

Protein Fold Recognition Using Genetic Algorithm Optimized Voting Scheme and Profile Bigram

Harsh Saini^{1*}, Gaurav Raicar¹, Sunil Lal¹, Abdollah Dehzangi², Seiya Imoto³, Alok Sharma^{1,2}

¹The University of the South Pacific, Fiji.

²Griffith University, Brisbane, Australia.

³University of Tokyo, Japan.

* Corresponding author; Tel.: +679 32 32946; email: harshsaini90@gmail.com

Manuscript submitted January 13, 2016; accepted March 20, 2016.

doi: 10.17706/jsw.11.8.756-767

Abstract: In biology, identifying the tertiary structure of a protein helps determine its functions. A step towards tertiary structure identification is predicting a protein's fold. Computational methods have been applied to determine a protein's fold by assembling information from its structural, physicochemical and/or evolutionary properties. It has been shown that evolutionary information helps improve prediction accuracy. In this study, a scheme is proposed that uses the genetic algorithm (GA) to optimize a weighted voting scheme to improve protein fold recognition. This scheme incorporates k-separated bigram transition probabilities for feature extraction, which are based on the Position Specific Scoring Matrix (PSSM). A set of SVM classifiers are used for initial classification, whereupon their predictions are consolidated using the optimized weighted voting scheme. This scheme has been demonstrated on the Ding and Dubchak (DD), Extended Ding and Dubchak (EDD) and Taguchi and Gromhia (TG) datasets benchmarked data sets.

Key words: K-separated bigrams, protein fold recognition, SCOP, PSSM, genetic algorithm, support vector machines.

1. Introduction

In the field of biological science, protein fold recognition refers to assigning a protein to one of a finite number of folds. This is considered a crucial transitional step in identifying a protein's tertiary structure [1]. Recognition of protein folds requires development of feature extraction and classification techniques. In literature, several feature extraction and classification techniques have been proposed. Dubchak *et al.* [2] have used structural and physicochemical based features for protein fold recognition. Taguchi and Gromiha [3] proposed structural features (occurrence) for protein fold recognition. Ghanty and Pal [3] have employed pairwise frequencies of amino acids, both for adjacent (PF2) and separated by one residue (PF1). Yang *et al.* 2011 [4] have used PF1 and PF2 in an augmented form in their study. Their features, subsequently, were of higher dimensions compared to Ghanty and Pal [3]. Nonetheless, high dimensionality can be controlled by utilizing either feature selection methods (eg. [5]–[10]) or dimensionality reduction methods (eg. [11]–[15]).

Recently in protein fold recognition, the use of evolutionary features have been showing good performance [16], [17].

Evolutionary features are extracted from Position Specific Scoring Matrix (PSSM) and are basically a relative measure of amino acid occurrence at a particular position in the protein sequence. Several researchers have used PSSM for improving protein fold recognition and some of these include auto cross-covariance [18], bi-gram [19],

tri-gram [20], sequence alignment via dynamic time warping [21] and ensemble features [22], [23].

Table 1. Summary of Ding and Dubchak Dataset

Fold	Number of training vectors	Number of test vectors
α		
Globin-like	13	6
Cytochromec	7	9
DNA-binding 3-helical-bundle	12	20
4-Helical up-and-down bundle	7	8
4-Helical cytokines	9	9
α EF-hand	6	9
β		
Immunoglobulin-like -sandwich	30	44
Cupredoxins	9	12
Viral coat and capsid proteins	16	13
ConA-like lectins/glucanases	7	6
SH3-like barrel	8	8
OB-fold	13	19
Trefoil	8	4
Trypsin-like serineproteases	9	4
Lipocalins	9	7
α/β		
(TIM)-barrel	29	48
FAD (also NAD)-binding motif	11	12
Flavodoxin-like	11	13
NAD(P)-binding Rossmann-fold	13	27
P-loop containing nucleotide	10	12
Thioredoxin-like	9	8
Ribonuclease H-like motif	10	12
Hydrolases	11	7
Periplasmic binding protein-like	11	4
$\alpha+\beta$		
β -Grasp	7	8
Ferredoxin-like	13	27
Small inhibitors, toxins, lectins	13	27

Furthermore, some of the classification techniques that have been explored include Linear Discriminant Analysis [24], K-Nearest Neighbors [25], Bayesian Classifiers [26]–[28], Support Vector Machines (SVM) [19]–[21], [28]–[30], Artificial Neural Networks (ANN) [31]–[33] and ensemble classifiers [34], [35]. Out of these mentioned classification techniques, SVM has showed promising results in protein fold recognition problem. However, it is shown in the literature that to further improve the protein folding accuracy, a good combination of features extraction technique as well as classification technique is needed [36], [37]. In this paper, a scheme is proposed where the Genetic Algorithm (GA) is utilized to determine the optimal weights for a weighted voting scheme. In this scheme, features are extracted using k-separated bigrams and classification is carried by utilizing multiple instances of SVM classifiers. The extracted information for each k-separated bigrams is used for classification by a separate SVM classifier whose individual predictions are combined using a weight voting system that has been optimized using the GA [38].

