

QSAR and Virtual Screening with Support Vector Machines

Jean-Philippe Vert

Jean-Philippe.Vert@ensmp.fr

Center for Computational Biology
Ecole des Mines de Paris

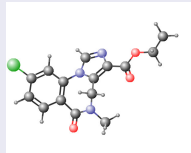
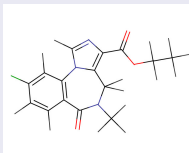
Chimiométrie 2006, Paris, France, December 1st, 2006

Ligand-Based Virtual Screening

Objective

Build models to **predict biochemical properties** of small molecules **from their structures**.

Structures



Properties

- binding to a therapeutic target,
- pharmacokinetics (ADME),
- toxicity...

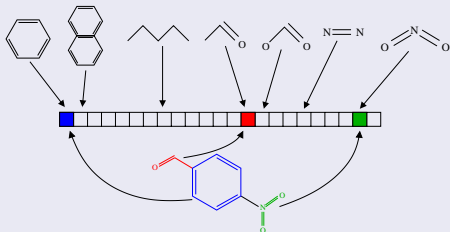
Issues and solution

Two important steps

- 1 Map each molecule to a **vector of fixed dimension**.
- 2 Apply an algorithm for **regression or pattern recognition**.

Example: 2D structural keys

A vector indexed by a limited set of **informative** structures



+ NN, PLS, decision tree, ...

Classical approaches

Difficulties

- **Expressivity** of the features (which features are relevant?)
- **Dimension** of the vector (memory storage, speed, statistical issues)

Our approach

Work **implicitly** in large (potentially infinite!) dimensions:

- Allows to consider a large number of **potentially important features**.
- **No need to store explicitly the vectors** (no problem of memory storage or hash clashes)
- Use of **regularized statistical algorithm** to handle the problem of large dimension

Classical approaches

Difficulties

- **Expressivity** of the features (which features are relevant?)
- **Dimension** of the vector (memory storage, speed, statistical issues)

Our approach

Work **implicitly** in large (potentially infinite!) dimensions:

- Allows to consider a large number of **potentially important features**.
- **No need to store explicitly the vectors** (no problem of memory storage or hash clashes)
- Use of **regularized statistical algorithm** to handle the problem of large dimension

- 1 Support Vector Machines and kernels
- 2 2D Kernel
- 3 3D Pharmacophore Kernel
- 4 Conclusion

- 1 Support Vector Machines and kernels
- 2 2D Kernel
- 3 3D Pharmacophore Kernel
- 4 Conclusion

The Machine Learning Paradigm

Objective

Predict a property y for objects x

- $x =$ **molecule**, gene sequence, picture, ...
- y is continuous (**regression**) or discrete (**pattern recognition**)

A two-step approach

- 1 **Training**: observe a set

$$\mathcal{S} = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

of **labeled objects**, and learn a function $f : \mathcal{X} \rightarrow \mathcal{Y}$

- 2 **Test**: Given a new object x , predict its label by $f(x)$.

The Machine Learning Paradigm

Objective

Predict a property y for objects x

- x = **molecule**, gene sequence, picture, ...
- y is continuous (**regression**) or discrete (**pattern recognition**)

A two-step approach

- 1 **Training**: observe a set

$$\mathcal{S} = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

of **labeled objects**, and learn a function $f : \mathcal{X} \rightarrow \mathcal{Y}$

- 2 **Test**: Given a new object x , predict its label by $f(x)$.

In biomedical research..

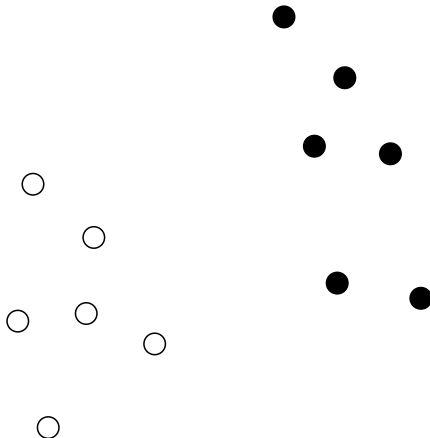
- **Virtual screening** : x is the description of a molecule, y is the activity / toxicity / drugability ...
- **Medical diagnosis and prognosis**: x is a set of features (age, weight, transcriptome...), y is the risk / type of tumor / expected evolution of disease.
- **Functional genomics** : x is a set of gene features (sequence, expression...), y is the function of the gene
- ...

