

See discussions, stats, and author profiles for this publication at: <http://www.researchgate.net/publication/5419532>

# Using support vector machine to predict $\beta$ - and $\gamma$ -turns in proteins

ARTICLE *in* JOURNAL OF COMPUTATIONAL CHEMISTRY · MAY 2008

Impact Factor: 3.6 · DOI: 10.1002/jcc.20929 · Source: PubMed

---

CITATIONS

26

---

DOWNLOADS

8

---

VIEWS

55

2 AUTHORS, INCLUDING:



Qianzhong Li

Inner Mongolia University

57 PUBLICATIONS 1,030 CITATIONS

SEE PROFILE

# Using Support Vector Machine to Predict $\beta$ - and $\gamma$ -Turns in Proteins

XIUZHEN HU,<sup>1,2</sup> QIANZHONG LI<sup>1</sup>

<sup>1</sup>Laboratory of Theoretical Biophysics, Department of Physics, College of Sciences and Technology, Inner Mongolia University, Hohhot, People's Republic of China

<sup>2</sup>Department of Physics, College of Sciences, Inner Mongolia University of Technology, Hohhot, People's Republic of China

Received 17 September 2007; Revised 22 December 2007; Accepted 27 December 2007

DOI 10.1002/jcc.20929

Published online 23 April 2008 in Wiley InterScience (www.interscience.wiley.com).

**Abstract:** By using the composite vector with increment of diversity, position conservation scoring function, and predictive secondary structures to express the information of sequence, a support vector machine (SVM) algorithm for predicting  $\beta$ - and  $\gamma$ -turns in the proteins is proposed. The 426 and 320 nonhomologous protein chains described by Guruprasad and Rajkumar (Guruprasad and Rajkumar J. Biosci 2000, 25,143) are used for training and testing the predictive model of the  $\beta$ - and  $\gamma$ -turns, respectively. The overall prediction accuracy and the Matthews correlation coefficient in 7-fold cross-validation are 79.8% and 0.47, respectively, for the  $\beta$ -turns. The overall prediction accuracy in 5-fold cross-validation is 61.0% for the  $\gamma$ -turns. These results are significantly higher than the other algorithms in the prediction of  $\beta$ - and  $\gamma$ -turns using the same datasets. In addition, the 547 and 823 nonhomologous protein chains described by Fuchs and Alix (Fuchs and Alix Proteins: Struct Funct Bioinform 2005, 59, 828) are used for training and testing the predictive model of the  $\beta$ - and  $\gamma$ -turns, and better results are obtained. This algorithm may be helpful to improve the performance of protein turns' prediction. To ensure the ability of the SVM method to correctly classify  $\beta$ -turn and non- $\beta$ -turn ( $\gamma$ -turn and non- $\gamma$ -turn), the receiver operating characteristic threshold independent measure curves are provided.

© 2008 Wiley Periodicals, Inc. J Comput Chem 29: 1867–1875, 2008

**Key words:** increment of diversity;  $\beta$ -turn;  $\gamma$ -turn; position conservation scoring function; support vector machine

## Introduction

Protein secondary structure prediction is a key step for predicting tertiary structure of proteins. In the past three decades, a large number of methods have been developed for predicting the regular secondary structures ( $\alpha$ -helix,  $\beta$ -strand) and coil in proteins.<sup>1–6</sup> However, they could not provide a directional change for the polypeptide chain. Therefore, the prediction of tight turns in proteins is as important as helix and strand prediction. The tight turns play an important role in protein, such as folding stability,<sup>7,8</sup> recognition,<sup>9,10</sup> and structure assembly.<sup>11</sup>

Tight turns can be divided into  $\delta$ -,  $\gamma$ -,  $\beta$ -,  $\alpha$ -, and  $\pi$ -turns according to the number of residues involved.<sup>12</sup> The  $\beta$ -turns are the most common turns in proteins. A  $\beta$ -turn comprises four consecutive residues, which does not form  $\alpha$ -helix, and the distance between  $C_{\alpha}(i)$  and  $C_{\alpha}(i + 3)$  is less than 7 Å. According to backbone dihedral angles in the inner residues,  $i + 1$  and  $i + 2$  will define different types of  $\beta$ -turns.