Table 2. Summary of Extended Ding and Dubchak Dataset

Folds	Number of Samples
A	
Globin-like	41
Cytochromec	35
DNA-binding 3-helical-bundle	322
4-Helical up-and-down bundle	69
4-Helical cytokines	30
Alpha; EF-hand	59
β	
Immunoglobulin-like β -sandwich	391
Cupredoxins	47
Viral coat and capsid proteins	60
ConA-like lectins/glucanases	57
SH3-like barrel	129
OB-fold	156
Trefoil	45
Trypsin-like serineproteases	45
Lipocalins	37
α/β	
(TIM)-barrel	336
FAD (also NAD)-binding motif	73
Flavodoxin-like	130
NAD(P)-binding Rossmann-fold	195
P-loop containing nucleotide	239
Thioredoxin-like	111
Ribonuclease H-like motif	128
Hydrolases	83
Periplasmic binding protein-like	16
$\alpha + \beta$	
β -Grasp	121
Ferredoxin-like	339
Small inhibitors, toxins, lectins	124

Table 3. Summary of Taguchi and Gromhia Dataset

Folds	Number of Samples
α	
Cytochrome C	25
DNA/RNA binding 3-helical bundle	103
Four helical up and down bundle	26
EF hand-like fold	25
SAM domain-like	26
$\alpha - \alpha$ super helix	47
β	
Immunoglobulin-like β -sandwich	173
Common fold of diphtheria toxin/transcription factors/cytochrome	28
Cupredoxin-like	30

Galactose-binding domain-like	25
Concanavalin A-like lectins/glucanases	26
SH3-like barrel	42
OB-fold	78
Double-stranded α -helix	34
Nucleoplasmin-like	42
α/β	
TIM α/β -barrel	145
NAD(P)-binding Rossmann-fold domains	77
FAD/NAD(P)-binding domain	31
Flavodoxin-like	55
Adenine nucleotide a hydrolase-like	34
P-loop containing nucleoside triphosphate hydrolases	95
Thioredoxin fold	32
Ribonuclease H-like motif	49
S-adenosyl-L-methionine-dependent methyltransferases	34
α/β -Hydrolases	37
$\alpha + \beta$	
β -Grasp, ubiquitin-like	42
Cystatin-like	25
Ferredoxin-like	118
Knottins	80
Rubredoxin-like	28

2. Dataset

In this research, data has been used from the benchmarked Ding and Dubchak (DD) protein sequence dataset. The dataset consists of a training set for the creation of the model and an independent test set for testing queries against the model. The data set belong to 27 SCOP folds which further represent the four major structural classes – α , β , $\alpha+\beta$, α/β . A brief listing of the DD dataset has been illustrated in Table 1. The training dataset consists of 311 protein sequences where any given pair of sequences do not have more than 35% sequence identity for aligned subsequences longer than 80 residues and the test set consists of 383 protein sequences where the sequence identity between any two given proteins is less than 40% [2]. Moreover, further benchmarking has been conducted using the Extend Ding and Dubchak (EDD) and Taguchi and Gromhia (TG) datasets. The EDD dataset has 27 SCOP folds with proteins having less than 40% sequence similarity while the TG dataset has 1612 proteins with less than 25% sequence similarity and comprises of 30 SCOP folds. A summary of these datasets is provided in Table 2 and Table 3 respectively. It should be noted that EDD and TG datasets do not have data distributed into predefined train and test sets, therefore, data was randomly divided in the ratio 3:2 for the purposes of training.

3. Procedure

3.1. Overview

The scheme, initially, commences with the extraction of PSSM from protein sequences using PSI-BLAST. This is succeeded by the calculation the k-separated bigrams and the features are extracted for $k=1, \dots, 11$. Finally, the SVM classifiers provide individual sets of predictions that are consolidated using a GA optimized weighted voting system. A flow-diagram of the classification procedure is illustrated in Fig 1.

