

Symbolic Learning vs. Graph Kernels: An Experimental Comparison in a Chemical Application

Luc Brun¹, Donatello Conte³, Pasquale Foggia³, Mario Vento³, and Didier Villemin²

¹ GREYC UMR CNRS 6072,

² LCMT UMR CNRS 6507

ENSICAEN-Université de Caen Basse-Normandie,
14050 Caen, France

`luc.brun@greyc.ensicaen.fr`, `didier.villemin@ensicaen.fr`

³ Dipartimento di Ingegneria dell'Informazione e di Ingegneria Elettrica,
Università di Salerno, Via Ponte Don Melillo, 1 I-84084 Fisciano (SA), Italy
`{dconte, pfoggia, mvento}@unisa.it`

Abstract. In this paper we present a quantitative comparison between two approaches, Graph Kernels and Symbolic Learning, within a classification scheme. The experimental case-study is the predictive toxicology evaluation, that is the inference of the toxic characteristics of chemical compounds from their structure. The results demonstrate that both approaches are comparables in terms of accuracy, but present pros and cons that are discussed in the last part of the paper.

1 Introduction

Many scientific areas are populated by applications dealing with structured information, i.e., complex information which can be seen as made of simpler parts suitably interconnected. Structured information is widely used in many applicative domains as software engineering, pattern recognition, chemistry, medicine, linguistics, etc. Usually structured information is represented by means of graphs because of their high expressive power.

Despite their attractiveness in terms of representational power, structural methods (i.e., methods dealing with structured information) imply complex procedures both in the recognition and in the learning processes. The difficulty of defining effective algorithms for facing this task is so high that the problem is considered still open.

In this paper we focus our attention on two different approaches using graphs for classification tasks. The first one is based on the idea of facing the learning problem directly in the representation space of the structured data (we will call it Symbolic Learning). The second is based on the notion of Graph Kernels that gives a measure of similarity between graphs which corresponds to a metric thus making possible the adoption of well-known statistical/neural paradigms. Here

we present several Graph Kernels based on the notion of bag of patterns and graph edit distance that both preserve most of the structural information of graphs.

The chosen application domain is the Predictive Toxicology Evaluation, that is the characterization of the toxic nature of chemical compounds from their chemical structure. Namely, we used data from the Predictive Toxicology Challenge⁴, a competition based on the prediction of cancerogenic molecules. This domain is appropriate for the analysis because of the structured nature of chemical compounds.

The paper is organized as follows: in section 2 a brief survey of the two approaches is presented, while details on the two methods will be given in the section 3; section 4 reports the results of an experimental case-study of both the approaches. Finally, in the last section, conclusions are outlined.

2 Related Works

The Symbolic Learning can be summarized as follows: the goal of the system is to build a structural description for each class (prototypes) so that the classification is conducted by mean of a structural matching between the prototypes and the structured representation of the object to classify.

Papers dealing with this approach can be ascribed to two rather different categories. The first category collects methods based on the assumption that the prototypical descriptions are built by interacting, more or less deeply, with an expert of the domain. To this concern, some authors [16] start from some ideal prototypes and refine manually, in successive approximations, a model of the possible deformations. Others [13] assume a deformational model and use a small training set for tuning the parameters contained in it. Despite the advantage of facing the problem without a training set or at least with a small training set, this approach is effective only in simple problems. The inadequacy of human knowledge to find a set of prototypes really representative of a given class significantly increases the risk of errors, especially in domains containing either many data or many different classes.

The methods ascribed to the second category face the problem in a formal way ([10]) by considering it as a symbolic machine learning problem, so formulated: "Given a suitably chosen set of input data, whose class is known and possibly some background domain knowledge, find out a set of optimal prototypical descriptions for each class". The literature in the field is populated by methods ([11, 14]) which mainly differ in the adopted formalism sometimes more general than that implied by the graphs.

Graph kernels methods are based on an implicit embedding of graphs within a vector space of large dimension. Several graph kernels are introduced in literature but graph kernels built on the notion of bag of patterns seem to be very promising. The design of such kernels, implies the following choices or heuristics:

⁴ <http://www.predictive-toxicology.org>

1. Selection of a bag of patterns from a graph,
2. Definition of a kernel between bags based on a kernel between patterns,
3. Definition of a kernel between patterns.