What is a SVM?

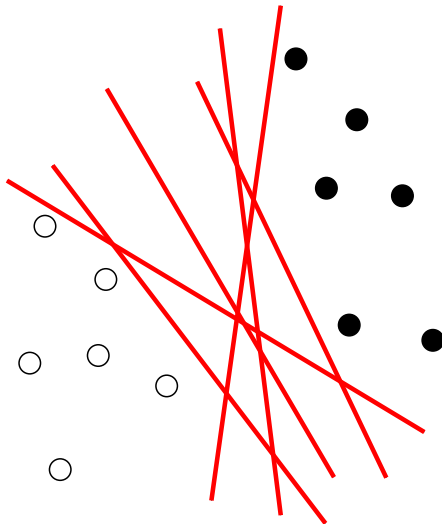
Main features

- an algorithm for **pattern recognition** and **regression**
- robust in **high dimension** (e.g., images, texts, microarrays, fingerprints)
- handles vectorial or **structured data** (e.g., sequences, graphs)
- allows easy integration of **heterogeneous data** (e.g., gene sequence and expression, docking score and molecule structure...)
- **state-of-the-art performance** on many real-world applications.

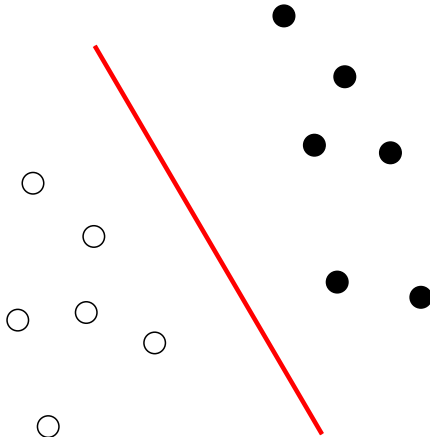
Simplest SVM



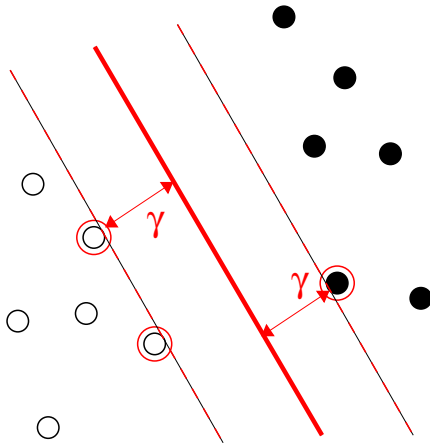
Simplest SVM



Simplest SVM



Simplest SVM



Linear SVM: implementation

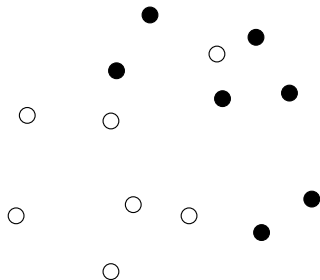
- After some algebra it is obtained by solving in $\alpha \in \mathbb{R}^n$ the following **quadratic program**:

$$\begin{aligned} \text{minimize} \quad & \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j x_i^\top x_j - \sum_{i=1}^n \alpha_i \\ \text{subject to} \quad & \alpha_i \geq 0, \quad i = 1, \dots, n, \\ & \sum_{i=1}^n \alpha_i y_i = 0. \end{aligned}$$

- Once α is found, the **classification function** is the sign of :

$$f(x) = \sum_{i=1}^n \alpha_i x_i^\top x + b.$$

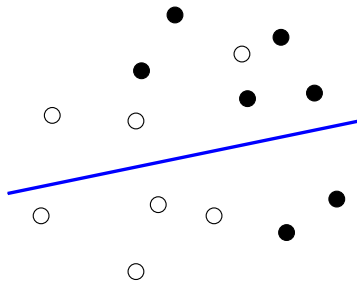
Linear SVM: non-separable case



Implementation

- Solution: find a **trade-off** between large margin and few misclassification
- Simple and elegant mathematical translation: replace $0 \leq \alpha_i$ by $0 \leq \alpha_i \leq C$, for some constant $C > 0$, in the optimization problem.

