

Prediction of protein structural class based on symmetrical recurrence quantification analysis

Ines Abdennaji, Mourad Zaied, Jean-Marc Girault

▶ To cite this version:

Ines Abdennaji, Mourad Zaied, Jean-Marc Girault. Prediction of protein structural class based on symmetrical recurrence quantification analysis. Computational Biology and Chemistry, 2021, 92, pp.107450. 10.1016/j.compbiolchem.2021.107450. hal-03676706

HAL Id: hal-03676706 https://hal.science/hal-03676706

Submitted on 10 Mar 2023

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.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

Prediction of Protein Structural Class Based on Symmetrical Recurrence Quantification Analysis

Ines Abdennaji , Mourad Zaied , Jean-Marc Girault

Research Team in Intelligent Machines, National School of Engineers of Gabes, B.P. W 6072 Gabes, Tunisia GSII ESEO -LAUM UMR CNRS 6613, 49000 Angers, France

Abstract

Protein structural class prediction for low similarity sequences is a significant challenge and one of the deeply explored subjects. This plays a important role in drug design, folding recognition of protein, functional analysis and several other biology applications. In this paper, we worked with two benchmark databases existing in the literature 1) 25PDB and 2) 1189 to apply our proposed method for predicting protein structural class. Initially, we transformed protein sequences into DNA sequences and then into binary sequences. Furthermore, we applied symmetrical recurrence quantification analysis (the new approach), where we got 8 features from each symmetry plot computation. Moreover, the machine learning algorithms such as Linear Discriminant Analysis, Random Forest and Support Vector Machine are used. In addition, comparison was made to find the best classifier for protein structural class prediction. Results show that symmetrical recurrence quantification as feature extraction method with RF classifier outperformed existing methods with an overall accuracy of 100% without overfitting.

Keywords: Protein structural classes, Symmetry, Symmetrical recurrence quantification analysis, Recurrence plot, Machine learning, SVM, LDA, Random Forest

1 1. Introduction

Today, the structural classes in four levels (quaternary, ternary, secondary and primary) play a significant role in theoretical and experimental studies of protein science. The protein quaternary and the tertiary structures are

Preprint submitted to Neural Networks

May 15, 2021

determined via the process of protein folding. Protein secondary structure is 5 the three-dimensional form of local segments of proteins whose amino acids 6 linear sequence (in a peptide or protein) forms the protein primary structure. 7 As mentioned by Chou et Zhang in 1995[1], it is important and helpful to 8 predict higher proteinic classes from primary proteinic sequences for two 9 reasons. Firstly, if the structural class of the protein under study is known 10 then the searching scope of conformation can be reduced [2]. Secondly, the 11 structural class is related to various protein properties [3]. Since there is 12 no simple and direct way for the protein tertiary structure prediction from 13 its primary structure, four secondary structural classes of proteins based on 14 the types and arrangement of their secondary structural class are proposed 15 by Levitt and Chothi [4]. These classes are the α , the β and those with a 16 mixture of α and β shapes called the α/β and the $\alpha + \beta$. 17

These four protein structural classes can be used to 1) implement a heuristic method for deciding tertiary structure [5], 2) reduce search space of probable conformations of tertiary structure [6, 2], 3) improve prediction of secondary structure accuracy and 4) predict function from amino acid sequence information. Protein structural class prediction plays an essential role in functional analysis, protein structures, drug designs and a lot of other similar applications in biology [7].

For the last 10 years, prediction of protein structural class for low similar-25 ity sequences [8, 9] is a tough challenge for the scientific community. There-26 fore, an automated and accurate protein structural class prediction for newly 27 established proteins is required. In order to extract the feature sequences 28 from protein, various feature extraction techniques are used in the recent 29 studies which can be later useful for classification of the structural classes. 30 Most of used techniques include Amino Acid Composition(AAC) [10, 11, 12], 31 Average Chemical Shift (ACS) [13, 14, 15], Pseudo Amino Acid (PSeAA) [16], 32 polypeptides composition [17], PsiBlast [18, 19] and etc. These techniques 33 do not facilitate to reach 70% of classification results individually therefore, 34 extracted features from different feature extraction techniques are fused. Fur-35 thermore, to classify the structural classes, various classification methods are 36 applied such as Fisher's Linear Discriminant Algorithm (LDA) [20], Support 37 Vector machines (SVM) [11, 21, 22, 23], Artificial Neural Network (ANN) 38 [24] and Bayesian Classifier [25]. 39

From the studies presented above, it is noticed that there is a great disparity in the protein sequences encoding and feature extraction. Furthermore, classification performance can be improved by using the fused feature engineering technique and machine learning methods. The need to introducenew simple methods with high performance is expected.

The proposed work is a continuation of the previously undertaken stud-45 ies [21, 26] on the use of recurrences and the recent work done by Girault 46 [27] based on the link between recurrences and symmetries. In addition, the 47 presence of symmetry in the tertiary structures of proteins [28] suggests that 48 symmetry can be an important property which has to be explored. Conse-49 quently, it is appropriate to investigate the consideration of symmetries for 50 the classification of proteins. The major contributions of this study is to 51 present: 52

- ⁵³ 1. a simpler protein sequences encoding;
- ⁵⁴ 2. an easy to use method;
- 3. new feature vectors based on symmetry concept and recurrence;
- 4. the best classifier by comparing different protein structural class pre diction models such as SVM, LDA and RF.

