

Structured Kernels for Automatic Detection of Protein Active Sites

Elisa Cilia¹, Alessandro Moschitti¹, Sergio Ammendola², and Roberto Basili¹

¹ Department of Computer Science, System and Production
University of Rome, Tor Vergata,
Via Della Ricerca Scientifica s.n.c., 00133, Roma, ITALY
`elisa.cilia@gmail.com`
`{moschitti,basili}@info.uniroma2.it`

² Ambiotec sas
SME
field: biotechnology
Via F. Acqua Mariana 125, 00040, Roma, ITALY
PI 07635910636
`sergio.ammendola@fastwebnet.it`

Abstract. In this paper, we design novel models based on Support Vector Machines and Kernel Methods for the automatic protein active site classification. We devise innovative attribute-value and tree substructure representations derived from biological and spatial information of proteins. We experimented such models with the Protein Data Bank adequately pre-processed to make explicit the active site information. Our results show that structural kernels used in combination with polynomial kernels can be effectively applied to discriminate an active site from other regions of a protein. Such finding is very important since it firstly shows the successful identification of catalytic sites of a very large family of catalytic proteins belonging to a broad classes of enzymes.

1 Introduction

Recent research in Bioinformatics has been devoted to the production and understanding of genomic data. One important step in this direction is the study of the relation between molecular structures and their functions, which in turn depends on the discovering of the protein active sites. As there is a large number of synthesized proteins which have no associated function yet, i.e. whose function remains unknown, automatic approaches for active site detection are critical.

Currently, the general strategy used to identify a protein active site involves the expertise of researchers and biologists accumulated in years of study on the target protein as for example in [1]. This manual approach is conducted essentially using homology based strategies, i.e. inferring the function of a new protein based on a close similarity to already annotated proteins. Sometimes proteins with the same overall tertiary structure can have different active sites, i.e. different functions and proteins with different overall tertiary structure can

show the same function and similar active sites. In these cases homology based approaches are inadequate. Moreover, there is no general automated approach to protein active site detection, although it is evident its usefulness to restrict the number of candidate sites and also to automatically learn rules characterizing an active site [2].

In this paper we define the problem of determining protein active sites in terms of a classification problem. We modeled protein active site based on both attribute/value and structural representations. The former representation is a set of standard linear features whereas the latter is constituted by tree structures extracted from graphs associated with proteins or their candidate sites. The graph nodes (or vertexes) represent amino acids (or better residues) and edges represent distances in the three-dimensional space between these residues.

We applied these representations to SVMs using polynomial kernels, tree kernels and some combinations of them. To experimentally evaluate our approach, we created a data set, using the protein structures retrieved from the *Protein Data Bank* (PDB) [3] maintained by the *Research Collaboratory for Structural Bioinformatics* (RCBS) at <http://www.rcbs.org>. The combined kernels show the highest F_1 measure, i.e. 68%, in the detection of active sites. This is an important and promising result considering that the baseline based on a random selection of active sites has an upperbound of only 2%.

In the remainder of this article Section 2 describes the faced problem. Section 3 describes the proposed linear and structural features. Section 4 describes the experimental evaluation and reports the results of the classification experiments. Finally, in Section 5, we summarize the results of the previous sections and propose other interesting future research lines.

2 Protein active site classification

An active site in a protein is a topological region which defines the protein function, in other words it is a functional domain in the protein three-dimensional structure (see also [2]). In a cell there are many types of proteins which carry out different functions. The enzymes are those proteins able to accelerate chemical processes inside a cell. This type of proteins are distinguished from structural and supplying proteins for their catalytical action on the large part of molecules constituting the living world. We limit our research to a particular class of enzymes, the hydrolases.

Hydrolases are maybe the most studied and known type of enzymes. They catalyze hydrolysis reactions, generically consisting in the cleavage of a biochemical compound thanks to the addition of a water molecule (H_2O). The characteristic of some hydrolases to catalyze reactions in the presence of a water molecule motivates our model: as an hydrolase active site, we choose a sphere in a three-dimensional space centered in the coordinates of the oxygen atom of a water molecule. This sphere includes a portion of the protein within its volume, that is a number of amino acids which could reciprocally interact with other amino acids in the surrounding space, or with water molecules. In this first analysis, we

consider a sphere with a radius of 8 \AA , which is the maximum distance needed for the water-residue interaction.

Figure 1(a), shows the active site of 1A2O protein structure and its representation according to our model. The protein residues are colored in light gray whereas the particular catalytic residues are in dark gray. The center of the sphere is the black colored oxygen atom of a water molecule.