Because of  $\beta$ -turn's ability to reverse the direction of a protein chain to 180°, it is responsible for the compact globular shape of proteins. The  $\beta$ -turn is an important component of  $\beta$ -hairpin structure and plays a vital role in protein folds. Enhanc-

ing  $\beta$ -turn prediction can have a direct effect on molecular recognition studies and the identification of important structural motifs, such as  $\beta$ -hairpins. It also contributes indirectly to the overall prediction of protein tertiary structures.

Some methods have been developed for prediction of  $\beta$ -turns based on the statistical model and machine learning technique.<sup>13–22</sup> The typical  $\beta$ -turns prediction was made by Kaur and Raghava.<sup>13</sup> They used the same dataset with 426 protein chains constructed by Guruprasad and Rajkumar<sup>23</sup> and the same measures, compared with some other  $\beta$ -turn prediction methods, such as Chou-Fasman,<sup>24</sup> the 1–4 and 2–3 correlation model,<sup>18</sup> the sequence coupled model,<sup>17</sup> GORBTURN (v3.0),<sup>25,26</sup> and

**Correspondence to:** Q.-Z. Li; e-mail: qzli@imu.edu.cn

Contract/grant sponsor: National Natural Science Foundation of China; contract/grant number: 30560039

Contract/grant sponsor: Natural Science Foundation of the Inner Mongolia of China; contract/grant numbers: 200508010509, 200607010101.

Contract/grant sponsor: Project for Excellent Subject-Directors of Inner Mongolia Autonomous Region

BTPRED.<sup>16</sup> In addition, an improved neural network method, BetaTPred2, was developed by Kaur and Raghava.<sup>14</sup> In this method, a great improvement in prediction performance was that the Matthews correlation coefficient (Mcc) = 0.43 had been achieved by using the multiple sequence alignment as input instead of the single amino acid sequence. Farther, Cai et al.<sup>19</sup> and Lin et al.<sup>20</sup> used the support vector machine (SVM) and the Markov Chains theory to predict  $\beta$ -turns, and obtained significant results. Fuchs and Alix<sup>22</sup> predicted  $\beta$ -turns in the dataset of 426 proteins using propensities and multiple alignments. The obtained Mcc and overall prediction accuracy were 0.42 and 74.8%, respectively. The better prediction accuracy based on SVM was obtained using input parameters of the predicted secondary structure and multiple alignment information among these methods.<sup>21,27</sup> The overall prediction accuracy and Mcc in the 7-fold cross-validation were 77.3% and 0.45,<sup>21</sup> and 79.8% and 0.45,<sup>27</sup> respectively, for 426 nonhomologous protein chains described by Guruprasad and Rajkumar.<sup>23</sup>

Recently, the different kinds of turns were studied by Street et al.<sup>28</sup> Their results provided a molecular explanation for the observation that reverse turns between elements of regular secondary structure can be classified into a small number of discrete conformations. And Bornot and de Brevern<sup>29</sup> analyzed the distributions of  $\beta$ -turns according to different secondary structure assignment methods.

The  $\gamma$ -turn is the second most characterized and commonly found turn, after the  $\beta$ -turn. A  $\gamma$ -turn is defined as a three-residue turn with a hydrogen bond between the carbonyl oxygen of residue  $i$  and the hydrogen of the amide group of residue  $i + 2$ . There are two types of  $\gamma$ -turns: inverse and classic.<sup>30</sup>

$\gamma$ -Turns also play an important role in protein folding and recognition. Experimentation indicated that some folds are achieved by  $\gamma$ -turn reverse,<sup>31</sup> and some acceptor combining sites are also in the  $\gamma$ -turn.<sup>32</sup> Although  $\gamma$ -turn content is fewer, it contained important information on the molecule recognition. Therefore, based on the amino acid sequence, predicting  $\gamma$ -turns is significant.

