



Protein fold recognition by alignment of amino acid residues using kernelized dynamic time warping

James Lyons ^a, Neela Biswas ^e, Alok Sharma ^{b,c,*}, Abdollah Dehzangi ^{c,d}, Kuldip K. Paliwal ^a



^a School of Engineering, Griffith University, Australia

^b School of Engineering and Physics, University of the South Pacific, Fiji

^c Institute for Integrated and Intelligent Systems (IIIS), Griffith University, Brisbane, Australia

^d National ICT Australia (NICTA), Brisbane, Australia

^e Royal Brisbane and Women's Hospital, Brisbane, Australia

HIGHLIGHTS

- Amino acid alignment method is developed to extract important features from protein sequences.
- The extraction of features is done by computing dissimilarity distances between proteins.
- This method shows significant improvement in protein fold recognition.
- Overall fold recognition was 4.3–7.6% higher than the existing methods.

ARTICLE INFO

Article history:

Received 26 December 2013

Received in revised form

5 March 2014

Accepted 21 March 2014

Available online 31 March 2014

Keywords:

Protein sequence

Fold recognition

Alignment method

Feature extraction

Classification

ABSTRACT

In protein fold recognition, a protein is classified into one of its folds. The recognition of a protein fold can be done by employing feature extraction methods to extract relevant information from protein sequences and then by using a classifier to accurately recognize novel protein sequences. In the past, several feature extraction methods have been developed but with limited recognition accuracy only.

Protein sequences of varying lengths share the same fold and therefore they are very similar (in a fold) if aligned properly. To this, we develop an amino acid alignment method to extract important features from protein sequences by computing dissimilarity distances between proteins. This is done by measuring distance between two respective position specific scoring matrices of protein sequences which is used in a support vector machine framework. We demonstrated the effectiveness of the proposed method on several benchmark datasets. The method shows significant improvement in the fold recognition performance which is in the range of 4.3–7.6% compared to several other existing feature extraction methods.

© 2014 Elsevier Ltd. All rights reserved.

1. Introduction

In biological sciences, deciphering the tertiary structures of proteins is considered to be an important and challenging task. The identification of tertiary structures provides information about protein functions which helps in understanding protein heterogeneity, protein–protein interactions and protein–peptide interactions. The computational ways of determining protein structures has gained considerable attention since it is normally very time

consuming to identify protein structures by crystallography methods. Protein fold recognition could help in the process of recognizing tertiary structure. The objective of protein fold recognition is to associate a fold to a novel protein sequence.

Protein fold recognition broadly covers feature extraction and classification tasks. The brief description of the work conducted in the past has been depicted in the Related Work Section. It has been shown in the literature that feature extraction methods using evolutionary information performs quite well in the fold recognition process (Altschul et al., 1997; Dong et al., 2009; Sharma et al., 2013b). In this work, we have used this information to build a feature extraction method for protein–protein alignment. For this, we extract position specific scoring matrices (PSSMs) using PSI-BLAST and build dissimilarity matrix between two protein

* Corresponding author at: School of Engineering and Physics, University of the South Pacific, Fiji. Tel.: +679 3232223, fax: +679 3231538.

E-mail addresses: alok.sharma@griffith.edu.au, sharma.al@usp.ac.fj (A. Sharma).

sequences and conduct dynamic time warping to find the alignment path. Many proteins share the same fold in spite of the variation in sequence lengths. Therefore, features extracted from aligned homologous proteins give discriminant features for protein fold recognition. In order to illustrate this, we picked 7 protein sequences and extracted their corresponding PSSMs for comparison. Out of 7 protein sequences, 4 protein sequences (Proteins A, B1, B2 and B3 in Fig. 1) belong to a particular fold and the remaining 3 protein sequences (Protein C1, C2 and C3 in Fig. 1) belong to different folds. We then used Protein A (see Fig. 1) and found dissimilarity matrices by comparing it with all the 6 remaining protein sequences. In the first three dissimilarity matrices (i, ii and iii), PSSMs from protein sequences in the same fold are compared and the next three dissimilarity matrices (iv, v and vi), protein sequences of mutually different folds are compared. We can observe that in the first 3 figures (i, ii and iii), a diagonal path can be seen (we call an alignment path), however, in the next 3 figures, this alignment path is not clearly observed. This alignment path (which shows the dissimilarity between two proteins) can be used to distinguish between proteins of one fold with that of another fold. This is a typical example, there could be variations depending upon different proteins. Nonetheless, dissimilarity distance could be used as a measure to observe dissimilarity between proteins. From biological perspective, proteins in the same fold often have amino acid subsequences that are highly

conserved. The alignment path (i.e., the dissimilarity distance) characterizes the subsequences of amino acids in these conserved regions via their PSSMs. If a certain subsequence is conserved in a fold, then each protein in that fold would have a low dissimilarity distance from that conserved region. This can help in discriminating folds that do not have the same amino acid subsequences. The details of the proposed scheme are described later. The proposed scheme provides promising results (in terms of recognition performance) when experimented on 3 benchmark datasets: Ding and Dubchak (DD) (Ding and Dubchak, 2001), Taguchi and Gromiha (TG) (Taguchi and Gromiha, 2007) and extended DD (EDD) (Dong et al., 2009). The 10-fold cross-validation recognition performance on DD dataset is 74.7%, on TG dataset is 74.0% and on EDD dataset is 90.2% which is very promising when compared with other existing feature extraction methods.

2. Related work

The development of protein fold recognition research can be broadly categorized into two main tasks: feature extraction and classification. For the former task, several feature extraction techniques have been developed using structural, physicochemical and evolutionary information. Dubchak et al. (1997) have proposed syntactical and physicochemical-based features for protein