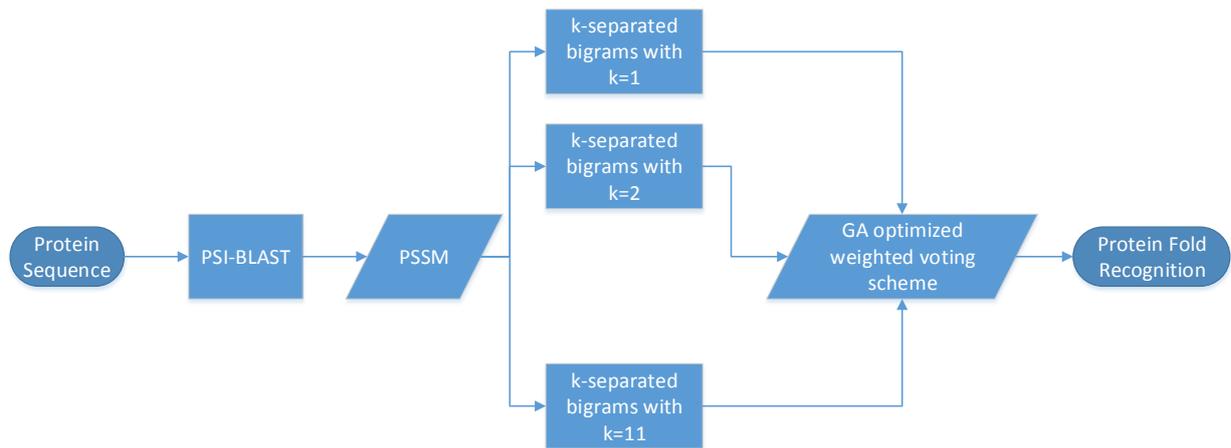


Fig. 1. A flow diagram of the proposed scheme.

3.2. Feature Extraction

The technique proposed in this study attempts to model relationships between amino acids that are non-adjacent in the protein sequence. To accomplish this, amino acid bigram probabilities are extracted from the sequential evolution probabilities in PSSM. In this study, we also explore k -separated bigrams that are non-adjacent. i.e., they are separated by other amino acids in the sequence whereby k determines the positional distance between the bigrams under consideration. This technique can be summarized mathematically as shown in Equation 1. If we let N be the PSSM matrix representation for a given protein, N will have L rows, where L is the length of the primary sequence and 20 columns since there are only 20 unique amino acids. The transition of the m th amino acid to the n th amino acid is:

$$T_{m,n}(k) = \sum_{i=1}^{L-k} N_{i,m} N_{i+k,n} \quad (1)$$

where $1 \leq m \leq 20, 1 \leq n \leq 20$ and $1 \leq k \leq K$

$$T(k) = [T_{1,1}(k), T_{1,2}(k), \dots, T_{1,20}(k), T_{2,1}(k), \dots, T_{2,20}(k), \dots, T_{20,1}(k), \dots, T_{20,20}(k)] \quad (2)$$

The equation stated in (1) constructs a matrix $T(k)$ or $F(k)$ (for clarity) (2) which contains 400 elements representing 400 amino acid transitions for just a single value of k . As stated previously, k represents the distance between the amino acid positions that are used to compute the transition probabilities. For $k=1$, the transition probabilities are computed between neighboring amino acids whereas, for $k=2$, the transition probabilities are computed between amino acids that are separated by 1 amino acid in between. Therefore, for $k=11$, the amino acids used to calculate the transition probabilities are separated by 10 amino acids. We consider k -separated bigrams only until $k=1,2,\dots,11$ due to the computational constraints and also due to classifier performance discussed in later sections. The k -separated bigrams have been visualized in Fig 2 for a sample protein ARTARA.

We calculate k -separated bigrams only up to $k=1,\dots,11$ due to the fact that the classification accuracy gradually decreases as the k gets larger and the weights assigned by the genetic algorithm to each k (shown in Fig 4) become too small to provide any significant input to the classification. Additionally, the accuracy for each featureset $F(k)$ gradually drops as the value of k increases (highlighted in Fig. 3).

Upon extracting all features with $k=1,2,\dots,11$, we get a feature set of 4400 features ($400 * 11$ features). Instead of considering all these features a single feature vector of length 4400, we consider features extracted for a particular value of k as separate and independent feature vector. This is partly due to the high computational requirements of processing the concatenated feature vector and the increased difficulty in classification with a feature vector of high dimensionality. Moreover, this approach will also clearly highlight the contribution by

individual k -separated bigrams. Therefore, for each value of k , we consider the extracted features to be representing the bigram transition probabilities with varying k -separation between the amino acid pairs. Each of these 11 different feature sets provides input for the set of SVM classifiers discussed in the later sections.