Compared with  $\beta$ -turns, however,  $\gamma$ -turns are seldom investigated. This is because of the lower occurrence of  $\gamma$ -turns in proteins. On the basis of multiple alignment and predicted secondary structure information, the  $\gamma$ -turns in proteins were predicted by Kaur and Raghava using neural network method.<sup>33</sup> The Mcc was 0.17 in the 5-fold cross-validation and the corresponding area under the receiver operating characteristic (ROC) curves was 0.73 for 320 nonhomologous protein chains described by Guruprasad and Rajkumar.<sup>23</sup> Guruprasad et al.<sup>23,34</sup> used the Markov Chains theory to predict  $\gamma$ -turns, and obtained significant results. Pham et al.<sup>27</sup> used SVM method at the residue level and turn level to predict  $\gamma$ -turns. The Mcc was 0.13 in the 5-fold cross-validation, and the overall prediction accuracy was 79.9% for 320 nonhomologous protein chains.

In this study, based on the datasets described by Guruprasad and Rajkumar<sup>23</sup> and by Fuchs and Alix,<sup>22</sup> by using the composite vector at the residue level, including the ID'PCSF values and predictive secondary structure information as inputting parameters of the SVM, the LIBSVM program packages are applied to predict  $\beta$ - and  $\gamma$ -turns, respectively. More accurately predicted results are obtained.

## Materials and Methods

### Materials

The two datasets described in the work of Guruprasad and Rajkumar<sup>23</sup> are used. A total of 426 and 320 nonhomologous protein chains with resolution  $<2.0$  Å and sequence identity  $<25\%$  are used in our method. Three hundred and fifteen out of 320 protein chains are contained in the above 426 protein chains. Four hundred and twenty six and 320 nonhomologous protein chains are respectively used for the prediction of  $\beta$ - and  $\gamma$ -turns.

In addition, the two datasets described in the work of Fuchs and Alix<sup>22</sup> are used. Among a total of 547 and 823 nonhomologous protein chains with a resolution of  $<2.0$  Å and a sequence identity of  $<25\%$ , each chain contains one minimum  $\beta$ -turn, 543 and 819 proteins that contained respectively 7912 and 11257  $\beta$ -turns, and are used for  $\beta$ -turns prediction. The 346 and 536 nonhomologous protein chains in two datasets, each chain containing one minimum  $\gamma$ -turn, are used for  $\gamma$ -turns prediction. They contain 873 and 1303  $\gamma$ -turns, respectively.

Two hundred and ten out of 543 protein chains are contained in the above 819 protein chains. One hundred and ninety one out of 543 protein chains are contained in the above 426 protein chains. Ninety out of 819 protein chains are contained in the above 426 protein chains.

The secondary structure was assigned to each amino acid of two protein datasets by using DSSP.<sup>35</sup> The program PROMOTIF<sup>36</sup> was implemented to identify the observed  $\beta$ - and  $\gamma$ -turn motifs.

## Methods

### Position Conservation Scoring Function

The position conservation scoring function (PCSF) method had been widely used in the prediction of transcription factor binding sites in genomes.<sup>37-41</sup> To consider the effect of position in  $\beta$ - or  $\gamma$ -turn sequence segments, PCSF will be constructed by calculating the position probability matrix (PPM) and the conservation index of position.