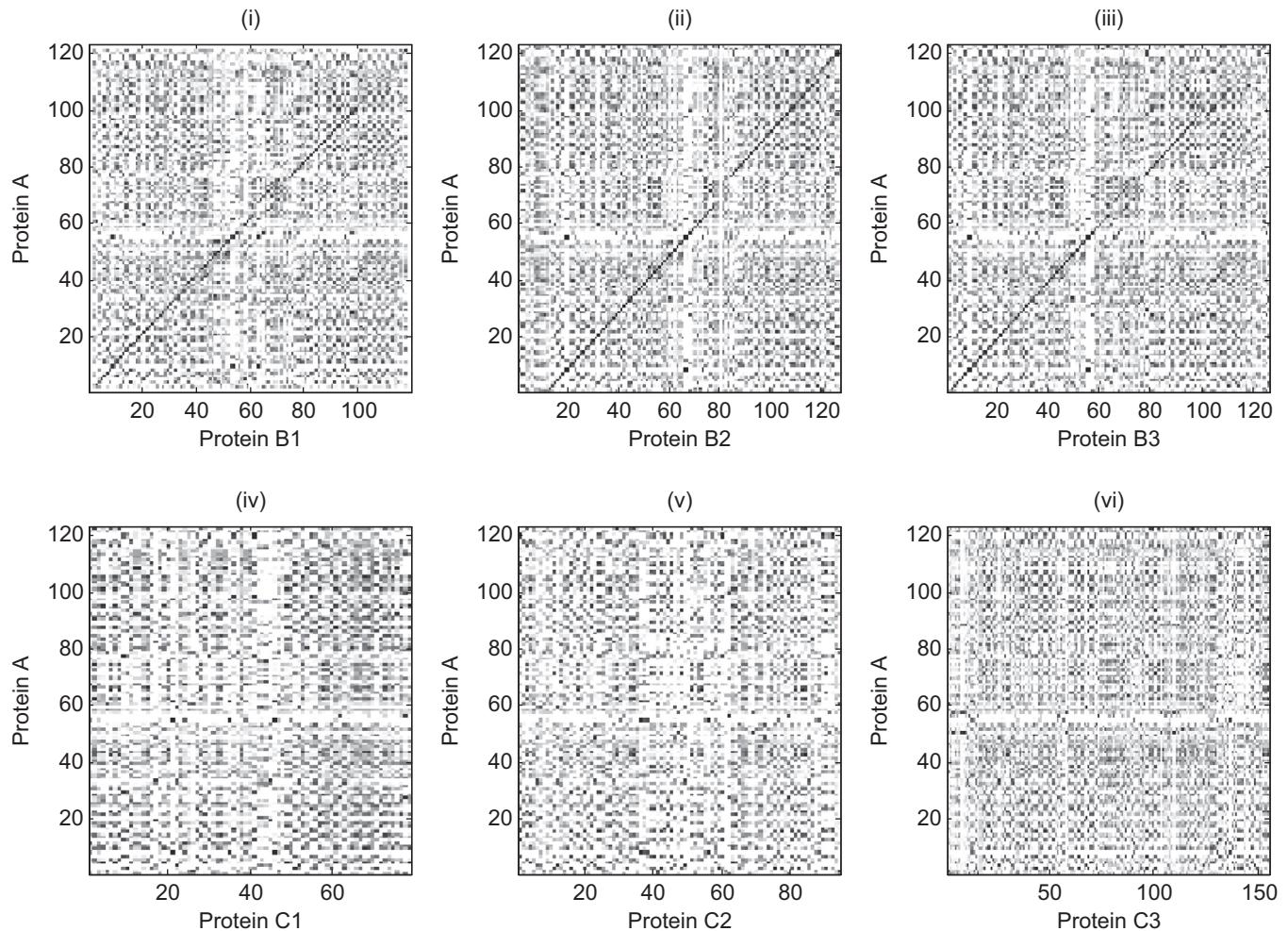


Fig. 1. An illustration using dissimilarity matrix of protein sequences. The pictures above represent similarities computed between PSSMs. Dark pixels indicate corresponding rows of each PSSM are very similar. Long sequences of similar PSSM rows manifest as dark lines in the pictures. The top 3 pictures are from proteins in the same fold, the bottom three all proteins are from different folds. The contrast of these pictures has been increased for clarity in viewing.

fold recognition. They used amino acids' composition (AAC) as syntactical-based features and the 5 following attributes of amino acids for deriving physicochemical-based features namely, hydrophobicity (H), predicted secondary structure based on normalized frequency of α -helix (X), polarity (P), polarizability (Z) and van der Waals volume (V). They used three descriptors (composition, transition and distribution) to compute the features. The AAC features comprise of 20 features and physicochemical-based features comprise of 105 features (21 features for each of the attributes used). The features proposed by Dubchak et al. (1997) have been widely used in the field of protein fold recognition (Chinnasamy et al., 2005; Krishnaraj and Reddy, 2008; Valavanis et al., 2010; Ding and Dubchak, 2001; Dehzangi et al., 2009, 2013a, 2013b, 2013c; Kecman and Yang, 2009; Kavousi et al., 2011; Dehzangi and Amnuaisuk, 2011; Chmielnicki and Stapor, 2012). Apart from the above mentioned 5 attributes used by Dubchak et al. (1997), features have also been extracted by incorporating other attributes of amino acids. Some of the other attributes used are: solvent accessibility (Zhang et al., 2012), flexibility (Najmanovich et al., 2000), bulkiness (Huang and Tian, 2006), first and second order entropy (Zhang et al., 2008a), size of the side chain of the amino acids (Dehzangi and Amnuaisuk, 2011). These physicochemical attributes are selected in an arbitrary way and recently a systematic way of selecting physicochemical attributes was proposed by Sharma et al. (2013a, 2012b). Ohlson et al. (2004) proposed a profile-profile alignment method to improve protein fold recognition. Taguchi and Gromiha (2007) proposed features which are based on amino acids' occurrence; Shamim et al. (2007) have extracted features from the structural information of amino acid residues and amino acid residue pairs; Ghanty and Pal (2009) proposed pairwise frequencies of amino acids separated by one residue (PF1) and pairwise frequencies of adjacent amino acid residues (PF2). There are 400 features each in PF1 and PF2. These pairwise frequency features (PF) are concatenated in the study conducted by Yang et al. (2011), thereby, having 800 features. If the dimensionality of a feature vector is very large then a few important features can be selected for further processing using feature selection or dimension reduction schemes (Sharma et al., 2006, 2011, 2012a, 2012c, 2012d, 2013; Sharma and Paliwal, 2007, 2008, 2010, 2012a, 2012b, 2012c; Paliwal and Sharma, 2011, 2012). To avoid completely losing the sequence-order information, the pseudo amino acid composition (Chou, 2001, 2005) or Chou's PseAAC (Lin and Lapointe, 2013) was proposed to replace the simple amino acid composition (AAC) for representing the sample of a protein. Since the concept of PseAAC was proposed in 2001, it has been widely used to study various attributes of proteins (see, e.g., Esmaeili et al., 2010; Nanni et al., 2012; Zhang et al., 2008b; Mohabatkar et al., 2011, 2013; Mohammad Beigi et al., 2011; Sahu and Panda, 2010; Nanni and Lumini, 2008), among many others (see a long list of papers cited in the References section of Chou (2011)). Recently, the concept of PseAAC was further extended to represent the feature vectors of DNA and nucleotides (Chen et al., 2012, 2013). Because it has been widely and increasingly used, recently two powerful soft-wares, called 'PseAAC-Builder' (Du et al., 2012) and 'propy' (Cao et al., 2013), were established for generating various special Chou's pseudo-amino acid compositions, in addition to the web-server 'PseAAC' (Shen and Chou, 2008) built in 2008. Dong et al. (2009) have shown autocross-covariance (ACC) transformation for protein fold recognition. Shen and Chou (2006), Kurgan et al. (2008) and Liu et al. (2012) have shown autocorrelation features for protein sequence, and Dehzangi and Amnuaisuk (2011) derived features by considering more physicochemical properties. Chou and Cai (2004) have used functional domain composition for predicting protein structural classes. Sharma et al. (2013b) have derived bi-gram features using evolutionary information (PSSM).

Evolutionary information (using PSSM) has been utilized in other work as well (Lin et al., 2013; Chou and Shen, 2007, 2008; Xiao et al., 2011a, 2011b; Lin et al., 2012; Chou et al., 2011, 2012; Wu et al., 2011, 2012; Shen and Chou, 2007). Paliwal et al. (2014) have proposed tri-gram features using evolutionary information. For the latter task case, several classifiers have been developed or used including linear discriminant analysis (Klein, 1986), Bayesian classifiers (Chinnasamy et al., 2005), Bayesian decision rule (Wang and Yuan, 2000), k -nearest neighbor (Shen and Chou, 2006; Ding and Zhang, 2008), hidden Markov model (Bouchaffra and Tan, 2006; Deschavanne and Tuffery, 2009), artificial neural network (Chen et al., 2007; Ying et al., 2009), support vector machine (SVM) (Ding and Dubchak, 2001; Shamim et al., 2007; Ghanty and Pal, 2009) and ensemble classifiers (Dehzangi et al., 2009, 2010; Yang et al., 2011; Dehzangi and Karamizadeh, 2011). Among these classifiers, SVM (or SVM-based for ensemble strategy) classifier exhibits quite promising results (Liu et al., 2012; Kurgan et al., 2008; Ghanty and Pal, 2009).

The extraction of relevant and informative features from protein sequences is a crucial step in identifying protein folds. In order to improve protein fold recognition, we focus on carefully developing the feature extraction method. Since SVM classifier (Vapnik, 1995) provides high recognition accuracy, we use SVM classifier to compare the performance of our feature extraction method with other feature extraction methods. SVM classifiers are often employed with the Radial Basis Function (RBF) kernel. The RBF kernel (along with other common SVM kernels such as the linear and polynomial kernel) requires fixed length feature vectors. This has motivated many previous works to try and extract fixed length representations of proteins so that they can then be efficiently compared. In this work we define a kernel designed to work with variable length data. This allows us to directly compare PSSM matrices, instead of first transforming the matrix into a fixed length vector prior to comparison.

As demonstrated by a series of recent publications (Chen et al., 2012, 2013; Min et al., 2013; Xiao et al., 2013; Xu et al., 2013; Qiu et al., 2014) and summarized in a comprehensive review (Chou, 2011), to develop a really useful predictor for a protein system, one needs to go through the following five steps: (i) select or construct a valid benchmark dataset to train and test the predictor; (ii) represent the samples with an effective formulation that can truly reflect their intrinsic correlation with the target to be predicted; (iii) introduce or develop a powerful algorithm to conduct the prediction; (iv) properly perform cross-validation tests to objectively evaluate the anticipated prediction accuracy; (v) establish a user-friendly web-server for the predictor that is accessible to the public. These steps are elaborated in the following sections.