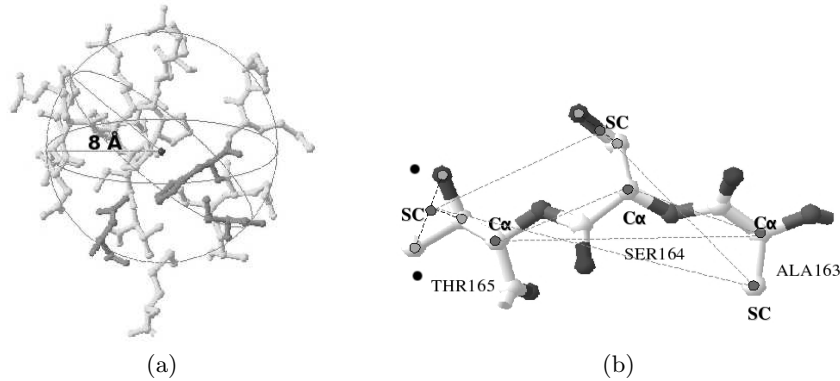


Fig. 1. (a) A sphere (positive example). (b) Distances.

2.1 The computational model

As stated in the previous section, we defined the functional site identification as a classification problem, where the objects we want to classify are protein active sites. We represent the portion of the protein contained in a spherical three-dimensional region with a completely connected graph. Each vertex of this graph is a residue and each edge represents the distance in the three-dimensional space between a pair of vertices.

Every amino acid is represented by two points in the three-dimensional space: one which represents an amino acid main-chain (the α -carbon atom of the amino acid, $C\alpha$) and one which represents an amino acid side-chain (the centroid between the coordinates of the atoms belonging to the amino acid side-chain, SC) (see Figure 1(b)). The same kind of approximation has been described in [4] because it seems to provide a good balance between fuzziness and specificity in these kind of applications.

In Figure 1(b) the three-dimensional SC-SC distances and $C\alpha$ - $C\alpha$ distances are indicated between the represented chain of three residues. An object (modeled by a graph centered on a water molecule) can be classified as being an active site or not with a binary classifier. Thus, we consider as a positive example, a graph whose set of vertex includes all the catalytic amino acids and as a negative example a graph which contain no catalytic amino acid. Moreover, to reduce the

task complexity, we extract, from the initial completely connected graph, some spanning substructures which preserve the edges within the maximum interaction distance of 5 Å between the side-chains of the residues.

The next section shows how the above representation model can be used along with Support Vector Machines to design an automatic active site classifier.

3 Automatic classification of active sites

Previous section has shown that the active site representation is based on graphs. To design the computational model of these latter, we have two possibilities: (1) we extract scalar features able to capture the most important properties of the graph and (2) we can use graph based kernels [5] in kernel-based machines such as Support Vector Machines [6]. Point (2) often leads to high computational complexity. We approached such problem by extracting a tree forest from the target graph and applying efficient tree kernels [7].

3.1 Scalar features

Scalar features refer to typical chemical values of the molecules described in the target graph. We defined 5 different types of such features (see Table 1):

The first class of linear features (C1) encodes chemical and physical properties of the graph. This class represents properties such as hydrophobicity, polarity, polarizability and Van der Waals volume of the amino acids composing the sphere. The encoding is the same used in [8] where the features were used to classify the function of proteins.

The second class of linear features (C2) encodes the amino acid composition of a spherical region. There is a feature associated with every labeled vertex (amino acid) in the graph, weighted with the inverse of the distance from the oxygen atom of the water molecule which is the center of the sphere. This group of features emphasizes the importance of the interaction distance of a residue with respect to a water molecule.

The third class of features (C3) represents charge or neutrality of a spherical region. This is measured by counting the number of positively or negatively charged amino acids.

Another group of linear features (C5) encodes the quantity of water in a sphere. This is measured by counting the number of water molecules within the sphere radius. This group of features is motivated by the fact that biologists observed that an active site is usually located in a hydrophobic core of the protein while on the surface the quantity of water is higher and the residues exposed to the solvent are not hydrophobic.

Finally, the last class of linear features (C6) is the one which measures the atomic density of the sphere calculated as the total number of atoms in the sphere.

It should be noted that (a) the last two classes of linear features are made discrete using a different number of value intervals. A feature is associated with

Table 1. Representation: feature classes

Linear Features	Description
1st Class	Physical and chemical properties (amino acid attributes)
2nd Class	Amino acidic Composition
3rd Class	Charge/Neutrality
5th Class	Water molecule quantity
6th Class	Atomic density
Structural Features	Description
4th Class	Tree substructures from tertiary structure

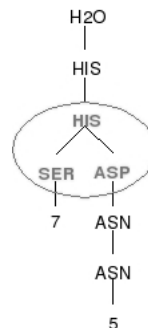
an example if the measured value of a certain property falls in the correspondent range. (b) These features are often used to describe protein structures in similar tasks of Bioinformatics [8] and to develop software for protein structure prediction like Modeler 7v7.