Kernels between graphs then correspond to kernels between bags of patterns. The implicit assumption being here that bags of patterns capture most of the structural information of graphs. Several kernels have been proposed based on different types of patterns: Vert [9] proposed to compare the set of sub-trees of two graphs; Graphlets kernels [17] are based on the number of common sub-graphs of two graphs. However, the comparison of complex patterns requires important processing times for an insight which not always clearly demonstrated [9]. As a consequence most of graph kernels are based on simpler patterns such as walks [8], trails [4] or paths. Bags of patterns may be finite or infinite. Infinite bags, such as bag of walks are computed implicitly using random walks defined on both graphs being compared [8]. Bags of patterns may also be computed explicitly [4] by enumerating all patterns of a given graph. Given two bags of patterns the similarity between both bags is often measured using convolution kernels [6]. Both bags may alternatively be considered as two sets in the feature space induced by the kernel between patterns. This point of view allows to measure the similarity between bags using general kernels between sets [3]. Finally, the similarity between patterns should be designed according both to the feature available on each element of the pattern and to the structural noise which may corrupt patterns. Using walks on labeled graphs the well known kernel [8] being equal to 1 if both walks are similar, an 0 otherwise, encodes the fact that labels are either equal or not comparable. The construction of a bag of patterns from a graph removes connections between elementary patterns within input graphs which are then only represented by their bags. As previously mentioned, assumption of bag of pattern kernels is that most of the relevant information about graphs is contained in the pattern and thus that most of the information is preserved by the transformation from a graph to a bag. Bunke [1] proposed several kernels based on the graph edit distance which do not rely on such an assumption.

3 Methods Description

3.1 Symbolic Learning

The Symbolic Learning Method presented in this section (hereafter called SLM) is inspired by Quinlan’s FOIL algorithm particularizing it for Attributed Relational Graph (ARG) representation. An extension of this representation has been presented for describing prototypes of Attributed Relational Graphs; these graphs, called Generalized Attributed Relational Graphs (GARGs), contain generalized nodes, edges and attributes. Then, a learning algorithm has been formulated to build such prototypes by means of a set of operations directly defined on graphs. The algorithm preserves the generality of the prototypes generated by the classical machine learning algorithms and moreover, similarly to most

machine learning systems ([10]), the prototypes obtained by this system are consistent, i.e. each prototype covers samples of a same class.

An ARG can be defined as a 6-tuple $(N, E, A_N, A_E, a_N, a_E)$, where N and $E \subset N \times N$ are, respectively, the sets of nodes and edges of the ARG, A_N and A_E the sets of nodes and edge attributes and, finally, a_N and a_E the functions which associate to each node or edge of the graph the corresponding attributes.

A GARG is used for representing a prototype of a set of ARGs. In order to allow a GARG to match a set of possibly different ARGs, in the SLM the attribute definition has been extended. First of all, the set of types of node and edge attributes is extended with the special type Φ , carrying no parameter and allowed to match any attribute type, ignoring the attribute parameters. We say that a GARG $G^* = (N, E, A_N, A_E, a_N, a_E)$ covers a sample G and we use the notation $G^* \triangleright G$ (the symbol \triangleright denotes the relation called covering) if there is a mapping $\mu : N^* \rightarrow N$ such that μ is a *monomorphism* and the attributes of the nodes and of the edges of G^* are *compatible* with the corresponding ones of G . The first condition requires that each primitive and each relation in the prototype is present also in the sample; note that the converse condition does not hold. The latter condition constrains the monomorphism to be consistent with the attributes of the prototype and of the sample in the sense that the type of the attribute of the prototype must be either equal to the special type Φ or to the type of the corresponding attribute of the sample. In the latter case, all the parameters of the attribute, which are actually sets of values, must contain the value of the corresponding parameter of the sample. The goal of the learning algorithm can be stated as follows: there is a (possibly infinite) set S^* of all the patterns that may occur, partitioned into C different classes S_1^*, \dots, S_C^* , with $S_i^* \cap S_j^* = \phi$ for $i \neq j$; to the algorithm is given a finite subset $S \subset S^*$ (training set) of labeled patterns ($S = S_1 \cup \dots \cup S_C$ with $S_i = S \cap S_i^*$), from which it tries to find a sequence of prototype graphs $G_1^*, G_2^*, \dots, G_p^*$, each labeled with a class identifier, such that:

$$\forall G \in S^* \exists i : G_i^* \triangleright G (\text{completeness of the prototype set}) \quad (1)$$

$$\forall G \in S^*, G_i^* \triangleright G \Rightarrow \text{class}(G) = \text{class}(G_i^*) (\text{consistency of the prototype set}) \quad (2)$$

where $\text{class}(G)$ and $\text{class}(G^*)$ refer to the class associated with sample G and prototype G^* , respectively. Of course, this is an ideal goal since only a finite subset of S^* is available to the algorithm; in practice, the algorithm can only demonstrate that completeness and consistency hold for the samples in S . On the other hand, (1) dictates that, in order to get as close as possible to the ideal case, the prototypes generated should be able to model samples also not found in S . However, they should not be too general otherwise (2) will not be satisfied.