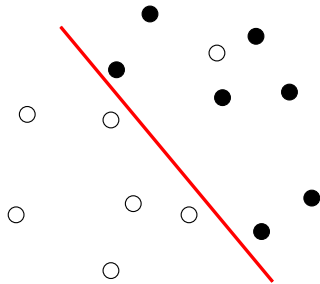
Linear SVM: non-separable case



Implementation

- Solution: find a **trade-off** between large margin and few misclassification
- Simple and elegant mathematical translation: replace $0 \leq \alpha_i$ by $0 \leq \alpha_i \leq C$, for some constant $C > 0$, in the optimization problem.

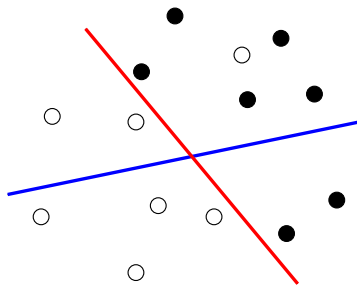
Linear SVM: non-separable case



Implementation

- Solution: find a **trade-off** between large margin and few misclassification
- Simple and elegant mathematical translation: replace $0 \leq \alpha_i$ by $0 \leq \alpha_i \leq C$, for some constant $C > 0$, in the optimization problem.

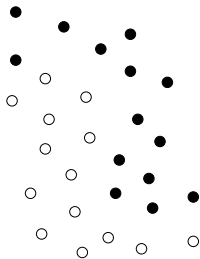
Linear SVM: non-separable case



Implementation

- Solution: find a **trade-off** between large margin and few misclassification
- Simple and elegant mathematical translation: replace $0 \leq \alpha_j$ by $0 \leq \alpha_j \leq C$, for some constant $C > 0$, in the optimization problem.

Nonlinear SVM



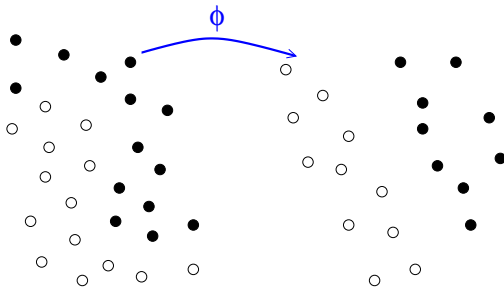
The idea

- Define a (nonlinear) mapping

$$\phi : \mathcal{X} \rightarrow \mathcal{F} \subset \mathbb{R}^p .$$

- Run a linear SVM in the feature space.

Nonlinear SVM



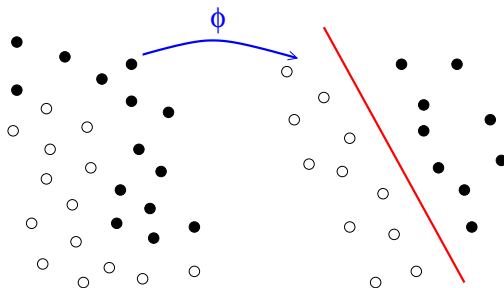
The idea

- Define a (nonlinear) mapping

$$\phi : \mathcal{X} \rightarrow \mathcal{F} \subset \mathbb{R}^p .$$

- Run a linear SVM in the feature space.