3.2 Structural features

We designed a class of structural features to encode the three-dimensional structure (tertiary structure) or better, the spatial configuration characterizing a spherical region, i.e. the set of amino acids composing it with their 3D distances. As previously mentioned this representation results in a completely connected graph since every vertex is connected to any other vertex in the sphere graph through an edge labeled with the 3D distance of the pair.

Starting from this completely connected graph, we extract some tree substructures using heuristics: for example, the one which preserves the maximum interaction distances to 5 Å between the side-chains of the residues and a minimum spanning tree algorithm. Such heuristics is motivated by the observation that to perform the catalytic function it is necessary that the side-chains of the catalytic residues can interact with each other and with the substrate. The maximum interaction distance between atoms in different residue side-chains is usually of about 3-4 Å. We chose a cut-off distance of 5 Å to take into consideration our approximation in the representation of residues (Figure 1(b)).

The applied heuristics lead possibly to the separation in disconnected components of the initial graph. From each of these components, using the Prim algorithm [9], we extract the spanning tree which minimizes the cost function $c(T)$, i.e. the interaction distances d_{xy} between the side-chains of the residues x and y . Note that as some graphs contain more than one connected component,

**Fig. 2.** Graphical representation of a tree of a sphere

the Prim algorithm is applied to each of them. This leads to the extraction of a tree forest.

We add the water molecule (center of the sphere) as root node to the obtained spanning tree. In Figure 2, we show a tree which can represent the spherical region in Figure 1(a). In bold light gray, we highlight the nodes which represent catalytic amino acids.

The tree substructures generated for each example constitute the features analyzed by our tree kernel function. If two examples are described by two tree forests, we can use as a kernel function the summation of a tree kernels applied to all possible pairs coming from such forests.

4 Experiments

In the subsequent subsections, we describe our classification experiments carried out on the data set that we generated from the Protein Data Bank.

4.1 Experimental set-up

The evaluations were carried out using the SVM-light-TK software [10] (available at <http://ai-nlp.info.uniroma2.it/moschitti/>) which encodes tree kernels in SVM-light [11]. We used the polynomial kernels for the linear features and tree kernels for the structural feature processing. More precisely, we used the SST and the PT kernels described in [7] on a simple tree, i.e. the main tertiary structure³, or on a tree forest (see Section 3.2). The former kernels are indicated with SST_T and PT_T whereas the latter are called SST_F and PT_F. The kernel for tree forest is simply the summation of all possible pairs of trees contained in two examples.

We experimented our models with the protein structures downloaded from the Protein Data Bank (PDB). We adequately pre-processed PDB files to obtain all the information of interest for this task. In particular, we created a data set of 14,688 examples from 48 hydrolases from the PDB structures. The data set is composed of 171 positive examples and 14,571 negative examples, which means a $\frac{1}{125}$ ratio between positive and negative examples.

The results were evaluated by applying a 5-fold cross validation⁴ on this data set measuring the performance with the F_1 measure⁵.

A noticeable attention was devoted to parameterization (cost factor, decay factor, etc.)

³ The single tree structure is the most relevant one in the forest, that is, the tree which contains at least a catalytic amino acid and the two nearest residue side-chains of the sphere

⁴ We separated the data set into five parts, each one composed of examples belonging to a set of nine or ten protein structures randomly assigned to this set.

⁵ F_1 assigns equal importance to Precision P and Recall R i.e. $F_1 = \frac{2P \cdot R}{P+R}$

Table 2. (a) Linear features performance. (b) Combined kernel performance.

(a)				(b)			
Linear	Precision	Recall	F_1		Precision	Recall	$F_1 \pm Std.Dev.$
C1	5.5%	66.7%	10.2%	L	62.3%	55.4%	56.2% ± 6.8
C2	55.9%	63.3%	59.4%	SST_F	66.2%	31.8%	39.9% ± 13.7
C3	20%	3.3%	5.7%	L+SST_F	82.9%	58.6%	68.3% ± 14.5
C5	2.2%	30%	4.1%				
C6	5.5%	13.3%	7.8%				

4.2 Experiment results

Table 2(a) reports the results on the 5 types of linear features using the polynomial kernel (degree 3). These results are only indicative as we did not run a cross validation procedure. We note that most linear features cannot discriminate between active and non-active site. Only, the second class, which encodes the structural information, shows a meaningful F_1 . The general low results of linear features is caused by the remarkable complexity of the task as suggested by the F_1 upperbound of the random selection, i.e. $\simeq 1.6\%$.