The remaining paper is arranged as follows. Material and methodology is presented in Section 2. Results illustrate in section 3 accompanied by discussion in section 4. Finally, Section 5 concludes the paper.

⁶¹ 2. Materials and Methods

62 2.1. The Framework

The framework diagram of this study is shown in Fig.1. First, data 63 set is split up into training and test sets with a ratio of 80:20. Second, 64 the training and test sets are preprocessed through a coding phase. Then, 65 symmetrical recurrence plots (SRP) are calculated and the feature extraction 66 step is performed by applying symmetrical recurrence quantification analysis 67 (SRQA). In total, three different features data sets are calculated: 8-SRQA-68 R, 8- SRQA-I, 16-SRQA, their definition will be presented in subsection 2.4. 69 Third, the machine learning models such as the RF, SVM and LDA are used 70 to training data set for training. The model parameters iteratively tune to 71 improve the performance in the training process. Lastly, test data set is used 72 to evaluate the the trained models. 73



Figure 1: Framework diagram.

74 2.2. Database

In this work, we used two benchmark databases containing low similarity proteins which are widely used for predicting protein structural classes: the database 25PDB includes 1673 protein sequences with 40% sequence homology and the database 1189 contains 1092 protein sequences with 25% sequence homology. Table 1 gives more details about the two databases [8] and the distribution of the four secondary structural classes.

Dataset	α	β	α/β	$\alpha + \beta$	Total
25PDB	443	443	346	441	1673
1189	223	294	334	241	1092

Table 1: Structure of the two data sets used in our study.

81 2.3. Reverse Encoding & DNA Codification

Each protein is formed with a linear sequence of Amino Acids (AAs). In addition, there are 20 standard genetic codes and multi-coded methods. So, each one protein could be expressed by different kinds of nucleotide sequences. The reverse encoding goes in inverse from protein to DNA sequence. As there is no uniqueness in the universal code of translating DNA into AAs, we used the codon (see in Table 2) as presented by Deschavanne and Tuffery [29]. In their study, the authors prove that this encoding gives the best results for

Proteii	ı Seque	nce:											
			D	P	F	L	V		H	H			
DNA S	Sequenc	e:					ļ		F	Cev e	rse E n	coding	
	GAC	CCA	TΊ	С	ст	A	GI	G			CAC	TCA	
Binary Representation													
11-1	-1-1	1 -1	-1 -]	L	-1-]	11	1	-11				-11-1	-11-1

Figure 2: Representation of Protein 1A6M by a binary sequence.

⁸⁹ the prediction of protein structural class. Furthermore, the authors guarantee

 $_{\tt 90}$ $\,$ the balance in base composition to maximize the difference between the AAs

91 codes.

Table 2: Reverse Encoding

There are a lot of representations of DNA sequences used in the biology field like: numerical representation [30], Chaos Game representation[31], binary representation[32] and, Etc. For the sake of simplicity, we used a unique DNA representation performed by Elio Conte et al. [33] which is based on attributing:

• (+1) to the purine: Adenine (A) and Guanine (G);

• (-1) to the pyrimidine: Cytosine (C) and Thymine (T);

The simple reverse binary encoding (reverse encoding + binary DNA encoding) constitutes the first contribution of our proposed approach. It permits the transformation of one protein sequence into a binary sequence, one example is shown in Fig2. This will help to visualize, extract and identify characteristics from the sequences such as symmetries and recurrences.

104 2.4. Proposed Approach

Our second contribution is an improvement of previously undertaken studies [20, 26] that are based on the use of recurrences. The improvement extracts four kinds of symmetrical recurrences as proposed initially in the recent work done by Girault [27]. These extracted symmetrical recurrences have two advantages: they use symmetry properties that have not been used
currently and the symmetrical nature of recurrences does not require an embedding phase, therefore, making it much simpler. To further explore the
symmetrical recurrences and to make the paper autonomous, we recalled the
concept of standard recurrences plot in the appendix.

114 2.4.1. Symmetrical Recurrence Plot

As proposed by Girault in [27], taking symmetrical properties of recur-115 rences in consideration make processes understandable and detect invisible 116 transitions effectively. The present work is an application of this new concept 117 to biological discrete sequences. From the concept of symmetrical recurrence 118 plot, four novel recurrence matrices are proposed. In [27], it is seen that re-119 spective matrices are sensitive to the occurrence of diverse symmetry types. 120 Four types of transformation are performed i.e. Translation, Reflection, In-121 version and Glide (TRIG). Furthermore, corresponding components of the 122 two-dimensional matrix M_k (a new matrix) can be presented in the general-123 ized framework as below: 124

$$M_k(j,i) = \ominus[\varepsilon - \parallel X(j) - G_k X(i) \parallel]$$
(1)

with ε a gauge and $k \in \{T, R, I, G\}$.

