

Kernel Methods for QSAR and Virtual Screening

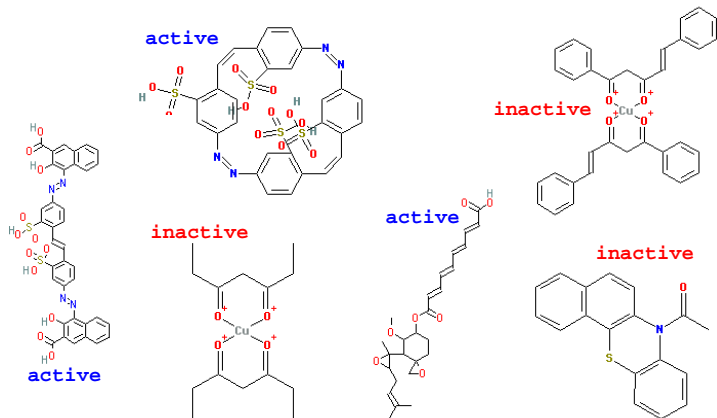
Jean-Philippe Vert

Jean-Philippe.Vert@ensmp.fr

Center for Computational Biology
Ecole des Mines de Paris, ParisTech

10th European Symposium on Statistical Methods for the Food
Industry, Louvain-la-Neuve, Belgium, January 23rd, 2008

Ligand-Based Virtual Screening and QSAR



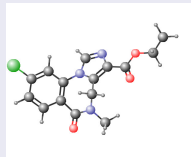
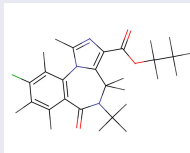
NCI AIDS screen results (from <http://cactus.nci.nih.gov>).

More formally...

Objective

Build models to **predict biochemical properties Y** of small molecules **from their structures X** , using a training set of (X, Y) pairs.

Structures X



Properties Y

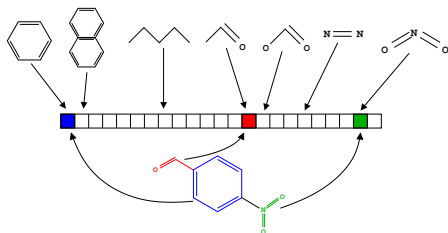
- binding to a therapeutic target,
- pharmacokinetics (ADME),
- toxicity...

Classical approaches

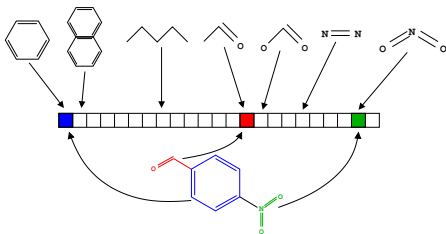
Two steps

- 1 Map each molecule to a **vector of fixed dimension** using **molecular descriptors**
 - Global properties of the molecules (mass, logP...)
 - 2D and 3D descriptors (substructures, fragments,)
- 2 Apply an algorithm for **regression or pattern recognition**.
 - PLS, ANN, ...

Example: 2D structural keys



Which descriptors?



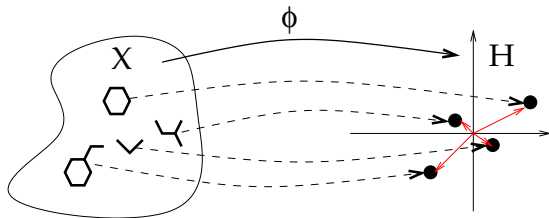
Difficulties

- **Many** descriptors are **needed** to characterize various features (in particular for 2D and 3D descriptors)
- But **too many** descriptors are **harmful** for memory storage, computation speed, statistical estimation