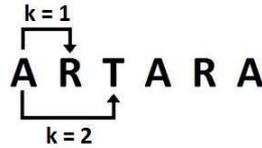


Fig. 2. A sample protein with illustrations of the concept of k -separated bigrams.

In order to illustrate this feature extraction technique, let us consider a fictional protein sequence ARTARA of length $L=5$, and also assume that there are only 3 amino acids for the purpose of illustration. Therefore, the length of the feature vector for each value of k will be 9 (since $3*3$). Assuming that the PSSM matrix for this sequence is given as shown in Table 4, the technique described in this paper will be utilized to extract features for $k=1$ and $k=2$. If we apply the proposed feature extraction technique with $k=1$, we can compute the bigram transition probabilities as shown in Table 5. Similarly, the bigram transition probabilities for $k=2$ can be shown in Table 6. As described previously, each of these feature sets are considered as independent and they provide input to specific classifier in the ensemble classification procedure.

3.3. Classification on Individual Features

SVM was used for classifying each of the feature sets represented by $F(k)$. SVM is a supervised learning model linked to machine learning algorithms that is used for pattern recognition. It is widely used in classification and regression analysis. In its simplest form, SVM accepts a set of inputs and then predicts for each input which of the two possible classes it falls under. For multi-class problems, SVM can still be used by reducing the problem into multiple binary classification problems. SVM aims to construct a hyper-plane in infinite-dimensional space such that a good level of separation is achieved between the classes thus lowering the generalization error of the classifier.

A different instance of SVM was used to build the classification model for every different $F(k)$. Therefore, in this scheme we have 11 different instances of SVM each of which is trained with a particular feature set, $F(k)$. As stated earlier, this is done so to mitigate the problem of high dimensionality with the concatenated feature set. The parameters for each of the SVM classifiers were tuned such that the classifiers yielded the highest training accuracies.

3.4. Fusion on Classifier Outputs

Since the classification model consisted of several instances of SVM classifiers, all of which provided their own sets of predictions, it was important to formulate a scheme to assimilate the predictions from the classifiers efficiently to produce one consolidated set of predictions. Considerable experimentation was carried out to determine the optimal scheme, which included simple majority voting, selection of best performing feature sets from $F(k)$ and genetic algorithm optimized weighted voting.

Initially, a simple majority voting system was used for determining the final prediction using all 11 SVM models. The classification accuracy of this scheme had a slight improvement over the highest individual performance of the SVM classifier. However, it was determined that not all classifiers were required equally for improving classification accuracy. Therefore, a new scheme was created whereby a subset of 5 best performing classifiers were selected based on their training accuracies and then simple majority voting was used to determine the final set of predictions, which resulted in further improvement in the classification accuracy.

However, this approach of selecting a subset of classifiers had its shortcomings. It was difficult to determine the optimal number of classifiers that need to be selected to yield the best results. Additionally, by removing certain classifiers from the final prediction calculation, we were effectively removing the input from the entire

feature set $F(k)$ for that particular value of k . Therefore, we pursued a more holistic approach whereby an evolutionary machine learning technique, Genetic Algorithm, was used to assist in consolidating the predictions.

Table 4. PSSM of Protein ARTARA

Amino Acid	A	R	T
A	0.10	0.60	0.30
R	0.45	0.35	0.20
T	0.20	0.56	0.24
A	0.31	0.42	0.27
R	0.66	0.17	0.17
A	0.13	0.71	0.16

Table 5. Bigram Transition Probabilities for Protein ARTARA with k=1

Amino Acid	A	R	T
A	0.4874	0.8923	0.3403
R	0.8129	0.8333	0.4538
T	0.4497	0.4844	0.2459

Table 6. Bigram Transition Probabilities for Protein ARTARA with k=2

Amino Acid	A	R	T
A	0.3318	0.4991	0.2291
R	0.6527	0.8764	0.4009
T	0.3155	0.4845	0.2100

The evolutionary approach to machine learning is based on computational models which includes natural selection and genetics. This is known as evolutionary computation which simulates evolution on a computer and encompasses genetic algorithms, evolutionary strategies and genetic programming. The latter techniques simulate evolution using selection, mutation and reproduction processes. Genetic Algorithm (GA) is basically an optimization algorithm which iteratively improves the quality of a solution until is optimal using a stochastic approach.