3. Dataset

In this study, three protein sequence datasets have been used: (1) DD-dataset (Ding and Dubchak, 2001), (2) TG-dataset (Taguchi and Gromiha, 2007) and (3) EDD-dataset (Dong et al., 2009). The DD-dataset that we have used consists of 311 protein sequences in the training set where two proteins have no more than 35% of sequence identity for aligned subsequence longer than 80 residues. The test set consists of 383 protein sequences where sequence identity is less than 40%. Both the sets belong to 27 Structural Classification of Proteins (SCOP) folds which represent all major structural classes: α , β , α/β , and $\alpha+\beta$ (Ding and Dubchak, 2001). These sets were divided originally by the donor of this dataset. The training set and test set have been merged as a single set of data in order to perform the k -fold cross-validation process.

The TG-dataset consists of 1612 protein sequences belonging to 30 different folding types of globular proteins from SCOP. The

names of the number of protein sequences in each of 30 folds have been described in Taguchi and Gromiha (2007). The sequence similarity of protein of TG datasets is no more than 25%.

The EDD-dataset consists of 3418 proteins with less than 40% sequential similarity belonging to the 27 folds that originally used in DD-dataset. We extracted the EDD-dataset from SCOP in similar manner to Dong et al. (2009) in order to study our proposed method using a larger number of samples.

4. Amino acid alignment method

In this section, we present the proposed feature extraction method based on the alignment of proteins. To present the overview, a flow diagram of the proposed scheme has been shown in Fig. 2. The model can be subdivided into the training phase and test phase. In the training phase a set of protein sequences is used to estimate the model parameters and in the test phase, the fold of a novel protein sequence is identified. During the training of the model, PSSM matrices of protein sequences are computed by using PSI-BLAST. In the pairwise analysis step, row vectors of two PSSM matrices are used to measure pairwise distance. By comparing all the row vectors in two PSSMs we get a dissimilarity matrix. This dissimilarity matrix is then used in dynamic time warping (DTW) stage to compute dissimilarity distance between two PSSM matrices of the corresponding proteins. The obtained dissimilarity distance is then used in the kernelization stage to compute kernel distance. A protein is compared progressively with all other proteins to form a kernel matrix. This kernel matrix will then be

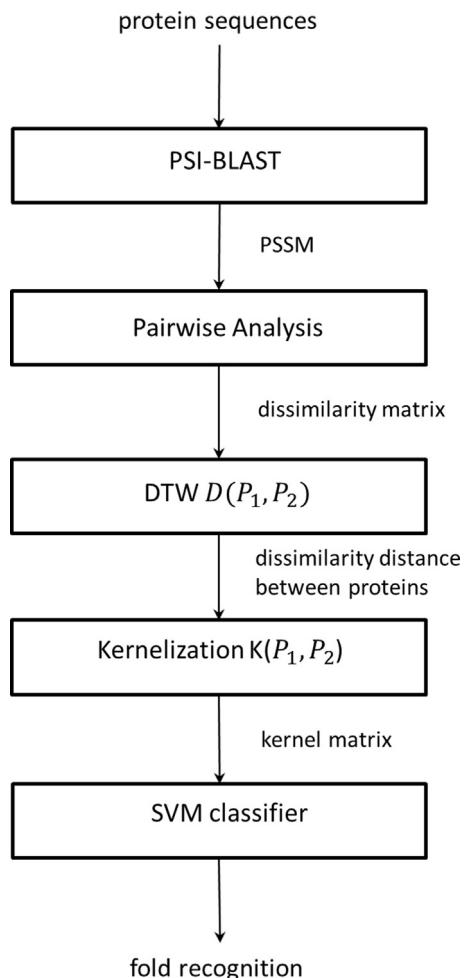


Fig. 2. A flow-diagram of protein sequence classification using alignment method.

used to train SVM parameters. Once the model parameters are estimated then the system can determine the fold of a novel protein sequence.

Let P and Q be the matrices representing PSSM (log probabilities) of two protein sequences of length L_1 and L_2 , respectively. PSSM matrix can be interpreted as the relative probability of substitution of amino acids. The matrix P will have L_1 rows and 20 columns and the matrix Q will have L_2 rows and 20 columns. Let p_i (for $i = 1, 2, \dots, L_1$) and q_j (for $j = 1, 2, \dots, L_2$) be the row vectors of P and Q , respectively. The dissimilarity cosine distance between p_i and q_j can be given as

$$d(p_i, q_j) = 1 - \frac{p_i q_j^T}{\sqrt{p_i p_i^T q_j q_j^T}}, \text{ for } i = 1, 2, \dots, L_1 \text{ and } j = 1, 2, \dots, L_2 \quad (1)$$

Calculating distance d for all L_1 rows and L_2 rows would give a $L_1 \times L_2$ dissimilarity matrix S . We then employ dynamic time warping to find the minimum cost path through the dissimilarity matrix S . This would give cumulative dissimilarity matrix D . The matrix D defines the total cost of alignment between (p_1, q_1) and (p_i, q_j) . Lower cost implies a better alignment, which indicates that the proteins are more similar. The computation of cumulative dissimilarity matrix D can be done in the following way

$$\begin{aligned} D_{ij} &= \min(D_{i-1,j}, D_{i,j-1}, D_{i-1,j-1}) + S_{ij}, \text{ for } i = 1, 2, \dots, L_1 \text{ and} \\ j &= 1, 2, \dots, L_2 \end{aligned} \quad (2)$$

where $D_{ij} = []$ (empty set) for $i \leq 0$ and/or $j \leq 0$ and $S_{ij} = d(p_i, q_j)$.

We define the distance between two PSSM matrices P and Q , as $D_{dtw}(P, Q)$. This can be expressed as $D_{dtw}(P, Q) = D_{L_1, L_2}$. The distance D_{dtw} represents dissimilarity between the aligned proteins. The kernel distance between PSSM matrices P and Q , can be represented as $K(P, Q)$, where γ is a kernel parameter (chosen by performing cross-validation on the training set). The kernel function $K(P, Q)$ is defined by $\exp(-D_{dtw}(P, Q)^2/\gamma^2)$. We then compute the kernel distance between all the pairs of proteins in the training set. This gives a kernel matrix K having n rows and n columns, where n is the number of training samples. The kernel matrix K is then further processed through the SVM classifier for parameter estimation and classification.

5. Support vector machine as a classifier

In this paper we used SVM (Vapnik, 1995) as a classifier. SVM is considered to be the state-of-the-art machine learning and pattern classification algorithm. It has been extensively applied in classification and regression tasks. SVM aims to find maximum margin hyper-plane (MMH) to minimize classification error. In SVM a function called the kernel K is used to project the data from input space to a new feature space, and if this projection is non-linear it allows non-linear decision boundaries (Bishop, 2006). This function K is usually considered as RBF kernel, polynomial kernel or linear kernel. These kernels require fixed length feature vectors. Since the protein sequences are of varying lengths, we can not use these kernels. However, in this work we have defined a kernel function that can cater for this varying length (of proteins) problem without limiting the proteins to a fixed length vector. This would provide SVM more relevant and useful information for protein fold recognition.

In order to find a decision boundary between two folds, SVM attempts to maximize the margin between the folds, and choose linear separations in a feature space. The classification of some known point in input space \mathbf{x}_i is y_i which is defined to be either -1 or $+1$. If \mathbf{x}' is a point in input space with unknown

classification then

$$y' = \text{sign} \left(\sum_{i=1}^n \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}') + b \right) \quad (3)$$

where y' is the predicted class of point \mathbf{x}' . The function $K()$ is the kernel; n is the number of support vectors; α_i are adjustable weights and b is a bias. We use libsvm (Chang and Lin, 2011) for training and testing with our kernel function.

6. An illustration of alignment method using a toy problem

In order to illustrate the alignment method, let us consider a toy example of two protein sequences $P = \text{VARA}$ and $Q = \text{VVARA}$ of corresponding length $L_1 = 4$ and $L_2 = 5$, respectively. Note that we assume that the toy proteins are made of 3 amino acids A, R and V. Tables 1a and 1b show the PSSM of these proteins.