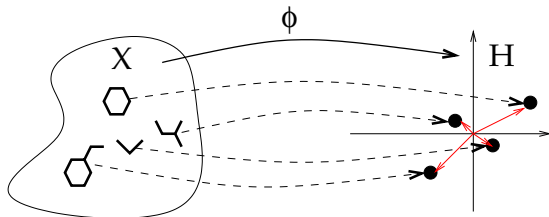
Definition

- Let $\Phi(x) = (\Phi_1(x), \dots, \Phi_p(x))$ be a vector representation of the molecule x
- The **kernel** between two molecules is defined by:

$$K(x, x') = \Phi(x)^\top \Phi(x') = \sum_{i=1}^p \Phi_i(x) \Phi_i(x').$$



The kernel trick



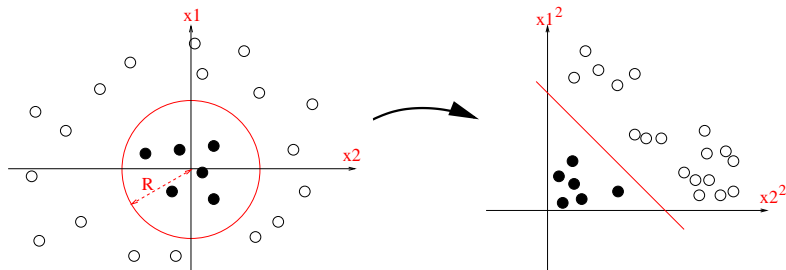
The trick

- 1 Computing the kernel $K(x, x')$ is often **more efficient** than computing $\Phi(x)$, especially in high or infinite dimensions! Ex:

$$K(x, x') = \exp\left(-\gamma\|x - x'\|^2\right).$$

- 2 Many linear algorithms for regression or pattern recognition can be **expressed only in terms of kernels**.

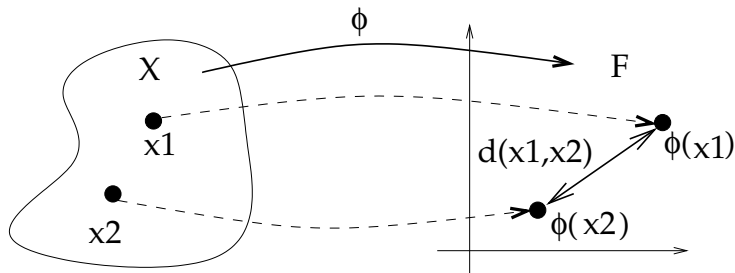
Kernel example: polynomial kernel



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

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

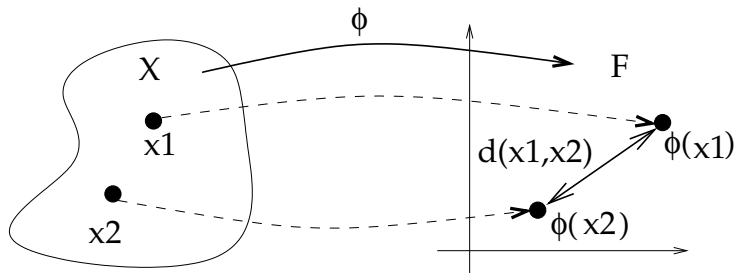
Kernel trick example: computing distances in the feature space



$$\begin{aligned}d_K(\mathbf{x}_1, \mathbf{x}_2)^2 &= \|\Phi(\mathbf{x}_1) - \Phi(\mathbf{x}_2)\|_{\mathcal{H}}^2 \\ &= \langle \Phi(\mathbf{x}_1) - \Phi(\mathbf{x}_2), \Phi(\mathbf{x}_1) - \Phi(\mathbf{x}_2) \rangle_{\mathcal{H}} \\ &= \langle \Phi(\mathbf{x}_1), \Phi(\mathbf{x}_1) \rangle_{\mathcal{H}} + \langle \Phi(\mathbf{x}_2), \Phi(\mathbf{x}_2) \rangle_{\mathcal{H}} - 2 \langle \Phi(\mathbf{x}_1), \Phi(\mathbf{x}_2) \rangle_{\mathcal{H}}\end{aligned}$$

$$d_K(\mathbf{x}_1, \mathbf{x}_2)^2 = K(\mathbf{x}_1, \mathbf{x}_1) + K(\mathbf{x}_2, \mathbf{x}_2) - 2K(\mathbf{x}_1, \mathbf{x}_2)$$

