



## Graph kernels in chemoinformatics

Benoit Gaüzère, Luc Brun, Didier Villemin

### ► To cite this version:

Benoit Gaüzère, Luc Brun, Didier Villemin. Graph kernels in chemoinformatics. Matthias Dehmer and Frank Emmert-Streib. Quantitative Graph Theory Mathematical Foundations and Applications, CRC Press, pp.425-470, 2015, 978-1-4665-8451-8. hal-01201933

**HAL Id: hal-01201933**

**<https://hal.science/hal-01201933>**

Submitted on 18 Sep 2015

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

## *Dedication*

This is the dedication.





## Part I

This is a Part



# Chapter 1

## Graph Kernels in chemoinformatics

**Benoit Gaüzère**

*GREYC CNRS UMR 6072, Université de Caen, France*

**Luc Brun**

*GREYC CNRS UMR 6072, Université de Caen, France*

**Didier Villemin**

*LCMT CNRS UMR 6507, ENSICAEN, France*

1.1	Introduction .....	4
1.1.1	General definitions .....	5
1.2	Graph similarity measures .....	8
1.2.1	Maximum common sub graphs .....	9
1.2.2	Frequent graphs .....	10
1.2.3	Graph edit distance .....	11
1.2.3.1	Graph edit distance approximation .....	11
1.2.3.2	Graph edit distance and machine learning methods ....	12
1.3	Graph embedding methods .....	14
1.3.1	Empirical kernel map .....	14
1.3.2	Isometric embedding .....	15
1.3.3	Local approaches .....	16
1.3.4	Global approaches .....	16
1.3.5	Conclusion .....	17
1.4	Kernel theory .....	17
1.4.1	Definitions .....	17
1.4.2	Kernel trick .....	19
1.4.3	Kernels and similarity measures .....	20
1.4.4	Kernel methods .....	20
1.4.4.1	Support vector machines .....	20
1.4.4.2	Support vectors machines for regression .....	22
1.4.4.3	Kernel ridge regression .....	23
1.4.5	Graph kernels and chemoinformatics .....	24
1.5	Graph kernels based on bags of patterns .....	24
1.5.1	Complete graph kernel .....	24
1.5.2	Linear pattern kernels .....	25
1.5.2.1	Random walks .....	25
1.5.2.2	Tottering .....	25
1.5.2.3	Paths .....	26
1.5.3	Non linear pattern kernels .....	27
1.5.3.1	Graphlets .....	27
1.5.3.2	Tree patterns .....	29
1.5.3.3	Weisfeiler-Lehman .....	29
1.5.3.4	Graph fragments .....	30
1.5.3.5	Treelets .....	31

1.5.4	3D pattern kernel .....	32
1.5.5	Cyclic pattern kernels .....	33
1.5.5.1	Cyclic patterns .....	33
1.5.5.2	Relevant cycle graph .....	34
1.5.5.3	Relevant cycle hypergraph .....	36
1.5.6	Pattern Weighting .....	37
1.5.6.1	A priori methods .....	38
1.5.6.2	Branching cardinality .....	38
1.5.6.3	Ratio depth/size .....	38
1.5.6.4	Multiple kernel learning .....	39
1.5.6.5	Simple MKL .....	39
1.5.6.6	Generalized MKL .....	40
1.5.6.7	Infinite MKL .....	40
1.6	Experiments .....	40
1.6.1	Boiling point prediction .....	41
1.6.2	Classification .....	42
1.6.2.1	AIDS dataset .....	42
1.6.2.2	PTC dataset .....	43
1.7	Conclusion .....	44

## 1.1 Introduction

Graphs provide a generic data structure widely used in chemo and bioinformatics to represent complex structures such as chemical compounds or complex interactions between proteins. However, the high flexibility of this data structure does not allow to readily combine it with usual machine learning algorithms based on a vectorial representation of input data.

Indeed, algorithms restricted to the graph domain are essentially restricted to k-nearest neighbors or k-medians algorithms. These algorithms use different measures of similarity based on the set of frequent sub graphs extracted from graphs or from the size of the maximum common sub graph of two graphs. A widely used graph similarity measure is based on a measure of the distortion required to transform one graph into another. This measure of dissimilarity, called graph edit distance, is NP-hard to compute. However, this complexity can be reduced by computing a sub optimal, but usually effective, graph edit distance. Nevertheless, graph edit distance does not usually fit all the requirements of an euclidean distance and hence can not be readily applied in conjunction with many machine learning methods.

Graph embedding methods aims to tackle this limitation by embedding graphs into explicit vectorial representations which allow the use of any machine learning algorithm defined on vectorial representations. However, encoding graphs as explicit vectors having a limited size induces a loss of information which may reduce prediction accuracy.

Graph kernels are defined as similarity measures between graphs. Under mild conditions, graph kernels correspond to scalar products between possibly implicit graph embeddings into an Hilbert space. Thanks to this graph

embedding, machine learning methods which may be rewritten so as to use only scalar products between input data, such as SVM, can be applied on graphs. Graph kernels thus provide a natural connection between graph space and machine learning.

A large family of kernel methods is based on a decomposition of graphs into bags of patterns. The prototypical example of this last family, is the complete graph kernel which is based on a decomposition of a graph into all its sub structures (sub graphs). Such a kernel being NP-Hard to compute, kernels based on bags of patterns must be restricted to the extraction of different family of sub structures. Given a particular type of bag of patterns, each sub structure of a bag encodes a different structural information. This information may be weighted according to the property to predict. Such a weighting scheme allows to highlight relevant patterns having an high influence on a particular property.

### 1.1.1 General definitions

This first section aims to introduce some basics of graph theory required to define graph similarity measures used in chemoinformatics.

**Definition 1 (Graph)** *An unlabeled graph is a pair  $G = (V, E)$  such that  $V$  corresponds to a set of nodes and  $E \subset V \times V$  corresponds to a set of edges connecting nodes. The size of the graph is defined by  $|V|$ . If  $(u, v) \in E$ ,  $u$  is said to be adjacent to  $v$ .*

**Definition 2 (Neighbourhood)** *Neighbourhood relationship is encoded by the function  $\Gamma : V \rightarrow \mathcal{P}(V)$  with  $\Gamma(v) = \{u \in V \mid (u, v) \in E\}$  where  $\mathcal{P}(V)$  denotes all subsets of  $V$ .*

**Definition 3 (Degree)** *The degree of a node  $v \in V$  is defined as  $|\Gamma(v)|$ .*

**Definition 4 (Non oriented graph)** *A graph is non oriented if for any pair of nodes  $(u, v) \in E$ ,  $(v, u) \in E$ . In this case  $(u, v)$  denotes indifferently the oriented edges  $(u, v)$  and  $(v, u)$ .*

In the following we will only use non oriented graphs. Hence unless otherwise stated, graphs are supposed to be non oriented.

**Definition 5 (Sub graph)** *A graph  $G' = (V', E')$  is a sub graph of  $G = (V, E)$ , denoted  $G' \subseteq G$ , if  $V' \subseteq V$  and  $E' \subseteq E$ . If  $E' = E \cap (V' \times V')$ ,  $G'$  is called a vertex induced sub graph of  $G$ .*

**Definition 6 (Walks, Trails, Paths)** *A walk of a graph  $G = (V, E)$  is a sequence of nodes  $W = (v_1 \dots v_n)$  connected by edges:  $(v_i, v_{i+1}) \in E$  for any  $i \in \{1, \dots, n-1\}$ . If each edge  $(v_i, v_{i+1})$  appears only once in  $W$ ,  $W$  is called a trail. If each vertex (and thus each edge) appears only once,  $W$  is called a path. The length of a walk is defined by its number of nodes.*



**Definition 7 (Distance between nodes)** The distance  $d_G(u, v)$  between two nodes  $u$  and  $v$  of a graph  $G = (V, E)$  is defined as the length of the shortest path between  $u$  and  $v$  in  $G$ .

**Definition 8 (Cycle)** A cycle is a path whose first node is equal to the last one. This node is the only one appearing twice in the sequence.

**Definition 9 (Connected graph)** Given a graph  $G = (V, E)$ , a set  $U \subset V$  is said to be connected if it exists a path between any pair of distinct nodes in  $U$ . The set  $U$  is called a connected component of  $G$  if it is not included into a larger connected set. The graph  $G$  is said to be connected if all its connected components are equal to  $V$ .

**Definition 10 (Tree)** A tree is a connected graph without cycles.

**Definition 11 (Bridge)** A bridge is an edge whose removal increases the number of connected components of the graph. The set of bridges of a graph  $G$  is denoted  $\mathcal{B}(G)$ .

**Definition 12 (Labeled graph)** A non oriented labeled graph  $G = (V, E, \mu, \nu)$  is a non oriented unlabeled graph  $G = (V, E)$  associated to a node labeling function  $\mu : V \rightarrow L_V$  and an edge labeling function  $\nu : E \rightarrow L_E$ , where  $L_V$  and  $L_E$  denote respectively sets of node and edge labels.

Two graphs are considered as equal if it exists a bijection between the nodes of both graphs which respect adjacency relationships. In this case both graphs are said to be isomorphic.

**Definition 13 (Graph isomorphism)** Two graphs  $G = (V, E)$  and  $G' = (V', E')$  are structurally isomorphic, denoted  $G \simeq_s G'$ , if and only if it exists a bijection

$$f : V \rightarrow V'$$

such that:

$$(u, v) \in E \Leftrightarrow (f(u), f(v)) \in E'.$$

If  $G$  and  $G'$  correspond to labeled graphs, i.e. if  $G = (V, E, \mu, \nu)$  and  $G' = (V', E', \mu', \nu')$ . A structural isomorphism between  $G$  and  $G'$  is called a graph isomorphism, denoted by  $G \simeq G'$ , if

$$\begin{aligned} \forall v \in V, \quad \mu(v) &= \mu'(f(v)) \quad \text{and} \\ \forall (u, v) \in E, \quad \nu(u, v) &= \nu'(f(u), f(v)). \end{aligned}$$

**Definition 14 (Partial sub graph isomorphism)** Let  $G = (V, E, \mu, \nu)$  and  $G' = (V', E', \mu', \nu')$  denote two graphs such that  $|V| \leq |V'|$ . It exists a partial structural sub graph isomorphism between  $G$  and  $G'$  if and only if it exists an injection:

$$f : V \rightarrow V'$$

such that:

$$\forall (u, v) \in V^2 \ (u, v) \in E \Rightarrow (f(u), f(v)) \in E'.$$

A partial structural sub graph isomorphism is a partial sub graph isomorphism, denoted  $G \subseteq_p G'$ , if:

$$\begin{aligned} \forall v \in V, \quad \mu(v) &= \mu'(f(v)) \quad \text{and} \\ \forall (u, v) \in E, \quad \nu(u, v) &= \nu'(f(u), f(v)). \end{aligned}$$

Note that the partial sub graph isomorphism relationship may also be denoted as homomorphism relationship in other references.

**Definition 15 (Sub graph isomorphism)** Let  $G = (V, E, \mu, \nu)$  and  $G' = (V', E', \mu', \nu')$  denote two graphs such that  $|V| \leq |V'|$ . It exists a structural sub graph isomorphism if it exists a partial structural sub graph isomorphism between  $G$  and  $G'$  and:

$$\forall (u, v) \in V'^2, (u, v) \in E' \Leftrightarrow (u, v) \in E$$

A structural sub graph isomorphism is a sub graph isomorphism denoted  $G \subseteq G'$ , if:

$$\begin{aligned} \forall v \in V, \quad \mu(v) &= \mu'(f(v)) \quad \text{and} \\ \forall (u, v) \in E, \quad \nu(u, v) &= \nu'(f(u), f(v)). \end{aligned}$$

Hypergraphs correspond to an extension of graphs which allows to define an adjacency relationship between more than two nodes. Hypergraphs have been introduced by Claude Berge [3].

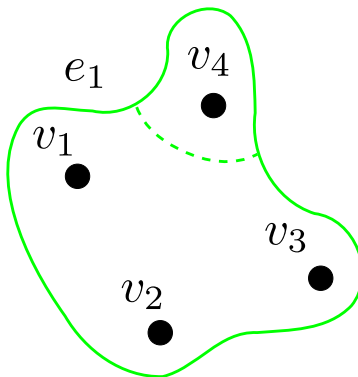
**Definition 16 (Hypergraph)** An hypergraph  $H = (V, E)$  is a pair of sets  $V$ , encoding hypergraph's nodes, and  $E = (e_i)_{i \in I} \subseteq \mathcal{P}(V)$  encoding hyperedges:

- $\forall i \in \{1, \dots, |E|\}, e_i \neq \emptyset,$
- $\cup_{i \in I} e_i = V$

As for graphs, the size of an hypergraph  $H = (V, E)$  is defined as  $|V|$  and two nodes  $u$  and  $v$  are adjacent if it exists  $e \in E$  such that  $\{u, v\} \subseteq e$ .

Oriented hypergraph [12] is an hypergraph where hyperedges connect two sets of nodes (figure 1.1):

**Definition 17 (Oriented hypergraph)** An oriented hypergraph  $H = (V, \vec{E})$  is a pair of sets  $V$ , encoding hypergraph's nodes, and  $\vec{E} = (e_i)_{i \in I} \subseteq \mathcal{P}(V) \times \mathcal{P}(V)$  encoding oriented hyperedges. An oriented hyperedge  $e = (s_u, s_v) \in \vec{E}$  with  $s_u = \{u_1, \dots, u_i\} \subset \mathcal{P}(V)$  and  $s_v = \{v_1, \dots, v_j\} \subset \mathcal{P}(V)$  encodes an adjacency relationship between the two sets of node  $s_u$  and  $s_v$  (figure 1.1).



**FIGURE 1.1:** Oriented hypergraph with oriented hyperedge  $e_1 = (\{v_1, v_2\}, \{v_3\})$ .

In the remaining part of this chapter we consider that an hyperedge  $(s_v, s_u)$  exists in  $\vec{E}$  for each  $(s_u, s_v) \in \vec{E}$  and we thus do not differentiate  $(s_u, s_v)$  from  $(s_v, s_u)$ . Note anyway that an oriented hypergraph remains different from an usual hypergraph since within an oriented hypergraph an edge connects two sets of vertices while within the usual hypergraph framework an edge is composed of a single set.

An usual molecular representation is defined by the molecular graph [15, 49] (figure 1.2). A molecular graph is a non oriented labeled graph encoding adjacency relationships between atoms of a molecule. The set of nodes encodes the set of atoms of a molecule while the set of edges encodes atomic bonds connecting these atoms. Each vertex is labeled by the chemical element of the corresponding atom and each edge is labeled by the type of the atomic bond (simple, double, triple or aromatic) connecting two atoms. From a chemical point of view, an atom may share pairs of valence electrons with a maximum of 8 atoms. Therefore, a molecule is encoded by a molecular graph having a degree bounded by 8. An usual convention consists in not encoding hydrogen atoms within molecular graph since they are implicitly encoded by the valency of other atoms. An other convention is to not explicitly represent carbon atoms within graphical representation of a molecular graph (figure 1.2).