¹²⁶ The theoretical framework proposed is similar to the one proposed in [34]:

$$|| X(j) - G_k X(i) || \le \varepsilon.$$
(2)

¹²⁷ Four types of operations are considered:

• $G_T[X(j)] = X(j+n)$ represents a translation of n samples, k = T;

• $G_R[X(j)] = X(-j+n)$ represents a reflection at the position n, k = R;

•
$$G_I[X(j)] = -X(-j+n)$$
 represents an inversion at the position n, k =
I;

• $G_G[X(j)] = -X(j+n)$ represents a glide reflection of n samples, k = 134 G.

An interesting properties of symmetrical recurrence plots are 1) not useful to embed and 2) sojourn points are naturally removed. This means that standard settings are fix to d = 1 (embedding dimension), $\tau = 0$ (time

delay). Also, the gauge is null ($\varepsilon = 0$) since we are working on binary 138 sequences. In the particular case of binary data, $M_T = M_R$ and $M_I = M_G$. 139 This is observed in Fig.3 where the four symmetrical recurrence plots (SRP) 140 were computed using Equ.1 by considering a protein sequence. We clearly 141 noticed that the Translation and Reflection presented the same plot. In 142 addition, Glide and Inversion gave the identical plot. Owing to these two 143 matching results, we will consider just the **Reflection** and the **Inversion** in 144 the rest of the paper. Finally, The quantification step is very significant and 145 useful to investigate the difference between local and global symmetries in 146 the symmetrical recurrence analysis. 147

148 2.4.2. Symmetry Recurrence Quantification analysis

In order to quantify the different types of recurrences, it is recommended to extend the current recurrence descriptors to other forms of recurrence such as symmetrical recurrences. Therefore, Symmetry Recurrence Quantification Analysis (SRQA) is proposed based on Recurrence Quantification Analysis (RQA) [35, 36, 37].

Eight descriptors are calculated for each recurrence matrix $M_R(j, i)$ and $M_I(j, i)$. Therefore, a total of sixteen descriptors were calculated with $k \in \{R, I\}$ (see (equation A.3 to equation A.11) in appendix): Recurrence Rate (RR_k) , Determinism (DET_k) , Entropy $(ENTR_k)$, Laminarity (LAM_k) , Maxline $(Lmax_k)$, Meanline (L_k) , Trapping Time (TT_k) and Trend $(TREND_k)$. Finally, we can define 3 sets of features as input's classifiers simply:

• 8-SRQA-R (RR_R , DET_R , $ENTR_R$, LAM_R , $Lmax_R$, L_R , TT_R , $TREND_R$);

• 8-SRQA-I (RR_I , DET_I , $ENTR_I$, LAM_I , $Lmax_I$, L_I , TT_I , $TREND_I$);

• 16-SRQA $(RR_R, DET_R, ENTR_R, LAM_R, Lmax_R, L_R, TT_R, TREND_R, RR_I, DET_I, ENTR_I, LAM_I, Lmax_I, L_I, TT_I, TREND_I).$

¹⁶⁴ 2.5. Prediction Model and Performance Metrics

As discussed in section I, the purpose of the study is to predict the protein structural classes such as α , β , α/β and $\alpha + \beta$. The framework for classification is presented and described in Fig 1. In our study, 3 sets of features (8-SRQA-R, 8-SRQA-I, 16-SRQA) are fed into machine learning classifiers. Furthermore, machine learning classifiers such as SVM, LDA are used as suggested in [18, 19, 20] to predict the protein structural class. Besides, the



Figure 3: (a) Translation Recurrence Plot, (b) Reflection Recurrence Plot, (c) Inversion Recurrence Plot, (d) Glide Recurrence Plot, for the Time series derived from protein 1A6M. The parameters used: $\varepsilon = 0$, d=1, $\tau = 0$.

ensemble technique such as RF is also considered. In order to compare each
classifier and validate the accuracy of classification models, performance metrics are utilized. We decide to use the performance metrics in line with the
recent studies such as overall accuracy and sensitivity. These measures are
calculated as below:

176

$$OA = \frac{TN + TP}{TP + FP + TN + FN} \tag{3}$$

$$sensitivity = \frac{TP}{FN + TP} \tag{4}$$

where TP and TN are # True Positive and # True Negative respectively. In addition, FP and FN are # False Positive and # False Negatives accordingly.

180 **3. Result**

Sensitivity (%) and Overall Accuracy (%) were calculated considering two
benchmark datasets (25PDB and 1189). For the sake of clarity, a synthesis
of results obtained with the three classifiers (SVM, LDA, RF) is presented
below in Tables 3-5. More details are presented in the appendix in Tables
A.7-A.9.

¹⁸⁶ 3.1. Support Vector Machine(SVM) Classifier

During the training process, three hyper-parameters were tuned such as the kernel coefficient gamma (auto mode), the polynomial kernel function degree (set to 3) and on/off probability estimates (set to TRUE). Finally, a test set was used to evaluate the model.