In the proposed scheme, GA was used to optimize weights for voting assigned to each classifier that ranged from 0 to 1, with 0 indicating no input for the final prediction and 1 indicating maximum input towards the final prediction. The final prediction was simply based on selecting the output variable with the highest weight. This has been summarized in Algorithm 1 shown below.

Loop through each class Sum the weighted votes by each classifier for that class End Loop Predicted Class := Class with greatest votes

Algorithm 1. Weighted voting scheme.

To determine the optimal weights for each classifier, real value genetic algorithm was used. The Optimization Toolbox from MATLAB provided such an implementation of the genetic algorithm for minimization problems. Since the aim was to maximize the classification accuracy, the fitness function had to be modified accordingly to represent the problem as a minimization function. The fitness function (Equation 3) and the genetic algorithm parameters are provided in Table 7. The genetic algorithm was allowed to terminate after the fitness values remained constant for 100 generations. On an average, the genetic algorithm ran for 150 generations.

Table 7. Classification Accuracies for k-Separated Bigrams for F(k)

F(k)	Training Acc.	Testing Acc.
1	65.1	68.9
2	62.1	68.9
3	66.3	68.9
4	63.2	65.0
5	64.0	66.1
6	63.7	67.4
7	61.9	65.5
8	62.4	65.3
9	62.9	65.5
10	61.9	66.3
11	60.8	65.0

Table 8. Classification Accuracies for k-Separated Bigrams for F(k)

Method	Training Accuracy (%)
All 11 classifiers, equal weights	68.8
5 top classifiers, equal weights	70.4
All 11 classifiers, GA optimized weights	72.7

Table 9. Classification Accuracies for k-Separated Bigrams for F(k)

Feature set	Accuracy (%)
ACC+HXPZV (Ding and Dubchak, 2001)	56.0
Shamim et al., (2007)	60.5
Ghanty and Pal (2009)	59.2
Chmielnicki and Stapor (2012)	62.8
AHVPZ (Yang et al., 2011)	44.7
AX (Yang et al., 2011)	40.3
AHXPZV (Yang et al., 2011)	49.4
PF (Yang et al., 2011)	60.8
AHVPZ+PF (Yang et al., 2011)	51.2
AHXPZV+PF (Yang et al., 2011)	52.7
Monogram (Sharma et al., 2013)	62.1
Bigram (Sharma et al., 2013)	69.5
This paper	71.5

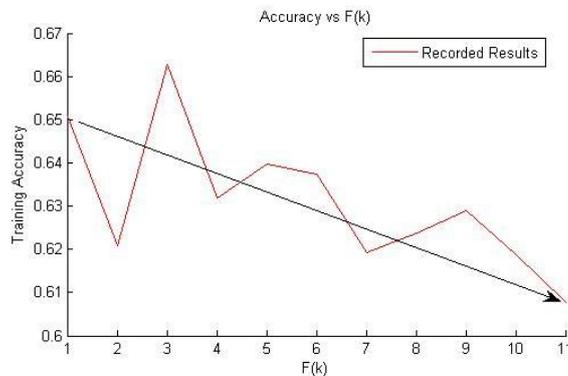


Fig. 3. A graph of training accuracies of F(k) for k=1,2,...,11.

4. Results and Discussion

The experimentation was performed on the benchmarked DD dataset to evaluate the performance of the classification scheme described previously. As previously stated, the DD dataset, which been widely adopted by many researchers, has divided the samples into two sets, training and testing datasets. The classification model described in this study was evaluated using independent testing using the test dataset and k-fold cross-validation testing. This strategy of performance evaluation is widely employed by researchers in literature. However, the jackknife test was not performed due to high computational demands.

Table 10. Summary of 10-Fold Cross Validation on DD, EDD and TG Datasets

Feature sets	DD	EDD	TG
PF1	50.6	50.8	38.8
PF2	48.2	49.9	38.8
PF	53.4	55.6	43.1
O	51.0	46.9	36.3
AAC	45.1	40.9	32.0
AAC+HXPZV	47.2	40.9	36.3
PSSM+PF1	64.6	85.9	66.4
PSSM+PF2	64.7	75.2	52.7
PSSM+PF	67.5	74.9	51.1
PSSM+O	63.5	79.3	58.8
PSSM+AAC	59.2	68.5	46.7
PSSM+AAC+HXPZV	58.2	61.9	44.0
Bigram	74.1	67.9	46.6
Trigram	73.8	84.5	68.1
Alignment	74.7	86.2	72.5
This paper	75.7	87.7	75.8