Kernel trick example: computing distances in the feature space



$$\begin{aligned}d_K(\mathbf{x}_1, \mathbf{x}_2)^2 &= \|\Phi(\mathbf{x}_1) - \Phi(\mathbf{x}_2)\|_{\mathcal{H}}^2 \\ &= \langle \Phi(\mathbf{x}_1) - \Phi(\mathbf{x}_2), \Phi(\mathbf{x}_1) - \Phi(\mathbf{x}_2) \rangle_{\mathcal{H}} \\ &= \langle \Phi(\mathbf{x}_1), \Phi(\mathbf{x}_1) \rangle_{\mathcal{H}} + \langle \Phi(\mathbf{x}_2), \Phi(\mathbf{x}_2) \rangle_{\mathcal{H}} - 2 \langle \Phi(\mathbf{x}_1), \Phi(\mathbf{x}_2) \rangle_{\mathcal{H}}\end{aligned}$$

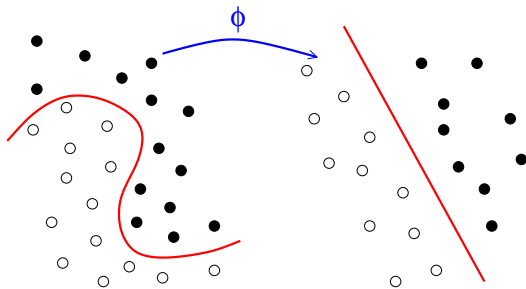
$$d_K(\mathbf{x}_1, \mathbf{x}_2)^2 = K(\mathbf{x}_1, \mathbf{x}_1) + K(\mathbf{x}_2, \mathbf{x}_2) - 2K(\mathbf{x}_1, \mathbf{x}_2)$$

Kernel methods

You don't like nearest-neighbor classification, or your problem is not binary classification, but you would like to benefit from the kernel trick (nonlinearity, structured data etc...)? Try other **kernel methods** that extend your favorite algorithm to handle kernels:

- **Support Vector Machines**,
- kernel PLS,
- kernel PCA,
- kriging,
- kernel perceptron,
- kernel logistic regression,
- and **many more!**

Example: Support Vector Machine



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

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

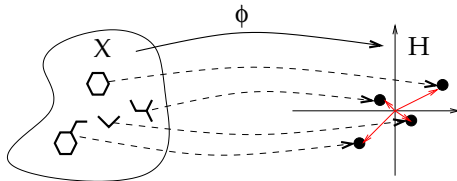
$$\sum_{i=1}^n \alpha_i y_i = 0.$$

Making kernels for molecules

- **Strategy 1:** use **well-known molecular descriptors** to represent molecules m as vectors $\Phi(m)$, and then use kernels for vectors, e.g.:

$$K(m_1, m_2) = \Phi(m_1)^\top \Phi(m_2).$$

- **Strategy 2:** invent **new kernels** to do things you can not do with strategy 1, such as using an infinite number of descriptors. We will now see two examples of this strategy, extending 2D and 3D molecular descriptors.



Summary

The problem

- **Regression** and **pattern recognition** over molecules
- Classical **vector representation** is both statistically and computationally **challenging**

The kernel approach

By defining a **kernel for molecules** we can 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) thanks to the **kernel trick**
- Use of **regularized statistical algorithm** (SVM, kernel PLS, kernel perceptron...) to handle the statistical problem of large dimension

Summary

The problem

- **Regression** and **pattern recognition** over molecules
- Classical **vector representation** is both statistically and computationally **challenging**

The kernel approach

By defining a **kernel for molecules** we can work **implicitly** in large (potentially infinite!) dimensions:

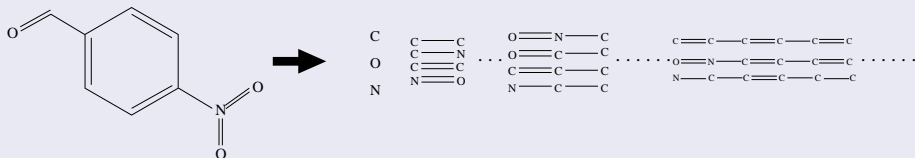
- Allows to consider a **large number** of **potentially important features**.
- **No need to store explicitly the vectors** (no problem of memory storage or hash clashes) thanks to the **kernel trick**
- Use of **regularized statistical algorithm** (SVM, kernel PLS, kernel perceptron...) to handle the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