Nonlinear SVM



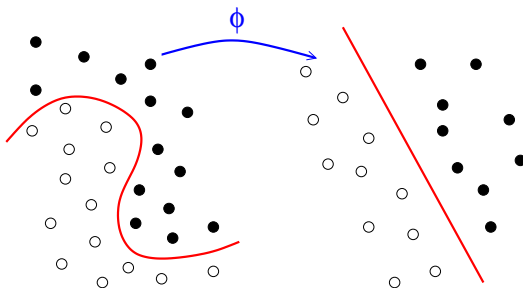
The idea

- Define a (nonlinear) mapping

$$\phi : \mathcal{X} \rightarrow \mathcal{F} \subset \mathbb{R}^p .$$

- Run a linear SVM in the feature space.

Nonlinear SVM



The idea

- Define a (nonlinear) mapping

$$\phi : \mathcal{X} \rightarrow \mathcal{F} \subset \mathbb{R}^p .$$

- Run a linear SVM in the feature space.

Nonlinear SVM: implementation

- Solve in $\alpha \in \mathbb{R}^n$:

$$\text{minimize} \quad \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \Phi(x_i)^\top \Phi(x_j) - \sum_{i=1}^n \alpha_i$$

$$\text{subject to} \quad 0 \leq \alpha_i \leq C, \quad i = 1, \dots, n,$$

$$\sum_{i=1}^n \alpha_i y_i = 0.$$

- Once α is found, the classification function is the sign of :

$$f(x) = \sum_{i=1}^n \alpha_i \Phi(x_i)^\top \Phi(x) + b$$

Important idea!

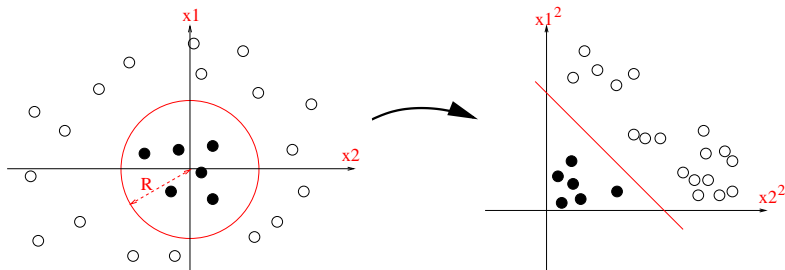
- To any mapping $\Phi : \mathcal{X} \rightarrow \mathcal{F}$ corresponds a **kernel function** K :

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

- SVM only need K , rather than Φ :

$$\begin{aligned} \text{minimize} \quad & \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(x_i, x_j) - \sum_{i=1}^n \alpha_i \\ \text{subject to} \quad & 0 \leq \alpha_j \leq C, \quad i = 1, \dots, n, \\ & \sum_{i=1}^n \alpha_i y_i = 0. \end{aligned}$$

Example: polynomial kernel



For $x = (x_1, x_2)^T \in \mathbb{R}^2$, let $\Phi(x) = (x_1^2, \sqrt{2}x_1x_2, x_2^2) \in \mathbb{R}^3$:

$$\begin{aligned} K(x, x') &= x_1^2 x_1'^2 + 2x_1x_2x_1'x_2' + x_2^2 x_2'^2 \\ &= (x_1x_1' + x_2x_2')^2 \\ &= (x^T x')^2. \end{aligned}$$

For vectors

- The **linear kernel**

$$K_{lin}(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}' .$$

- The **polynomial kernel**

$$K_{poly}(\mathbf{x}, \mathbf{x}') = \left(\mathbf{x}^\top \mathbf{x}' + a \right)^d .$$

- The **Gaussian RBF kernel**:

$$K_{Gaussian}(\mathbf{x}, \mathbf{x}') = \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2} \right) .$$

Main features

- 1 There exist **conditions** to ensure that a function $K(x, x')$ is a valid kernel (symmetry, positive definiteness).
- 2 No need to compute the corresponding Φ .
- 3 A kernel K can be thought of as a **measure of similarity** (inner products) between the data points.
- 4 The kernel trick allows to work **implicitly** in a (possibly large-dimensional) feature space, in particular:
 - to obtain **non-linear** versions of linear methods (nonlinear kernels)
 - to extend these methods to **non-vector data** (kernels for general objects)
- 5 SVM are designed **not to overfit** the training data even in infinite dimension.
- 6 **Kernel engineering** for complex objects is a hot topic!