A description of the learning algorithm is presented in the following: the algorithm starts with an empty list L of prototypes and tries to cover the training set by successively adding consistent prototypes. When a new prototype is found, the samples covered by it are eliminated and the process continues on the remaining samples of the training set. Then a sample is compared sequentially against the prototypes in the same order in which they have been generated, and

it is attributed to the class of the first prototype that covers it. In this way, each prototype implicitly entails the condition that the sample is not covered by any previous prototype. The algorithm fails if no consistent prototype covering the remaining samples can be found. It is worth noting that the test of consistency in the algorithm actually checks whether the prototype is almost consistent, i.e. almost all the samples covered by G belongs to the same class:

$$\text{Consistent}(G^*) \Leftrightarrow \max_i (|S_i(G^*)| / |S(G^*)|) \geq \theta \quad (3)$$

where $S(G^*)$ denotes the sets of all the samples of the training set covered by a prototype G^* , and $S_i(G^*)$ the samples of the class i covered by G^* and θ is a threshold close to 1. For further details about the representation and the algorithm you can refer to [5].

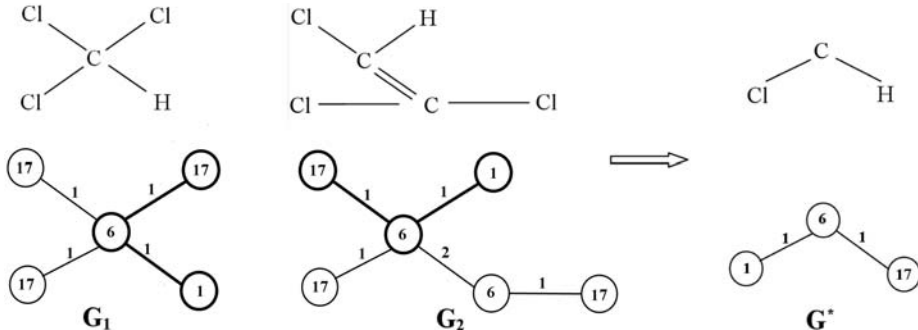


Fig. 1. An example of GARG covering two graphs. The example is extracted from the application domain used in the experimental phase. In bold the parts of the two graphs covered by G^* .

An heuristic function H is introduced for evaluating how promising is the provisional prototype. The heuristic function used in the algorithm is given by the product of H_{compl} and H_{cons} :

$$H = H_{compl}(S, G^*) \cdot H_{cons}(S, G^*) \quad (4)$$

and it is based on estimation of consistency and completeness of prototype. The use of H_{cons} will drive the algorithm towards consistent prototypes while the term H_{compl} is introduced in order to privilege general prototypes with respect to prototypes which, albeit consistent, cover only a small number of samples.

See Figure 1 for an example of a GARG G^* that covers the graphs G_1 and G_2 representing two chemical compounds.

3.2 Graph Kernels

In this section we will introduce two kind of Graph Kernels: Bag of trails Kernel and Laplacian Kernel.

Bag of trails Kernel Let us consider a *graph* $G = (V, E)$ where V denotes the set of vertices and $E \subset V \times V$ the set of edges. We define a *simple-graph* as a graph with no multiple edges between two nodes and no loop (an edge linking a node with itself). We define a *trail* as an alternating sequence of nodes and edges with distinct edges and a *path* as a trail with distinct nodes. A *bag of trails* T associated to G is defined as a set of trails of G whose cardinality is denoted by $|T|$. The kernel between trails is denoted K_{trail} .

Suard [20] has proposed numerous kernels for bags of paths which can be easily extended to bags of trails. We proposed to extend the mean kernel to bag of trails. By considering the bag of trails as set and a trail kernel K_{trail} , we can design a convolution kernel called *mean kernel* as follows:

$$K_{mean}(T_1, T_2) = \frac{1}{|T_1|} \frac{1}{|T_2|} \sum_{t \in T_1} \sum_{t' \in T_2} K_{trail}(t, t'). \quad (5)$$

where T_1 and T_2 denote two bags of trails.