The performance of SVM on each of the individual feature sets $F(k)$ with $k=1,2,\dots,11$ was encouraging. The results, shown in Table 8, clearly indicate that there is discriminatory information present in k-separated bigrams even when the bigrams are separated by 10 amino acids in between. In order to consolidate the individual classifier predictions, we explored simple majority voting, subset selection with majority voting and weighted voting approaches, the results of which are reported in Table 9. As shown by the results, each of these schemes lead to an improvement on the highest individual performance by SVM, however, GA optimized weighted voting scheme leads to the highest improvement in performance. The optimized weights, shown in Figure 4, yielded 72.7% training accuracy. It should be noted that all evaluations up till this stage were carried out by using the training dataset only unless otherwise stated and these results were analyzed to determine the various parameters. Upon finalizing optimal parameters and weights, the model was now ready for evaluations using the standardized test dataset.

The experiment consisted of two parts; in the first part, the parameters for the model were optimized using the training set and the separate test set was employed to determine the classification accuracy of the model. In the second part, the training and testing datasets were combined and n-fold cross validation was used to evaluate the performance of the model. For this method of performance evaluation, the parameter values determined previously in the independent dataset test were used.

For the first part (using the training set and the independent test set), the performance of the proposed technique showed improvement compared to previous works in literature. It was noted that the highest accuracy achieved was 71.5% (this paper) which showed a 2% improvement compared to existing feature extraction techniques. The results are summarized in Table 10.

For the second part, n-fold cross validation procedure was performed on the merged training and test dataset for $n=10$. The same experimental procedure as the first part was adopted in this part which included the

ensemble of classifiers with optimal parameters as illustrated in Table 11.

The proposed scheme was compared against various other schemes that use information structural and evolutionary information for fold recognition. These techniques included PF1 and PF2 [3], PF [39], Occurrence (O) [40], AAC and AAC+HXPZV [2], which compute feature sets from the original protein sequences. In addition, ACC [18], Bi-gram [19], Tri-gram [20] and Alignment [21] are also included since they compute features directly from the evolutionary information present in PSSM. Moreover, features have been computed from the consensus sequences for PF1, PF2, O, AAC and AAC+HXPZV to obtain additional feature sets for comparison. A prefix of PSSM+ indicates that the features have been computed on the consensus sequence.

The optimal weights that were achieved using GA in the first part were used to combine the classification outputs. The accuracies were compared to past literature works and it was seen that there was an improvement of 2.0%. The highest classification accuracy recorded was 76.7% (this paper).

In order to provide a more comprehensive evaluation of the proposed scheme, the results for EDD and TG datasets are also shown in Table 11. During the experiment, $k=8$ and $k=7$ were used for EDD and TG datasets respectively. These values had been determined similar to the approach described for DD dataset previously.

5. Conclusion

In this study, a feature extraction technique has been employed, which is based on sequential evolution probabilities. This technique uses varying distances of amino acid transitions within the sequence to compute the features. A variety of classifiers were used on each transition distance and the predictions of the classifiers were fused using a weighted voted scheme. Genetic Algorithm was further used to optimize the weights so that the best weight distribution can be determined which would give the optimal classification accuracy.

The proposed technique gave promising results, and the accuracy noted was **76.7%** via n-fold cross validation for the DD dataset. The highest recorded accuracies for the EDD and TG datasets were **87.7%** and **75.8%** respectively. Since user-friendly and publicly accessible web-servers represent the future direction for developing practically more useful models, simulated methods, or predictors we shall make efforts in our future work to provide a web-server for the method presented in this study.

Reference

- [1] Chmielnicki, W. (2012). A hybrid discriminative/generative approach to protein fold recognition. *Neurocomputing*, 75(1), 194-198.
- [2] Dubchak, I., & Ding, C. H. (2001). Multi-class protein fold recognition using support vector machines and neural networks. *Bioinformatics*, 17(4), 349-358.
- [3] Pal, N. R., & Ghanty, P. (2009). Prediction of protein folds: Extraction of new features, dimensionality reduction, and fusion of heterogeneous classifiers. *NanoBioscience, IEEE Trans*, 8(1), 100-110.
- [4] Chen, X., & Yang, J. Y. (2011). Improving taxonomy based protein fold recognition by using global and local features. *Proteins Struct. Funct. Bioinforma*, 79(7), 2053-2064.
- [5] Miyano, S., Imoto, S., Paliwal, K. K., & Sharma, A. (2014). A feature selection method using improved regularized linear discriminant analysis. *Mach. Vis. Appl.*, 25(3), 775-786.
- [6] Sharma, A., Imoto, S., Miyano, S., & Sharma, V. (2012). Null space based feature selection method for gene expression data. *Int. J. Mach. Learn. Cybern.*, 3(4), 269-276.
- [7] Sharma, A., Imoto, S., & Miyano, S. (2012). A between-class overlapping filter-based method for transcriptome data analysis. *J. Bioinform. Comput. Biol.*, 10(5).
- [8] Sharma, A., Imoto, S., & Miyano, S. (2012). A filter based feature selection algorithm using null space of covariance matrix for DNA microarray gene expression data. *Curr. Bioinform.*, 7(3), 289-294.
- [9] Sharma, A., Imoto, S., & Miyano, S. (2012). A top-r feature selection algorithm for microarray gene expression data. *IEEE/ACM Trans. Comput. Biol. Bioinforma.*, 9(3), 754-764.
- [10] Sharma, A., Koh, C. H., Imoto, S., & Miyano, S. (2011). Strategy of finding optimal number of features on gene expression data. *Int. J. Mach. Learn. Cybern.*, 47(8), 480-482.