Dataset	Scenarios	Sensitivity						
		All- α	All- β	α/β	$\alpha + \beta$	OA		
25PDB	8-SRQA-I	100	82	70	71	81.2		
	16-SRQA	100	81	82	79	86.0		
1189	8-SRQA-I	47	90	100	48	74		
	16-SRQA	62	100	94	57	80.4		

Table 3: Sensitivity (%) of our method using SVM on the two benchmark datasets. Scenarios correspond to the two best feature sets.

In Table 3, the best result is obtained with SVM[16-SRQA] in the both benchmark datasets for example All- α : 100%, All- β : 81%, α/β : 82% and $\alpha+\beta$: 79% and with 86% overall accuracy for the database 25PDB, and All- α : 62%, All- β : 100%, α/β : 94% and $\alpha+\beta$: 57% and with 80.4% overall accuracy for the database 1189. According to Table 3 SVM classifier performs better with the database 25PDB as compared to the database 1189 considering sensitivity.

¹⁹⁸ 3.2. Linear Discriminant Analysis (LDA) classifier

During the training process, default hyper-parameters were used with a 199 dimensionality reduction. Finally, a test set was used to evaluate the model. 200 In Table 4, the best result is obtained with LDA[16-SRQA] in the both 201 benchmark datasets for example All- α : 99%, All- β : 95%, α/β : 99% and 202 $\alpha + \beta$: 97% and with 97% overall accuracy for the database 25PDB, and All- α : 203 99%, All- β : 97%, α/β : 96% and $\alpha+\beta$: 100% and with 98.2% overall accuracy 204 for the database 1189. According to Table 4, LDA classifier performs better 205 with the database 1189 as compared to the database 25PDB considering 206 sensitivity. 207

Dataset	Scenarios	Sensitivity						
		All- α	All- β	α/β	$\alpha + \beta$	OA		
25PDB	8-SRQA-R	98	92	100	96	96.4		
	16-SRQA	99	95	99	97	97		
1189	8-SRQA-I	97	95	98	100	98		
	16-SRQA	99	97	96	100	98.2		

Table 4: Sensitivity (%) of our method using LDA on the two benchmark datasets. Scenarios correspond to the two best feature sets.

208 3.3. Random Forest (RF) classifier

During the training process, hyper-parameter such as the number of estimators and the maximum depth were tuned. These parameters were selected as 9 (for the number of estimators) and 6 (for the maximum depth). Finally, a test set was used to evaluate the model.

In Table 5, the best result is obtained with RF [8-SQRA-I] in both benchmark datasets for example All- α : 100%, All- β : 100%, α/β : 100% and $\alpha + \beta$: 100% and with 100%. According to Table 5, RF classifiers performs in a similar way to whatever the dataset based on sensitivity.

217 3.4. Classifier Comparison:

From Tables 3, 4 and 5, it can be claimed that the best combination between classifier input features and the classifier is RF[8-SRQA-I] with on overall of 100% without overfitting on both benchmark data sets with the same encoding. The second best combination is obtained with LDA[16-SRQA] with an overall of 97%. The worst combination is obtained with SVM[16-SRQA] with an overall of 80.5%. Consequently, we recommend using RF[8-SRQA-I].

	а ·		n	• . • •				
Dataset	Scenarios	Sensitivity						
		All- α	All- β	α/β	$\alpha + \beta$	OA		
25PDB	8-SRQA-I	100	100	100	100	100		
	16-SRQA	78	81	75	94	82		
1189	8-SRQA-I	100	100	100	100	100		
	16-SRQA	91	93	97	87	92.2		

Table 5: Sensitivity (%) of our method using RF on the two benchmark data sets. Scenarios correspond to the two best feature sets.

4. Discussion:

In this study, we showed the possibility to classify the 4 protein structural 226 classes: All- α , All- β , α/β , $\alpha + \beta$ without error by considering: 1) a binary 227 encoding of protein sequences, 2) the calculation of symmetrical recurrences 228 and its 8 associated descriptors/features and 3) a classifier. In our study, 229 the best combination of classifiers and their inputs is RF [8-SRQA-I]. From 230 our point of view, the joint use of 1) a simple encoding, 2) taking into ac-231 count descriptors based on symmetrical recurrences and 3) use the ensemble 232 classifier is proved very significant for better results. 233

In Table 6 a comparison is made between our method (8-SRQA-I) and 234 existing methods (RQA) obtained in [20] and [26] on the same data sets. 235 In [20], the protein sequences are encoded in two time series via the Chaos-236 Game-Representation (CGR) approach, 8-RQA and a LDA classifier were 237 applied. In [20], the data are embedded in a space with d = 8 dimensions 238 and with a delay $\tau=2$ and $\varepsilon=0.3$. The results obtained having sensitivity % 239 64.3(All- α), 65(All- β), 61.7(α/β) and 65($\alpha + \beta$) with an overall prediction 240 accuracy of 64% for 25PDB dataset. Similar behavior is seen in [26] with 241 overall prediction accuracy 90% and LDA was applied. 242