Let p_i (for $i = 1, \dots, 4$) and q_j (for $j = 1, \dots, 5$) are the row vectors of PSSMs of P and Q , respectively. To compute the dissimilarity distance between row 1 of Table 1a ($p_1 = [1, 5, 6]$) and row 1 of Table 1b ($q_1 = [2, 3, 2]$), we employ Eq. (1) as follows:

$$\begin{aligned} d(p_1, q_1) &= 1 - \frac{p_1 q_1^\top}{\sqrt{p_1 p_1^\top q_1 q_1^\top}} \\ d(p_1, q_1) &= 1 - 0.8933 = 0.1067 \\ (\text{since } p_1 q_1^\top &= 29; p_1 p_1^\top = 62 \text{ and } q_1 q_1^\top = 17) \end{aligned} \quad (4)$$

In a similar way, dissimilarity distance can be computed between all the rows of Tables 1a and 1b. This would give similarity matrix S as follows:

$$S = \begin{bmatrix} 0.1067 & 1.0618 & 0.1171 & 0.1991 & 1.2272 \\ 0.3789 & 0.1484 & 0.6206 & 0.3479 & 0.3701 \\ 0.0541 & 0.3301 & 0.1372 & 0.2767 & 0.4870 \\ 0.3031 & 0.0677 & 0.4029 & 0.5490 & 0.1685 \end{bmatrix} \quad (5)$$

Dissimilarity matrix S is used in computing cumulative dissimilarity matrix D using dynamic programming (Eq. (2)) to find the minimum cost path (alignment path) as follows:

$$\begin{aligned} D_{11} &= \min(D_{01}, D_{10}, D_{00}) + S_{11} \\ &= S_{11} = 0.1067 \\ (\text{since, and do not exist and considered as empty}) \quad (6) \end{aligned}$$

In a similar way, we can compute $D_{21} = D_{11} + S_{21} = 0.4856$; $D_{12} = D_{11} + S_{12} = 1.1685$ and $D_{22} = \min(D_{12}, D_{21}, D_{11}) + S_{22} = 0.1067 +$

Table 1a
PSSM of the protein P .

| Amino acids | A | R | V |
|-------------|---|---|----|
| V | 1 | 5 | 6 |
| A | 4 | 6 | -3 |
| R | 3 | 3 | 1 |
| A | 5 | 3 | -1 |

Table 1b
PSSM of the protein Q .

| Amino acids | A | R | V |
|-------------|---|---|----|
| V | 2 | 3 | 2 |
| V | 5 | 2 | -3 |
| A | 5 | 4 | 6 |
| R | 0 | 4 | 1 |
| A | 2 | 0 | -1 |

$0.1484 = 0.2552$. The computed matrix D is given as follows:

$$D = \begin{bmatrix} 0.1067 & 1.1685 & 1.2856 & 1.4847 & 2.7119 \\ 0.4856 & 0.2551 & 0.8757 & 1.2236 & 1.5937 \\ 0.5397 & 0.5852 & 0.3923 & 0.6690 & 1.1560 \\ 0.8428 & 0.6074 & 0.7952 & 0.9413 & 0.8375 \end{bmatrix} \quad (7)$$

By using matrix D , the distance between two proteins can be computed which is simply given by $D_{dtw}(P, Q) = D(4, 5) = 0.8375$. Suppose the kernel parameter $\gamma = 10$ (evaluated by doing cross-validation on the training set) then kernel distance would be $K(P, Q) = \exp(-D_{dtw}(P, Q)^2/\gamma^2) = \exp(-0.8375^2/10^2) = 0.9930$. If $K(P, Q) = 1$ then it translates that proteins P and Q are very similar to each other. Further, if there are n training data then it will give $n \times n$ kernel matrix K which will be processed through SVM classifier for its parameter estimation.

7. Results and discussions

We carried out experiments on 3 benchmark datasets: DD, TG and EDD, to show the effectiveness of our proposed feature extraction method. We employ SVM classifier from libsvm (Chang and Lin, 2011) to find the accuracy of protein fold recognition where the accuracy is defined as the percentage of correctly recognized proteins to all the proteins of the test set. The SVM classifier is widely used in classification task. It finds maximum margin hyper-plane to minimize classification error. For the SVM classifier, kernel K is used. The kernel and SVM parameters, gamma and C , are optimized using grid search. In statistical prediction, the following three procedures are often used to examine a predictor for its effectiveness in practical application: independent dataset test, subsampling test, and jackknife test. However, of the three test procedures, the jackknife test is deemed the least arbitrary that can always yield a unique result for a given benchmark dataset as elaborated in Chou and Shen (2010). However, to reduce the computational time, we adopted the k -fold cross-validation in this study as done by many investigators with SVM as the prediction engine. We use datasets to perform k -fold cross-validation for $k = 5, 6, 7, 8, 9$ and 10. For statistical stability we performed 50 times k -fold cross-validation in this paper.¹

The proposed feature extraction method has been compared with several other feature extraction methods and the results have been shown in Tables 2–4. The following feature sets are considered for the experiment: PF1, PF2 (Ghanty and Pal, 2009), PF (Yang et al., 2011), Occurrence (O) (Taguchi and Gromiha, 2007), AAC, AAC+HXPZV (Ding and Dubchak, 2001), ACC (Dong et al., 2009), mono-gram and bi-gram (Sharma et al., 2013b). We have also updated the protein sequences to get the consensus sequence by using their corresponding PSSMs; i.e., each amino acid of a protein sequence is replaced by the amino acid that has the highest probability in PSSM. After this updating procedure, we have used the same feature extraction techniques (PF1, PF2, PF, O, AAC and AAC+HXPZV) again to obtain the recognition

¹ In statistical prediction, the following three cross-validation methods are often used to examine a predictor for its effectiveness in practical application: independent dataset test, subsampling or k -fold crossover test, and jackknife test. However, of the three test methods, the jackknife test is deemed the least arbitrary that can always yield a unique result for a given benchmark dataset as elaborated in Chou and Shen (2010) and demonstrated by Eqs. (28)–(30) in Chou, 2011. Accordingly, the jackknife test has been increasingly used and widely recognized by investigators to examine the quality of various predictors (see, e.g., Esmaeili et al., 2010; Chen et al., 2012, 2013; Feng et al., 2013; Hajisharifi et al., 2014; Chou et al., 2012). However, to reduce the computational time, we adopted the k -fold cross-validation in this study as done by many investigators with SVM as the prediction engine.

Table 2

Recognition accuracy by k -fold cross validation procedure for various feature extraction techniques using SVM classifier on DD-dataset.

| Feature sets | $k=5$ | $k=6$ | $k=7$ | $k=8$ | $k=9$ | $k=10$ |
|---|--------------|--------------|--------------|--------------|--------------|--------------|
| PF1 (Ghanty and Pal, 2009) | 48.6 46.3 | 49.1 47.0 | 49.5 47.5 | 50.1 47.7 | 50.5 47.9 | 50.6 48.2 |
| PF (Yang et al., 2011) | 51.2 | 52.2 | 52.6 | 52.9 | 53.4 | 53.4 |
| O (Taguchi and Gromiha, 2007) | 49.7 | 50.4 | 50.8 | 50.8 | 51.1 | 51.0 |
| AAC (Ding and Dubchak, 2001) | 43.6 | 43.9 | 44.2 | 44.8 | 44.6 | 45.1 |
| AAC+HXPZV ⁺ (Ding and Dubchak, 2001) | 45.1 | 46.2 | 46.5 | 46.8 | 46.9 | 47.2 |
| ACC (Dong et al., 2009) | 65.7 | 66.6 | 66.8 | 67.5 | 67.7 | 68.0 |
| PSSM+PF1 | 62.5 | 63.2 | 63.7 | 64.2 | 64.5 | 64.6 |
| PSSM+PF2 | 62.7 | 63.3 | 64.1 | 64.2 | 64.6 | 64.7 |
| PSSM+PF | 65.5 | 66.2 | 66.5 | 66.9 | 67.1 | 67.5 |
| PSSM+O | 62.5 | 62.1 | 62.5 | 62.9 | 63.4 | 63.5 |
| PSSM+AAC | 57.5 | 58.1 | 58.4 | 58.7 | 59.1 | 59.2 |
| PSSM+AAC+HXPZV | 55.9 | 56.9 | 57.1 | 57.7 | 58.0 | 58.2 |
| Mono-gram (Sharma et al., 2013b) | 67.7 | 68.4 | 68.6 | 69.1 | 69.4 | 69.6 |
| Bi-gram (Sharma et al., 2013b) | 72.6 | 73.1 | 73.7 | 73.7 | 74.1 | 74.1 |
| Alignment method (this paper) | 72.6 | 73.5 | 73.8 | 74.2 | 74.7 | 74.7 |

Table 3

Recognition accuracy (in percentage) by k -fold cross validation procedure for various feature extraction techniques using SVM classifier on TG dataset.