### PPM and Conservation Index of Position

The PPM includes  $20 \times L$  elements ( $L$  is the length of the sequence segments, 20 denotes the 20 native amino acids). Each element in the matrix represents probability at a corresponding position, which is defined as:

$$P_{i,x} = \frac{n_{i,x} + s_x}{N_i + \sum_{x=1}^{20} s_x} \quad (1)$$

where,  $n_{i,x}$  and  $s_x$ , respectively, denote the real counts and pseudo counts for amino acid  $x$  at the  $i$ -th position of the sequence segments.  $s_x$  is calculated by<sup>37,38</sup>:

$$s_x = \frac{\sqrt{N_i}}{20} \quad (2)$$

where  $N_i$  is the total number of the sequences.

To reflect the action of position information in the sequence segment, the conservation index at the  $i$ -th position may be defined by the following expressions<sup>39,40</sup>:

$$c_i = \frac{100}{\log 20} \left( \sum_{x=1}^{20} p_{i,x} \log p_{i,x} + \log 20 \right) \quad (3)$$

$p_{i,x}$  is defined by equation (1),  $c_i$  value equals 100 for full conservation at the  $i$ -th position,  $c_i$  equals 0 for random amino acids at the  $i$ -th position. The conservation index of position in the amino acid sequence reflects the difference of amino acid compositions in same position between different datasets (i.e. between  $\beta$ -turn and non- $\beta$ -turn datasets, or between  $\gamma$ -turn and non- $\gamma$ -turn datasets).

### Scoring Function

For an arbitrary sequence segment  $S$  with  $N$  amino acids (i.e.  $S = (x_1, x_2, \dots, x_N)$ , where  $x_i$  is the amino acid at position  $i$  in segment  $S$ ), the score of segment  $S$  can be defined as<sup>39,41</sup>:

$$F(S) = F(x_1, x_2, \dots, x_N) = \frac{\sum_{i=1}^N c_i (p_{i,x} - p_{i,x}^{\min})}{\sum_{i=1}^N c_i (p_{i,x}^{\max} - p_{i,x}^{\min})} \quad (4)$$

It is easily proven that  $0 \leq F(S) \leq 1$ . Here,  $p_{i,x}^{\min} = \min(p_{i,x})$  and  $p_{i,x}^{\max} = \max(p_{i,x})$  are the minimal and maximal values of amino acid probabilities at position  $i$ , respectively.

For an arbitrary sequence segment, the class of this segment may be predicted by the maximum among  $F(S)^\beta$  and  $F(S)^{\text{non-}\beta}$  ( $F(S)^\gamma$  and  $F(S)^{\text{non-}\gamma}$ ), and can be formulated as follows:

$$\text{For prediction } \beta\text{-turns: } F(S)^\xi = \text{Max} \{F(S)^\beta, F(S)^{\text{non-}\beta}\}$$

$$\xi \in \beta\text{-turn or non-}\beta\text{-turn.}$$

$$\text{For prediction } \gamma\text{-turns: } F(S)^\xi = \text{Max} \{F(S)^\gamma, F(S)^{\text{non-}\gamma}\}$$

$$\xi \in \gamma\text{-turn or non-}\gamma\text{-turn.}$$

The operator Max means taking the maximum value among those in the parentheses, and then the  $\xi$  will give the segment class to which the predicted segment should belong (i.e. center amino acid of segment is should belong).

### Increment of Diversity

The increment of diversity (ID) algorithm, a whole uncertain measure and total information of a system on state space, is essentially a measure of the composition similarity level for two systems.<sup>42</sup> The ID algorithm has been applied in the recognition of protein structural class,<sup>43</sup> the exon-intron splice site prediction,<sup>44</sup> and the prediction of subcellular location of proteins.<sup>45</sup>