In Table 6 a comparison is made between our best results obtained with 243 RF[8-SRQA-I] and other existing methods [20, 38, 26, 39, 40, 41]. From 244 Table 6, it is clearly shown that our best configuration i.e. RF[8-SRQA-I] 245 outperformed the recent results of Wang et al [39]. In addition, from Ta-246 bles A.7-A.9 we see that RF[8-SRQA-I], RF[8-SRQA-R] overpass the recent 247 results of Wang et al [39] for any type of symmetry. Moreover, the percent-248 age accuracy and sensitivity on training and testing are the almost same. 249 Besides, the RF model is tuned with 6 as maximum depth and results ap-250 proached to 100%. Cross-validation shows the same results. Therefore, our 251 classification model is not overfitting. 252

It is often tricky to find the right parameters. In our approach, the data are binary coded therefore, $\varepsilon = 0$. In order to choose the embedded dimension, a value greater than 1 eliminates false recurrences (sojourn point). With the combined use of binary data and symmetric recurrences, there are no more sojourn points. Therefore, search for an embedding dimension or a delay is not required. Consequently, the right parameters are d=1, $\tau=0$ and $\varepsilon=0$.

²⁶⁰ 5. Conclusion

In this paper, we have shown that the judicious combination of 1) a sim-261 ple reverse encoding followed by a binary coding, 2) the calculation of sym-262 metrical recurrences features and, 3) a classifier like RF, provide the best 263 classifications of 4 structural classes of proteins such as All- α , All- β , α/β 264 and $\alpha + \beta$ without overfitting. The simple recurrences settings $(d = 1, \tau = 0)$ 265 and $\varepsilon = 0$) are proved useful to calculate the recurrences. The consideration of 266 symmetry suggested by the presence of symmetric tertiary structures of pro-267 teins results in 100% classification without error. This proposed classification 268 method can be used for other applications having binary or quaternary data. 269 Furthermore, our proposed method will help in improving the drug design, 270 folding recognition of protein, functional analysis and several other biology 271 applications. 272

273 Appendix A. Appendix

274 Appendix A.1. Recurrence Plot

Eckmann et al. [42] proposed the concept of recurrence plot initially to identify the presence of identical neighboring points in a time series such as

Dataset	Methods	Classifier		Sensi	tivity		OA
			All- α	All- β	α/β	$\alpha + \beta$	
25PDB	AAD-CGR $[20]$	LDA	64.3	65	61.7	65	64
	SCPRED[40]	SVM	92.6	80.1	74	71	79.7
	Zhang et Al.[41]	SVM	96.7	80.8	82.4	75.5	83.7
	H. Olyaee[26]	LDA	95.6	89.5	88.1	87	90
	WD PseAAC[38]	SVM	95.7	97.7	94.8	84.4	93.1
	Wang $[39]$	KNN	98	98.9	98	97.5	98.1
	Our Method	LDA	99	95	99	97	97
		SVM	100	81	82	79	86
		RF	100	100	100	100	100
1189	AAD-CGR [20]	LDA	62.3	67.7	63.1	66.5	65.2
	SCPRED $[40]$	SVM	89.1	86.7	89.6	53.8	80.6
	Zhang et Al[41]	SVM	92.4	84.4	84.4	73.4	83.6
	H. Olyaee[26]	LDA	92.3	90.1	86.5	75.2	-
	WD PseAAC[38]	SVM	98.7	99	94	68.9	90.8
	Wang et $Al[39]$	KNN	98.2	99.3	99.1	91.3	97.3
	Our Method	LDA	99	97	96	100	98.2
		SVM	62	100	94	57	80.4
		RF	100	100	100	100	100

Table 6: Comparison of our method (8-SRQA-I) with other studies.

 x_{1} x_{1} x_{2} , x_{2} , x_{N} . This time series is embedded into a phase space with an embedding dimension d and a time delay τ . Two points such as X(i)and X(j) in the d-dimensional space are considered recurrent if they satisfy the following test [42]:

$$|| X(i) - X(j) || \le \varepsilon.$$
(A.1)

A two-dimensional matrix N x N, M (recurrence Matrix) can be calculated as followed:

$$M(i,j) = \ominus [\varepsilon - \parallel X(i) - X(j) \parallel].$$
(A.2)

The recurrence matrix M is a binary matrix composed of zeros and ones where zero components present the same state and non-zero components exhibit different states. Besides, the right selection of d is very important, as incorrect selection will lead towards recurrences contamination [35, 43] (false recurrences/sojourn points). In addition, The embedding dimension is usually set with $d \ge 2$ to avoid the presence of sojourn points. The increase in dimensionality usually reduces the number of false recurrences; however, other approaches are proposed by Zaylaa et al.[44].