---

## 1.2 Graph similarity measures

Most of existing methods used in chemoinformatics are based on the similarity principle which states that similar compounds have similar properties [24]. If we consider the molecular representation given by the molecular

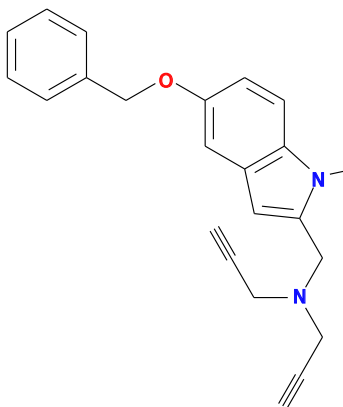


FIGURE 1.2: A molecular graph

graph, comparing molecular compounds thus relies on comparing molecular graphs. Therefore, in order to be able to predict physical or biological properties of molecular compounds, we have to define molecular graph similarity measures.

### 1.2.1 Maximum common sub graphs

A first similarity measure between graphs based on graph theory is defined from the maximum common sub graph. Intuitively, two graphs are considered as similar if they share a large common structure and dissimilar otherwise. In order to define the notion of maximum common sub graph, let us first define the notion of common sub graph:

**Definition 18 (Common sub graph)** *A graph  $G$  is a common sub graph of two graphs  $G_1$  and  $G_2$  if it exists two sub graphs  $\hat{G}_1$  and  $\hat{G}_2$  of  $G_1$  and  $G_2$  such that:*

$$G \simeq \hat{G}_1 \simeq \hat{G}_2$$

Common sub graphs correspond thus to sets of nodes and edges common to two graphs. In order to formally define maximum common sub graphs, let us additionally define the notion of maximal common sub graph:

**Definition 19 (Maximal common sub graph)** *A common sub graph of two graphs  $G_1$  and  $G_2$  is maximal if it is a common sub graph of  $G_1$  and  $G_2$  and if it is not itself a sub graph of another common sub graph of  $G_1$  and  $G_2$ .*

The notion of maximal common sub graph corresponds to set of nodes and edges common to two graphs which can not be enlarged. Using maximal common sub graphs, we can now define maximum common sub graph:

**Definition 20 (Maximum common sub graph)** *A common sub graph  $G$  is maximum if it is a maximal common sub graph of  $G_1$  and  $G_2$  and it does not exist a larger common sub graph of  $G_1$  and  $G_2$ .*

A maximum common sub graph corresponds thus to the largest maximal common sub graph.

The size of a maximum common sub graph of two graphs may be seen as a similarity measure between these two graphs. Methods based on maximum common sub graph consider two graphs as similar if they share a large common structure. Conversely, graphs sharing a lot of small common structures but not a large one are not considered as similar. This particularity may alter the accuracy of methods based on maximum common sub graphs on chemical prediction problems where molecular activity is due to a large number of small structures.

### 1.2.2 Frequent graphs

A widely used approach in chemoinformatics consists in finding sub structures which are responsible of a particular activity. This family of methods is thus more dedicated to activity prediction problems since from a chemical point of view, a molecular compound may have a particular property if it contains a sub structure, called a pharmacophore. This chemical consideration has sustained the emergence of a family of methods based on the discovery of frequent sub graphs [10, 56, 38]. These methods are based on the following assumptions:

- A sub structure is considered as frequent if its number of occurrences is greater than a frequency threshold  $\sigma$  within a set of positive molecular compounds and is insignificant within a set of negative molecular compounds;
- A sub structure of size  $k$  may be frequent if it is composed of at least two frequent sub structures of size  $k-1$ .

Among the set of methods using such an approach, methods described in [10, 56] are based on an iterative algorithm (algorithm 1) which aims to find frequent sub graphs of size  $k$  from the set of frequent sub graphs of size  $k-1$ .

These approaches obtain a good prediction accuracy on many datasets within a reasonable computational time. However, the main drawback of these approaches is the frequent sub graphs hypothesis. Indeed, the hypothesis that one pharmacophore is responsible of an activity is no longer valid when the activity of molecular compounds is induced by a set of different pharmacophores. In this case, each pharmacophore is associated to few molecular

---

**Algorithm 1:** Generic algorithm used by frequent sub graphs methods
 

---

```

1 Find a set of sub graphs  $\mathcal{S}_2$  having a size  $k = 2$  and a frequency greater
  than  $\sigma$ . ;
2 for  $k = 3 \rightarrow k\text{-max}$  do
3   Build a set of candidates  $\mathcal{C}_k$  of size  $k$  from  $\mathcal{S}_{k-1}$ . ;
4    $\mathcal{S}_k \leftarrow$  frequent sub graphs in  $\mathcal{C}_k$ ;
  
```

---

compounds and each pharmacophore does not reach the minimal frequency threshold in order to be considered as a frequent sub structure. Such configurations induce thus a poor prediction model.

### 1.2.3 Graph edit distance

Graph edit distance aims to define a distance between graphs. This distance is based on a measure of distortion induced by a transformation of one graph into another. Graph edit distance between two graphs is defined by the minimal cost associated to an edit path which transforms the first graph into the second one. This distance is defined as a sequence of edit operations, each of these elementary operations being defined as:

- a node/edge insertion: adding one node/edge to the graph;
- a node/edge deletion: removing one node/edge to the graph;
- a node/edge substitution: replace one node/edge label.

Each of these edit operations is associated to a cost  $c(.) \in \mathbb{R}_+$ . Considering these edit costs, graph edit distance is defined as:

$$d_{edit}(G, G') = \min_{(e_1, \dots, e_k) \in \mathcal{C}(G, G')} \sum_{i=1}^k c(e_i) \quad (1.1)$$

where  $\mathcal{C}(G, G')$  encodes all edit paths transforming  $G$  into  $G'$ . Each edit path encodes a sequence of edit operations  $e_1, \dots, e_k$ , each operation  $e_i$  being associated to an edit cost  $c(e_i)$ .

Graph edit distance may be computed using  $A^*$  algorithm [20]. This algorithm consists in building a rooted tree where each path from the root to a leaf encodes a possible edit path. Then, the root to leaf path corresponding to the lowest sum of edit operation costs is defined as the optimal edit path. However, the computational complexity induced by this method is exponential according to the number of graphs' nodes which limits its application to very small graphs.

### 1.2.3.1 Graph edit distance approximation

In order to reduce the complexity required by the computation of graph edit distance, Riesen proposed a polynomial algorithm [44] which reduces graph edit distance complexity. Given two graphs  $G = (V, E)$  of size  $n = |V|$  and  $G' = (V', E')$  of size  $m = |V'|$ , this algorithm is based on a complete bipartite graph  $G_a = (V_a, E_a)$  where  $V_a = V_G \cup V_{G'}$ ,  $V_G = \{V \cup \{\varepsilon_1, \dots, \varepsilon_m\}\}$ ,  $V_{G'} = \{V' \cup \{\varepsilon'_1, \dots, \varepsilon'_n\}\}$  and  $E_a = \{(u, v) \mid u \in V_G \text{ and } v \in V_{G'}\}$ . Bipartite graph  $G_a$  encodes a node to node matching between the two sets  $V_G$  and  $V_{G'}$ . A matching  $u \in V \rightarrow v \in V'$  encodes a substitution whereas a matching  $u \in V \rightarrow \varepsilon'_i$  (resp.  $\varepsilon_i \rightarrow v \in V'$ ) encodes a node deletion (resp. a node insertion). Each edge  $(u, v) \in V_G \times V_{G'}$  is weighted by the cost associated to the edit operation encoded by matching  $u \rightarrow v$ . Therefore, each matching  $\varepsilon_i \rightarrow \varepsilon'_j$  is associated to zero cost since they do not encode any edit operation.

Optimal matching, i.e. matching from  $V$  to  $V'$  having the lowest sum of edge cost, is computed by means of Munkre's algorithm (or Hungarian algorithm) [36]. This optimal matching is then associated to an edit path, itself being associated to a cost. However, the computed cost does not correspond necessarily to graph edit distance since there is no guaranty that the computed edit path is the optimal one. Therefore, on one hand the main drawback of this algorithm is that the computed dissimilarity measure does not correspond to exact graph edit distance. On the other hand, this dissimilarity measure can be computed in polynomial time which allows the use of this approach on chemoinformatics problem.

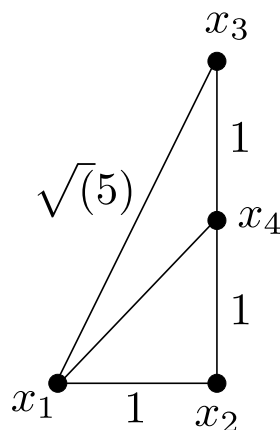
### 1.2.3.2 Graph edit distance and machine learning methods

If we consider that costs associated to elementary edit operations define a distance (proposition 21), graph edit distance fulfills the four properties of a metric and thus defines a distance [37].

**Definition 21 (Distance)** *A function  $d : \mathcal{X}^2 \rightarrow \mathbb{R}_+$  is a distance over  $\mathcal{X}$  if, for any  $(x_i, x_k) \in \mathcal{X}^2$ , it fulfills the four following properties:*

- *non negativity:*  $d(x_i, x_j) \geq 0$ ;
- *self distance:*  $d(x_i, x_j) = 0 \Leftrightarrow x_i \simeq x_j$ ;
- *symmetry:*  $d(x_i, x_j) = d(x_j, x_i)$ ;
- *triangle inequality:* for any  $x_k \in \mathcal{X}$ ,  $d(x_i, x_j) \leq d(x_i, x_k) + d(x_k, x_j)$ .

However, graph edit distance does not fulfill the fifth property of an euclidean distance (Definition 22) [9] and hence does not correspond to a distance in an euclidean space. From a practical point of view, this limitation induces that the metric space defined by the edit distance on the set of graphs can not be isometrically embedded (neither explicitly nor implicitly) into an Hilbert space. This last point restricts drastically the set of machine learning algorithm which may be used in conjunction with the graph edit distance.



**FIGURE 1.3:** Four point polyhedron and point to point distances. Distance between  $x_1$  and  $x_4$  can not be retrieved from the four first properties of a distance (example taken from [9]).

The need of a fifth property to define euclidean distances can be illustrated using a distance matrix completion problem. Let us consider four points  $x_1, x_2, x_3$  and  $x_4$  forming a polyhedron (figure 1.3). Distance matrix  $D \in \mathbb{R}^{4,4}$  encodes distances between the four points of this polyhedron. Let us consider the following incomplete distance matrix  $D$ :

$$D = \begin{bmatrix} 0 & 1 & 5 & d_{14} \\ 1 & 0 & 4 & 1 \\ 5 & 4 & 0 & 1 \\ d_{41} & 1 & 1 & 0 \end{bmatrix} \quad (1.2)$$

The first property of a distance induces that  $d_{14} \geq 0$  and  $d_{41} \geq 0$ . The third property induces that  $d_{14} = d_{41}$  and finally, the fourth property of a distance allows us to bound  $d_{14}$  by  $\sqrt{5} - 1 \leq d_{14} \leq 2$ . However, only one value corresponds to an euclidean embedding since  $x_4$  is at a distance of 1 from  $x_2$  and  $x_3$ . Moreover,  $x_2$  and  $x_3$  are separated by a distance of 2 which induces that  $x_2, x_3$  and  $x_4$  must be collinear. All others distance values fulfill the first four distance properties but they do not correspond to an embedding regardless of the number of dimensions [9]. Therefore, the four properties of a distance are not sufficient to ensure that a distance corresponds to a distance in an euclidean space.

**Definition 22 (Relative angle inequality [9])** For any  $i, j, l \neq k$ ,  $k = \{1, \dots, N\}$ ,  $i < j < l$ . For  $N \geq 4$ , let consider a set of distinct points  $\mathcal{X}_k = \{x_1, \dots, x_N\}$ . Relative angle inequality property is fulfilled if all points  $x_k \in \mathcal{X}_k$  fulfill the following inequalities:



$$\begin{aligned} \cos(\theta_{ikl} + \theta_{lkj}) &\leq \cos(\theta_{ikj}) \leq \cos(\theta_{ikl} - \theta_{lkj}) \\ 0 &\leq \theta_{ikl}, \theta_{lkj}, \theta_{ikj} \leq \pi \end{aligned} \quad (1.3)$$

where  $\theta_{ikj} = \theta_{jki}$  encodes angle between vectors  $\vec{x_k x_i}$  and  $\vec{x_k x_j}$ .

**Proposition 1 (Euclidean distance)** *A function  $d : \mathcal{X}^2 \rightarrow \mathbb{R}_+$  corresponds to a distance in an euclidean space  $\mathbb{R}^N$  if  $d$  is a distance (proposition 21) and  $d$  fulfills the relative angle inequality (definition 22).*

An alternative characterization of an euclidean distance function is provided by the following proposition [9]:

**Proposition 2 (Euclidean distance matrix)** *The matrix  $D$  corresponds to an euclidean distance iff:*

$$D \in \mathbb{S}_h^N \text{ and } \forall c \in \mathbb{R}^N, s.t. \begin{cases} e^t c = 0 \\ \|c\| = 1 \end{cases} \quad c^t D c \leq 0 \quad (1.4)$$

where  $e = (1, \dots, 1)^t$  and  $\mathbb{S}_h^N$  corresponds to the set of  $N \times N$  symmetric matrices having 0 main diagonal.

Since graph edit distance does not define a distance in an euclidean space it can not be directly used in machine learning methods where an implicit or explicit vectorial embedding is mandatory. Despite this, graph edit distance encodes a dissimilarity measure between two molecular graphs which may be used in conjunction with some algorithms such as k-nearest neighbors or k-median.

### 1.3 Graph embedding methods

Methods presented in section 1.1 are defined within graph space. Despite the fact that this space allows to encode complex structures, it suffers from a lack of mathematical properties. Therefore, the use of many well know machine learning methods is not possible within such a space. In order to use conjointly graphs and usual machine learning methods, one possibility consists in defining an explicit vectorial space where each graph is embedded. Using such a scheme, each graph is encoded by a vector which can be used conjointly with any machine learning method defined on vectorial representations. Such vectorial representations of graphs are called graph embeddings.

### 1.3.1 Empirical kernel map

In order to combine graph edit distance and machine learning algorithms, some methods aim to define a graph embedding based on graph edit distance. Riesen has proposed to encode a graph by an explicit vector which encodes the distances between an input graph and a set of graph prototypes [42]. Given a set of prototype graphs denoted by  $\mathcal{G}_{ref} = \{G_1, \dots, G_N\}$ , each graph is embedded into an explicit euclidean space by the following embedding function  $\Phi_{EKM}^{\mathcal{G}} : \mathcal{G} \rightarrow \mathbb{R}^N$ :

$$\Phi_{EKM}^{\mathcal{G}_{ref}}(G) = (d(G, G_1), \dots, d(G, G_N)) \quad (1.5)$$

This method thus encodes each graph by a vector which may be used in any machine learning method defined on euclidean spaces. Although this method uses a widely used graph dissimilarity measure, the choice of prototype graphs strongly impacts the embedding quality. Riesen [42] addressed this last problem through different heuristics providing “good” sets of prototype graphs.