| Feature sets | $k=5$ | $k=6$ | $k=7$ | $k=8$ | $k=9$ | $k=10$ |
|------------------------------------|-------------|-------------|-------------|-------------|-------------|-------------|
| PF1 (Ghanty and Pal, 2009) | 38.1 | 38.4 | 38.6 | 38.7 | 38.8 | 38.8 |
| PF2(Ghanty and Pal, 2009) | 38.0 | 38.4 | 38.5 | 38.6 | 38.7 | 38.8 |
| PF (Yang et al., 2011) | 42.3 | 42.6 | 42.7 | 43.0 | 43.0 | 43.1 |
| O (Taguchi and Gromiha, 2007) | 35.8 | 36.1 | 36.2 | 36.1 | 36.3 | 36.3 |
| AAC (Ding and Dubchak, 2001) | 31.5 | 31.5 | 31.7 | 31.8 | 31.9 | 32.0 |
| AAC+HXPZV (Ding and Dubchak, 2001) | 35.7 | 36.0 | 36.1 | 36.2 | 36.3 | 36.3 |
| ACC (Dong et al., 2009) | 64.9 | 65.4 | 65.9 | 66.2 | 66.4 | 66.4 |
| PSSM+PF1 | 51.1 | 51.5 | 52.0 | 52.3 | 52.4 | 52.7 |
| PSSM+PF2 | 50.2 | 50.4 | 50.7 | 50.8 | 51.0 | 51.1 |
| PSSM+PF | 57.2 | 57.8 | 58.0 | 58.3 | 58.5 | 58.8 |
| PSSM+O | 46.0 | 46.3 | 46.5 | 46.5 | 46.7 | 46.7 |
| PSSM+AAC | 43.2 | 43.5 | 43.6 | 43.8 | 43.8 | 44.0 |
| PSSM+AAC+HXPZV | 45.6 | 45.9 | 46.0 | 46.2 | 46.3 | 46.6 |
| Mono-gram (Sharma et al., 2013b) | 57.2 | 57.3 | 58.2 | 58.4 | 58.8 | 58.8 |
| Bi-gram (Sharma et al., 2013b) | 67.1 | 67.5 | 67.6 | 67.8 | 68.1 | 68.1 |
| Alignment method (this paper) | 72.0 | 72.7 | 73.0 | 73.5 | 73.6 | 74.0 |

Table 4

Recognition accuracy (in percentage) by k -fold cross validation procedure for various feature extraction techniques using SVM classifier on EDD dataset.

| Feature sets | $k=5$ | $k=6$ | $k=7$ | $k=8$ | $k=9$ | $k=10$ |
|------------------------------------|-------------|-------------|-------------|-------------|-------------|-------------|
| PF1 (Ghanty and Pal, 2009) | 50.2 | 50.5 | 50.5 | 50.7 | 50.8 | 50.8 |
| PF2 (Ghanty and Pal, 2009) | 49.3 | 49.5 | 49.7 | 49.8 | 49.8 | 49.9 |
| PF (Yang et al., 2011) | 54.7 | 55.0 | 55.2 | 55.4 | 55.5 | 55.6 |
| O (Taguchi and Gromiha, 2007) | 46.4 | 46.6 | 46.6 | 46.7 | 46.7 | 46.9 |
| AAC (Ding and Dubchak, 2001) | 40.3 | 40.6 | 40.7 | 40.7 | 40.9 | 40.9 |
| AAC+HXPZV (Ding and Dubchak, 2001) | 40.2 | 40.4 | 40.6 | 40.7 | 40.9 | 40.9 |
| ACC (Dong et al., 2009) | 84.9 | 85.2 | 85.4 | 85.6 | 85.8 | 85.9 |
| PSSM+PF1 | 74.1 | 74.5 | 74.7 | 75.0 | 75.1 | 75.2 |
| PSSM+PF2 | 73.7 | 74.1 | 74.5 | 74.6 | 74.7 | 74.9 |
| PSSM+PF | 78.2 | 78.6 | 78.8 | 79.0 | 79.1 | 79.3 |
| PSSM+O | 67.6 | 68.0 | 68.1 | 68.3 | 68.3 | 68.5 |
| PSSM+AAC | 60.9 | 61.3 | 61.5 | 61.6 | 61.7 | 61.9 |
| PSSM+AAC+HXPZV | 66.7 | 67.2 | 67.4 | 67.7 | 67.8 | 67.9 |
| Mono-gram (Sharma et al., 2013b) | 76.2 | 76.3 | 76.6 | 76.8 | 77.0 | 76.9 |
| Bi-gram (Sharma et al., 2013b) | 83.6 | 84.0 | 84.1 | 84.3 | 84.3 | 84.5 |
| Alignment method (this paper) | 89.4 | 89.7 | 89.9 | 90.0 | 90.1 | 90.2 |

performance. In Tables 2–4, we have placed the results for PSSM updated protein sequences (or the consensus sequence) in the columns 2–7 of the row of PSSM+FEAT, where FEAT is any feature extraction technique. The highest recognition accuracy of a particular k -fold cross-validation is mentioned in bold face. It can be observed from Table 2 (on DD dataset) that the highest accuracy of protein fold recognition is 74.7% which is obtained by alignment

method (when $k=9$ and $k=10$) followed by bi-gram method which is 74.1% (when $k=10$). Besides the enhancement achieved compared to bi-gram and mono-gram methods that we have recently proposed in our previous study, we achieved an improvement of 7% prediction accuracy compared to ACC method (which has been proposed by Dong et al. (2009) and remained unbeaten ever since). In general, the protein fold prediction accuracy by

alignment method is around 0.6% to 29% higher than other methods.

Table 3 shows accuracy on TG dataset. It can be observed from the table that the highest accuracy of protein fold recognition is by alignment method. For the first time, we have enhanced the prediction accuracy to over 70% when the sequential similarity is less than 25%. We report 74.0% (when $k = 10$) prediction accuracy for TG benchmark followed by bi-gram method which is 68.1% (Sharma et al., 2013b). In general, the accuracy is around 5.9% to 40.5% higher than other feature extraction methods.

Next, **Table 4** depicts protein fold recognition accuracy on EDD dataset. It can be seen from the table that the highest accuracy is again obtained by alignment method. For the first time, we have enhanced the protein fold prediction accuracy to over 90% when the sequential similarity rate is less than 40%. We report 90.2% (when $k = 10$) prediction accuracy for the EDD benchmark followed by ACC which is 85.9% (Dong et al., 2009). In general, the protein fold prediction enhancement achieved by alignment method compared to previously reported results for the EDD benchmark is from 4.3% to 49.2%.

It can be observed that for TG and EDD datasets the performance by the alignment method was comparatively better than the performance on DD dataset. This could be because TG and EDD are large datasets which enough samples which lead to better parameter estimation on training data.

In order to study the statistical significance of the prediction accuracy enhancement reported in this study, we conduct the paired *t*-test on our achieved results compared to the highest results reported in the literature. Associated probability value achieved for the paired *t*-test is $p = 0.03$ which confirms the statistical significance of our reported enhancement in this study compared to the state-of-the-art results found in the literature for protein fold recognition.

Furthermore, we have conducted precision, sensitivity and specificity analysis of all the features used in this paper over 3 datasets to provide more information about the statistical significant of our achieved results (Kurgan and Homaeian, 2006; Dehzangi et al., 2014). Sensitivity measures the ratio of correctly classified samples to the whole number of test samples for each class which are classified as correct samples and calculated as follows:

$$\text{Sensitivity} = \frac{TP}{TP+FN} \times 100, \quad (8)$$

while TP represents true positive and FN represents false negative samples. Specificity, as other evaluation criterion used in this study measures the ratio of correctly rejected samples to the whole number of rejected test samples and is calculated as follows:

$$\text{Specificity} = \frac{TN}{TN+FP} \times 100, \quad (9)$$

where TN represents true negative and FP represents false positive. The third evaluation criteria used in this study is precision which represents how relevant the number of TP is to the whole number of positive prediction and is calculated as follows:

$$\text{Precision} = \frac{TP}{TP+FP} \times 100. \quad (10)$$

The sensitivity, specificity and precision are computed for each class and then average over all the classes are computed and reported in Figs. 3–5.