Following Haussler [6], this kernel is positive definite on the bag of trails domain if and only if K_{trail} is positive definite on the trail domain.

However, the mean kernel compares all the paths of both bags without considering their relevance in both bags. As a consequence, irrelevant paths may corrupt the overall value of the kernel due to the "meaning" effect. Recently Dupé et al. [4], proposed to limit this effect by using a relevance measure of trails. Let T_1 and T_2 two bags of trails, the *weighted mean kernel* between these two bags is then defined as:

$$K_{weighted}(T_1, T_2) = \frac{1}{|T_1|} \frac{1}{|T_2|} \sum_{t \in T_1} \sum_{t' \in T_2} \langle R_{T_1}(t), R_{T_2}(t') \rangle^d K_{trail}(t, t'). \quad (6)$$

where $d \in \mathbb{R}^+$ allows to weight the importance of the relevance measure and $R_T(t)$ denote the relevance of trail t according to the bag T . This relevance measure is defined as the similarity (the kernel) between the input trail and the mean trail of the bag ($R_T(t) = \frac{1}{|T|} \sum_{t' \in T} K_{trail}(t, t')$).

Kernel defined by equation 6 is a convolution kernel and so is positive definite if K_{trail} is positive definite.

The above bag of trail kernels are based on a trail kernel K_{trail} . A kernel between two trails $t = (v_1, e_1, \dots, v_n)$ and $t' = (v'_1, e'_1, \dots, v'_p)$ can be derived from the path kernel proposed by Kashima [8]. This kernel is defined as 0 if both trails have not the same size and as follows otherwise:

$$K_{trail}(t, t') = K_v(\varphi(v_1), \varphi(v'_1)) \prod_{i=2}^{|t|} K_e(\psi(e_{v_{i-1}v_i}), \psi(e_{v'_{i-1}v'_i})) K_v(\varphi(v_i), \varphi(v'_i)), \quad (7)$$

where $\varphi(v)$ and $\psi(e)$ denote respectively the vectors of features associated to the node v and the edge e . The terms K_v and K_e denote two kernels for respectively node's and edge's features. Using labeled graph K_v and K_e are defined as Dirac functions comparing the labels of respectively vertices and edges. Since K_{trail} is a tensor product kernel, it is positive definite since K_e and K_v are positive definite kernels.

Laplacian Kernel The graph edit distance between two graphs is defined as the minimal number of vertices and edge removal/addition/relabeling required to transform one graph into an other [12]. The computational complexity of the exact computation of the edit distance is exponential in the number of vertices. Fortunately, Riesen and Bunke [15] proposed an heuristic to efficiently approximate the edit distance between two graphs.

Unfortunately, the edit distance does not define a metric and, as a consequence, kernels directly based on the edit distance does not define a valid kernel.

Let us consider a set of input graphs $\{G_1, \dots, G_n\}$ defining our graph test database. Given a kernel k , the gram matrix K associated to this database is an $n \times n$ matrix defined by $K_{i,j} = k(G_i, G_j)$. As denoted by Steinke [19], the inverse of K (or its pseudo inverse if K is not invertible) may be considered as a regularisation operator on the set of vector of dimension n . Conversely, the inverse (or pseudo inverse) of any definite positive regularization operator may be considered as a kernel. More precisely, within an interpolation or classification scheme, the optimal mapping of a real value f^* to each graph $\{G_1, \dots, G_n\}$ which interpolates a vector of know values y while minimising the norm induced by K in \mathbb{R}^n is achieved by solving the following minimisation problem:

$$f^* = \arg \min_{f \in \mathbb{R}^n} CLoss(f, y, K) + f^t K^{-1} f$$

where $CLoss(\cdot, \cdot, \cdot)$ denotes a given loss function encoding the distance between vector f and the vector of known values y .

One scheme to design a kernel consists thus to first build a definite positive regularization operator and then to take its inverse (or its pseudo inverse) to obtain a kernel. Having an interpolation problem in mind, we would like that our regularization operator penalizes different values of f^* mapped onto graphs with a small edit distance. We thus consider the Laplacian operator defined as follows: given the set of graphs, $\{G_1, \dots, G_n\}$, we first consider the $n \times n$ adjacency matrix $W_{i,j} = e^{-\frac{d(G_i, G_j)}{\sigma}}$ where σ is a tuning variable and $d(\cdot, \cdot)$ denotes the edit distance [15]. The normalized Laplacian of $\{G_1, \dots, G_n\}$ is then defined as $L = I - D^{-\frac{1}{2}} W D^{-\frac{1}{2}}$ where D is a diagonal matrix defined by $D_{i,i} = \sum_{j=1}^n W_{i,j}$.