²⁹¹ Appendix A.2. Recurrence Quantification Analysis

The Recurrence Quantification Analysis (RQA) extracts quantitative fea-292 tures (descriptors) from the binary matrix M. It permits to measure differ-293 ently appearing recurrence plots (RPs) with the help of small-scale structures 294 present inside it. A main advantage of RQA is to give effective information 295 for non-stationary and short data while other techniques fail to do so. It 296 may be applied to versatile types of data. Moreover, to quantify the com-297 plexity, different measures of RQA introduced heuristically in [35, 36, 37] as 298 described below. 299

Recurrence Rate (RR): it measures of the density of recurrent points present in the matrix M. RR ranges between 0 to 100% where 100% reflect that all the points are recurrent.

$$RR = \frac{1}{N^2} \sum_{i,j=1}^{N} M(i,j) \quad \forall i \neq j$$
(A.3)

Determinism (DET): it measures the presence of temporal correlation and appears through the presence of diagonal/anti diagonal. DET is the percentage of recurrence points assembled to build diagonal lines.

$$DET = \frac{\sum_{l=l_{min}}^{N} lp(l)}{N^2(RR)} \tag{A.4}$$

where p(l) represents the probability of finding diagonal line/ anti-diagonal of l. $l_m in$ is the segment which is shortest and considered often as 2.

Entropy: it measures the deterministic structures complexity within the system. It depends on the bin-number sensitively.

$$ENTR = -\sum_{l=l_{min}}^{N} p(l)ln(p(l))$$
(A.5)

where p(l) represents the chances of occurrence is that diagonal segment is of exact length(l) which is calculated based on frequency distribution P(l).

$$p(l) = \frac{p(l)}{\sum_{l=l_{min}}^{N} p(l)}$$
(A.6)

Laminarity: it measures of the total number of recurrence points which combine to form a vertical line.

$$LAM = \frac{\sum_{v=v_{min}}^{N} vp(v)}{N^2(RR)}$$
(A.7)

where p(v) represents the probability of finding vertical lines of v which has at least $v_m in$ as length.

316 Maxline is the longest length of the diagonal line.

$$L_{MAX} = max(l_i; i = 1, \dots, N_l) \tag{A.8}$$

Meanline: the vertical and diagonal line's length can be measured. Therefore, the average diagonal line length is called the meanline which is associated with the predictability interval of the dynamic system.

$$L = \frac{\sum_{l=l_{min}}^{N} l(p(l))}{\sum_{l=l_{min}}^{N} p(l)}$$
(A.9)

Trapping Time (TT): it measures the average length of the vertical lines, which is directly connected to the laminarity interval of the dynamic system i.e. how long the dynamic system will remain in some specific state.

$$TT = \frac{\sum_{v=v_{min}}^{N} v(p(v))}{\sum_{v=v_{min}}^{N} p(v)}$$
(A.10)

Trend (TREND): it is the regression coefficient of the linear association among the density of recurrence points in a line parallel to the line of Identity and its distance to the line of Identity. In addition, the trend gives significant information about the system's stationarity.

$$TREND = \frac{\sum_{l=1}^{\overline{N}} (i - \frac{\overline{N}}{2})(RR - \langle RR \rangle)}{\sum_{i=1}^{\overline{N}} (i - \frac{\overline{N}}{2})^2}$$
(A.11)

 $_{327}$ where \overline{N} is the Maximal number of diagonals parallel to the LOI.

328 Appendix A.3. Tables

Tables A.7-A.9 present sensitivity (%) obtained with the three different classifiers (SVM, LDA and RF). In each table, the two benchmark data sets 25PDB and 1189 are analyzed. For each row, All- α , All- β , α/β and $\alpha + \beta$ are tested.

Firstly, considering the SVM classifier, it is observed in Table A.7 that the best performances are globally obtained with 16-SRQA (Fusion) in 25PDB benchmark data set. However, with data set 1189, the best outcome comes from 8-SRQA-R.

Secondly, considering the LDA classifier, it is noticed in Table A.8 that the best performances are globally obtained with 16-SRQA (Fusion) in both benchmark data sets.

Thirdly, considering the RF classifier, it can be seen in Table A.9 that the best performances are globally obtained with 8-SQRA-I in both benchmark data sets. Note that outcomes obtained with 8-SRQA-R and 16-SRQA are fairly close to those obtained with 8-SRQA-I.

Dataset	Methods	Sensitivity					
		All- α	All- β	α/β	$\alpha + \beta$		
25PDB	8-SRQA-R	100	76	66	79		
	8-SRQA-I	100	82	70	71		
	16-SRQA	100	81	82	79		
1189	8-SRQA-R	44	100	100	98		
	8-SRQA-I	47	90	100	48		
	16-SRQA	62	100	94	57		

Table A.7: Sensitivity of our proposed method using SVM on the two benchmark datasets. Classifier input features are 8-SRQA-R (Reflection), 8-SRQA-I (Inversion) and 16-SRQA.

344 References

- [1] K.-C. Chou, C. Zhang, Prediction of protein structural classes, Crit ical reviews in biochemistry and molecular biology 30 (1995) 275–349.
 doi:10.3109/10409239509083488.
- [2] I. Bahar, A. R. Atilgan, R. L. Jernigan, B. Erman, Understanding
 the recognition of protein structural classes by amino acid composition,
 Proteins: Structure, Function, and Bioinformatics 29 (1997) 172–185.