In order to boost the classification performance, we experimented with the structural kernels. Table 2(b) summarizes the cross validation results: Row 2 reports the results with polynomial kernels on all the linear features (L), Row 3 shows the outcomes of the SST kernels on the tree forest (SST_F) and Row 4 illustrates the performance of the polynomial kernel summed to the SST kernel on the tree forest (L+SST_F). The \pm sign precedes the standard deviation evaluated on the 5 folds. It is worth to note that the F_1 obtained with the linear features (56.24%) improves by 12 absolute points if we use the combined model (L+SST_F), i.e. 68%.

We also experimented different variants of Tree Kernels, i.e. based on PTs. The results of the cross validation experiments are summarized in Table 3: Row 2 reports the results with polynomial kernel plus SST_F (applied to linear features and a forest structure), Row 3 reports the cross validation results of polynomial kernel plus SST_T (applied to linear features and a tree structure) and finally Row 4 illustrates the performance of the additive combination of polynomial with the PT kernel (PT_T) (on linear features and a tree structure).

The results show that the highest F_1 measure can be achieved with the SST_F but quite similar performance can be obtained representing examples with only a tree structure in the forest, i.e. SST_T. In contrast to our expectations the PT kernel, which may be considered the one most suitable for this task, shows the lowest F_1 . The most plausible explanation is the highest complexity on deriving its correct parameterization.

L+TK	Precision	Recall	$F_1 \pm Std.Dev.$
SST_F	82.9%	58.6%	68.3% ± 14.5
SST_T	79.7%	51.7%	62.3% ± 10.4
PT_T	80.4%	41.2%	54.4% ± 9.1

Table 3. Tree kernel impact

Overall, the very good F_1 of our best model suggests that our classification system can be a useful tool to help the biology researcher to study the protein functions.

5 Conclusions

In this paper, we have studied the problem of the identification of protein functional sites. We have defined a novel computational representation based on biological and spatial considerations and several classes of linear and structural features.

The experiments with SVMs using polynomial and tree kernels and their combinations show that the highest F_1 , i.e. 68%, is achieved with the combined model. Such finding is very important since it firstly shows the successful identification of catalytic sites of a very large family of catalytic proteins belonging to a broad classes of enzymes. Moreover, our work highlights the importance of structural information in the detection of protein active sites. This result motivates the need of structural representations which we efficiently modeled by means of tree kernels.

References

1. Cilia, E., Fabbri, A., Uriani, M., Scialdone, G.G., Ammendola, S.: The signature amidase from *sulfolobus solfataricus* belongs to the cx3c subgroup of enzymes cleaving both amides and nitriles: Ser195 and cys145 are predicted to be the active sites nucleophiles. *The FEBS Journal* **272** (2005) 4716–4724
2. Tramontano, A.: *The ten most wanted solutions in Protein Bioinformatics*. Chapman & Hall/CRC Mathematical Biology and Medicine Series (2005)
3. Berman, H.M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T.N., Weissig, H., Shindyalov, I.N., Bourne, P.E.: The protein data bank. *Nucleic Acids Res* **28**(1) (2000) 235–242
4. Meng, E.C., Polacco, B.J., Babbitt, P.C.: Superfamily active site templates. *PROTEINS: Structure, Function, and Bioinformatics* **55** (2004) 962–976
5. Gärtner, T.: A survey of kernels for structured data. *Multi Relational Data Mining (MRDM)* **5** (2003) 49–58
6. Vapnik, V.: *The Nature of Statistical Learning Theory*. Springer (1995)
7. Moschitti, A.: Efficient convolution kernels for dependency and constituent syntactic trees. In: *Proceedings of The 17th European Conference on Machine Learning*, Berlin, Germany, 2006, Berlin, Germany (2006)
8. Borgwardt, K.: *Graph-based Functional Classification of Proteins using Kernel Methods*. Ludwig Maximilians University of Monaco (2004)
9. Prim, R.C.: Shortest connection networks and some generalizations. *Bell Syst. Tech. Journal* **36** (1957) 1389–1401
10. Moschitti, A.: A study on convolution kernel for shallow semantic parsing. In: *Proceedings of the 42th Conference on Association for Computational Linguistic (ACL-2004)*, Barcelona, Spain (2004)
11. Joachims, T.: Making large-scale svm learning practical. In: *Advances in Kernel Methods - Support Vector Learning*, B. Schölkopf, C. Burges and A. Smola editors (1999)