- [11] Sharma, A., Paliwal, K. K., Imoto, S., & Miyano, S. (2013). Principal component analysis using QR decomposition. *Int. J. Mach. Learn. Cybern.*, 4(6), 679-683.
- [12] Paliwal, K. K., & Sharma, A. (2010). Improved direct LDA and its application to DNA microarray gene expression data. *Pattern Recognit. Lett.*, 31(16), 2489-2492.
- [13] Sharma, A., & Paliwal, K. K. (2008). Rotational linear discriminant analysis technique for dimensionality reduction. *Knowl. Data Eng. IEEE Trans.*, 202(10), 1336-1347.
- [14] Paliwal, K. K., & Sharma, A. (2008). Cancer classification by gradient LDA technique using microarray gene expression data. *Data Knowl. Eng.*, 66(2), 338-347.
- [15] Paliwal, K. K., & Sharma, A. (2007). Fast principal component analysis using fixed-point algorithm. *Pattern Recognit. Lett.*, 28(10), 1151-1155.
- [16] Liu, T., Geng, X., Zheng, X., Li, R., & Wang, J. (2012). Accurate prediction of protein structural class using auto covariance transformation of PSI-BLAST profiles. *Amino Acids*, 42(6), 2243-2249.
- [17] Liu, T., & Jia, C. (2010). high-accuracy protein structural class prediction algorithm using predicted secondary structural information. *J. Theor. Biol.*, 267(3), 272-275.
- [18] Dong, Q., Zhou, S., & Guan, J. (2009). A new taxonomy-based protein fold recognition approach based on autocross-covariance transformation. *Bioinformatics*, 2655-2662.
- [19] Paliwal, K. K., Dehzangi, A., Lyons, J., & Sharma, A. (2013). A feature extraction technique using bi-gram probabilities of position specific scoring matrix for protein fold recognition. *J. Theor. Biol.*, 320, 41-46.
- [20] Dehzangi, A., Lyons, J., Sharma, A., & Paliwal, K. K. (2014). A tri-gram based feature extraction technique using linear probabilities of position specific scoring matrix for protein fold recognition. *NanoBioscience, IEEE Trans.*, 13(1), 44-50.
- [21] Lyons, J., Biswas, N., Sharma, A., Dehzangi, A., & Paliwal, K. K. (2014). Protein fold recognition by alignment of amino acid residues using kernelized dynamic time warping. *J. Theor. Biol.*, 354, 137-145.
- [22] Sattar, A., & Dehzangi, A. (2013). Ensemble of diversely trained support vector machines for protein fold recognition. *Intelligent Information and Database Systems*, 335-344.
- [23] Dehzangi, A., Paliwal, K., Lyons, J., Sharma, A., & Sattar, A. (2013). Enhancing protein fold prediction accuracy using evolutionary and structural features. *Pattern Recognition in Bioinformatics*, 196-207.
- [24] Klein, P. (1986). Prediction of protein structural class by discriminant analysis. *Biochim. Biophys. Acta (BBA)-Protein Struct. Mol. Enzymol.*, 874(2), 205-215.
- [25] Zhang, T. L., & Ding, Y. S. (2007). Using pseudo amino acid composition and binary-tree support vector machines to predict protein structural classes. *Amino Acids*, 33(4), 623-629.
- [26] Chinnasamy, A., Sung, W. K., & Mittal, A. (2005). Protein structure and fold prediction using tree-augmented naive bayesian classifier. *J. Bioinform. Comput. Biol.*, 3(4), 803-819.
- [27] Anand, A., Pugalenth, G., & Suganthan, P. N. (2008). Predicting protein structural class by SVM with class-wise optimized features and decision probabilities. *J. Theor. Biol.*, 253(2), 375-380.
- [28] Sharma, A., Dehzangi, A., Lyons, J., Imoto, S., Miyano, S., Nakai, K., et al. (2014). Evaluation of Sequence Features from Intrinsically Disordered Regions for the Estimation of Protein Function. *PLoS One*, 9(2), e89890.
- [29] Sharma, A., Paliwal, K. K., Dehzangi, A., Lyons, J., Imoto, S., & Miyano, S. (2013). A strategy to select suitable physicochemical attributes of amino acids for protein fold recognition. *BMC Bioinformatics*, 14(1), 233.
- [30] Saini, H., Raicar, G., Sharma, A., Lal, S., Dehzangi, A., Rajeshkannan, A., et al. (2014). Protein structural class prediction via k-separated bigrams using position specific scoring matrix. *J. Adv. Comput. Intell. Intell. Informatics*, 8(4), 474-479.
- [31] Cai, Y. D., & Zhou, G. P. (2000). Prediction of protein structural classes by neural network. *Biochimie*, 8(4), 783-785.
- [32] Jahandideh, S., Abdolmaleki, P., Jahandideh, M., & Asadabadi, E. B. (2007). Novel two-stage hybrid neural discriminant model for predicting proteins structural classes. *Biophys. Chem.*, 128(1), 87-93.

