- Combined with a hinge loss leads to a **fused SVM** (Rapaport et al., 2008);

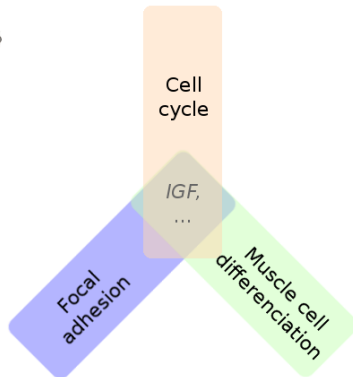
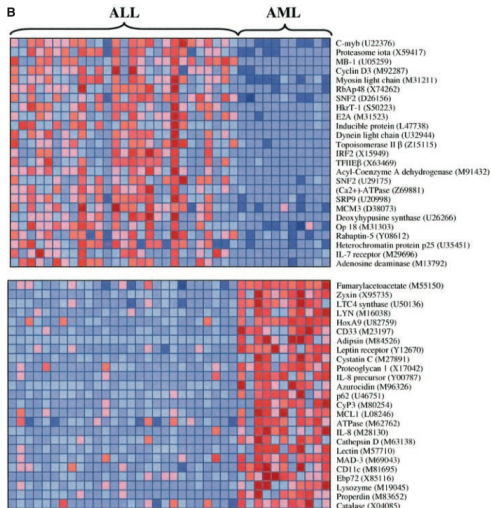
# Application: metastasis prognosis in melanoma



- 1 Explicit computation of features : the case of graph features
- 2 Using kernels
  - Introduction to kernels
  - Graph kernels
  - Kernels for gene expression data using gene networks
- 3 Using sparsity-inducing shrinkage estimators
  - Feature selection for all subgraph indexation
  - Classification of array CGH data with piecewise-linear models
  - **Structured gene selection for microarray classification**
- 4 Conclusion



# How to select jointly genes belonging to the same pathways?

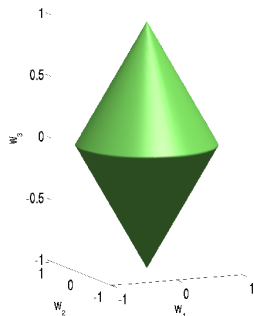


# Selecting pre-defined groups of variables

## 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*:

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

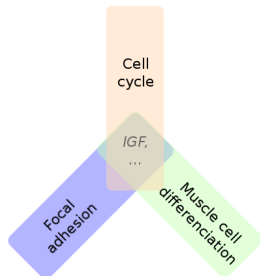


$$\Omega(w_1, w_2, w_3) = \|(w_1, w_2)\|_2 + \|w_3\|_2$$

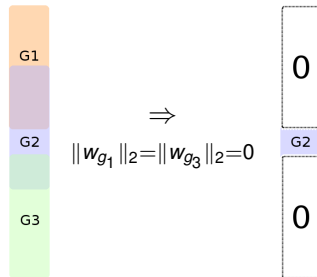
# What if a gene belongs to several groups?

## Issue of using the group-lasso

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

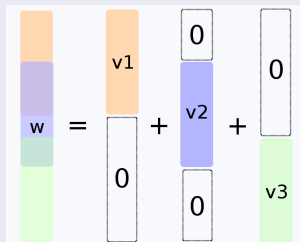


Removal of *any* group containing a gene  $\Rightarrow$  the weight of the gene is 0.

## An idea

Introduce latent variables  $v_g$ :

$$\begin{cases} \min_{w,v} L(w) + \lambda \sum_{g \in \mathcal{G}} \|v_g\|_2 \\ w = \sum_{g \in \mathcal{G}} v_g \\ \text{supp}(v_g) \subseteq g. \end{cases}$$



## Properties

- Resulting support is a *union* of groups in  $\mathcal{G}$ .
- Possible to select one variable without selecting all the groups containing it.
- Setting one  $v_g$  to 0 doesn't necessarily set to 0 all its variables in  $w$ .