In the state space of  $t$  dimension, the diversity measure for diversity source  $S: \{m_1, m_2, \dots, m_t\}$  is defined as<sup>42-45</sup>:

$$D(S) = M \log M - \sum_i m_i \log m_i \quad (5)$$

In the same state space, ID between the source of diversity  $X(n_1, n_2, \dots, n_t)$  and  $S(m_1, m_2, \dots, m_t)$  is defined as:

$$\text{ID}(X, S) = D(X + S) - D(X) - D(S) \quad (6)$$

It is easily proven that the ID can be written as:

$$\text{ID}(X, S) = D(M, N) - \sum_i D(m_i, n_i) \quad (7)$$

$$D(M, N) = (M + N) \log(M + N) - M \log M - N \log N \quad (8)$$

$$D(m_i, n_i) = (m_i + n_i) \log(m_i + n_i) - m_i \log m_i - n_i \log n_i \quad (9)$$

here  $N = \sum_i n_i, M = \sum_i m_i$ . If  $m_i$  or  $n_i$  equals zero, then  $D(m_i, n_i) = 0$ .

An arbitrarily sequence segment may be predicted by the minimum among  $\text{ID}^\beta$  and  $\text{ID}^{\text{non-}\beta}$  ( $\text{ID}^\gamma$  and  $\text{ID}^{\text{non-}\gamma}$ ), and can be formulated as follows:

$$\text{For prediction } \beta\text{-turns: } \text{ID}^\xi = \text{Min} \{\text{ID}^\beta, \text{ID}^{\text{non-}\beta}\}$$

$$\xi \in \beta\text{-turn or non-}\beta\text{-turn.}$$

$$\text{For prediction } \gamma\text{-turns: } \text{ID}^\xi = \text{Min} \{\text{ID}^\gamma, \text{ID}^{\text{non-}\gamma}\}$$

$$\xi \in \gamma\text{-turn or non-}\gamma\text{-turn.}$$

The operator Min means taking the minimum value among those in the parentheses, and then the  $\xi$  will give the segment class to which the predicted segment should belong (i.e. center amino acid of segment is should belong). This method is also applied to predict  $\gamma$ -turns in the proteins.

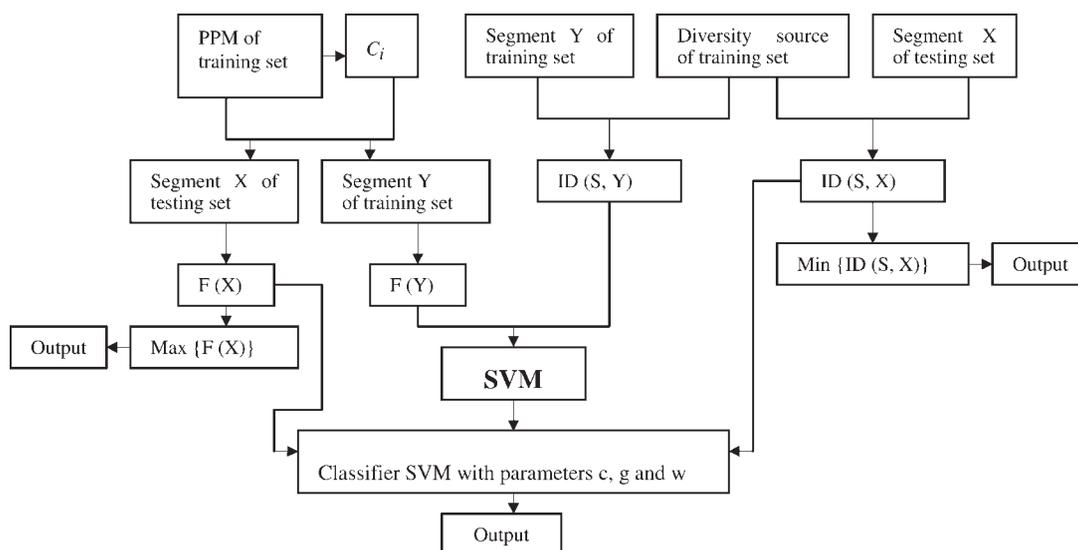
### Support Vector Machine

The SVM is an extremely successful learning system. SVM has been widely used in pattern recognition<sup>46,47</sup> and prediction of the secondary structure in proteins.<sup>5,48</sup> In this article, LIBSVM software packages are used.<sup>49</sup>

### Parameter Selection

The probabilities of 20 amino acids at each position are important parameters for predicting  $\beta$ - and  $\gamma$ -turns and have been introduced by the previous investigators.<sup>13,14,18,27</sup> They are selected as parameters ( $A_p$ ) of PCSF, using the training dataset to construct PPMs (contained  $20 \times L$  elements). For every sequence segment, the 2 score can be obtained for the  $\beta$ -turn and non- $\beta$ -turn ( $\gamma$ -turn and non- $\gamma$ -turn).