### 1.3.2 Isometric embedding

Another approach aims to define a graph embedding such that this embedding encodes a dissimilarity measure [25]. Let consider a set of graphs  $\mathcal{G} = \{G_1, \dots, G_n\}$  and a function encoding a graph dissimilarity  $d : \mathcal{G} \times \mathcal{G} \rightarrow \mathbb{R}$ . We define a dissimilarity matrix  $D$  by  $D_{i,j} = d(G_i, G_j)^2 \in \mathbb{R}^{n \times n}$ . The method described in [25] consists in computing  $n$  vectors  $x_i$  having  $p$ -dimensions such that the distance between  $x_i$  and  $x_j$  is as close as possible from the dissimilarity measure between  $G_i$  and  $G_j$  encoded by  $D_{i,j}$ .

In order to define vectorial representations  $(x_i)_{i \in \{1, \dots, n\}}$ , a method consists in defining a matrix  $S$  which encodes scalar products between vectorial representations of graphs. Considering a distance matrix  $D$ , matrix  $S$  is defined by  $S = (I - ee^t)D(I - ee^t)$  with  $e = (1, \dots, 1)^t$  and satisfies:  $D_{i,j} = S_{i,i} + S_{j,j} - 2S_{i,j}$ . This last equation corresponds to the well known relationship between euclidean distance and scalar products:  $\|x_i - x_j\|^2 = \langle x_i, x_i \rangle + \langle x_j, x_j \rangle - 2\langle x_i, x_j \rangle$ . If  $S$  is a semi definite positive matrix, then its spectral decomposition is given by  $S = V\Lambda V^t$  where  $\Lambda$  is a diagonal matrix containing the (positive) eigenvalues of  $S$ . Matrix  $S$  can then be written as  $XX^t$  with  $X = V(\Lambda^{\frac{1}{2}})$ . Each entry  $S_{i,j}$  of  $S$  is then defined as a scalar product between two rows of matrix  $X$ . Taking  $x_i$  as an embedding of graph  $G_i$  we indeed obtain:  $D_{i,j} = d^2(G_i, G_j) = \|x_i - x_j\|^2$ .

However, if the dissimilarity matrix  $D$  does not correspond to a distance in an euclidean space (definition 22), matrix  $S$  is not necessarily semi definite positive and we have to regularize it. One possibility is to use the constant shift embedding method which consists in removing negative eigenvalues of  $S$  by subtracting its lowest (negative) eigenvalue. Regularized matrix  $S'$  is then defined by:  $S' = S - \lambda_n(S)I$  where  $\lambda_n(S)$  denotes the lowest negative

eigenvalue of  $S$ . This regularization step alters the initial distances according to the following equation:  $D' = D - 2\lambda_n(S)(ee^t - I)$ . The magnitude of this alteration depends linearly on the lowest eigenvalue of  $S$ . Such a regularization step should be performed only for small values of  $|\lambda_n(S)|$ .

### 1.3.3 Local approaches

Other graph embedding methods which are not based on graph edit distance can be defined using different approaches. A first family consists in encoding each graph by a set of information describing local adjacency relationships and labeling information. Attribute statistics based embedding [19] defines a graph embedding where each coordinate of the vectorial representation encodes either the number of occurrences of a node label, or the number of edges incident to two node labels. For  $|L_V|$  node labels and  $|L_E|$  edge labels, the size of a vector encoding a non oriented graph is thus equal to  $\frac{1}{2}|L_E|(|L_V| + 1)$ . This type of embedding limits the size of encoded paths and thus the amount of encoded information. Indeed, including paths of size 3 induces a huge increase of the size of vector in order to encode all possible labeled paths of size 3.

In order to encode more structural information, topological embedding [48] encodes the number of occurrences of a set of unlabeled structures corresponding to the set of isomorphic graphs having at most  $N$  nodes. Labeling information is included within a matrix representation by means of histograms: For each structure, an histogram encodes the distribution of node and edge labels among all labeled subgraphs corresponding to the unlabeled structure. This method allows to encode larger structures which encodes more structural information than paths of size 2. However, the multiple configurations of labels within all subgraphs isomorphic to a given unlabeled structure is reduced to a simple histogram of node and edge labels hence losing an important part of the label information. This loss of information is induced by the limit on the size of vectors used to encode graphs.

### 1.3.4 Global approaches

Instead of using local approaches, some methods aim to define an embedding which globally encodes a graph rather than using a concatenation of local characteristics. Fuzzy multilevel graph embedding [29] defines a vectorial representation encoding both local characteristics, such as node degrees or node and edge labels, and simple global information such as number of nodes and edges of graphs. Vectorial representation provided by a fuzzy multilevel graph embedding includes thus different levels of analysis of the graph.

Spectral embedding methods [28, 6, 27] base the vectorial representation of a graph on different characteristics of the spectrum of the graph laplacian matrix. The basic idea of this family of methods is that the spectrum of the graph laplacian matrix is insensitive to any rotation of the adjacency matrix

and may thus be considered as a characteristic of the graph rather than a characteristic of its adjacency matrix. The resulting graph embedding are based on global characteristics of the graphs which are however often difficult to interpret in terms of graph properties.

### 1.3.5 Conclusion

Graph embedding methods define graph vectorial representations. Conversely to methods defined within graph space, these vectorial representations can be used in any machine learning methods defined on vectors. However, the encoding of a graph by an explicit vector of limited size induces necessarily a loss of information. This loss may alter prediction models since some relevant information may not be encoded by a vector of limited size.

---

## 1.4 Kernel theory

Methods presented in previous sections are based on two different approaches. The first one consists in defining similarity measures between graphs without reference to any embedding. Such similarity measures may be combined with few machine learning algorithms such as the k-nearest neighbors or the k-median algorithms. The second one consists in defining vectorial representations of graphs. These methods may be combined with any machine learning algorithm but the transformation from graphs to vectors may induce important loss of information due to the limited vector size.

In order to combine a graph similarity measure with many machine learning algorithms we have thus to transform our graphs into vectors. In the same time, the dimension of the embedding space must be sufficiently large (even infinite) in order to avoid to loss most of the graph information by this transformation. From this point of view, graph kernels provide a nice mathematical framework which allows to define a similarity measure between graphs which corresponds to a scalar product in some Hilbert space. The key point, known as the kernel trick, is that many machine learning algorithms may be rewritten solely in terms of scalar products between input data. The substitution of scalar products by graph kernels within these methods allows then to avoid the explicit transformation of graphs into vectors, hence allowing to work in Hilbert spaces of arbitrary dimension.

### 1.4.1 Definitions

Intuitively, a kernel  $k : \mathcal{X}^2 \rightarrow \mathbb{R}$  between two objects  $x$  and  $x'$  corresponds to a similarity measure defined as a scalar product between two projections

$\Phi_{\mathcal{H}}(x)$  and  $\Phi_{\mathcal{H}}(x')$  of  $x$  and  $x'$  into an Hilbert space  $\mathcal{H}$ :

$$\forall (x, x') \in \mathcal{X}^2, k(x, x') = \langle \Phi_{\mathcal{H}}(x), \Phi_{\mathcal{H}}(x') \rangle. \quad (1.6)$$

In order to define a valid kernel  $k$ , it is not mandatory to explicitly define the embedding function  $\Phi_{\mathcal{H}} : \mathcal{X} \rightarrow \mathcal{H}$ . However, to ensure that such an embedding does exist, the kernel  $k$  has to fulfill some properties.

**Definition 23 (Positive definite kernel)** *A positive definite kernel on  $\mathcal{X}^2$  is a function  $k : \mathcal{X}^2 \rightarrow \mathbb{R}$  symmetric:*

$$\forall (x, x') \in \mathcal{X}^2, k(x, x') = k(x', x) \quad (1.7)$$

*and semi definite positive (Mercer's condition [34]):*

$$\forall \{x_1, \dots, x_n\} \in \mathcal{X}^n, \forall c \in \mathbb{R}^n, \sum_i^n \sum_j^n c_i k(x_i, x_j) c_j \geq 0. \quad (1.8)$$

Some usual kernels defined on vectorial representations are listed in table 1.1.

**TABLE 1.1:** Usual definite positive kernels between vectors.

Linear	$k(x, y) = x^T y$
Gaussian	$k(x, y) = \exp(-\frac{\ x-y\ ^2}{2\sigma^2})$
Polynomial	$k(x, y) = (x^T y)^d + c, c \in \mathbb{R}, d \in \mathbb{N}$
Cosine	$k(x, y) = \frac{(x^T y)}{\ x\  \ y\ }$
Intersection	$k(x, y) = \sum_{i=1}^N \min(x_i, y_i)$

**Definition 24 (Gram matrix)** *A Gram matrix  $K$  associated to a kernel  $k$  on a finite set  $X = \{x_1, \dots, x_N\}$  is a  $N \times N$  matrix defined by:*

$$\forall (i, j) \in \{1, \dots, N\}^2, K_{i,j} = k(x_i, x_j). \quad (1.9)$$

For any finite set of objects  $X = \{x_1, \dots, x_N\}$ , the Gram matrix associated to a positive definite kernel  $k$  is semi definite positive. Conversely, if for any set  $X = \{x_1, \dots, x_N\}$ , the Gram matrix  $K$  associated to kernel  $k$  is semi definite positive, then  $k$  is a positive definite kernel.

A positive definite kernel  $k$  can be built from a combination of positive definite kernels  $k_1, \dots, k_n$  (proposition 3).

**Proposition 3 (Kernel combination [2])** *Let  $k_1$  and  $k_2$  denote two definite positive kernels defined on  $\mathcal{X}^2$ ,  $\mathcal{X}$  corresponding to a non empty space. We have:*

1. *The set of positive definite kernels is a closed convex cone. Therefore:*

- let  $w_1, w_2 \geq 0$ , kernel  $k_3 := w_1 k_1 + w_2 k_2$  is positive definite;
  - let  $k_n$  a sequence of positive definite kernels and  $k(x, x') := \lim_{n \rightarrow \infty} k_n(x, x')$ , then  $k$  is a positive definite kernel.
2. The product of two definite positive kernels is a positive definite kernel.
3. Let us consider that for  $i = 1, 2$ ,  $k_i$  is a positive definite kernel on  $\mathcal{X}_i^2$ ,  $\mathcal{X}_i$  being defined as a non empty space and let us consider  $(x_1, y_1) \in \mathcal{X}_1^2$  and  $(x_2, y_2) \in \mathcal{X}_2^2$ . Then, the tensor product  $k_1 \otimes k_2((x_1, x_2), (y_1, y_2)) = k_1(x_1, y_1)k_2(x_2, y_2)$  and the direct sum  $k_1 \oplus k_2((x_1, x_2), (y_1, y_2)) = k_1(x_1, y_1) + k_2(x_2, y_2)$  correspond to positive definite kernels on  $(\mathcal{X}_1 \times \mathcal{X}_2)^2$ .

Following [1], if  $k$  is a positive definite kernel on  $\mathcal{X}$ , then it exists an Hilbert space  $\mathcal{H}$ , having scalar product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ , and an embedding  $\Phi : \mathcal{X} \rightarrow \mathcal{H}$  such that:

$$\forall (x, x') \in \mathcal{X}^2, k(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathcal{H}} \quad (1.10)$$

The Hilbert space  $\mathcal{H}$  is called the Reproducing Kernel Hilbert Space (RKHS) of  $k$  or more usually the feature space.

#### 1.4.2 Kernel trick

Let us consider a polynomial kernel  $k(x, y) = \langle x, y \rangle^2$  with  $x = (x_1, x_2)$  and  $y = (y_1, y_2) \in \mathbb{R}^2$ . Kernel value  $k(x, y)$  is thus equals to:

$$k(x, y) = x_1^2 y_1^2 + x_2^2 y_2^2 + \sqrt{2}(x_1 x_2) \sqrt{2}(y_1 y_2) \quad (1.11)$$

Despite the fact that this last equation does not correspond to the usual equation of a scalar product, it indeed corresponds to a scalar product between a mapping of  $x$  and  $y$  through the following function:

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \xrightarrow{\Phi} \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1 x_2 \end{pmatrix} \quad (1.12)$$

Note that, using a polynomial kernel of degree 2, we implicitly work in a space of dimension 3, while our original data are of dimension 2. Let us now suppose that we wish to combine our kernel with a k-nearest neighbors algorithm. We have thus to compute distances between vectors in the feature space of dimension 3. Such distances may be computed as follows:

$$\begin{aligned} d_k^2(x, y) &= \|\Phi(x) - \Phi(y)\|^2 = \langle \Phi(x), \Phi(x) \rangle + \langle \Phi(y), \Phi(y) \rangle \\ &\quad - 2 \langle \Phi(x), \Phi(y) \rangle \\ &= k(x, x) + k(y, y) - 2k(x, y) \end{aligned} \quad (1.13)$$

Therefore, the k-nearest neighbors algorithm, may be applied in the feature

space of dimension 3 without computing the transformation  $\Phi$  but by using solely our kernel between data of dimension 2. This trick, known as the kernel trick allows us to avoid the explicit transformation of our original input data in many machine learning algorithms. This last point allows us to avoid to compute and to store vectors of very large or even infinite dimension, e.g. for the gaussian kernel (table 1.1).

Applied to graphs, the kernel trick allows us to avoid the explicit transformation of graphs into vectors and to focus our attention on the definition of a similarity measure. The fact that the implicit Hilbert space may be of very large dimension, allows to reduce the loss of information induced by the transformation from graphs to vectors.

### 1.4.3 Kernels and similarity measures

From a mathematical point of view, a kernel is defined as a scalar product between objects embedded into an Hilbert space. However, kernels are generally considered as similarity measures. This relationship between scalar product and similarity measure may be explained by the relationship between scalar product and euclidean distance (equation 1.13):

$$\begin{aligned} \|\Phi(x) - \Phi(y)\|^2 &= k(x, x) + k(y, y) - 2k(x, y) \\ \Rightarrow k(x, y) &= \frac{1}{2}(\|\Phi(x)\|^2 + \|\Phi(y)\|^2 - \|\Phi(x) - \Phi(y)\|^2) \end{aligned} \quad (1.14)$$

If we consider normalized vectorial representations, i.e.  $\|\Phi(x)\| = 1 \forall x \in \mathcal{X}$ , we obtain:

$$k(x, y) = 1 - \frac{1}{2}d_k^2(x, y)$$

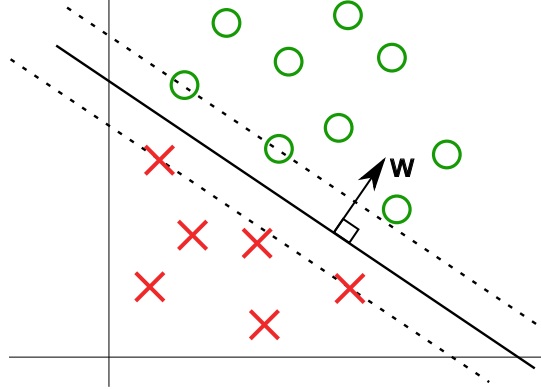
where  $d_k^2(x, y) = \|\Phi(x) - \Phi(y)\|^2$ .