Dataset	Methods	Sensitivity					
		All- α	All- β	α/β	$\alpha + \beta$		
25PDB	8-SRQA-R	98	92	100	96		
	8-SRQA-I	99	92	99	95		
	16-SRQA	99	95	99	97		
1189	8-SRQA-R	100	93	98	100		
	8-SRQA-I	97	95	98	100		
	16-SRQA	99	97	96	100		

Table A.8: Sensitivity of our method using LDA on the two benchmark datasets. Classifier input features are 8-SRQA-R (Reflection), 8-SRQA-I (Inversion) and 16-SRQA.

Dataset	Methods	Sensitivity					
		All- α	All- β	α/β	$\alpha + \beta$		
25PDB	8-SRQA-R	99	99	100	100		
	8-SRQA-I	100	100	100	100		
	16-SRQA	78	81	75	94		
1189	8-SRQA-R	100	100	100	100		
	8-SRQA-I	100	100	100	100		
	16-SRQA	91	93	97	87		

Table A.9: sensitivity of our method using RF on the two benchmark datasets. Classifier input features are 8-SRQA-R (Reflection), 8-SRQA-I (Inversion) and 16-SRQA.

- [3] K. Nishikawa, T. Ooi, Correlation of the amino acid composition of a protein to its structural and biological characters, The Journal of Biochemistry 91 (1982) 1821–1824.
- ³⁵⁴ [4] M. Levitt, C. Chothia, Structural patterns in globular proteins, Nature ³⁵⁵ 261 (1976) 552–558.
- L. Carlacci, K. C. Chou, G. M. Maggiora, A heuristic approach to
 predicting the tertiary structure of bovine somatotropin, Biochemistry
 30 (1991) 4389–4398.
- [6] K.-C. Chou, Energy-optimized structure of antifreeze protein and its
 binding mechanism, Journal of molecular biology 223 (1992) 509-517.
- [7] M. M. Gromiha, S. Selvaraj, Protein secondary structure prediction in
 different structural classes., Protein engineering 11 (1998) 249–251.

- [8] L. A. Kurgan, L. Homaeian, Prediction of structural classes for protein sequences and domains—impact of prediction algorithms, sequence representation and homology, and test procedures on accuracy, Pattern Recognition 39 (2006) 2323–2343.
- [9] K.-C. Chou, Progress in protein structural class prediction and its im pact to bioinformatics and proteomics, Current Protein and Peptide
 Science 6 (2005) 423-436.
- [10] K.-C. Chou, A novel approach to predicting protein structural classes in
 a (20-1)-d amino acid composition space, Proteins: Structure, Function,
 and Bioinformatics 21 (1995) 319–344.
- J. K. Kim, S.-Y. Bang, S. Choi, Sequence-driven features for prediction
 of subcellular localization of proteins, Pattern recognition 39 (2006)
 2301–2311.
- ³⁷⁶ [12] Y.-D. Cai, J. Hu, X. Liu, K.-C. Chou, Prediction of protein structural classes by neural network method, J Mol Des 1 (2002) 332–338.
- [13] S. Fias, S. Van Damme, P. Bultinck, Multidimensionality of delocalization indices and nucleus independent chemical shifts in polycyclic aromatic hydrocarbons, Journal of computational chemistry 29 (2008) 381 358-366.
- [14] H. Lin, C. Ding, Q. Song, P. Yang, H. Ding, K.-J. Deng, W. Chen,
 The prediction of protein structural class using averaged chemical shifts,
 Journal of Biomolecular Structure and Dynamics 29 (2012) 1147–1153.
- ³⁸⁵ [15] Z. Feng, X. Hu, Z. Jiang, H. Song, M. A. Ashraf, The recognition of ³⁸⁶ multi-class protein folds by adding average chemical shifts of secondary ³⁸⁷ structure elements, Saudi Journal of Biological Sciences 23 (2016) 189– ³⁸⁸ 197.
- [16] Y. Liang, S. Zhang, Predict protein structural class by incorporating two
 different modes of evolutionary information into chou's general pseudo
 amino acid composition, Journal of Molecular Graphics and Modelling
 78 (2017) 110–117.
- ³⁹³ [17] X.-D. Sun, R.-B. Huang, Prediction of protein structural classes using
 ³⁹⁴ support vector machines, Amino acids 30 (2006) 469–475.