- 1 2D Kernel
- 2 3D Pharmacophore Kernel
- 3 Conclusion

Motivation: 2D Fingerprints

Features

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



Pros

- Many features
- Easy to detect

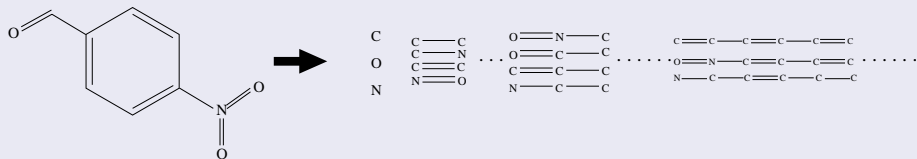
Cons

- Too many features?
- Hashing \implies clashes

Motivation: 2D Fingerprints

Features

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



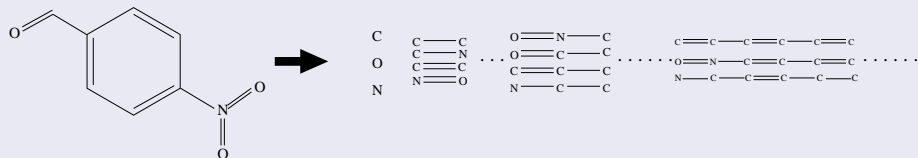
Pros

- 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

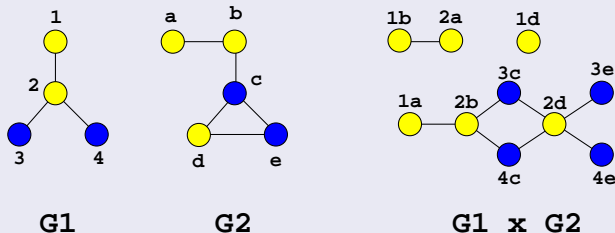
The 2D fingerprint kernel between two molecules x and x' can be computed with a **worst-case complexity** $O(|x| \times |x'|^3)$ (much faster in practice).

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 trick

- Rephrase the kernel computation as that as counting the number of walks on a graph (the product graph)

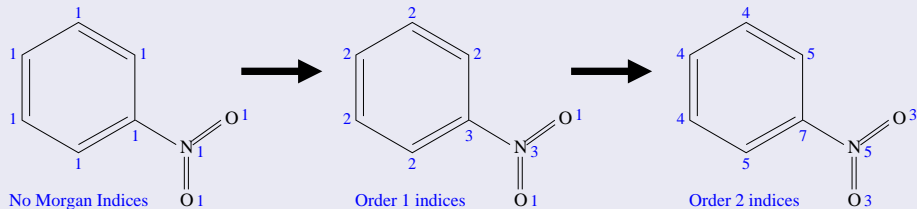


- The infinite counting can be factorized

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

Extensions 1: label enrichment

Atom relabeling with the Morgan index

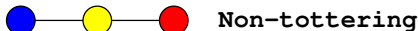
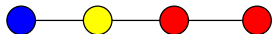


- **Compromise** between **fingerprints** and **structural keys features**.
- Other **relabeling** schemes are possible.
- **Faster computation with more labels** (less matches implies a smaller product graph).

Extension 2: Non-tottering walk kernel

Tottering walks

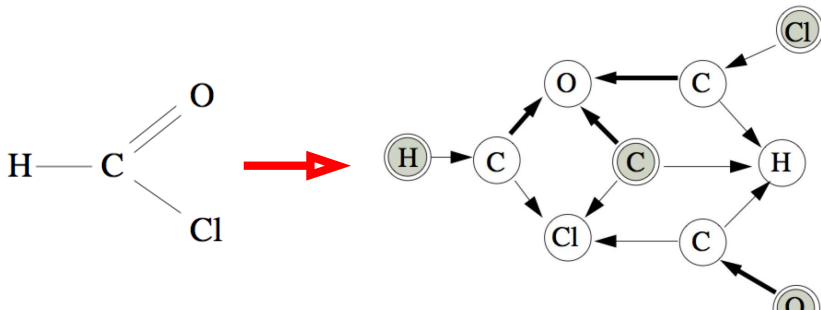
A **tottering walk** is a walk $w = v_1 \dots v_n$ with $v_i = v_{i+2}$ for some i .