In this case, kernels are defined as the opposite of the euclidean distance between vectors in the feature space. Intuitively, euclidean distance encodes a dissimilarity between objects. Therefore, a kernel function defined as a decreasing function of the distance encodes a similarity measure between objects. An high value corresponds to an high similarity whereas a low value, close to 0, encodes a high dissimilarity between objects.

### 1.4.4 Kernel methods

#### 1.4.4.1 Support vector machines

**Problem definition** Support Vector Machines [5] (SVM) corresponds to a machine learning algorithm. The classification problem addressed by SVM is the following: Given a set of objects labeled by a class  $\{x_i, y_i\}_{i=1}^n$ ,  $x_i \in \mathbb{R}^d$  and  $y_i \in \mathcal{Y} = \{-1, +1\}$ , learn a function  $f : \mathbb{R}^d \rightarrow \mathcal{Y}$  such that  $f(x_i) = y_i$ . SVM algorithm consists in finding an hyperplane having  $d - 1$  dimensions which separates data points according to their classes (figure 1.4). The optimal hyperplane is the one which maximizes the distance between the hyperplane



**FIGURE 1.4:** Hyperplane splitting data according to their classes (green circles and red crosses). Dashed lines encode margins.

and the closest data points of each class. Such a distance is called the margin. Considering the hyperplane defined by  $\langle w, x \rangle + b = 0$ , this distance is inversely proportional to the norm of vector  $w$ . Finding the best hyperplane thus relies to solve:

$$\begin{aligned} & \underset{w}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 \\ & \text{subject to:} \\ & y_i(\langle w, x \rangle + b) \geq 1, \forall i \in \{1, \dots, n\} \end{aligned} \quad (1.15)$$

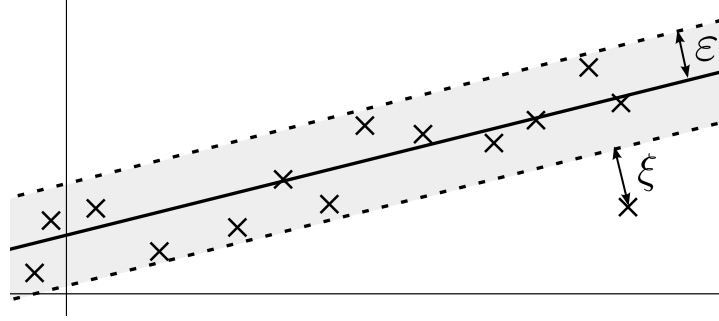
Considering a non linear separable case, i.e. where it does not exist an hyperplane having  $d - 1$  dimensions which separates data, Cortes and Vapnik proposed to include slack variables  $\xi_i \in \mathbb{R}^+$  [8]. These variables allows to take into account classification errors during training. The addressed problem is then defined by:

$$\begin{aligned} & \underset{w}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ & \text{subject to:} \\ & y_i(\langle w, x \rangle + b) \geq 1 - \xi_i, \forall i \in \{1, \dots, n\} \\ & \xi_i \geq 0, \forall i \in \{1, \dots, n\} \end{aligned} \quad (1.16)$$

where  $C \in \mathbb{R}_+$  is a regularization parameter which allows to weight the influence of errors made during training. An high  $C$  value favors a learning without errors, and thus potentially an over learning. Conversely, a low  $C$  value allows more errors during training and thus a greater generalization.

Using a kernel  $k$  together with SVM allows to find a linear hyperplane which solves equation 1.16 in the feature space associated to kernel  $k$ . Thanks to the kernel trick, this hyperplane may be a non linear separator in original data space. For example, let consider the kernel defined by equation 1.12. The hyperplane equation  $\langle w, x \rangle + b = 0$  is computed within the kernel feature space



FIGURE 1.5: Hyperplane associated to an  $\varepsilon$ -tube (shaded).

and corresponds to a non linear equation in  $\mathbb{R}^2$ :

$$w_1x_1^2 + w_2x_2^2 + w_3\sqrt{2}x_1x_2 + b = 0 \quad (1.17)$$

Note that such an equation corresponds to a quadric.

#### 1.4.4.2 Support vectors machines for regression

Regression problems consist in predicting a continuous value, conversely to classification problems which aim to predict a discrete value encoding a class. More formally, a regression problem is defined as: Given a learning set  $\{x_i, y_i\}_{i=1}^n$ , composed of a set of  $n$  objects  $\mathcal{X} = \{x_1, \dots, x_n\}$  with  $x_i \in \mathbb{R}^d$ , each object being associated to a value  $y_i \in \mathbb{R}$ , learn a prediction function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  such that  $\hat{y}_i \simeq f(x_i)$ .

Support vector machines are initially defined as a solver for classification problems. In order to handle regression problems, SVM have been adapted [11] by including an  $\varepsilon$ -insensitive cost function such as defined in [51]. Instead of computing an hyperplane splitting data according to their classes, SVM for regression consists in computing an hyperplane  $w$  associated to an  $\varepsilon$ -tube which includes data points to predict.

Formally, the minimization problem addressed by SVM for regression is defined by:

$$\begin{aligned} & \underset{w}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 \\ & \text{subject to:} \\ & \begin{cases} y_i - \langle w, x_i \rangle - b \leq \varepsilon, \forall i \in \{1, \dots, n\} \\ \langle w, x_i \rangle + b - y_i \geq \varepsilon, \forall i \in \{1, \dots, n\} \end{cases} \end{aligned} \quad (1.18)$$

Minimizing equation 1.18 relies to compute a linear function which approximates  $y$  values with an accuracy  $\varepsilon$ . In order to allow errors during training, slack variables for regression have been introduced into minimization problem. Similarly to SVM classification, these slack variables allow to weight, according to  $C$ , prediction errors made during the learning step. The minimization

problem is then defined as:

$$\begin{aligned} & \underset{w}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) \\ & \text{subject to:} \quad \forall i \in \{1, \dots, n\} \begin{cases} y_i - \langle w, x_i \rangle - b \leq \varepsilon + \xi_i \\ \langle w, x_i \rangle + b - y_i \leq \varepsilon + \xi_i^* \\ \xi_i, \xi_i^* \geq 0 \end{cases} \end{aligned} \quad (1.19)$$

This regression algorithm allows to compute a prediction function up to an accuracy defined by  $\varepsilon$ . This parameter is particularly interesting in chemoinformatics since it may encode inaccuracies of physical properties induced by experiments which are not predictable using molecular graphs.

#### 1.4.4.3 Kernel ridge regression

Considering ridge regression framework, a regression problem can be addressed by computing a linear function which encodes relationships between data  $x$  and responses  $y$ . This function may be computed by minimizing the following objective function:

$$\begin{aligned} & \underset{w}{\text{minimize}} \quad J(w) \\ & \text{with:} \quad J(w) = \|y - Xw\|^2 + \lambda \|w\|^2 \end{aligned} \quad (1.20)$$

Minimizing the first term of equation 1.20 relies to solve a least squares problem corresponding to the minimization of prediction errors. The second term ( $\lambda \|w\|^2$ ) corresponds to a regularization term which aims to penalize vectors  $w$  having an high norm. The factor  $\lambda$  allows to weight the influence of the regularization term into the minimization problem. Therefore, an high  $\lambda$  value may prevent prediction model from over learning which may induce a better generalization of the prediction model when applied on test sets.

The objective function to minimize, defined by equation 1.20, corresponds to a sum of  $L_2$  norms which thus defines a convex function. Global minimum is reached for a value  $w$  satisfying  $\frac{\partial J}{\partial w} = 0$ . This optimal  $w$  is obtained by the following analytic form:

$$w^* = X^t (X X^t + \lambda I)^{-1} y \quad (1.21)$$

Predicting property value of a new data  $x'$  relies thus on computing  $\hat{y} = w^{*t} x'$ .

Let us consider that it exists a vector  $\alpha \in \mathbb{R}^n$  such that  $w = X^t \alpha$ . We have thus  $\alpha = (\lambda I + X X^t)^{-1} y$  and prediction of a new data is given by  $\hat{y} = X X^t \alpha$ . We can notice that the access to data  $X$  is only performed through  $X X^t$ . Therefore, minimizing objective function can be performed without accessing directly to raw data  $\{x_1, \dots, x_n\}$  but only through the scalar product defined on this data since the matrix  $X X^t \in \mathbb{R}^{n \times n}$  encodes scalar product  $\langle x_i, x_j \rangle$  for each pair of objects  $(x_i, x_j) \in \mathcal{X}^2$ . Within the kernel framework,  $X X^t$  corresponds to the Gram matrix  $K$  associated to the kernel  $k(x_i, x_j) = \langle x_i, x_j \rangle$ .

Given a non linear kernel  $k_{\mathcal{H}}$ , associated to a feature space  $\mathcal{H}$  and an embedding function  $\Phi_{\mathcal{H}} : \mathcal{X} \rightarrow \mathcal{H}$ , ridge regression using  $k_{\mathcal{H}}$  relies to compute a vector  $w$  lying on feature space  $\mathcal{H}$ , instead of original data space  $\mathcal{X}$ .

### 1.4.5 Graph kernels and chemoinformatics

Considering conclusions of sections 1.2 and 1.3, graph kernels provide a solution to both limitations involved by methods based on graph embeddings and methods defined in graph space. On one hand, graph kernels allow to define a similarity measure which is not limited by vector sizes since graph vectorial representations are not required to be explicitly and exhaustively defined. On the other hand, graph kernels guarantee that a vectorial representation exists hence allowing the conjoint use of kernels with many machine learning methods. Therefore, graph kernels define a natural connection between molecular graphs and machine learning algorithms. From this point of view, many graph kernels have been defined to solve QSAR/QSPR problems in chemoinformatics.

---

## 1.5 Graph kernels based on bags of patterns

A graph kernel family widely used in chemoinformatics consists in defining similarity between graphs from a set of sub structures extracted from graphs to be compared. Formally, given a set of sub structures  $\mathcal{P}$  and a sub kernel  $k' : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}$ , a convolution kernel [21]  $K : \mathcal{G} \times \mathcal{G} \rightarrow \mathbb{R}$  is defined as a sum of sub kernel values. For each pair of sub structures  $(p, p') \in \mathcal{P}^2$  extracted from graphs by a decomposition function  $D : \mathcal{G} \rightarrow \mathcal{P}$ , sub kernel  $k$  encodes the similarity between two molecular graphs according to  $p$  and  $p'$ . Kernel  $k$  is thus equal to:

$$k(G, G') = \sum_{p \subseteq D(G)} \sum_{p' \subseteq D(G')} k'(p, p') \quad (1.22)$$

Kernels presented in this section are defined as convolution kernels. Although these kernels use a same approach, they mainly differ on the set of sub structures used to define them.

### 1.5.1 Complete graph kernel

A first trivial approach consists in defining the set of sub structures by all possible sub graphs.

**Definition 25 (Complete graph kernel)** *Let  $G$  and  $G'$  two graphs. Com-*

plete graph kernel is defined by:

$$k_{complete}(G, G') = \sum_{p \subseteq G} \sum_{p' \subseteq G'} k_{iso}(p, p') \quad (1.23)$$

with

$$k_{iso}(p, p') = \begin{cases} 1 & \text{iff } p \simeq p' \\ 0 & \text{otherwise} . \end{cases} \quad (1.24)$$

Complete graph kernel [14] relies on embedding graphs into an Hilbert space where each dimension encodes the number of occurrences of a sub graph of the considered graph. However, computing this kernel is NP-hard since it relies on determining if two graphs are isomorphic. Therefore, this kernel can not be efficiently computed which avoid using such an approach to chemoinformatics problems.

## 1.5.2 Linear pattern kernels

### 1.5.2.1 Random walks

A less complex approach consists in deducing graph similarity from the number of random walks common to both graphs to be compared. Two methods [14, 26] have been proposed both based on the bag of random walks  $\mathcal{W}(G)$  of a graph  $G$ . Kernels defined using random walks are based on the following formulation:

**Definition 26 (Marginalized kernel)** *Given two graphs  $G$  and  $G'$ ,  $\mathcal{W}(G)$  and  $\mathcal{W}(G')$  encode the sets of random walks in  $G$  and  $G'$  and  $k_{\mathcal{W}} : \mathcal{W} \times \mathcal{W} \rightarrow \mathbb{R}$  is defined as a kernel between walks. Marginalized kernel is then defined as:*

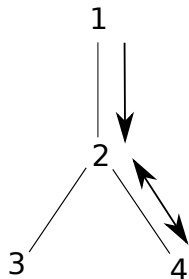
$$k_{rw}(G, G') = \sum_{w \in \mathcal{W}(G)} \sum_{w' \in \mathcal{W}(G')} p_G(w) p_{G'}(w') k_{\mathcal{W}}(w, w') \quad (1.25)$$

where  $p_G(w)$  corresponds to the probability of traversing  $w$  in  $G$ .

Kernel  $k_{\mathcal{W}}$  corresponds to a kernel encoding a similarity between walks according to their labeling. Usually, this kernel is binary, i.e.  $k_{\mathcal{W}}$  is equal to 1 if label sequences are similar and 0 otherwise. Method defined in [14] uses a different approach based on a direct product graph but is based on the same feature space and thus encodes the same information. Different implementations have been proposed [26, 53, 14, 31] and this kernel may be computed in polynomial time with the number of nodes of both graphs.

### 1.5.2.2 Tottering

A random walk corresponds to a node sequence which may oscillate between two connected neighbours (figure 1.6). This phenomenon, called tottering, induces random walks which are not representative of the corresponding



**FIGURE 1.6:** A path traversal altered by tottering corresponds to a label sequence equals to 1, 2, 4, 2, 4, 2, ...

molecular graph structure since a same information may be repeated indefinitely. In order to limit the influence of this non representative structures, Mahé and al. [31] proposed to transform a molecular graph into an oriented graph in  $\mathcal{O}(|V|^2)$ . Then, kernel is computed on this new molecular representation which allows to avoid oscillations between two connected nodes and thus to reduce tottering phenomenon.

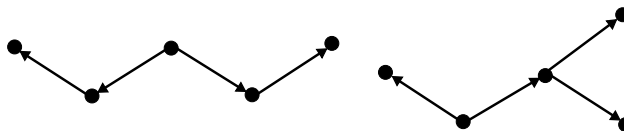
### 1.5.2.3 Paths

Instead of using random walks which may suffer from tottering, others kernels base their similarity measure on paths extracted from graphs to be compared. Although enumerating all paths of a graph remains a NP-Hard problem [4], a kernel can be defined using all paths composed of at most  $n$  nodes [40]. This path enumeration is performed using a depth first traversal from each node of the graph. Computing the kernel on two graphs  $G = (V, E)$  and  $G' = (V', E')$  relies on extracting paths from  $G$  and  $G'$  which induces  $\mathcal{O}(n(|V||E| + |V'||E'|))$  operations. Enumerated paths are encoded by a vector which can be seen as a molecular descriptor where each descriptor encodes the number of occurrences of a path. The similarity between these vectors may be computed using usual kernels such as Tanimoto kernel or MinMax kernel [40].