Kernel and kernel methods summary

Performance

- **State-of-the-art** in many real-world applications
- Resistant to **large** dimensions

Data representation

- Data do not need to be explicitly **vectors**
- A **similarity function** $K(x, x')$ between data is enough
- K must be symmetric and positive definite

Kernels in chemoinformatics

- We need **kernels for molecules!**
- **Inner products** of classical vector / fingerprint representations will work, but **we can do better.**

Kernel and kernel methods summary

Performance

- **State-of-the-art** in many real-world applications
- Resistant to **large** dimensions

Data representation

- Data do not need to be explicitly **vectors**
- A **similarity function** $K(x, x')$ between data is enough
- K must be symmetric and positive definite

Kernels in chemoinformatics

- We need **kernels for molecules!**
- **Inner products** of classical vector / fingerprint representations will work, but **we can do better.**

Kernel and kernel methods summary

Performance

- **State-of-the-art** in many real-world applications
- Resistant to **large** dimensions

Data representation

- Data do not need to be explicitly **vectors**
- A **similarity function** $K(x, x')$ between data is enough
- K must be symmetric and positive definite

Kernels in chemoinformatics

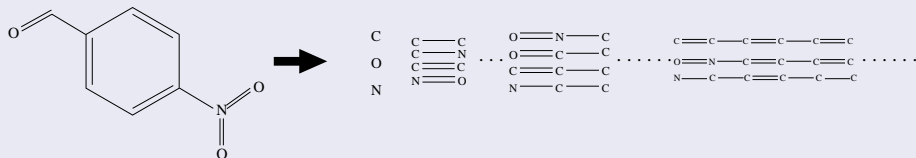
- We need **kernels for molecules!**
- **Inner products** of classical vector / fingerprint representations will work, but **we can do better.**

- 1 Support Vector Machines and kernels
- 2 2D Kernel**
- 3 3D Pharmacophore Kernel
- 4 Conclusion

Motivation: 2D Fingerprints

Features

A vector indexed by a large set of **molecular fragments**



Pros

- Many features
- Easy to detect

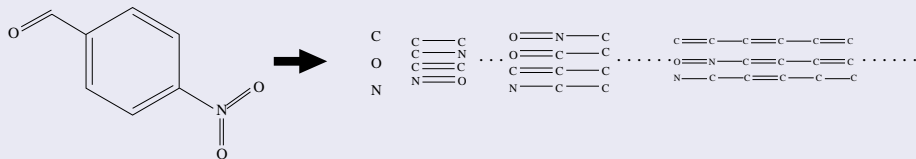
Cons

- Too many features?
- Hashing \implies clashes

Motivation: 2D Fingerprints

Features

A vector indexed by a large set of **molecular fragments**



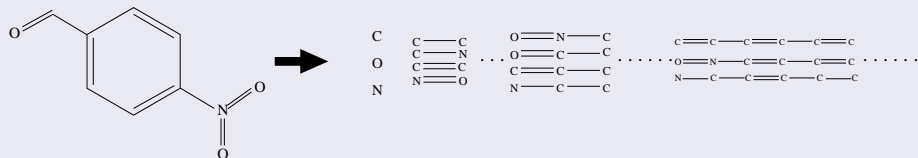
Pros

- Many features
- Easy to detect

Cons

- Too many features?
- Hashing \implies clashes

SVM approach



Let $\Phi(x)$ the vector of fragment counts:

- Long fragments lead to large dimensions :
SVM can learn in high dimension
- $\Phi(x)$ is too long to be stored, and hashes induce clashes:
SVM do not need $\Phi(x)$, they just need the kernel

$$K(x, x') = \phi(x)^T \phi(x') .$$

2D fingerprint kernel

Definition

- For any $d > 0$ let $\phi_d(x)$ be the vector of counts of **all fragments of length d** :

$$\phi_1(x) = (\#(C), \#(O), \#(N), \dots)^T$$

$$\phi_2(x) = (\#(C-C), \#(C=O), \#(C-N), \dots)^T \quad \text{etc...}$$

- The **2D fingerprint kernel** is defined, for $\lambda < 1$, by

$$K_{2D}(x, x') = \sum_{d=1}^{\infty} \lambda^d \phi_d(x)^T \phi_d(x').$$

- This is an **inner product** in the space of **2D fingerprints of infinite length**.