In the ID algorithm,  $D_0$ ,  $C_0$ , and  $H_2$  are, respectively, selected as the parameters of diversity source of the  $\beta$ -turns and non- $\beta$ -turns. The  $D_0$  denotes the frequencies of 400 dipeptide compositions from 20 amino acids;  $C_0$  denotes the frequencies of six hydrophathy characteristics at each position; and  $H_2$  denotes dipeptide frequencies of the residue  $i$  and reduces  $i + 3$  for the six hydrophathy characteristics. They constructed  $20 \times 20$ ,  $6 \times 7$ , and  $6 \times 6$  dimensions of state space, respectively. For example,  $D_0$  is selected as the parameters of diversity source in equation (5), based on the equation (7) definition, for all sequence segments, the 2 ID can be obtained for the  $\beta$ -turns and non- $\beta$ -turns.



**Figure 1.** The architecture of the turn prediction system. Note: The architecture contains the PCSF prediction system, the ID algorithmic prediction system, and the SVM prediction system. The  $C_i$  is the conservation index at the  $i$ -th position; the ID is the increment of diversity; the  $F$  is the scoring function.

$A_0$  and  $H_0$  are, respectively, selected as the parameters of diversity source of the  $\gamma$ -turns and non- $\gamma$ -turns.  $A_0$  denotes the frequencies of 20 native amino acids at each position;  $H_0$  denotes the frequencies of contiguous residue dipeptide compositions from six hydrophathy characteristics. They constructed the  $20 \times 5$  and  $6 \times 6$  dimensions of state space, respectively.

When using SVM algorithm to predict  $\beta$ -turn and non- $\beta$ -turn, for every sequence segment, the 2 scores can be obtained by the PCSF algorithm, and 6 ID (for  $D_0$ ,  $C_0$ , and  $H_2$  parameters) values can be obtained by ID algorithm. The 8 parameters are selected as input parameters of the SVM.

For  $\gamma$ -turn and non- $\gamma$ -turn prediction, for every sequence segment, the 2 scores can be obtained by the PCSF algorithm, and 4 ID (for  $A_0$  and  $H_0$  parameters) values can be obtained by ID algorithm. The 6 parameters are selected as input parameters of the SVM.

The architecture of the turn prediction system is shown in Figure 1.

#### Filtering and k-Fold Cross-Validation

In the case of predicting  $\beta$ -turns ( $\gamma$ -turns) at the residue level, the sliding windows contain 7 (5) amino acids, and the center amino acid of the segment is a predicted amino acid. To ensure prediction veracity, the structure characteristic of the  $\beta$ -turn ( $\gamma$ -turn) that contained 4 (or 3) consecutive residues is considered.<sup>16</sup>

To compare our method with other methods, the same 7-fold cross validation technique used by Kaur and Raghava<sup>13,14</sup> and Zhang et al.<sup>21</sup> for predicting  $\beta$ -turns is used. The 426 proteins are randomly divided into 7 subsets (6 subsets contained 61 chains; 1 subset contained 60 protein chains). Also, the 543 proteins are randomly divided into 7 subsets (3 subsets contained 77 chains; 4 subsets contained 78 protein chains). The 819 proteins are randomly divided into 7 subsets (7 subsets contained

117 chains). Each subset is an unbalanced set that retains the naturally occurring proportion of  $\beta$ -turns and non- $\beta$ -turns (1:3). The methods were trained on 6 subsets, and the performance is measured on the remaining seventh subset. This process is repeated 7 times so that each subset is tested.