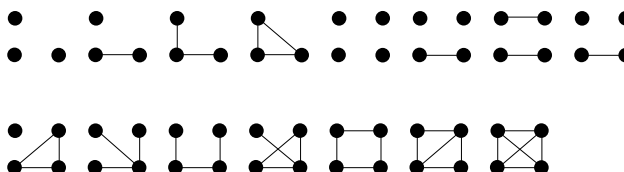
Another approach consists in computing similarity between graphs from the lengths of shortest paths between any pair of nodes. This method is based on Floyd's transformation [13] which consists in transforming a graph  $G = (V, E, \mu, \nu)$  into a complete graph  $G_F = (V, E_F, \mu, \nu_F)$  where each edge  $e_f = (v_i, v_j) \in V^2$  is labeled by the length of a shortest path between  $v_i$  and  $v_j$  in  $G$ . Considering the two graph transformations  $G_F = (V, E_F, \mu, \nu_F)$  and  $G'_F = (V', E'_F, \mu', \nu'_F)$  of two graphs  $G = (V, E, \mu, \nu)$  and  $G' = (V', E', \mu', \nu')$ , the shortest paths kernel is defined by:

$$k_{sp}(G, G') = \sum_{e \subseteq E_F} \sum_{e' \subseteq E'_F} k_{path}^1(e, e') \quad (1.26)$$

where  $k_{path}^1(e, e')$  is a kernel defined between edges. This kernel may be defined



**FIGURE 1.7:** Two different oriented graphs having a same representation in a feature space based on linear patterns.



**FIGURE 1.8:** All graphlets of size 3 and 4 [46].

as a Dirac function which considers two edges  $e = (u, v)$  and  $e' = (u', v')$  as similar if and only if  $\mu(u) = \mu'(u')$ ,  $\mu(v) = \mu'(v')$  and  $\nu_F(e) = \nu'_F(e')$ . Intuitively, two edges are similar if they connect two pairs of nodes having same labels and separated by the same distance in  $G$  and  $G'$ . This kernel can be computed in polynomial time ( $\mathcal{O}(|V|^4)$ ), most of the time being dedicated to the computation of Floyd's transformation and to the comparison of the  $|V|^2$  edges of transformed graphs. However, this approach only encodes the length of a shortest path connecting two nodes, but not the sequences of node and edge labels composing these shortest paths. This loss of information leads to a less accurate kernel.

### 1.5.3 Non linear pattern kernels

Although linear patterns allows to define a kernel which can be computed in polynomial time, most of the structural information included within molecular graphs can not be encoded using only linear patterns (figure 1.7).

#### 1.5.3.1 Graphlets

In order to increase the structural information encoded within bags of patterns, graphlet kernel [46] has been introduced. This kernel is based on bags of patterns defined as all unlabeled graphs having  $k$  nodes,  $k \in \{3, 4, 5\}$  (figure 1.8). This kernel may be efficiently computed for graphs having a bounded maximum degree (say by  $d$ ). Note that a molecular graph (section 1.1.1) has a maximal degree bounded by 8.

Graphlet enumeration is divided in two steps. The first one corresponds to the enumeration of connected graphlets and the second one to the enumeration of disconnected graphlets. Connected graphlets are divided into two classes. The first class corresponds to graphlets of size  $k$  containing at least a path of

size  $k - 1$  and the second class to others graphlets. The first class of graphlets is enumerated from paths of length  $k - 1$ , themselves enumerated using a depth first traversal. The second class of graphlets is enumerated from  $3$ -star graphlets which correspond to graphlets having 4 nodes with a central node having a degree equals to 3. Others graphlets are then enumerated by looking at the neighborhood of  $3$ -stars.

Concerning disconnected graphlets of size 3, their enumeration is based on pairs of nodes  $(v_1, v_2) \in E$ . If we consider a third node  $v_3$  distinct from  $v_1$  and  $v_2$ , three different configurations are possible (figure 1.9). These different configurations may be distinguished using the neighborhood of  $(v_1, v_2)$ . We have thus:

- $v_3 \notin \Gamma(v_1) \cup \Gamma(v_2) \rightarrow F_1$  configuration;
- $v_3 \in \Gamma(v_1) \setminus \Gamma(v_2) \vee v_3 \in \Gamma(v_2) \setminus \Gamma(v_1) \rightarrow F_2$  configuration;
- $v_3 \in \Gamma(v_1) \cap \Gamma(v_2) \rightarrow F_3$  configuration;

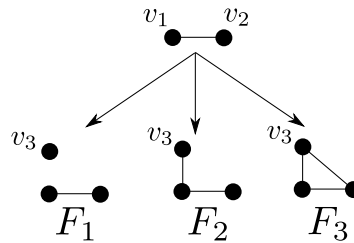
The method defined to enumerate graphlets of size 4 and 5 is based on the same scheme then the one used for graphlets of size  $k = 3$ : the set of 11 graphlets having 4 nodes is enumerated from the set of graphlets having 3 nodes. Similarly, the 34 graphlets having 5 nodes is enumerated from the ones composed of 4 nodes. This method allows to compute the distribution of graphlets of size  $k$  in  $\mathcal{O}(nd^{k-1})$  operations which corresponds to a linear complexity with the number of nodes of graphs if the maximum degree is bounded. Graphlet kernel is then defined as a scalar product between normalized vectors encoding the number of occurrences of each graphlet:

**Definition 27 (Graphlet kernel)** Let  $G = (V, E) \in \mathcal{G}$  and  $G' = (V', E') \in \mathcal{G}$  two graphs and  $D_G$  a vector encoding the number of occurrences of graphlets in  $G$ :

$$D_G(i) = \frac{\#\text{graphlet}(i) \subseteq G}{\sum_j \#\text{graphlet}(j) \subseteq G} \quad (1.27)$$

Graphlet kernel is then defined by:

$$k_{\text{graphlets}} = D_G^\top D_{G'} \quad (1.28)$$



**FIGURE 1.9:** Different possible topological configurations involving a pair of nodes and a third node.

Although this kernel allows to encode more structural information than kernels based on linear patterns, it does not encode the similarity involved by labeling. Therefore, this kernel is limited to unlabeled graphs which corresponds to a limited application domain in chemoinformatics. Indeed, atom's chemical elements are an important information and must be taken into account in order to define an accurate similarity measure.

### 1.5.3.2 Tree patterns

In order to include molecular graphs labeling into kernel computation, tree pattern kernel [41, 32] is defined as a kernel based on non linear and labeled patterns. This method deduces the similarity between two graphs from the number of their common tree patterns.

**Definition 28 (Tree pattern)** Let  $G(V, E) \in \mathcal{G}$ . If it exists a node  $r \in V$ , then  $r$  is a tree pattern of  $G$  rooted in  $r$  and having an height equals to 1. let  $t_1, t_2, \dots, t_n$   $n$  tree patterns respectively rooted in  $r_1, r_2, \dots, r_n$ , with  $r_i \neq r_j, \forall i \neq j$ . If  $(r, r_1), (r, r_2), \dots, (r, r_n) \in E$ , then  $r(t_1, t_2, \dots, t_n)$  is a tree pattern rooted in  $r$ .  $r$  is defined as the parent of each  $r_i$ .

Each tree pattern is associated to a dimension in the feature space. Each dimension encodes the number of tree patterns included within the two graphs to be compared.

**Definition 29 (Neighborhood matching set)** Let  $G = (V, E, \mu, \nu)$  and  $G' = (V', E', \mu', \nu')$  denote two graphs and let us consider two nodes  $r \in V$  and  $s \in V'$ . Neighborhood matching set  $M_{r,s}$  is defined by:

$$M_{r,s} = \{R \subseteq \Gamma(r) \times \Gamma(s) \mid \begin{array}{l} (\forall (a, b), (c, d) \in R : a \neq c \wedge b \neq d) \wedge \\ (\forall (a, b) \in R : \mu(a) = \mu'(b) \wedge \nu(r, a) = \nu'(s, b)) \end{array}\} \quad (1.29)$$

**Definition 30 (Kernel between tree patterns)** Let  $G = (V, E, \mu, \nu)$  and  $G' = (V', E', \mu', \nu')$  denote two graphs, the kernel value between two tree patterns of height  $h$  and rooted in  $r$  and  $s$  is defined by:

$$k(r, s, h) = \lambda_r \lambda_s \sum_{R \in M_{r,s}} \prod_{(r', s') \in R} k(r', s', h-1). \quad (1.30)$$

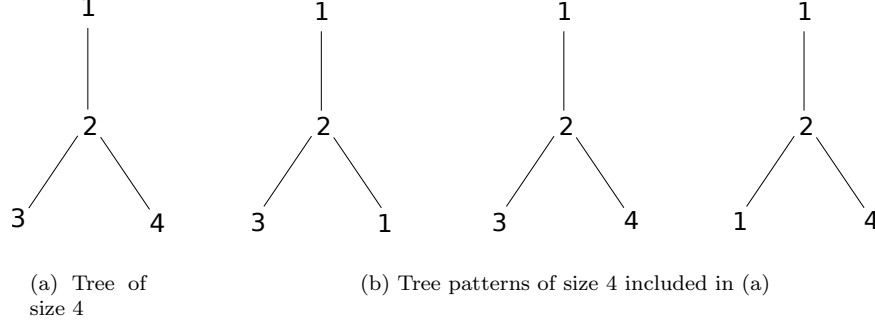
For  $h = 1$ ,  $k(r, s, 1)$  is equals to 1 if  $\mu(r) = \mu'(s)$ , 0 otherwise.

Kernel value is strictly greater than 0 if the two tree patterns are isomorphic. Parameters  $\lambda_r$  and  $\lambda_v$  are defined as positive real numbers lower than 1 such that large tree patterns have a low contribution. Tree pattern kernel is then defined as:

**Definition 31 (Tree pattern kernel)** Let  $G = (V, E, \mu, \nu) \in \mathcal{G}$  and  $G' = (V', E', \mu', \nu') \in \mathcal{G}$ . Tree pattern kernel for an height  $h$  is defined by:

$$k(G, G', h) = \sum_{r \in V} \sum_{s \in V'} k(r, s, h) \quad (1.31)$$





**FIGURE 1.10:** Difference between tree patterns and trees : graph (a) includes only one sub tree of size 4 rooted in 1 (itself) but 3 different tree patterns (b), all rooted in 1.

#### 1.5.3.3 Weisfeiler-Lehman

Weisfeiler-Lehman kernels [47, 45] correspond to graph kernels based on a subset of tree patterns corresponding to tree structures encoding all nodes within a given radius. This algorithm, based on Weisfeiler-Lehman isomorphism test, allows to compute a kernel based on a sub family of tree patterns with a linear complexity according to the number of nodes of both graphs. It thus provides faster computational times than the tree pattern kernel (section 1.5.3.2). Moreover, Weisfeiler-Lehman kernels computes an explicit representation of the feature space associated to kernels. This explicit enumeration allows to perform pattern enumeration only  $N$  times where  $N$  is the size of the training set. Using an implicit enumeration, pattern enumeration has to be performed for each pair of graph, i.e.  $N^2$  times.

One has to notice that tree patterns differ from trees (figure 1.10). Similarly to the difference between paths and walks, a sub tree pattern may include a same node twice. Therefore, information encoded by tree patterns may also be altered by tottering (section 1.5.2.2). However, graph transformation proposed by Mahé and al. to prevent tottering in random walks kernel can also be applied to kernels based on tree patterns.

#### 1.5.3.4 Graph fragments

Graph fragment kernel [55] corresponds to a kernel based on a set of patterns defined as all connected sub graphs having at most  $l$  edges. This set of subgraphs, denoted  $GF$ , can be divided into different subsets:

- *PF (Path Fragments)*: encodes all linear patterns;
- *TF (Tree Fragments)*: encodes all sub trees having at least one node  $v$  such that  $d(v) \geq 2$ .

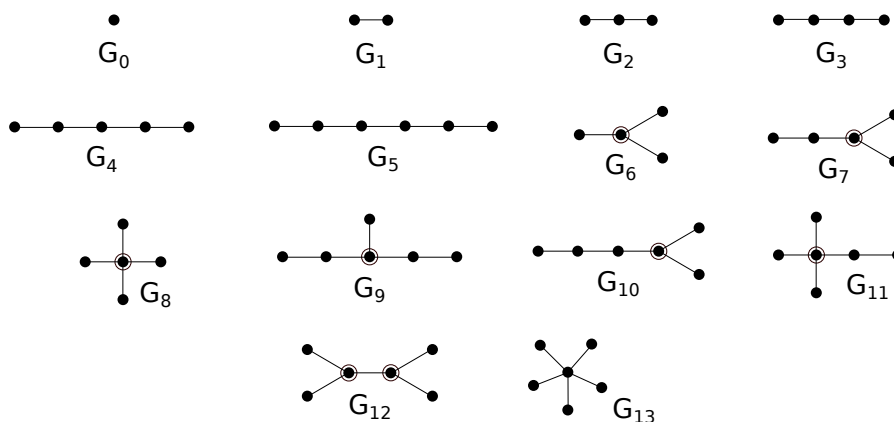


FIGURE 1.11: Set of treelet's structures

- $AF$  (*Acyclic Fragments*) is defined as  $AF = TF \cup PF$  which thus corresponds to all sub trees.
- The fourth sub set corresponds to the set difference  $GF \setminus AF$ . This sub set encodes all sub graphs containing at least one cycle.

This set of patterns defines the feature space associated to this kernel where each coordinate  $\Phi_i(G)$  encodes either the absence/presence or the number of occurrences of a given sub graph in  $G$ . A molecular graph  $G$  is thus encoded by a vector  $\Phi(G) = \{\Phi_1(G), \dots, \Phi_d(G)\}$ , where  $d$  encodes the number of different patterns enumerated for a given length  $l$ . The kernel between two graphs  $G$  and  $G'$  is then defined as a Min-Max kernel [40] between the two vectors  $\Phi(G)$  and  $\Phi(G')$ .

This method, based on an exhaustive enumeration of all sub graphs, relies to compute the complete graph kernel for a  $l$  big enough (section 1.5.1). Since this problem is a NP-hard problem and thus not feasible for large graphs, graph fragment is generally defined for  $l = 7$  in order to obtain reasonable computational times. Moreover, no significant accuracy gain is observed for length  $l \geq 5$ .

#### 1.5.3.5 Treelets

Conversely to methods based on tree patterns, treelet kernel [16] uses a bags of patterns defined as a sub set of strict sub trees, called treelets. Treelets are defined as the set of all labeled trees having at most 6 nodes (figure 1.11). Conversely to methods based on a non limited set of substructures, such as the tree pattern kernel, considering a set of predefined structures allows to define an efficient *ad hoc* linear enumeration. This enumeration is performed in two steps:

The first step consists to identify treelet's structure in a graph. This step

is based on the method used to enumerate connected graphlets [46] (section 1.5.3.1). Graphlet enumeration method is extended to enumerate tree structures up to 6 nodes while keeping a linear complexity when considering degree bounded graphs. This first step allows to associate a structure index  $(G_0, \dots, G_{13})$  (figure 1.11) to each treelet.