- T. Liu, X. Zheng, J. Wang, Prediction of protein structural class for low similarity sequences using support vector machine and psi-blast profile,
 Biochimie 92 (2010) 1330–1334.
- [19] L. Li, X. Cui, S. Yu, Y. Zhang, Z. Luo, H. Yang, Y. Zhou, X. Zheng,
 Pssp-rfe: accurate prediction of protein structural class by recursive
 feature extraction from psi-blast profile, physical-chemical property and
 functional annotations, PLoS One 9 (2014) e92863.
- [20] J.-Y. Yang, Z.-L. Peng, Z.-G. Yu, R.-J. Zhang, V. Anh, D. Wang, Prediction of protein structural classes by recurrence quantification analysis
 based on chaos game representation, Journal of Theoretical Biology 257
 (2009) 618–626.
- [21] P. Sudha, D. Ramyachitra, P. Manikandan, Enhanced artificial neural
 network for protein fold recognition and structural class prediction, Gene
 Reports 12 (2018) 261–275.
- [22] A. Anand, G. Pugalenthi, P. Suganthan, Predicting protein structural
 class by svm with class-wise optimized features and decision probabilities, Journal of theoretical biology 253 (2008) 375–380.
- [23] X.-J. Zhu, C.-Q. Feng, H.-Y. Lai, W. Chen, L. Hao, Predicting protein
 structural classes for low-similarity sequences by evaluating different features, Knowledge-Based Systems 163 (2019) 787–793.
- ⁴¹⁵ [24] W. Bao, Y. Chen, D. Wang, Prediction of protein structure classes with
 ⁴¹⁶ flexible neural tree, Bio-medical materials and engineering 24 (2014)
 ⁴¹⁷ 3797–3806.
- ⁴¹⁸ [25] Z. Aydin, A. Singh, J. Bilmes, W. S. Noble, Learning sparse models for
 ⁴¹⁹ a dynamic bayesian network classifier of protein secondary structure,
 ⁴²⁰ BMC bioinformatics 12 (2011) 154.
- ⁴²¹ [26] M. H. Olyaee, A. Yaghoubi, M. Yaghoobi, Predicting protein structural
 ⁴²² classes based on complex networks and recurrence analysis, Journal of
 ⁴²³ theoretical biology 404 (2016) 375–382.
- ⁴²⁴ [27] J.-M. Girault, Recurrence and symmetry of time series: Application to transition detection, Chaos, Solitons & Fractals 77 (2015) 11–28.

- ⁴²⁶ [28] R. Xu, M. Li, H. Chen, Y. Huang, Y. Xiao, A symmetry-related
 ⁴²⁷ sequence-structure relation of proteins, Chinese Science Bulletin 50
 ⁴²⁸ (2005) 536.
- [29] P. Deschavanne, P. Tuffery, Exploring an alignment free approach for
 protein classification and structural class prediction, Biochimie 90 (2008)
 615–625.
- [30] H. K. Kwan, S. B. Arniker, Numerical representation of dna sequences,
 in: 2009 IEEE International Conference on Electro/Information Technology, IEEE, 2009, pp. 307–310.
- [31] H. J. Jeffrey, Chaos game representation of gene structure, Nucleic acids
 research 18 (1990) 2163–2170.
- ⁴³⁷ [32] V. R. F, Evolution of long-range fractal correlations and 1/f noise in ⁴³⁸ dna base sequences, Rev. Lett 68 (1992) 3805–8.
- [33] S. Conte, A. Giuliani, Identification of possible differences in coding
 and non-coding fragments of dna sequences by using the method of
 the recurrence quantification analysis, arXiv preprint arXiv:0910.3516
 (2009).
- ⁴⁴³ [34] J.-M. Girault, A. Humeau-Heurtier, Centered and averaged fuzzy en-⁴⁴⁴ tropy to improve fuzzy entropy precision, Entropy 20 (2018) 287.
- [35] C. L. Webber Jr, J. P. Zbilut, Dynamical assessment of physiological
 systems and states using recurrence plot strategies, Journal of applied
 physiology 76 (1994) 965–973.
- [36] N. Marwan, N. Wessel, U. Meyerfeldt, A. Schirdewan, J. Kurths, Recurrence-plot-based measures of complexity and their application to heart-rate-variability data, Physical review E 66 (2002) 026702.
- [37] L. Trulla, A. Giuliani, J. Zbilut, C. Webber Jr, Recurrence quantification
 analysis of the logistic equation with transients, Physics Letters A 223
 (1996) 255–260.
- [38] B. Yu, L. Lou, S. Li, Y. Zhang, W. Qiu, X. Wu, M. Wang, B. Tian,
 Prediction of protein structural class for low-similarity sequences using
 chou's pseudo amino acid composition and wavelet denoising, Journal
 of Molecular Graphics and Modelling 76 (2017) 260–273.

- [39] S. Wang, X. Wang, Prediction of protein structural classes by different
 feature expressions based on 2-d wavelet denoising and fusion, BMC
 bioinformatics 20 (2019) 701.
- [40] L. Kurgan, K. Cios, K. Chen, Scpred: accurate prediction of protein
 structural class for sequences of twilight-zone similarity with predicting
 sequences, BMC bioinformatics 9 (2008) 226.
- ⁴⁶⁴ [41] L. Zhang, X. Zhao, L. Kong, A protein structural class prediction ⁴⁶⁵ method based on novel features, Biochimie 95 (2013) 1741–1744.
- [42] J. Eckmann, S. O. Kamphorst, D. Ruelle, et al., Recurrence plots of
 dynamical systems, World Scientific Series on Nonlinear Science Series
 A 16 (1995) 441–446.
- [43] T. March, S. Chapman, R. Dendy, Recurrence plot statistics and the
 effect of embedding, Physica D: Nonlinear Phenomena 200 (2005) 171–
 184.
- [44] A. Zaylaa, J. Charara, J.-M. Girault, Reducing sojourn points from
 recurrence plots to improve transition detection: Application to fetal
 heart rate transitions, Computers in biology and medicine 63 (2015)
 251–260.