Further information regarding these three evaluation criteria can be found in Kurgan and Homaeian (2006), and Dehzangi et al. (2013d). Fig. 3, depicts the analysis on DD dataset, Fig. 4 on EDD dataset and Fig. 5 on TG dataset. It can be observed from Figs. 3–5 that specificity is high for all the feature sets. However, precision

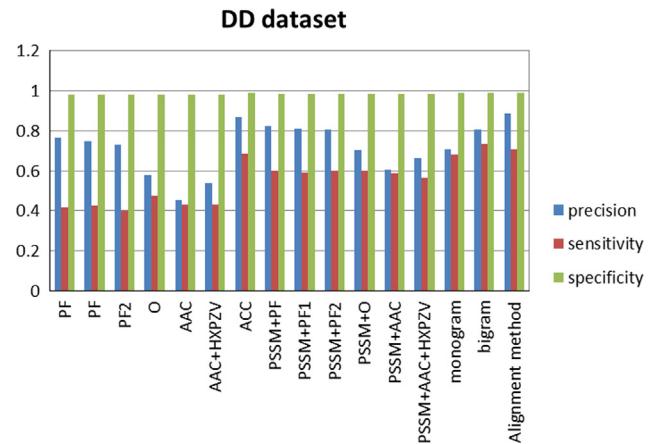


Fig. 3. Precision, sensitivity and specificity of all feature sets on DD dataset.

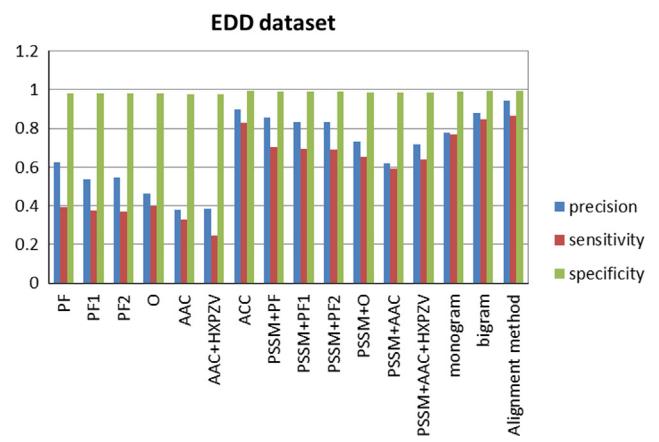


Fig. 4. Precision, sensitivity and specificity of all feature sets on EDD dataset.

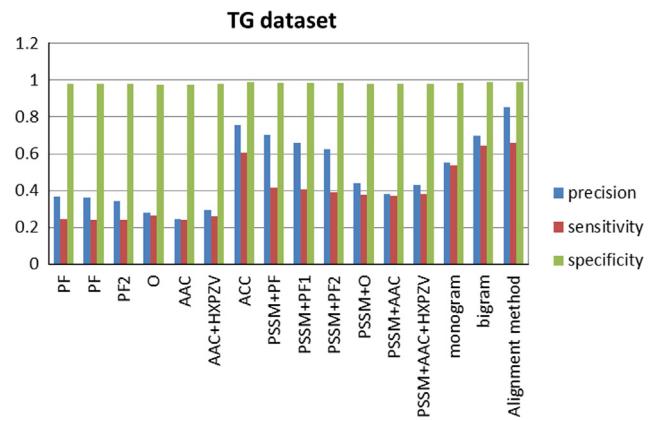


Fig. 5. Precision, sensitivity and specificity of all feature sets on TG dataset.

and sensitivity varies. For all the datasets, precision and sensitivity are quite promising for alignment method.

Since user-friendly and publicly accessible web-servers represent the future direction for developing practically more useful models, simulated methods, or predictors (Chou and Shen, 2009; Lin and Lapointe, 2013), we shall make efforts in our future work to provide a web-server for the method presented in this paper. However, since it is very useful to have accessible codes for developing practically more useful models, we have provided

Matlab based code for our method <http://www.staff.usp.ac.fj/~sharma_al/index.htm>.

8. Conclusion

In this work, we developed feature extraction method using amino acid alignment scheme. The technique used PSSM log probabilities of protein sequences, to determine the distance between two proteins. This method has been compared with several other existing feature extraction methods and very promising results have been obtained. It was noted that the proposed method outperformed existing methods for three commonly used benchmarks. We have reported 74.6% prediction accuracy on DD benchmark. For the first time, we have also achieved to over 70% and 90% prediction accuracies for protein fold recognition when the sequential similar rates are less than 25% and 40%, respectively. We observed 74.0% and 90.2% prediction accuracies for TG and EDD benchmarks, respectively. These reported results are over 5.9% and 4.3% better than the best results reported for these two benchmark datasets.