The second step consists to encode the labels of each instance of treelet found within a graph. This step is based on a tree traversal of each instance of treelet based on Morgan algorithm [35]. This traversal provides a string which, as shown in [16], identifies the label of each node and each edge of a treelet given the index of its structure. Note that, since the size of treelet's structure is bounded, the computation of the string encoding treelet's labels is performed in constant time. This last point is an important advantage over alternative encodings used for example by frequent subgraphs or graph fragments methods [10].

The concatenation of treelet's index and treelet's key defines an unique code for each treelet which allows to perform an explicit enumeration of all treelets composing a graph. Based on this enumeration, we define a function  $f$  which associates to each graph  $G$  a vector  $f(G)$  whose components encode the number of occurrences of each treelet  $t$  found in  $G$ :

$$f(G) = (f_t(G))_{t \in \mathcal{T}(G)} \text{ with } f_t(G) = |t \subseteq G| \quad (1.32)$$

where  $\mathcal{T}(G)$  denotes the set of treelets extracted from  $G$  and  $\subseteq$  the sub graph isomorphism relationship. Then, similarity between treelet distributions is computed using a sum of sub kernels between treelet's frequencies:

$$k_{\mathcal{T}}(G, G') = \sum_{t \in \mathcal{T}(G) \cap \mathcal{T}(G')} k(f_t(G), f_t(G')) \quad (1.33)$$

where  $k(.,.)$  defines any positive definite kernel between real numbers such as linear kernel, gaussian kernel or intersection kernel. Note that, conversely to tree pattern kernel (section 1.5.3.2), this kernel explicitly enumerates sub trees by computing the number of occurrences of each pattern.

#### 1.5.4 3D pattern kernel

Three-dimensional molecular information may be an important and useful information for some chemoinformatics problems such as docking or optical angle rotation. Among the few existing kernels including such information, 3D pharmacophore kernel [30] encodes a set of patterns corresponding to triplets of distinct atoms. Each pattern can thus be understood as a potential pharmacophore and encodes euclidean distances between each pair of nodes included within a corresponding triplet. These distances are measured on the most stable 3D conformation of a molecule. Pharmacophore kernel is defined as a convolution kernel based on a kernel  $k_{\mathcal{P}}$  defined on atom triplets. This last kernel is defined as a combination of two kernels:

- A first kernel which compares nodes labels of the two triplets to be compared;
- a second kernel which encodes spatial information by comparing distances between atoms by means of gaussian kernels. More formally, if  $p = (x_i)_{i=\{1,\dots,3\}}$  and  $p' = (x'_i)_{i=\{1,\dots,3\}}$  are two atom triplets. This kernel is defined as:

$$k_{spatial}(p, p') = \prod_{i=1}^3 k_{dist}(\|x_i - x_{i+1}\|, \|x'_i - x'_{i+1}\|) \quad (1.34)$$

Kernel  $k_{dist}$  may be defined as a gaussian kernel and aims to compare two distances.

This kernel encodes spatial information of molecular compounds and allows to define a relevant similarity measure for chemoinformatics problems including 3D information. However, distances used in this kernel are relative to a priori 3D conformations of both molecules while some alternative conformations may better explain the property to predict.

### 1.5.5 Cyclic pattern kernels

Similarly to 3D information, cyclic information may have a particular influence on molecular properties since cycles reduce the atom's degrees of freedom.

#### 1.5.5.1 Cyclic patterns

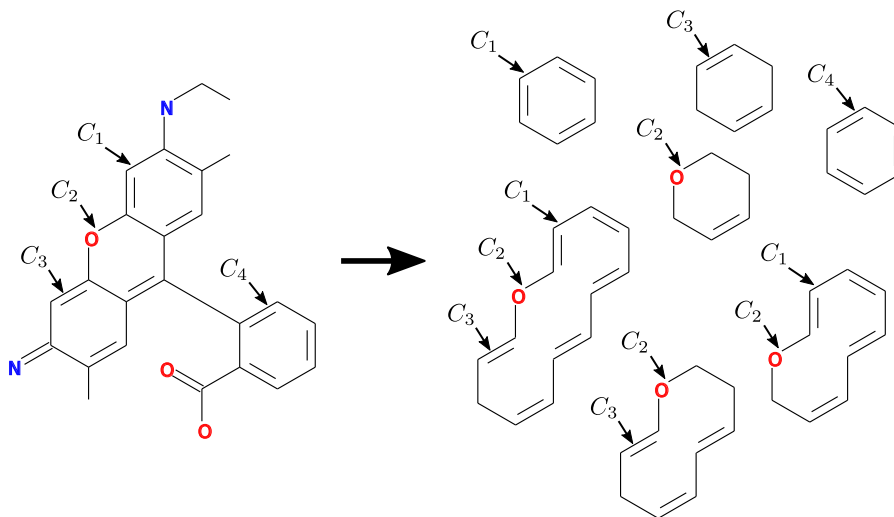
A first approach consists in defining a kernel based on the set of simple cycles of a molecular graph (figure 1.12). Considering a graph  $G$ , cyclic pattern kernel [23] is based on a decomposition of  $G$  into two sub sets:

- A first subset  $C(G)$  corresponding to the set of simple cycles included within  $G$  (figure 1.12);
- and a second subset  $\mathcal{B}(G)$  defined as the set of bridges of  $G$  and corresponding to atoms and edges not included in  $C(G)$ .  $\mathcal{B}(G)$  corresponds thus to a forest.

Considering these two sub sets, cyclic pattern kernel is defined as the sum of a tree pattern kernel applied on  $\mathcal{B}(G)$  together with a kernel applied on  $C(G)$ . Considering two graphs  $G$  and  $G'$ , cyclic pattern kernel is thus defined as:

$$k_C(G, G') = k_{TP}(\mathcal{B}(G), \mathcal{B}(G')) + k_{inter}(C(G), C(G')) \quad (1.35)$$

where  $k_{TP}(\mathcal{B}(G), \mathcal{B}(G'))$  corresponds to tree pattern kernel (section 1.5.3.2) and  $k_{inter}(C(G), C(G'))$  is defined as an intersection kernel which computes



**FIGURE 1.12:** Decomposition of a molecular graph into its set of simple cycles.

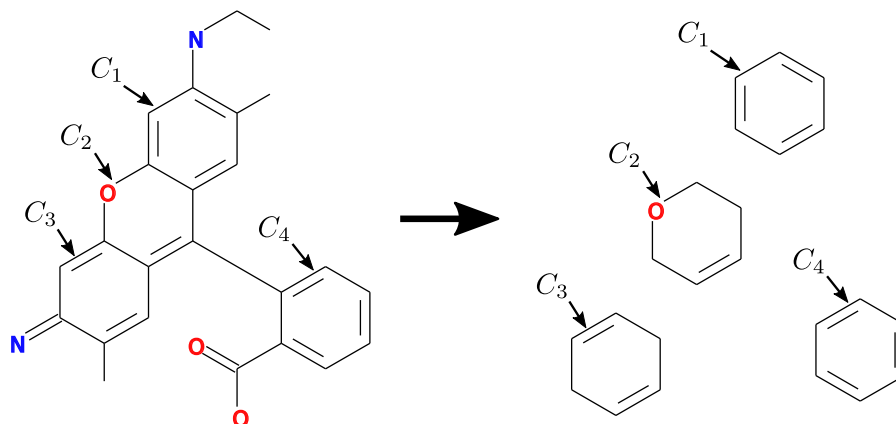
the number of common cycles in  $C(G)$  and  $C(G')$ . Similarity between molecular graphs is thus deduced from a sum of two kernels encoding respectively similarities between their cyclic and acyclic parts.

However, computing the set of simple cycles of a graph is a NP-hard problem which may however be addressed in chemoinformatics problem since molecular graphs generally include a low number of simple cycles. In order to tackle this complexity, cyclic pattern kernel has been defined using the set of relevant cycles (denoted  $\mathcal{C}_{\mathcal{R}}$ ) [22, 54] instead of the set of simple cycles. The enumeration of relevant cycles can be performed in polynomial time which reduces kernel computational complexity. Relevant cycles can be understood as a set of cycles which defines a basis for encoding all cycles included in a molecular graph (figure 1.13).

Cyclic pattern kernel allows to explicitly encode cyclic similarity into kernel computation which may provide an accurate similarity measure for chemoinformatics problems. However, although this approach allows to catch some cyclic information, it does not encode relationships between molecular cycles.

#### 1.5.5.2 Relevant cycle graph

In order to encode more cyclic information, relevant cycle relationships may be encoded using relevant cycle graph [18]. This graph representation aims to encode adjacency relationships between relevant cycles and thus to encode the cyclic system of a molecule. This graph is defined as  $G_{\mathcal{C}} = (\mathcal{C}_{\mathcal{R}}, E_{\mathcal{C}_{\mathcal{R}}}, \mu_{\mathcal{C}_{\mathcal{R}}}, \nu_{\mathcal{C}_{\mathcal{R}}})$  where each vertex encodes a relevant cycle and two vertices are connected by an edge if corresponding cycles share at least



**FIGURE 1.13:** Decomposition of a molecular graph into its set of relevant cycles.

one vertex of the initial graph (figure 1.14). Relevant cycle graph labeling functions are defined as follows:

- $\mu_{\mathcal{C}_{\mathcal{R}}}(C)$ : Each cycle  $C$  is defined by a sequence of edge and vertex labels encountered during the traversal of  $C$ . In order to obtain a sequence invariant to cyclic permutations,  $\mu_{\mathcal{C}_{\mathcal{R}}}(C)$  is defined as the sequence having the lowest lexicographic order.
- $\nu_{\mathcal{C}_{\mathcal{R}}}(e)$ : An edge  $e$  in  $G_{\mathcal{C}}$  encodes a path shared by two cycles and corresponds to a sequence of edge and node labels. Since such a path may be traversed from its two extremities, we define  $\nu_{\mathcal{C}_{\mathcal{R}}}(e)$  as the sequence of lowest lexicographic order.

In order to compute a cyclic similarity using relevant cycle graphs, treelet kernel (section 1.5.3.5) may be applied on the set of treelets extracted from relevant cycle graphs. This kernel is thus defined as follows:

$$k_{\mathcal{C}}(G, G') = \sum_{t_{\mathcal{C}} \in \mathcal{T}(G_{\mathcal{C}}) \cap \mathcal{T}(G'_{\mathcal{C}})} k(f_{G_{\mathcal{C}}}(t_{\mathcal{C}}), f_{G'_{\mathcal{C}}}(t_{\mathcal{C}})) \quad (1.36)$$

Where  $t_{\mathcal{C}}$  encodes a configuration of adjacent cycles present both in  $G$  and  $G'$ .

Then, using cyclic pattern kernel, a kernel based on relevant cycle graph may be defined as a sum of a kernel encoding acyclic similarity, such as treelet kernel  $k_{\mathcal{T}}$  applied on the original molecular graphs, and a kernel  $k_{\mathcal{C}}$  applied on their relevant cycle graphs:

$$k(G, G') = k_{\mathcal{T}}(G, G') + \lambda k_{\mathcal{C}}(G, G') \quad (1.37)$$

where  $\lambda \in \mathbb{R}_+$  allows to weight the contribution of cyclic kernel.

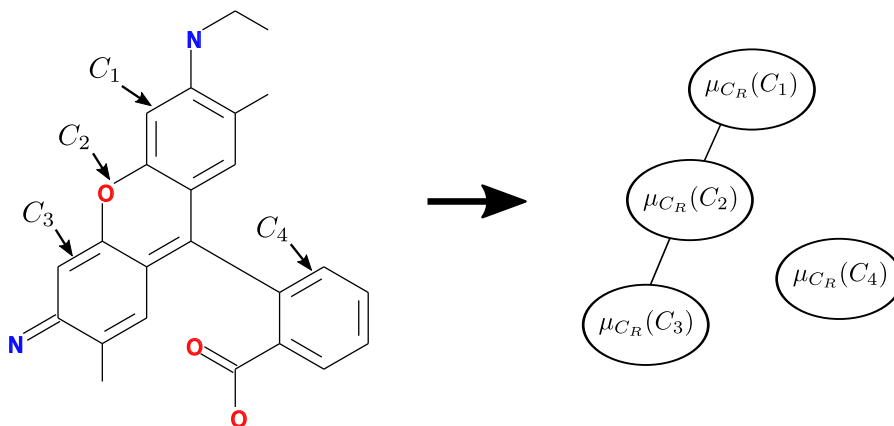


FIGURE 1.14: Relevant cycle graph computed from a molecular graph.

This kernel allows to take into account more cyclic information since, conversely to cyclic pattern kernel, relationships between relevant cycles are encoded in the relevant cycle graph. However, such as cyclic pattern kernel, global similarity measure is based on a sum which splits cyclic information and acyclic information. This last point constitutes a drawback since within the chemoinformatics framework, relationships between a cycle and its substituents may have an influence on molecular properties and must be encoding in kernel definition.

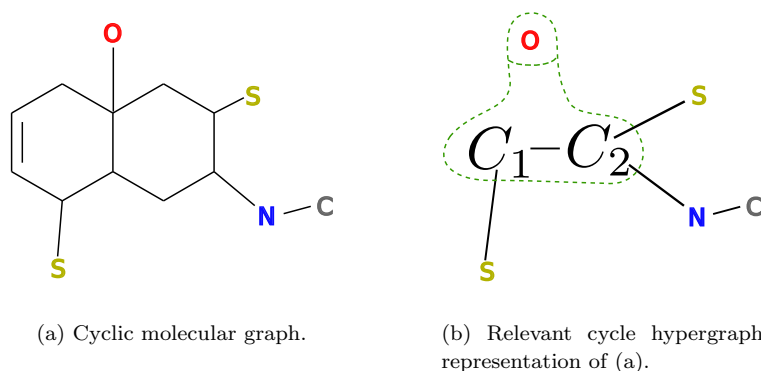
#### 1.5.5.3 Relevant cycle hypergraph

Relevant cycle hypergraph [17] (figure 1.15) corresponds to an oriented hypergraph representation [3, 12] which allows to encode both cyclic and acyclic informations into a single representation. This molecular representation is based on relevant cycle graph which encodes molecular cyclic systems. Relevant cycle graph may be augmented by adding acyclic parts in order to obtain the relevant cycle hypergraph. Given a graph  $G = (V, E, \mu, \nu)$ , its corresponding relevant cycle hypergraph  $H_{RC}(G) = (V_{RC}, E_{RC})$  is defined as follows:

- The set of nodes  $V_{RC}$  is defined as an union of two subsets:
  1. A first subset corresponding to the set of relevant cycles,
  2. and a second subset corresponding to the set of atoms not included within a cycle of  $G$ .