- 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).



MUTAG dataset

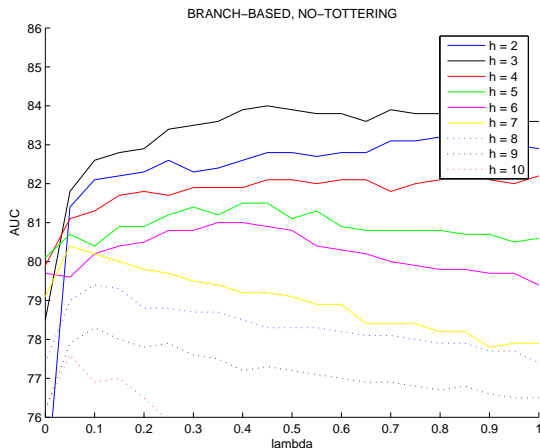
- aromatic/hetero-aromatic compounds
- high mutagenic activity /no mutagenic activity, assayed in *Salmonella typhimurium*.
- 188 compounds: 125 + / 63 -

Results

10-fold cross-validation accuracy

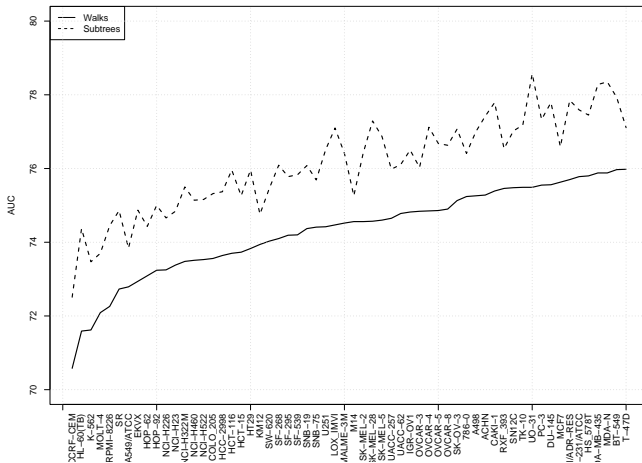
Method	Accuracy
Progol1	81.4%
2D kernel	91.2%

Subtree kernels



AUC as a function of the branching factors for different tree depths (from Mahé et al., 2007).

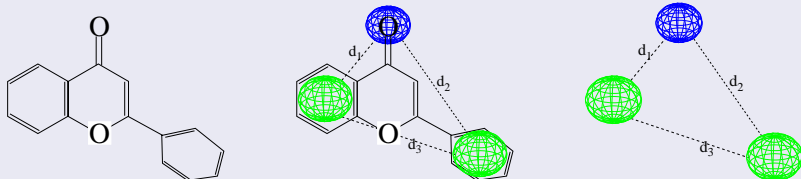
2D Subtree vs walk kernels



Screening of inhibitors for 60 cancer cell lines.

- 1 2D Kernel
- 2 3D Pharmacophore Kernel**
- 3 Conclusion

3-points pharmacophores



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

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

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**:

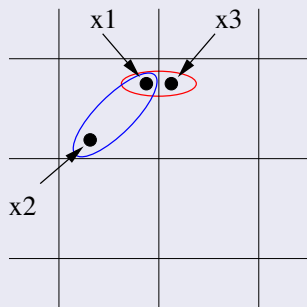
$$K(x, x') = \sum_{t \in \mathcal{T}_d} \phi_t(x) \phi_t(x') .$$

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 large**, then they are **too specific**

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



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

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

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 2D Kernel
- 2 3D Pharmacophore Kernel
- 3 Conclusion**

- 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 (Xerox, formerly at ParisTech)
- 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/>