References

- Altschul, S.F., Madden, T.L., Schaffer, A.A., Zhang, J.H., Zhang, Z., Miller, W., Lipman, D.J., 1997. Gapped blast and psi-blast: a new generation of protein database search programs. *Nucleic Acids Res.* 17, 3389–3402.
- Bishop, C.M., 2006. *Pattern Recognition and Machine Learning*. Springer Science, NY.
- Bouchaffra, D., Tan, J., 2006. Protein fold recognition using a structural Hidden Markov Model. In: Proceedings of the 18th International Conference on Pattern Recognition, pp. 186–189.
- Cao, D.S., Xu, Q.S., Liang, Y.Z., 2013. Propy: a tool to generate various modes of Chou's PseAAC. *Bioinformatics* 29, 960–962.
- Chang, C.-C., Lin, C.-J., 2011. LIBSVM: a library for support vector machines. *ACM Trans. Intell. Syst. Technol.* 2 (3), 27:1–27:27.
- Chen, K., Zhang, X., Yang, M.Q., Yang, J.Y., 2007. Ensemble of probabilistic neural networks for protein fold recognition. In: Proceedings of the 7th IEEE International Conference on Bioinformatics and Bioengineering (BIBE), pp. 66–70.
- Chen, W., Lin, H., Feng, P.M., Ding, C., Zuo, Y.C., et al., 2012. iNuc-physchem: a sequence-based predictor for identifying nucleosomes via physicochemical properties. *PLoS One* 7, e47843.
- Chen, W., Feng, P.M., Lin, H., et al., 2013. iRSpot-PseDNC: identify recombination spots with pseudo dinucleotide composition. *Nucleic Acids Res.* 41, e69.
- Chinnasamy, A., Sung, W.K., Mittal, A., 2005. Protein structure and fold prediction using tree-augmented naive Bayesian classifier. *J. Bioinf. Comput. Biol.* 3 (4), 803–819.
- Chmielnicki, W., Stapor, K., 2012. A hybrid discriminative-generative approach to protein fold recognition. *Neurocomputing* 75, 194–198.
- Chou, K.C., 2001. Prediction of protein cellular attributes using pseudo amino acid composition. *Proteins* 43, 246–255 (erratum: 2001, vol. 44, 60).
- Chou, K.C., Cai, Y.D., 2004. Predicting protein structural class by functional domain composition. *Biochem. Biophys. Res. Commun.* 321, 1007–1009 (Corrigendum: ibid, 2005, vol. 329, 1362).
- Chou, K.C., Shen, H.B., 2007. MemType-2L: a web server for predicting membrane proteins and their types by incorporating evolution information through Pse-PSSM. *Biochem. Biophys. Res. Commun.* 360, 339–345.
- Chou, K.C., Shen, H.B., 2008. ProtIdent: a web server for identifying proteases and their types by fusing functional domain and sequential evolution information. *Biochem. Biophys. Res. Commun.* 376, 321–325.
- Chou, K.C., Shen, H.B., 2010. Cell-PLoc: a package of web servers for predicting subcellular localization of proteins in various organisms (updated version: Cell-PLoc 2.0: an improved package of web-servers for predicting subcellular localization of proteins in various organisms). *Nat. Sci.* 2, 1090–1103, <http://dx.doi.org/10.4236/ns.2010.210136> (Nature Protocols, 2008, 3, 153–162).
- Chou, K.C., Shen, H.B., 2009. Review: recent advances in developing web-servers for predicting protein attributes. *Nat. Sci.* 2, 63–92.
- Chou, K.C., 2005. Using amphiphilic pseudo amino acid composition to predict enzyme subfamily classes. *Bioinformatics* 21, 10–19.
- Chou, K.C., 2011. Some remarks on protein attribute prediction and pseudo amino acid composition (50th Anniversary Year Review). *J. Theor. Biol.* 273, 236–247.
- Chou, K.C., Wu, Z.C., Xiao, X., 2011. iLoc-Euk: a multi-label classifier for predicting the subcellular localization of singleplex and multiplex eukaryotic proteins. *PLoS One* 6, e18258.
- Chou, K.C., Wu, Z.C., Xiao, X., 2012. iLoc-Hum: using accumulation-label scale to predict subcellular locations of human proteins with both single and multiple sites. *Mol. Biosyst.* 8, 629–641.
- Dehzangi, A., Ammuaisuk, S.P., 2011. Fold prediction problem: the application of new physical and physicochemical-based features. *Protein Pept. Lett.* 18, 174–185.
- Dehzangi, A., Ammuaisuk, S.P., Dehzangi, O., 2010. Enhancing protein fold prediction accuracy by using ensemble of different classifiers. *Aust. J. Intell. Inf. Process. Syst.* 26 (4), 32–40.
- Dehzangi, A., Ammuaisuk, S.P., Ng, K.H., Mohandes, E., 2009. Protein fold prediction problem using ensemble of classifiers. In: Proceedings of the 16th International Conference on Neural Information Processing, Part II, pp. 503–511.
- Dehzangi, A., Karamizadeh, 2011. Solving protein fold prediction problem using fusion of heterogeneous classifiers. *Inf. Int. Interdiscip.* J. 14 (11), 3611–3622.
- Dehzangi, A., Paliwal, K.K., Sharma, A., Dehzangi, O., Sattar, A., 2013a. A combination of feature extraction methods with an ensemble of different classifiers for protein structural class prediction problem. *IEEE/ACM Trans. Comput. Biol. Bioinf.* 10 (3), v564–v575.
- Dehzangi, A., Paliwal, K.K., Lyons, J., Sharma, A., Sattar, A., 2013b. Exploring potential discriminatory information embedded in psm to enhance protein structural class prediction accuracy. In: Proceeding of the Pattern Recognition in Bioinformatics. PRIB 2013, LNBI 7986, pp. 208–219.
- Dehzangi, A., Paliwal, K.K., Lyons, J., Sharma, A., Sattar, A., 2013c. Enhancing protein fold prediction accuracy using evolutionary and structural features. In: Proceeding of the Pattern Recognition in Bioinformatics. PRIB 2013, LNBI 7986, pp. 196–207.
- Dehzangi, A., Paliwal, K.K., Lyons, J., Sharma, A., Sattar, A., 2014. Proposing a highly accurate protein structural class predictor using segmentation-based features. *BMC Genomics* 15 (Suppl 1), S2.
- Deschavanne, P., Tuffery, P., 2009. Enhanced protein fold recognition using a structural alphabet. *Proteins: Struct. Funct. Bioinf.* 76, 129–137.
- Ding, C., Dubchak, I., 2001. Multi-class protein fold recognition using support vector machines and neural networks. *Bioinformatics* 17 (4), 349–358.
- Ding, Y.S., Zhang, T.L., 2008. Using Chou's pseudo amino acid composition to predict subcellular localization of apoptosis proteins: an approach with immune genetic algorithm-based ensemble classifier. *Patt. Recog. Lett.* 29, 1887–1892.
- Dong, Q., Zhou, S., Guan, J., 2009. A new taxonomy-based protein fold recognition approach based on autocross-covariance transformation. *Bioinformatics* 25 (20), 2655–2662.
- Du, P., Wang, X., Xu, C., Gao, Y., 2012. PseAAC-builder: a cross-platform stand-alone program for generating various special Chou's pseudo-amino acid compositions. *Anal. Biochem.* 425, 117–119.
- Dubchak, I., Muchnik, I., Kim, S.K., 1997. Protein folding class predictor for SCOP: approach based on global descriptors. In: Proceedings of 5th International Conference on Intelligent Systems for Molecular Biology, pp. 104–107.
- Esmaeili, M., Mohabatkar, H., Mohsenzadeh, S., 2010. Using the concept of Chou's pseudo amino acid composition for risk type prediction of human papilloma-viruses. *J. Theor. Biol.* 263, 203–209.
- Feng, P.M., Chen, W., Lin, H., et al., 2013. iHSP-PseRAAC: identifying the heat shock protein families using pseudo reduced amino acid alphabet composition. *Anal. Biochem.* 442, 118–125.
- Ghanty, P., Pal, N.R., 2009. Prediction of protein folds: extraction of new features, dimensionality reduction, and fusion of heterogeneous classifiers. *IEEE Trans. Nano Biosci.* 8, 100–110.
- Hajisharifi, Z., Pirayee, M., Mohammad Beigi, M., Behbahani, M., Mohabatkar, H., 2014. Predicting anticancer peptides with Chou's pseudo amino acid composition and investigating their mutagenicity via Ames test. *J. Theor. Biol.* 341, 34–40.
- Huang, J.T., Tian, J., 2006. Amino acid sequence predicts folding rate for middle-size two-state proteins. *Proteins: Struct. Funct. Bioinf.* 63 (3), 551–554.
- Kavousi, K., Moshiri, B., Sadeghi, M., Araabi, B.N., Moosavi-Movahedi, A.A., 2011. A protein fold classifier formed by fusing different modes of pseudo amino acid composition via PSSM. *Comput. Biol. Chem.* 35 (1), 1–9.
- Kecman, V., Yang, T., 2009. Protein fold recognition with adaptive local hyper plane Algorithm. *Computational Intelligence in Bioinformatics and Computational Biology, CIBCB '09. IEEE Symposium*, pp. 75–78.
- Klein, P., 1986. Prediction of protein structural class by discriminant analysis. *Biochim Biophys Acta* 874, 205–215.
- Krishnaraj, Y., Reddy, C.K., 2008. Boosting methods for protein fold recognition: an empirical comparison. *IEEE Int. Conf. Bioinf. Biomed.*, 393–396, <http://dx.doi.org/10.1109/BIBM.2008.83>
- Kurgan, L.A., Homaeian, L., 2006. Prediction of structural classes for protein sequences and domains—impact of predictions algorithms, sequence representation and homology, and test procedures on accuracy. *Pattern Recognit.* 39, 2323–2343.
- Kurgan, L.A., Zhang, T., Zhang, H., Shen, S., Ruan, J., 2008. Secondary structure-based assignment of the protein structural classes. *Amino Acids* 35, 551–564.
- Lin, S.X., Lapointe, J., 2013. Theoretical and experimental biology in one. *J. Biomed. Sci. Eng. (JBSE)* 6, 435–442.
- Lin, W.Z., Fang, J.A., Xiao, X., et al., 2012. Predicting secretory proteins of malaria parasite by incorporating sequence evolution information into pseudo amino acid composition via grey system model. *PLoS One* 7, e49040.
- Lin, W.Z., Fang, J.A., Xiao, X., et al., 2013. iLoc-Animal: a multi-label learning classifier for predicting subcellular localization of animal proteins. *Mol. BioSyst.* 9, 634–644.
- Liu, T., Geng, X., Zheng, X., Li, R., Wang, J., 2012. Accurate prediction of protein structural class using autocovariance transformation of PSI-BLAST profiles. *Amino Acids* 42, 2243–2249.
- Min, J.L., Xiao, X., Chou, K.C., 2013. iEzy-drug: a web server for identifying the interaction between enzymes and drugs in cellular networking. *BioMed. Res. Int.* 2013, 701317.
- Mohabatkar, H., Beigi, M.M., Abdolahi, K., Mohsenzadeh, S., 2013. Prediction of allergenic proteins by means of the concept of Chou's pseudo amino acid composition and a machine learning approach. *Med. Chem.* 9, 133–137.