All atoms of the original molecular graph are thus encoded in the relevant cycle hypergraph either by a vertex encoding a cycle or by the atom itself. In the same way as for vertices,

- the set of hyperedges  $E_{RC}$  is composed of two subsets:



**FIGURE 1.15:** Relevant cycle hypergraph representation of a molecule. We can note that hyperedge (in dashed green) encodes an adjacency between oxygen atom in one hand and two cycles in the other hand.

1. A set of edges composed of:
  - edges between relevant cycle vertices, corresponding to the set of edges of the relevant cycle graph,
  - edges of  $G$  which connect two acyclic atoms or connect a single relevant cycle to another single relevant cycle ( $C_1$  and  $C_2$  in figure 1.15) or an acyclic part of  $G$  ( $C_1$  and  $S$  in figure 1.15),
2. and a set of hyperedges which allows to encode special cases where an edge connects at least two distinct relevant cycles to another part of the molecule. For example, an hyperedge is required to encode adjacency relationship between  $C_1$  and  $C_2$  in one hand and an oxygen atom on the other hand in figure 1.15.

Similarly to vertices, each atomic bond of a molecular graph is encoded within its relevant cycle hypergraph. Therefore, a relevant cycle hypergraph encodes all information of a molecular graph.

Considering such a new hypergraph representation, similarity must be computed using an hypergraph kernel. In order to encode local relationships between cycles and substituents, treelet kernel has been adapted [17] to enumerate treelets included in relevant cycle hypergraphs. This kernel defines a similarity based both on an enumeration of relevant cycles of both graphs (pattern  $G_0$  in figure 1.11) and on an enumeration of more complex treelets encoding relationships between cycles and relationships between cyclic and acyclic parts of a molecule.



### 1.5.6 Pattern Weighting

Considering graph kernels based on bags on patterns, it may be interesting to weight each pattern according to its influence for a given property. Similarly to the notion of pharmacophore, some patterns may have a particular influence whereas some do not encode any information about the property to predict. Considering a kernel defined as a sum of sub kernels, where each sub kernel encodes a similarity according to a given pattern  $p \in \mathcal{P}$ , pattern weighting relies on computing a weight for each term of the sum. The weighted kernel is defined as:

$$k_{\text{weight}}(x, x') = \sum_{p \in \mathcal{P}} w(p) * k_p(x, x') \quad (1.38)$$

where  $w : \mathcal{P} \rightarrow \mathbb{R}_+$  encodes the influence of each pattern: an high value  $w(p)$  encodes an high influence of pattern  $p$  and thus corresponds to an high contribution to the similarity measure. Conversely, a  $w(p)$  close to or equals to 0 consists in removing the contribution of pattern  $p$  to the similarity measure.

#### 1.5.6.1 A priori methods

Tree pattern kernel as defined by Mahé [33] (Section 1.5.3.2) includes a tree pattern weighting. This kernel is defined on the same feature space than the original kernel but includes a weighting of each pattern:

**Definition 32 (Weighted tree pattern kernel)** *Let us consider two graphs  $G = (V, E) \in \mathcal{G}$  and  $G' = (V', E') \in \mathcal{G}$  together with the set of tree patterns of size  $h : \mathcal{T}_h$  a counting function of tree patterns:  $\phi_t : \mathcal{G} \rightarrow \mathbb{N}$  and a weighing function:  $w : \mathcal{T}_h \rightarrow \mathbb{R}^+$ . The weighted tree pattern kernel is defined by:*

$$k^h(G, G') = \sum_{t \in \mathcal{T}_h} w(t) \phi_t(G) \phi_t(G'). \quad (1.39)$$

This tree pattern kernel definition corresponds to an adaptation of equation 1.38 to tree pattern kernel. Two weighting functions based on structural information included in tree patterns have been proposed [33].

#### 1.5.6.2 Branching cardinality

The first one is based on branching cardinality. Branching cardinality is defined as the number of leaves  $-1$ . Given a tree pattern  $t = (V, E)$ , branching cardinality is defined by:

$$\text{branch}(t) = \sum_{v \in V} \mathbf{1}(\Gamma(v) = 1) - 1 \quad (1.40)$$

where  $\mathbf{1}(x)$  is equals to 1 if  $x$  is true. Weighting function is then defined by:

$$w_{\text{branch}}(t) = \lambda^{\text{branch}(t)} \quad (1.41)$$

where  $\lambda \in \mathbb{R}_+^*$ . This weighting function aims to favor linear patterns if  $\lambda < 1$ . Conversely, structurally complex patterns will be favored if  $\lambda > 1$ .

### 1.5.6.3 Ratio depth/size

The second weighting function is based on a ratio between the height and size of a tree pattern. For a given tree pattern  $t = (V, E)$  rooted in  $r \in V$ , the size-based weighting function is defined by:

$$w_{size}(t) = \lambda^{|V|-h} \quad (1.42)$$

Similarly to section 1.5.6.2,  $\lambda < 1$  favors tree patterns having similar sizes and heights, which correspond to linear patterns. Conversely,  $\lambda > 1$  favor more structurally complex tree patterns.

These two weighting functions allow to weight tree pattern influence according to their structural information. Parameter  $\lambda$  must be tuned using *a priori* chemical knowledge or cross validation. However, this weighting framework does not allow to weight specifically each tree pattern according to a property to predict.

### 1.5.6.4 Multiple kernel learning

Considering a pattern weighting problem as defined in equation 1.38, Multiple Kernel Learning (MKL) aims to compute an optimal weighting function  $w$  according to a given dataset. Considering a finite set of  $m$  patterns, weighting function  $w$  is encoded by a vector  $d \in \mathbb{R}_+^m$  where each coordinate  $d_i$  encodes a weight associated to a pattern. The weighted kernel is then defined as:

$$k_{MKL}(x, x') = \sum_{i=1}^m d_i * k_i(x, x') \quad (1.43)$$

Multiple kernel learning methods consist in computing an optimal vector  $d$  according to a prediction task.

### 1.5.6.5 Simple MKL

Simple MKL [39] consists in optimizing problems addressed by SVM (section 1.4.4.1) by computing an optimal weighting. Considering a kernel as defined in equation 1.43, simple MKL consists in solving:

$$\underset{d}{\text{minimize}} J(d) \text{ such that } \begin{cases} \sum_{i=1}^m d_i = 1 \\ d_i > 0, \forall i \in \{1, \dots, m\} \end{cases} \quad (1.44)$$

with

$$J(d) = \begin{cases} \min_{w, b, \xi} & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ \text{subject to:} & y_i(\langle w, x \rangle + b) \geq 1 - \xi_i, \forall i \in \{1, \dots, n\} \\ & \xi_i \geq 0, \forall i \in \{1, \dots, n\} \end{cases} \quad (1.45)$$

This minimization problem includes a constraint on the  $L_1$  norm of vector  $d$  which induces sparsity on vector  $d$ . This sparsity allows to keep only most relevant patterns into kernel computation and to remove irrelevant ones. Simple

MKL problem is resolved by alternating a classical SVM resolution together with a projected gradient descent according to vector  $d$ . This projected gradient step allows to minimize objective function while ensuring that constraints on vector  $d$  are fulfilled. Since problem defined in equation 1.44 is convex, iterating these two steps leads to the optimal vector  $d$ .

#### 1.5.6.6 Generalized MKL

In order to allow more flexibility on sparsity constraint, generalized MKL [52] defines an objective function which includes a Tikhonov regularization on vector  $d$  instead of an equality constraint in simple MKL. Corresponding minimization problem is defined as:

$$\underset{d}{\text{minimize}} J(d) \text{ such that } d_i > 0, \forall i \in \{1, \dots, m\} \quad (1.46)$$

with

$$J(d) = \begin{cases} \min_{w,b,\xi} & \frac{1}{2}\|w\|^2 + C \sum_{i=1}^n \xi_i + \sigma \|d\|_1 \\ \text{subject to :} & y_i(\langle w, x \rangle + b) \geq 1 - \xi_i, \forall i \in \{1, \dots, n\} \\ & \xi_i \geq 0, \forall i \in \{1, \dots, n\} \end{cases} \quad (1.47)$$

Parameter  $\sigma$  allows to weight the influence of  $L_1$  norm regularization into the minimization problem. An high  $\sigma$  favors vector  $d$  having a low  $L_1$  norm and thus sparse vectors. Conversely, a low  $\sigma$  relaxes sparsity constraint. The method used to resolve this problem is closely similar to the one proposed to solve simple MKL method. The main difference is that the gradient is no longer projected on simplex since equality constraint no longer exists.

#### 1.5.6.7 Infinite MKL

In order to be able to deal with thousands of patterns, one have to either compute each gram matrix for each sub kernel at each iteration or store each gram matrix. When considering thousands of patterns and graphs, these two options induce too much computational time or memory space to be applicable. In order to be able to handle such datasets, infinite MKL [57] defines a MKL method based on Simple MKL which only considers a sub set of sub kernels at each iteration.

This method only performs optimization using a sub set of active kernels, that is to say kernels having a weight  $d_i$  strictly greater than 0. At the end of each Simple MKL optimization, active kernel set is updated by removing kernels having a null weight and potentially active kernels are added to the set of active kernels. Potentially active kernels may be determined using an oracle based on KKT conditions and gradients associated to each sub kernel [57].

This approach allows to consider an high, possibly infinite, number of sub kernels and allows to select most relevant ones according to a particular property to predict.

**TABLE 1.2:** Boiling point prediction on acyclic molecule dataset using 90% of the dataset as train set and remaining 10% as test set. Execution times are displayed in seconds.

Method	RMSE( $^{\circ}C$ )	Gram matrix (s)	
		learning	prediction
(1) Edit Distance	10.27	1.35	0.05
(2) Graph embedding	10.19	2.74	0.01
(3) Path kernel	12.24	7.83	0.18
(4) Random walks kernel	18.72	19.10	0.57
(5) Tree pattern kernel	11.02	4.98	0.03
(6) Treelet kernel	8.10	0.07	0.01
(7) Treelet kernel with MKL	5.24	70	0.01

## 1.6 Experiments

Different kernels presented in section 1.5 have been tested on several chemoinformatics datasets. These datasets are divided in regression and classification problems<sup>1</sup>

### 1.6.1 Boiling point prediction

The first prediction problem used to test different methods is based on a dataset composed of 185 acyclic molecules [7]. This prediction problem consists in predicting boiling point of molecules. Table 1.2 shows prediction ability of different methods together with execution times. The first line of table 1.2 shows results obtained by a gaussian kernel applied on approximate edit distance (section 1.2.3.1). Note that since graph edit distance does not define an euclidean distance, Gram matrix is regularized in order to be semi-definite positive. The second line corresponds to the embedding method described in section 1.3.1. We can first note that methods based on edit distance obtains intermediate results on this dataset. These results may be due either to the approximation of edit distance or the loss of information induced by regularization or embedding. Next lines correspond to graph kernels defined in section 1.5.3. Line 3 corresponds to a kernel based on paths (section 1.5.2.3) and line 4 corresponds to a random walks kernel (section 1.5.2.1). These two kernels suffer from the low expressiveness of linear patterns and we can note that tottering phenomenon degrades prediction accuracy. Finally the third last lines correspond to kernels based on non linear patterns. As we can see line 7, best results are obtained thanks to the combination of multiple kernel learning (section 1.5.6.4) which allows to only consider relevant patterns.

<sup>1</sup>All these dataset are available on the IAPR TC15 Web page: <https://iapr-tc15.greyc.fr/links.html>

**TABLE 1.3:** Results on AIDS dataset.

Method	Classification Accuracy
(1) KNN using graph edit distance	97.3%
(2) Graph embedding	98.2%
(3) Edit distance	99.7%
(4) Path kernel	98.5%
(5) Random walks kernel	98.5%
(6) Treelet kernel	99.1%
(7) Treelet kernel with MKL	99.7%

Execution times show that explicit enumeration of treelets allows to compute Gram matrix in less time than tree pattern kernel or random walks kernels. Conversely, pattern weighting allows to define a more accurate similarity measure while requiring the highest computational time. Note that prediction time is not altered since weighting is only computed during learning step.

This first experiment highlights the gain obtained by using non linear patterns instead of linear patterns. Moreover, pattern weighting step allows to greatly increase prediction accuracy by removing irrelevant patterns from similarity measure computation. This weighting step allows to get best results among tested methods

### 1.6.2 Classification

The two next experiments correspond to binary classification problems which consist in predicting if a molecule has a particular activity or not.

#### 1.6.2.1 AIDS dataset

First classification experiment has been performed on a graph database provided by [43]. This dataset, defined from the AIDS Antiviral Screen Database of Active Compounds, is composed of 2000 chemical compounds some of them being disconnected. These chemical compounds have been screened as active or inactive against HIV and they are split into three sub sets:

- A train set composed of 250 compounds used to train SVM;
- a validation set composed of 250 compounds used to find parameters set giving the best accuracy result;
- A test set composed of the remaining 1500 compounds used to test the classification model.

This dataset is composed of a large set of different chemical compounds including both cyclic and acyclic molecules and composed of several heteroatoms.

Table 1.3 shows results obtained by different methods on this dataset. First line of table 1.3 corresponds to a classifier defined as a k-nearest neighbors algorithm using approximate graph edit distance. Line 2 corresponds to a graph embedding as described section 1.3.1 and line 3 corresponds to a gaussian kernel applied on graph edit distance and regularized. We can first note that the use of k-nearest neighbors or graph embedding approach does not lead to good results on this dataset. Conversely, gaussian kernel on edit distance (table 1.3, line 4) combined with SVM obtains the best results on this dataset. Note that the regularization added to Gram matrix does not alter prediction accuracy. This may be explained by the fact that classification results may not be altered by a reasonable distortion of the graph edit distance induced by the regularization step. Conversely regression problems which consists in predicting a real value instead of a binary class may be altered by any modification of the initial distance. Next lines correspond to kernels using different bags of patterns. Kernels corresponding to line 4 and line 5 are based on linear patterns (section 1.5.2). As seen in regression experiment, the low expressiveness of kernels based on linear patterns leads to poor results on this dataset. Conversely, methods based on non linear patterns (lines 6 and 7) obtain better results since they encode more structural information than kernels based on linear patterns. Moreover, as observed in regression experiment, multiple kernel learning step allows to reach the best results on this classification problem. Note that structures having biggest weights can be assimilated to pharmacophores. These relevant patterns according to MKL may be analyzed by chemical experts.

This experiment confirms the conclusions drawn in section 1.6.1 on a regression problem with only acyclic molecules.

#### 1.6.2.2 PTC dataset

The second classification experiment is a classification problem taken from the Predictive Toxicity Challenge [50] which aims to predict carcinogenicity of chemical compounds applied to female (F) and male (M) rats (R) and mice (M). This experiment is based on ten different datasets, each of them being composed of one train set and one test set. Table 1.4 shows the number of correctly classified molecules over the ten test sets for each method and for each class of animal.