Theorem

The 2D fingerprint kernel between two molecules x and x' can be computed with a **worst-case complexity** $O(|x| \times |x'|^3)$ (much faster in practice).

Remarks

- The complexity is not related to the **length** of the fragments considered (although faster computations are possible if the length is limited).
- Solves the problem of **clashes** and **memory storage**.
- Allows to work with **infinite-length fingerprints** without computing them!

Theorem

The 2D fingerprint kernel between two molecules x and x' can be computed with a **worst-case complexity** $O(|x| \times |x'|^3)$ (much faster in practice).

Remarks

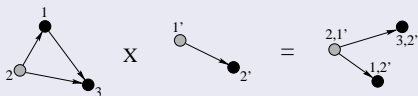
- The complexity is not related to the **length** of the fragments considered (although faster computations are possible if the length is limited).
- Solves the problem of **clashes** and **memory storage**.
- Allows to work with **infinite-length fingerprints** without computing them!

2D kernel computation: Sketch (1/2)

- Let $\mathcal{F}(x)$ be the set of fragments of the molecule x (with repeats). Let $l(f)$ be the label of fragment f (e.g., $c - c$), and $|f|$ its length. Then the kernel can be rewritten:

$$K_{2D}(x, x') = \sum_{f \in \mathcal{F}(x)} \sum_{f' \in \mathcal{F}(x')} \mathbf{1}(l(f) = l(f')) \lambda^{|f|}.$$

- For any two molecules (graphs) G_1 and G_2 , compute the **product graph** $G = G_1 \times G_2$:



- There is a **bijection** between:
 - each **fragments of G** ,
 - each **pair of fragments** in G_1 and G_2 with same label.

2D kernel computation: Sketch (2/2)

- Therefore the kernel can be rewritten:

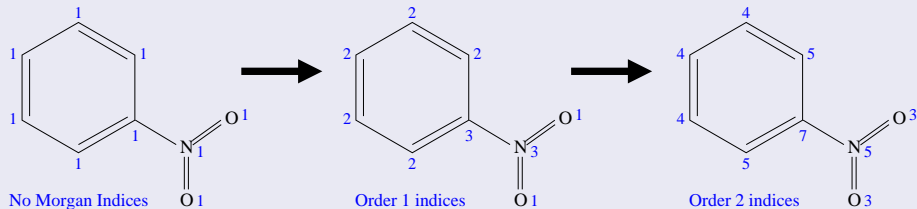
$$K_{2D}(x, x') = \sum_{f \in \mathcal{F}(G)} \lambda^{|f|}.$$

- Let A be the **adjacency matrix** of G . For any $d \geq 1$, $[A^d]_{i,j}$ is the **number of fragments of length d starting in i and ending in j** .
- Therefore the kernel is the sum of the elements of the matrices:

$$\lambda A + \lambda^2 A^2 + \lambda^3 A^3 + \dots = (I - \lambda A)^{-1} - I.$$

Extensions 1: label enrichment

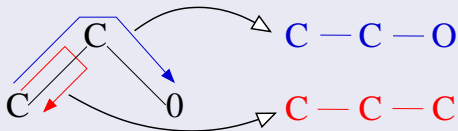
Atom relabeling with the Morgan index



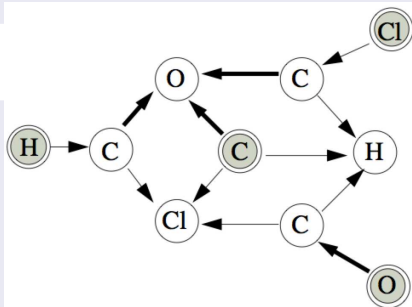
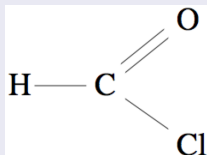
- **Compromise** between **fingerprints** and **structural keys features**.
- Other **relabeling** schemes are possible.
- **Faster computation with more labels** (less matches implies a smaller product graph).