## Overlap norm

$$\left\{ \begin{array}{l} \min_{w,v} L(w) + \lambda \sum_{g \in \mathcal{G}} \|v_g\|_2 \\ w = \sum_{g \in \mathcal{G}} v_g \\ \text{supp}(v_g) \subseteq g. \end{array} \right. = \min_w L(w) + \lambda \Omega_{\text{overlap}}(w)$$

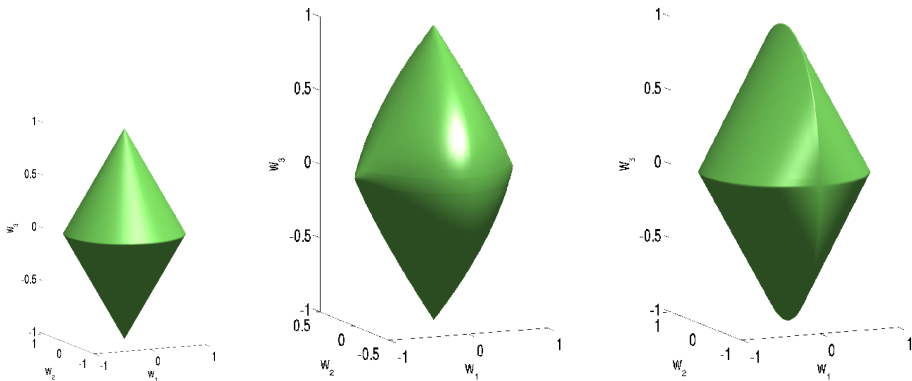
with

$$\Omega_{\text{overlap}}(w) \triangleq \left\{ \begin{array}{l} \min_v \sum_{g \in \mathcal{G}} \|v_g\|_2 \\ w = \sum_{g \in \mathcal{G}} v_g \\ \text{supp}(v_g) \subseteq g. \end{array} \right. \quad (*)$$

## Property

- $\Omega_{\text{overlap}}(w)$  is a norm of  $w$ .
- $\Omega_{\text{overlap}}(\cdot)$  associates to  $w$  a specific (not necessarily unique) decomposition  $(v_g)_{g \in \mathcal{G}}$  which is the argmin of  $(*)$ .

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

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

Then

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

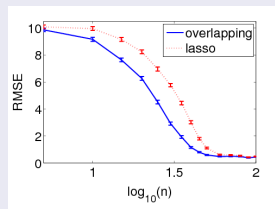
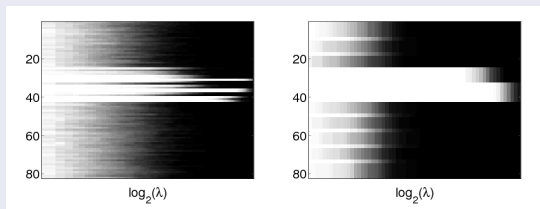
the optimal solution  $\hat{w}$  admits a unique decomposition  $(\hat{v}_g)_{g \in \mathcal{G}}$  such that

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



## Synthetic data: overlapping groups

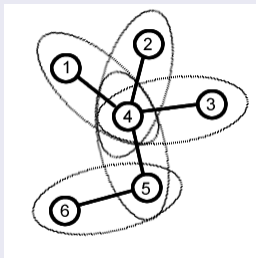
- 10 groups of 10 variables with 2 variables of overlap between two successive groups :  $\{1, \dots, 10\}, \{9, \dots, 18\}, \dots, \{73, \dots, 82\}$ .
- Support: union of 4<sup>th</sup> and 5<sup>th</sup> groups.
- Learn from 100 training points.



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

## Graph lasso

- Consider groups that are subgraphs whose union would give such connected components (e.g., edges  $E$ ).



- $\Omega_{\text{graph}}(\mathbf{w}) = \min_{\mathbf{v} \in \mathcal{V}_E} \sum_{e \in E} \|v_e\| \quad \text{s.t.} \quad \sum_{e \in E} v_e = \mathbf{w}, \text{supp}(v_e) = e.$

# Graph lasso vs kernel on graph

- 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	$\ell_1$	$\Omega_{\text{OVERLAP}}^G(\cdot)$
ERROR	$0.38 \pm 0.04$	$0.36 \pm 0.03$
# PATH.	148, 58, 183	6, 5, 78
PROP. PATH.	0.32, 0.14, 0.41	0.01, 0.01, 0.17

- 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.	1.1, 1, 1.0	1.3, 1.4, 1.2

- 1 Explicit computation of features : the case of graph features
- 2 Using kernels
  - Introduction to kernels
  - Graph kernels
  - Kernels for gene expression data using gene networks
- 3 Using sparsity-inducing shrinkage estimators
  - Feature selection for all subgraph indexation
  - Classification of array CGH data with piecewise-linear models
  - Structured gene selection for microarray classification
- 4 Conclusion

- Machine learning with complex and structured data becomes the rule
- We surveyed several ideas
  - Feature construction
  - Learning with kernels
  - Learning with sparsity
- Performance and interpretability are both important
- Many promising bridges between machine learning and data mining!