We know results from spectral graph theory [2], establish that L is a symmetric, semi definite positive matrix whose eigenvalues belongs to $[0, 2]$. Unfortunately, the Laplacian is also well know for being non invertible. The only semi definite property of the Laplacian matrix forbids a direct inversion of this matrix. Moreover, the pseudo inverse of the Laplacian induces numerical instabilities which does not lead to a reliable kernel. Therefore, following Smola [18], we rather regularize the spectrum of L . The regularized version of L , denoted as \tilde{L} , is defined as :

$$\tilde{L} = \sum_{i=1}^n r(\lambda_i) u_i u_i^t \quad (8)$$

where $\{\lambda_i, u_i\}$ denotes the eigensystem of L while $r(\cdot)$ is a regularization function. Many regularization functions may be used, however we restricted our choice to Regularized Laplacian ([18]): $r(\lambda) = 1 + \eta\lambda$ where η is a tuning variable.

The regularized laplacian \tilde{L} is invertible and its inverse $K = \tilde{L}^{-1}$ is taken as a kernel.

4 Experimental Evaluation

The application used for the comparison between the two approaches is the Predictive Toxicology Challenge.

4.1 The Application Domain and the Test Sets

Nowadays, almost every human being in industrialized societies is exposed to chemical compounds whose long-term effects are not evident. The carcinogenicity of such compounds, for instance, is often unknown. The US National Toxicology Program (NTP) ⁵ conducted several experimentations on the effects of many chemical compounds on mice and rats, but such empirical results are expensive and very slow to obtain. Hence the necessity to find carcinogenicity models able to generate reliable toxicity predictions for other compounds. The Predictive

Dataset		FM	MM	FR	MR
Training Set	Cancerogenic	78	67	62	70
	Non-Cancerogenic	124	123	140	124
Test Set		78	64	98	105

Table 1. Number of Chemical Compounds in the sets used for the experimentation.

Toxicology Challenge ⁶ was initiated in this aim. It enabled participants to submit machine learning programs which, from the large training datasets of the NTP, provide carcinogenicity predictions of compounds in test sets. According to the Structure Activity Relationship (SAR) hypothesis, which assume that the activity of a molecule can be entirely determined by its structure, the only informations used for toxicity prediction of a molecule are related to its chemical structure.

Four distinct training and test sets were provided (see Table 1), since mice and rats reactions to carcinogens compounds are significantly different. In this paper we consider all the provided datasets using, for both the methods, the following representation: each chemical compound is represented by a graph $G(V, E, \alpha, \beta)$ where V represents the set of atoms, E is the set of edges representing bonds, α is the nodes labelling function that associates an atom symbol to each node of the graph and β is the edges labelling function that associates a bond multiplicity to each edge of the graph.

⁵ National Toxicology Program: <http://ntp.niehs.nih.gov/>

⁶ <http://www.predictive-toxicology.org/>

4.2 Evaluation and Comparison

In Table 2 we report the comparative results obtained on FM, MM, FR and MR datasets by the two different approaches.

The three methods (two methods from the Graph Kernels) have comparable results. However the time required by Symbolic Learning method is very high.

We can conclude that in problems that have stricts time constraints Graph Kernels method is very effective while Symbolic Learning can be used only on smaller problems, or where there are less time constraints. On the other hand Symbolic Learning generates explicit, declarative knowledge on the classes in the form of GARGs that may be understood with little effort by a human expert: a direct inspection of such explicit prototypes can be sufficient to validate them or possibly to delete the badly formed prototypes that could be detrimental for classification performance. The Graph Kernel approach, instead, is able to perform a classification, but not to explain the criteria that guided its choice.

Method	FM	MM	FR	MR	Elapsed Time
Symbolic Learning	60.4%	52.1%	62.0%	62.0%	~ 2000 min
Bag of trails Kernel	61.2%	50.4%	62.8%	65.6%	~ 2 min
Laplacian Kernel	59.2%	55.2%	57.7%	60.9%	~ 1 min 30 sec

Table 2. Accuracy and Elapsed Time for the two methods on the FM, MM, FR and MR dataset.