- [33] Jahandideh, S., Abdolmalek, P., Jahandideh, M., & Hayatshahi, S. H. (2007). Novel hybrid method for the evaluation of parameters contributing in determination of protein structural classes. *J. Theor. Biol.*, 244(2), 275-281.
- [34] Dehzangi, A., Amnuaisuk, S. P., Ng, K. H., & Mohandesi, E. (2009). Protein fold prediction problem using ensemble of classifiers. *Neural Information Processing*, 503-511.
- [35] Kedarisetti, K. D., Kurgan, L., & Dick, S. (2006). Classifier ensembles for protein structural class prediction with varying homology. *Biochem. Biophys. Res. Commun.*, 348(3), 981-988.
- [36] Kurgan, L., Zhang, T., Zhang, H., Shen, S., & Ruan, J. (2008). Secondary structure-based assignment of the protein structural classes. *Amino Acids*, 35(3), 551-564.
- [37] Kurgan, L., & Chen, K. (2007). Prediction of protein structural class for the twilight zone sequences. *Biochem. Biophys. Res. Commun.*, 357(2), 453-460.
- [38] Goldberg, D. E., & Holland, J. H. (1988). Genetic algorithms and machine learning. *Mach. Learn.*, 3(2), 95-99.
- [39] Yang, T., Kecman, V., Cao, L., Zhang, C., & Huang, J. Z. (2011). Margin-based ensemble classifier for protein fold recognition. *Expert Syst. Appl.*, 38(10), 12348-12355.
- [40] Taguchi, Y. H., & Gromiha, M. M. (2007). Application of amino acid occurrence for discriminating different folding types of globular proteins. *BMC Bioinformatics*, 8(1), 404.



Harsh Saini is a software developer and an aspiring data scientist. He has completed his bachelor's majoring in computing and information sciences and a postgraduate diploma in computing science from the University of the South Pacific. Recently, he has completed Masters by research focusing on data mining and machine learning techniques in the field of bioinformatics. Primarily, the research involved developing feature extraction methods to extract discriminatory information from raw protein sequences to classify folds, structural classes and subcellular locations. Additionally, another component of the research required using machine learning techniques to develop prediction methodologies to improve classification accuracies.



Alok Sharma received the B.Tech. degree from the University of the South Pacific (USP), Suva, Fiji, in 2000 and the M.Eng. degree, with an academic excellence award, and the Ph.D. degree in the area of pattern recognition from Griffith University, Brisbane, Australia, in 2001 and 2006, respectively.

He was with the University of Tokyo, Japan (2010-2012) as a Research Fellow. He is an A/Prof. at the USP and an Adjunct A/Prof. at the Institute for Integrated and Intelligent Systems (IIIS), Griffith University. He participated in various projects carried out in conjunction with Motorola (Sydney), Auslog Pty., Ltd. (Brisbane), CRC Micro Technology (Brisbane), the French Embassy (Suva), and JSPS (Japan). His research interests include pattern recognition, computer security, human cancer classification and protein fold and structural class prediction problems. He reviewed several articles and is in the editorial board of several journals.

Dr. Alok is a member of the Institute of Electrical and Electronics Engineers (IEEE).