Extensions 2: filter out tottering fragments

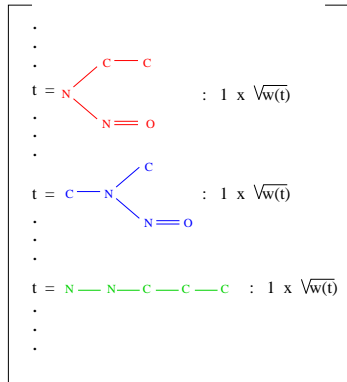
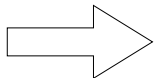
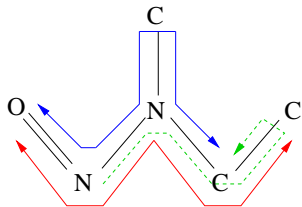
Tottering fragments



Solution: graph transform



Extensions 3: tree-like fragments



MUTAG dataset

- aromatic/hetero-aromatic compounds
- high mutagenic activity /no mutagenic activity
- 188 compounds: 125 + / 63 -

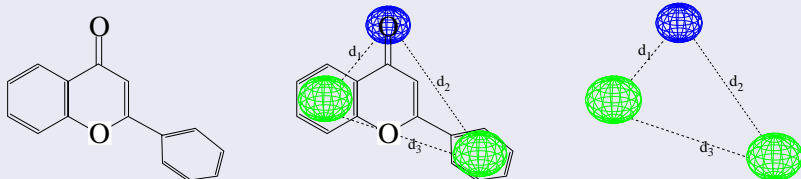
Results

10-fold cross-validation accuracy

Method	Accuracy
Progol1	81.4%
2D kernel	91.2%

- 1 Support Vector Machines and kernels
- 2 2D Kernel
- 3 3D Pharmacophore Kernel**
- 4 Conclusion

3-points pharmacophores



A set of 3 atoms, and 3 inter-atom distances:

$$\mathcal{T} = \{((x_1, x_2, x_3), (d_1, d_2, d_3)), x_i \in \{\text{atom types}\}; d_i \in \mathbb{R}\}$$

3D fingerprint kernel

Pharmacophore fingerprint

- 1 **Discretize** the space of pharmacophores \mathcal{T} (e.g., 6 atoms or groups of atoms, 6-7 distance bins) into a finite set \mathcal{T}_d
- 2 Count the number of occurrences $\phi_t(x)$ of each pharmacophore bin t in a given molecule x , to form a **pharmacophore fingerprint**.

3D kernel

A simple 3D kernel is the **inner product of pharmacophore fingerprints**:

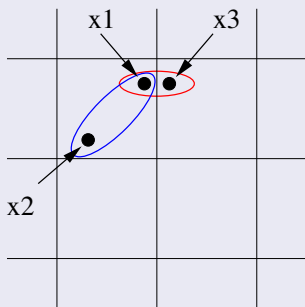
$$K(x, x') = \sum_{t \in \mathcal{T}_d} \phi_t(x) \phi_t(x') .$$

Discretization of the pharmacophore space

Common issues

- 1 If the bins are **too large**, then they are **not specific enough**
- 2 If the bins are **too small**, then they are **too specific**