- Mohammad Beigi, M., Behjati, M., Mohabatkar, H., 2011. Prediction of metalloprotease family based on the concept of Chou's pseudo amino acid composition using a machine learning approach. *J. Struct. Funct. Genomics* 12, 191–197.
- Mohabatkar, H., Mohammad Beigi, M., Esmaeili, A., 2011. Prediction of GABA (A) receptor proteins using the concept of Chou's pseudo-amino acid composition and support vector machine. *J. Theor. Biol.* 281, 18–23.
- Najmanovich, R., Kuttner, J., Sobolev, V., Edelman, M., 2000. Side-chain flexibility in proteins upon ligand binding. *Proteins: Struct. Funct. Bioinf.* 39 (3), 261–268.
- Nanni, L., Lumini, A., 2008. Genetic programming for creating Chou's pseudo amino acid based features for submitochondria localization. *Amino Acids* 34, 653–660.
- Nanni, L., Lumini, A., Gupta, D., Garg, A., 2012. Identifying bacterial virulent proteins by fusing a set of classifiers based on variants of Chou's pseudo amino acid composition and on evolutionary information. *IEEE/ACM Trans. Comput. Biol. Bioinf.* 9, 467–475.
- Ohlson, T., Wallner, B., Elofsson, A., 2004. Profile–profile methods provide improved fold-recognition: a study of different profile–profile alignment methods. *Proteins: Struct. Funct. Bioinf.* 57, 188–197.
- Paliwal, K.K., Sharma, A., Lyons, J., Dehzangi, A., 2014. A tri-gram based feature extraction technique using linear probabilities of position specific scoring matrix for protein fold recognition. *IEEE Trans. Nanobiosci.* 13 (1).
- Paliwal, K.K., Sharma, A., 2011. Approximate LDA technique for dimensionality reduction in the small sample size case. *J. Pattern Recognit. Res.* 6 (2), 298–306.
- Paliwal, K.K., Sharma, A., 2012. Improved pseudoinverse linear discriminant analysis method for dimensionality reduction. *Int. J. Pattern Recognit. Artif. Intell.* 26 (1), 1250002-1–1250002-9.
- Qiu, W.R., Xiao, X., et al., 2014. iRSpot-TNCPseAAC: identify recombination spots with trinucleotide composition and pseudo amino acid components. *Int. J. Mol. Sci.* 15, 1746–1766.
- Sahu, S.S., Panda, G., 2010. A novel feature representation method based on Chou's pseudo amino acid composition for protein structural class prediction. *Comput. Biol. Chem.* 34, 320–327.
- Sharma, A., Paliwal, K.K., Dehzangi, A., Lyons, J., Imoto, S., Miyano, S., 2013a. A strategy to select suitable physicochemical attributes of amino acids for protein fold recognition. *BMC Bioinf.* 14, 233.
- Sharma, A., Lyons, J., Dehzangi, A., Paliwal, K.K., 2013b. A feature extraction technique using bi-gram probabilities of position specific scoring matrix for protein fold recognition. *J. Theor. Biol.* 320 (7), 41–46.
- Sharma, A., Imoto, S., Miyano, S., Sharma, V., 2012a. Null space based feature selection method for gene expression data. *Int. J. Mach. Learn. Cybern.* 3 (4), 269–276.
- Sharma, A., Imoto, S., Miyano, S., 2012b. A top-r feature selection algorithm for microarray gene expression data. *IEEE/ACM Trans. Comput. Biol. Bioinf.* 9 (3), 754–764.
- Sharma, A., Paliwal, K.K., 2012a. A two-stage linear discriminant analysis for face-recognition. *Pattern Recognit. Lett.* 33 (9), 1157–1162.
- Sharma, A., Imoto, S., Miyano, S., 2012c. A filter based feature selection algorithm using null space of covariance matrix for DNA microarray gene expression data. *Curr. Bioinf.* 7 (3), 289–294.
- Sharma, A., Imoto, S., Miyano, S., 2012d. A between-class overlapping filter-based method for transcriptome data analysis. *J. Bioinf. Comput. Biol.* 10 (5), 1250010-1–1250010-20.
- Sharma, A., Paliwal, K.K., 2012b. A gene selection algorithm using Bayesian classification approach. *Am. J. Appl. Sci.* 9 (1), 127–131.
- Sharma, A., Paliwal, K.K., 2012c. A new perspective to null linear discriminant analysis method and its fast implementation using random matrix multiplication with scatter matrices. *Pattern Recognit.* 45 (6), 2205–2213.
- 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, <http://dx.doi.org/10.1007/s13042-012-0131-7>.
- Sharma, A., Paliwal, K.K., 2008. A gradient linear discriminant analysis for small sample sized problem. *Neural Process. Lett.* 27 (1), 17–24.
- Sharma, A., Koh, C.H., Imoto, S., Miyano, S., 2011. Strategy of finding optimal number of features on gene expression data. *Electron. Lett.* 47 (8), 480–482.
- Sharma, A., Paliwal, K.K., 2010. Regularisation of eigenfeatures by extrapolation of scatter-matrix in face-recognition problem. *IEE Electron. Lett.* 46 (10), 450–475.
- Sharma, A., Paliwal, K.K., 2007. Fast principal component analysis using fixed-point algorithm. *Pattern Recognit. Lett.* 28 (10), 1151–1155.
- Sharma, A., Paliwal, K.K., Onwubolu, G.C., 2006. Class-dependent PCA, LDA and MDC: a combined classifier for pattern classification. *Pattern Recognit.* 39 (7), 1215–1229.
- Shamim, M.T.A., Anwaruddin, M., Nagarajaram, H.A., 2007. Support vector machine-based classification of protein folds using the structural properties of amino acid residues and amino acid residue pairs. *Bioinformatics* 23 (24), 3320–3327.
- Shen, H.B., Chou, K.C., 2006. Ensemble classifier for protein fold pattern recognition. *Bioinformatics* 22, 1717–1722.
- Shen, H.B., Chou, K.C., 2007. EzyPred: A top-down approach for predicting enzyme functional classes and subclasses. *Biochem. Biophys. Res. Commun.* 364, 53–59.
- Shen, H.B., Chou, K.C., 2008. PseAAC: a flexible web-server for generating various kinds of protein pseudo amino acid composition. *Anal. Biochem.* 373, 386–388.
- Taguchi, Y.-H., Gromiha, M.M., 2007. Application of amino acid occurrence for discriminating different folding types of globular proteins. *BMC Bioinf.* 8, 404.
- Vapnik, V.N., 1995. *The Nature of Statistical Learning Theory*. Springer-Verlag, New York.
- Wang, Z.Z., Yuan, Z., 2000. How good is prediction of protein-structural class by the component-coupled method? *Proteins* 38, 165–175.
- Yang, T., Kecman, V., Cao, L., Zhang, C., Huang, J.Z., 2011. Margin-based ensemble classifier for protein fold recognition. *Expert Syst. Appl.* 38, 12348–12355.
- Ying, Y., Huang, K., Campbell, C., 2009. Enhanced protein fold recognition through a novel data integration approach. *BMC Bioinf.* 10 (1), 267.
- Valavanis, I.K., Spyrou, G.M., Nikita, K.S., 2010. A comparative study of multi-classification methods for protein fold recognition. *Int. J. Comput. Intell. Bioinf. Syst. Biol.* 1 (3), 332–346.
- Wu, Z.C., Xiao, X., et al., 2011. iLoc-Plant: a multi-label classifier for predicting the subcellular localization of plant proteins with both single and multiple sites. *Mol. BioSyst.* 7, 3287–3297.
- Wu, Z.C., Xiao, X., et al., 2012. iLoc-Gpos: a multi-layer classifier for predicting the subcellular localization of singleplex and multiplex gram-positive bacterial proteins. *Protein Pept. Lett.* 19, 4–14.
- Xiao, X., Min, J.L., Wang, P., et al., 2013. iCDI-PseFpt: identify the channel-drug interaction in cellular networking with PseAAC and molecular fingerprints. *J. Theor. Biol.* 337C, 71–79.
- Xu, Y., Shao, X.J., Wu, L.Y., Deng, N.Y., et al., 2013. iSNP-AAPair: incorporating amino acid pairwise coupling into PseAAC for predicting cysteine S-nitrosylation sites in proteins. *PeerJ* 1, e171.
- Xiao, X., Wu, Z.C., et al., 2011a. A multi-label classifier for predicting the subcellular localization of gram-negative bacterial proteins with both single and multiple sites. *PLoS One* 6, e20592.
- Xiao, X., Wu, Z.C., et al., 2011b. iLoc-Virus: a multi-label learning classifier for identifying the subcellular localization of virus proteins with both single and multiple sites. *J. Theor. Biol.* 284, 42–51.
- Zhang, S.W., Zhang, Y.L., Yang, H.F., Zhao, C.H., Pan, Q., 2008a. Using the concept of Chou's pseudo amino acid composition to predict protein subcellular localization: an approach by incorporating evolutionary information and von Neumann entropies. *Amino Acids* 34, 565–572.
- Zhang, H., Zhang, T., Gao, J., Ruan, J., Shen, S., Kurgan, L.A., 2012. Determination of protein folding kinetic types using sequence and predicted secondary structure and solvent accessibility. *Amino Acids* 42 (1), 271–283, <http://dx.doi.org/10.1007/s00726-010-0805-y>, Epub 2010 Nov 17.
- Zhang, T.L., Ding, Y.S., Chou, K.C., 2008b. Prediction protein structural classes with pseudo amino acid composition: approximate entropy and hydrophobicity pattern. *Theor. Biol.* 250, 186–193.