Table 1.4 shows results obtained by different kernels encoding cyclic information in different ways. The first line of this table corresponds to treelet kernel which does not encode any cyclic information since it is only based on acyclic patterns. Lines 2 to 4 correspond to methods encoding different levels of cyclic information. Line 2 corresponds to cyclic pattern kernel which simply compares common cycles (section 1.5.5.1). Line 3 corresponds to a treelet kernel applied on relevant cycle graph (section 1.5.5.2) which encodes cycle relationships. Line 4 corresponds to results obtained by treelet kernel adapted to relevant cycle hypergraph (section 1.5.5.3) comparison [17]. First, we can

**TABLE 1.4:** Classification accuracy on PTC dataset.

Method	# correct predictions			
	MM	FM	MR	FR
(1) Treelet kernel (TK)	208	205	209	212
(2) Cyclic pattern kernel	209	207	202	228
(3) TK on relevant cycle graph	211	210	203	232
(4) TK on relevant cycle hypergraph (TCH)	217	224	207	233
(5) TK + MKL	217	224	223	250
(6) TC + MKL	216	213	212	237
(7) TCH + MKL	225	229	215	239
(8) TK + $\lambda$ TCH	225	230	224	252

note that best results are obtained by kernel encoding both cyclic and acyclic relationships which validates the relevance of including cyclic information, and more particularly adjacency relationships between cyclic and acyclic parts of a molecule.

Then, lines 5 to 7 correspond to different kernels combined with multiple kernel learning. This weighting step shows that kernel based on cyclic information obtains best results on two datasets over four (line 7, datasets MM and FM) and kernel only based on acyclic patterns obtains best results over the two others datasets (line 5, datasets MR and FR). Note that pattern weighting step allows us to reduce the number of patterns included within kernel computation from about 3500 to 150, depending on dataset. Finally, since different kernels obtain best results over the four datasets, weighted combination of kernel encoding cyclic information and a kernel encoding acyclic information leads to the best results on this dataset (line 8). This combination shows the flexibility of kernel approaches by means of multiple kernel learning and linear combinations of kernels.

## 1.7 Conclusion

Graph kernel framework allows to define scalar products between implicit or explicit vectorial representations of graphs in a given feature space. On one hand, conversely to methods based on graph theory, graph kernels can be used in well known and widely used machine learning methods such as SVM. On the other hand, exemption of an explicit vectorial representation allows to encode more information than methods based on an explicit and fixed size vectorial representation. This characteristic allows to define accurate graph similarity measures which encode most of structural and labeling information encoded by

molecular graphs. Therefore, graph kernels allow to combine efficient machine learning methods with accurate and expressive similarity measures.

Defining a graph kernel consists in defining a graph similarity measure which encodes a maximum of useful information and which fulfills all properties required to define a positive definite kernel. Kernels based on bags of patterns deduce molecular graph similarities from similarities of bags of patterns extracted from these graphs. Kernels based on bags of patterns aim to encode a maximum of information while keeping an efficient computational time in order to be applicable to datasets encountered in chemoinformatics. Kernels based on non linear patterns encode more structural information than kernels based on linear patterns and can be computed in linear time when applied on molecular graphs. Moreover, some bags of patterns are defined such as they explicitly encode cyclic information into similarity measure. This information is particularly useful in chemoinformatics since molecular cycles have a great influence on molecular properties. Among kernels based on bags of patterns, relevant cycle hypergraph encodes both acyclic and cyclic parts and their relationships into a single representation which allows to explicitly encode adjacency relationships between a cycle and its substituents.

Kernel theory allows to define a kernel from a linear combination of sub kernels. Considering kernels based on bags of patterns, each sub kernel may be defined as a kernel encoding a molecular similarity according to a particular pattern. From this point of view, multiple kernel learning methods allow to compute an optimal weight for each sub kernel according to a property to predict. This weight corresponds to the contribution of each pattern to kernel computation and may be understood as a measure of the influence of each pattern. On one hand, this weighting step allows to increase accuracy of prediction models by removing irrelevant patterns from the kernel computation. On the other hand, patterns corresponding to high weights may be seen as pharmacophores. These pharmacophores obtained without *a priori* chemical knowledge may be analyzed by chemical experts to understand some chemical or biological properties.

In conclusion, graph kernels provide an useful framework which may obtain accurate prediction results by combining expressive similarity measures and powerful machine learning methods.





---

## Bibliography

- [1] N. Aronszajn. Theory of reproducing kernels. *Transactions of the American Mathematical Society*, 68(3):337–404, 1950.
- [2] C. Berg, J.P.R. Christensen, and P. Ressel. *Harmonic Analysis on Semi-groups: Theory of Positive Definite and Related Functions*. Applied Mathematical Sciences. Springer-Verlag, 1984.
- [3] C. Berge. *Graphs and hypergraphs*, volume 6. Elsevier, 1976.
- [4] K.M. Borgwardt and H. Kriegel. Shortest-Path kernels on graphs. *Fifth IEEE International Conference on Data Mining (ICDM'05)*, pages 74–81, 2005.
- [5] B.E. Boser, I.M. Guyon, and V.N. Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the fifth annual workshop on Computational learning theory*, pages 144–152. ACM, 1992.
- [6] T. Caelli and S. Kosinov. An eigenspace projection clustering method for inexact graph matching. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26:515–519, 2004.
- [7] D. Cherqaoui, D. Villemin, A. Mesbah, J. M. Cense, and V. Kvasnicka. Use of a Neural Network to Determine the Normal Boiling Points of Acyclic Ethers, Peroxides, Acetals and their Sulfur Analogues. *J. Chem. Soc. Faraday Trans.*, 90:2015–2019, 1994a.
- [8] C. Cortes and V. Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, September 1995.
- [9] J. Dattorro. *Convex optimization & Euclidean distance geometry*. Meboo Publishing USA, 2005.
- [10] M. Deshpande, M. Kuramochi, and G. Karypis. Frequent sub-structure-based approaches for classifying chemical compounds. In *Proceedings of the Third IEEE International Conference on Data Mining, ICDM '03*, pages 35–, Washington, DC, USA, 2003. IEEE Computer Society.
- [11] H. Drucker, C. J. C. Burges, L. Kaufman, A. Smola, and V. Vapnik. Support vector regression machines. *Advances in neural information processing systems*, pages 155–161, 1997.

- [12] A. Ducournau. *Hypergraphes: clustering, réduction et marches aléatoires orientées pour la segmentation d'images et de vidéo*. PhD thesis, École Nationale d'Ingénieurs de Saint-Étienne., 2012.
- [13] R. W. Floyd. Algorithm 97: Shortest path. *Commun. ACM*, 5(6):345, June 1962.
- [14] T. Gärtner, P. A. Flach, and S. Wrobel. On graph kernels: Hardness results and efficient alternatives. In *Proceedings of the 16th Annual Conference on Computational Learning Theory and the 7th Kernel Workshop*, 2003.
- [15] J. Gasteiger, editor. *Handbook of Chemoinformatics*. Wiley-VCH, 1. edition, 2003.
- [16] B. Gaüzère, L. Brun, and D. Villemin. Two New Graphs Kernels in Chemoinformatics. *Pattern Recognition Letters*, 33(15):2038–2047, 2012.
- [17] B. Gaüzère, L. Brun, and D. Villemin. Relevant cycle hypergraph representation for molecules. In W. G. Kropatsch, N. M. Artner, Y. Haxhimusa, and X. Jiang, editors, *Graph-Based Representations in Pattern Recognition*, volume 7877 of *Lecture Notes in Computer Science*, pages 111–120. Springer Berlin Heidelberg, 2013.
- [18] B. Gaüzère, L. Brun, D. Villemin, and M. Brun. Graph kernels based on relevant patterns and cycle information for chemoinformatics. In *Proceedings of ICPR 2012*, pages 1775–1778. IAPR, IEEE, November 2012.
- [19] J. Gibert, E. Valveny, and H. Bunke. Graph embedding in vector spaces by node attribute statistics. *Pattern Recognition*, 45(9):3072–3083, 2012.
- [20] P.E. Hart, N.J. Nilsson, and B. Raphael. A formal basis for the heuristic determination of minimum cost paths. *Systems Science and Cybernetics, IEEE Transactions on*, 4(2):100–107, 1968.
- [21] D. Haussler. Convolution kernels on discrete structures. Technical report, University of California at Santa Cruz, 1999.
- [22] T. Horváth. Cyclic pattern kernels revisited. In Springer-Verlag, editor, *Proceedings of the 9th Pacific-Asia Conference on Knowledge Discovery and Data Mining*, volume 3518, pages 791 – 801, 2005.
- [23] T. Horváth, T. Gärtner, and S. Wrobel. Cyclic pattern kernels for predictive graph mining. In *Proceedings of the 2004 ACM SIGKDD international conference on Knowledge discovery and data mining - KDD '04*, page 158. ACM Press, 2004.
- [24] M. A. Johnson and G. M. Maggiora, editors. *Concepts and Applications of Molecular Similarity*. Wiley, 1990.

- [25] S. Jouili and S. Tabbone. Graph embedding using constant shift embedding. In *Proceedings of the 20th International conference on Recognizing patterns in signals, speech, images, and videos, ICPR'10*, pages 83–92, Berlin, Heidelberg, 2010. Springer-Verlag.
- [26] H. Kashima, K. Tsuda, and A. Inokuchi. Marginalized kernels between labeled graphs. In *International Conference on Machine Learning*, volume 3, pages 321–328, 2003.
- [27] B. Luo, R. C. Wilson, and E. Hancock. A spectral approach to learning structural variations in graphs. *Pattern Recognition*, 39(6):1188–1198, 2006.
- [28] B. Luo, R. C. Wilson, and E. R. Hancock. Spectral embedding of graphs. *Pattern Recognition*, 36(10):2213 – 2230, 2003.
- [29] M. M. Luqman, J.-Y. Ramel, J. Lladós, and T. Brouard. Fuzzy multilevel graph embedding. *Pattern Recognition*, 46(2):551–565, 2013.
- [30] P. Mahé, L. Ralaivola, V. Stoven, and J.-P. Vert. The pharmacophore kernel for virtual screening with support vector machines. *Journal of chemical information and modeling*, 46(5):2003–14, 2006.
- [31] P. Mahé, N. Ueda, T. Akutsu, J.-L. Perret, and J.-P. Vert. Extensions of marginalized graph kernels. In *Twenty-first international conference on Machine learning - ICML '04*, page 70. ACM Press, 2004.
- [32] P. Mahé and J.-P. Vert. Graph kernels based on tree patterns for molecules. *Machine Learning*, 75(1):3–35, October 2008.
- [33] P. Mahé and J.-P. Vert. Graph kernels based on tree patterns for molecules. *Machine Learning*, 75(1)(September 2008):3–35, 2009.
- [34] J. Mercer. Functions of Positive and Negative Type, and Their Connection with the Theory of Integral Equations. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 83(559):69–70, November 1909.
- [35] H. L. Morgan. The generation of a unique machine description for chemical structures-a technique developed at chemical abstracts service. *J. of Chem. Doc.*, 5(2):107–113, 1965.
- [36] J. Munkres. Algorithms for the assignment and transportation problems. *Journal of the Society for Industrial and Applied Mathematics*, 5(1):32–38, 1957.
- [37] M. Neuhaus and H. Bunke. *Bridging the Gap Between Graph Edit Distance and Kernel Machines*. World Scientific Publishing Co., Inc., River Edge, NJ, USA, 2007.

- [38] G. Poezevara, B. Cuissart, and B. Crémilleux. Discovering emerging graph patterns from chemicals. In *Proceedings of the 18th International Symposium on Methodologies for Intelligent Systems (ISMIS 2009)*, pages 45–55, Prague, 2009. LNCS.
- [39] A. Rakotomamonjy, F. Bach, S. Canu, and Y. Grandvalet. SimpleMKL. *Journal of Machine Learning Research*, 9:2491–2521, 2008.
- [40] L. Ralaivola, S. J. Swamidass, H. Saigo, and P. Baldi. Graph kernels for chemical informatics. *Neural networks : the official journal of the International Neural Network Society*, 18(8):1093–1110, October 2005.
- [41] J. Ramon and T. Gärtner. Expressivity versus efficiency of graph kernels. In *First International Workshop on Mining Graphs, Trees and Sequences*, pages 65–74, 2003.
- [42] K. Riesen. *Classification and Clustering of Vector Space Embedded Graphs*. PhD thesis, Institut für Informatik und angewandte Mathematik Universität Bern, 2009.
- [43] K. Riesen and H. Bunke. Iam graph database repository for graph based pattern recognition and machine learning. In *Proceedings of the 2008 Joint IAPR International Workshop on Structural, Syntactic, and Statistical Pattern Recognition, SSPR & SPR '08*, pages 287–297, Berlin, Heidelberg, 2008. Springer-Verlag.
- [44] K. Riesen and H. Bunke. Approximate graph edit distance computation by means of bipartite graph matching. *Image and Vision Computing*, 27(7):950–959, 2009.
- [45] N. Shervashidze and K. M. Borgwardt. Fast subtree kernels on graphs. In *Advances in Neural Information Processing Systems*, pages 1660–1668, 2009.
- [46] N. Shervashidze, T. Petri, K. Mehlhorn, K. M. Borgwardt, and S. Viswanathan. Efficient graphlet kernels for large graph comparison. In *International Conference on Artificial Intelligence and Statistics*, pages 488–495, 2009.
- [47] N. Shervaszide. *Scalable Graph Kernels*. PhD thesis, Universität Tübingen, 2012.
- [48] N. Sidere, P. Héroux, and J.-Y. Ramel. Vector representation of graphs: Application to the classification of symbols and letters. In *ICDAR*, pages 681–685. IEEE Computer Society, 2009.
- [49] R. Todeschini and V. Consonni. *Molecular Descriptors for Chemoinformatics*, volume 41. WILEY-VCH, Weinheim (Allemagne), 2009.

- [50] H. Toivonen, A. Srinivasan, R. D. King, S. Kramer, and C. Helma. Statistical evaluation of the predictive toxicology challenge 2000-2001. *Bioinformatics*, 19(10):1183–1193, 2003.
- [51] V. Vapnik. *The Nature of Statistical Learning Theory*. Springer, 1995.
- [52] M. Varma and D. Ray. Learning the discriminative power-invariance trade-off. In *IEEE 11th International Conference on Computer Vision, 2007 (ICCV 2007)*, pages 1–8. IEEE, 2007.
- [53] S.V.N. Vishwanathan, K.M. Borgwardt, and N.N. Schraudolph. Fast computation of graph kernels. *Advances in Neural Information Processing Systems*, 19:1449, 2007.
- [54] P. Vismara. Union of all the minimum cycle bases of a graph. *The Electronic Journal of Combinatorics*, 4(1):73–87, 1997.
- [55] N. Wale, I. A. Watson, and G. Karypis. Comparison of descriptor spaces for chemical compound retrieval and classification. *Knowledge and Information Systems*, 14(3):347–375, 2008.
- [56] X. Yan and J. Han. gspan: Graph-based substructure pattern mining. In *Proceedings of the 2002 IEEE International Conference on Data Mining, ICDM '02*, pages 721–, Washington, DC, USA, 2002. IEEE Computer Society.
- [57] F. Yger and A. Rakotomamonjy. Wavelet kernel learning. *Pattern Recognition*, 44(10):2614–2629, 2011.