In all cases, the **arbitrary position of boundaries between bins** affects the comparison:



$$\rightarrow d(x_1, x_3) < d(x_1, x_2)$$

$$\text{BUT } \text{bin}(x_1) = \text{bin}(x_2) \neq \text{bin}(x_3)$$

Kernels between pharmacophores

A small trick

$$\begin{aligned}K(x, y) &= \sum_{t \in \mathcal{T}_d} \phi_t(x) \phi_t(y) \\&= \sum_{t \in \mathcal{T}_d} \left(\sum_{p_x \in \mathcal{P}(x)} \mathbf{1}(\text{bin}(p_x) = t) \right) \left(\sum_{p_y \in \mathcal{P}(y)} \mathbf{1}(\text{bin}(p_y) = t) \right) \\&= \sum_{p_x \in \mathcal{P}(x)} \sum_{p_y \in \mathcal{P}(y)} \mathbf{1}(\text{bin}(p_x) = \text{bin}(p_y))\end{aligned}$$

General pharmacophore kernel

$$K(x, y) = \sum_{p_x \in \mathcal{P}(x)} \sum_{p_y \in \mathcal{P}(y)} K_P(p_x, p_y)$$

New pharmacophore kernels

- Discretizing the pharmacophore space is equivalent to taking the following kernel between individual pharmacophores:

$$K_P(p_1, p_2) = \mathbf{1} (\text{bin}(\mathbf{p}_x) = \text{bin}(\mathbf{p}_y))$$

- For general kernels, there is **no need for discretization!**
- For example, if $d(p_1, p_2)$ is a Euclidean distance between pharmacophores, take:

$$K_P(p_1, p_2) = \exp(-\gamma d(p_1, p_2)) .$$

4 public datasets

- BZR: ligands for the benzodiazepine receptor
- COX: cyclooxygenase-2 inhibitors
- DHFR: dihydrofolate reductase inhibitors
- ER: estrogen receptor ligands

	TRAIN		TEST	
	Pos	Neg	Pos	Neg
BZR	94	87	63	62
COX	87	91	61	64
DHFR	84	149	42	118
ER	110	156	70	110

Results (accuracy)

Kernel	BZR	COX	DHFR	ER
2D (Tanimoto)	71.2	63.0	76.9	77.1
3D fingerprint	75.4	67.0	76.9	78.6
3D not discretized	76.4	69.8	81.9	79.8

1 Support Vector Machines and kernels

2 2D Kernel

3 3D Pharmacophore Kernel

4 Conclusion

Summary

- SVM is a **powerful and flexible machine learning algorithm**. The kernel trick allows the manipulation of **non-vectorial objects** at the cost of defining a kernel function.
- The 2D kernel for molecule extends classical fingerprint-based approaches. It solves the problem of **bit clashes**, allows **infinite fingerprints** and **various extensions**.
- The 3D kernel for molecule extends classical pharmacophore fingerprint-based approaches. It solves the problems of **bit clashes** and of **discretization**.
- Both kernels improve upon their classical counterparts, and provide **competitive results** on benchmark datasets.

Acknowledgements

- Pierre Mahé (CBIO)
- Tatsuya Akutsu, Nobuhisa Ueda, Jean-Luc Perret (Kyoto University)
- Liva Ralaivola (U Marseille)

- Kashima, H., Tsuda, K., and Inokuchi, A. *Marginalized kernels between labeled graphs*. Proceedings of the 20th ICML, 2003, pp. 321-328.
- P. Mahé, N. Ueda, T. Akutsu, J.-L. Perret, and J.-P. Vert. *Graph kernels for molecular structure-activity relationship analysis with SVM*. J. Chem. Inf. Model., 45(4):939-951, 2005.
- P. Mahé, L. Ralaivola, V. Stoven, and J-P Vert. *The pharmacophore kernel for virtual screening with SVM*. J. Chem. Inf. Model., 46(5):2003-2014, 2006.
- P. Mahé and J.-P. Vert. *Graph kernels based on tree patterns for molecules*. Technical report HAL:ccsd-00095488, 2006.
- P. Mahé. *Kernel design for virtual screening of small molecules with support vector machines*. PhD thesis, Ecole des Mines de Paris, 2006.
- **Open-source kernels for chemoinformatics:**
<http://chemcpp.sourceforge.net/>