In conclusion the results confirm that there is no way to decide what is the best method because both approaches have their peculiar virtues and deficiencies.

In our view, the two approaches should not be seen as competitors; instead, only their cooperative integration can provide us with more powerful tools for knowledge exploitation. The first step toward the integration of the two methodologies could be their combined use in a multiexpert system [7]. As widely demonstrated in the literature, an accurate selection of the combining criteria allows us to significantly improve the results obtained by the single experts by retaining their strengths, while overcoming their weaknesses.

5 Conclusions

In this paper we presented a comparison between two approaches, Graph Kernels and Symbolic Learning, within a classification scheme. Our experiments indicate that both approaches are comparables in terms of accuracy. Graph kernel methods require less execution times but does not provide an explanation of the classification criteria. Symbolic Learning, on the other hand have the reverse advantages and drawbacks. Both methods are thus clearly complementary.

Future steps, as we discussed in the last section, could be an integration of the two approaches in order to improve the results obtained by the application of the single approaches.

Acknowledgments. We would like to express our thanks to Alice Kijewski and David Lemaresquier for their support in the experimental phase.

References

1. Bunke, H., Riesen, K.: A family of novel graph kernels for structural pattern recognition. In: XII Iberoamerican Congress on Pattern Recognition. pp. 20–31 (2007)
2. Chung, F.R.K.: Spectral Graph Theory (CBMS Regional Conference Series in Mathematics, No. 92). American Mathematical Society (February 1997)
3. Dupé, F.X., Brun, L.: Hierarchical bag of paths for kernel based shape classification. In: SSPR. pp. 227–236 (2008)
4. Dupé, F.X., Brun, L.: Tree covering within a graph kernel framework for shape classification. In: XV ICIAP (2009)
5. Foggia, P., Genna, R., Vento, M.: Symbolic vs connectionist learning: an experimental comparison in a structured domain. IEEE Transaction on Knowledge and Data Engineering 13-2, 176–195 (2001)
6. Haussler, D.: Convolution kernels on discrete structures. Tech. rep., Department of Computer Science, University of California at Santa Cruz (1999)
7. Iannello, G., Percannella, G., Sansone, C., Soda, P.: On the use of classification reliability for improving performance of the one-per-class decomposition method. Data & Knowledge Engineering 68, 1398–1410 (2009)
8. Kashima, H., Tsuda, K., Inokuchi, A.: Marginalized kernel between labeled graphs. In: In Proc. of the Twentieth International conference on Machine Learning (2003)
9. Mahé, P., Vert, J.P.: Graph kernels based on tree patterns for molecules. Machine Learning 75(1), 3–35 (2008)
10. Michalski, R.: Pattern recognition as rule-guided inductive inference. IEEE Transaction on PAMI 2-4, 349–361 (1980)
11. Muggleton, S.: Inductive logic programming. New Generation Computing 8-4, 295–318 (1991)
12. Neuhaus, M., Bunke, H.: Bridging the Gap Between Graph Edit Distance and Kernel Machines. World Scientific Publishing Co., Inc., River Edge, NJ, USA (2007)
13. Nishida, H.: Shape recognition by integrating structural descriptions and geometrical/statistical transforms. CVIU 64, 248–262 (1996)
14. Quinlan, J.R.: Learning logical definitions from relations. ML 5-3, 239–266 (1990)
15. Riesen, K., Bunke, H.: Approximate graph edit distance computation by means of bipartite graph matching. Image Vision Computing 27(7), 950–959 (2009)
16. Rocha, J., Pavlidis, T.: A shape analysis model with applications to a character recognition system. IEEE Transaction on PAMI 16-4, 393–404 (1994)
17. Shervashidze, N., Vishwanathan, S.V., Petri, T.H., Mehlhorn, K., Borgwardt, K.M.: Efficient graphlet kernels for large graph comparison. In: Twelfth International Conference on Artificial Intelligence and Statistics (2009)
18. Smola, A.J., Kondor, R.I.: Kernels and regularization on graphs. In: XV Annual Conference on Learning Theory. pp. 144–158 (2003)
19. Steinke, F., Schölkopf, B.: Kernels, regularization and differential equations. Pattern Recognition 41(11), 3271 – 3286 (2008)
20. Suard, F., Rakotomamonjy, A., Benschrair, A.: Kernel on bag of paths for measuring similarity of shapes. In: European Symposium on Artificial Neural Networks. Bruges-Belgique (April 2007)