Similarly, for the  $\gamma$ -turns prediction, we have used the 5-fold cross validation technique, which is used by Kaur and Raghava.<sup>33</sup> The 320, 346, and 536 proteins are randomly divided into 5 subsets, respectively: 5 subsets containing 64 chains, 4 subsets containing 69 chains, 1 subset containing 70 protein chains, 4 subsets containing 107 chains, and 1 subset containing 108 protein chains. Each subset is an unbalanced set that retains the naturally occurring proportion of  $\gamma$ -turns and non- $\gamma$ -turns (1:28~30), (1:30~33), and (1:26~30).

#### Performance Measures and Threshold Independent Measures

The performance of  $\beta$ - and  $\gamma$ -turns prediction is estimated by four parameters: the overall prediction accuracies ( $Q_{\text{total}}$ ); percentages of correctly predicted  $\beta$ - or  $\gamma$ -turns ( $Q_{\text{pred}}$ ); percentages of observed  $\beta$ - or  $\gamma$ -turns that are correctly predicted ( $Q_{\text{obs}}$ ); and Mcc, calculated by:

$$Q_{\text{total}} = \frac{p + n}{p + n + u + o}$$

$$Q_{\text{pred}} = \frac{p}{p + o}$$

$$Q_{\text{obs}} = \frac{p}{p + u}$$

$$\text{Mcc} = \frac{(p \times n) - (u \times o)}{\sqrt{(p + o) \times (p + u) \times (n + o) \times (n + u)}}$$

**Table 1.** The Classifications of Amino Acids.

Classification	Amino Acids
Strongly hydrophilic or polar	<b>R, D, E, N, Q, K, H</b>
Strongly hydrophobic	<b>L, I, V, A, M, F</b>
Weakly hydrophilic or weakly hydrophobic	<b>S, T, Y, W</b>
Proline	<b>P</b>
Glycine	<b>G</b>
Cysteine	<b>C</b>

where  $p$  is the number of correctly classified  $\beta$ -turn (or  $\gamma$ -turn) residues,  $n$  is the number of correctly classified non- $\beta$ -turn (non- $\gamma$ -turn) residues,  $o$  is the number of non- $\beta$ -turn (non- $\gamma$ -turn) residues incorrectly classified as  $\beta$ -turn ( $\gamma$ -turn) residues, and  $u$  is the number of  $\beta$ -turn ( $\gamma$ -turn) residues incorrectly classified as non- $\beta$ -turn (non- $\gamma$ -turn) residues.

For a prediction method, it is important to know the prediction reliability. The area under the ROC curve used by Kaur and Raghava<sup>13,14,33</sup> is taken as a reliable index because it provides the effectiveness of  $\beta$ -turns ( $\gamma$ -turns) prediction. The measure of overall accuracy is not dependent on a particular threshold.<sup>50</sup> In the ROC plot, all sensitivity values (true positive fraction) for all available thresholds are displayed on the  $y$ -axis, and all (1-specificity) values (false-positive fraction) for all available thresholds are shown on the  $x$ -axis. Sensitivity ( $S_n$ ) and specificity ( $S_p$ ) are defined as:

$$S_n = p/(p + u) \quad S_p = n/(n + o)$$

**Table 2.** The Predictive Results of Different Methods for the  $\beta$ -Turns in the 426 Proteins Using the 7-Fold Cross-Validation.

Method (parameter)	$M$	$Q_{\text{total}}$ (%)	$Q_{\text{pred}}$ (%)	$Q_{\text{obs}}$ (%)	Mcc
PCSF ( $A_p$ )	$20 \times 7$	75.7	55.2	32.4	0.26
ID ( $D_0$ )	400	67.8	43.2	69.7	0.29
ID ( $C_0$ )	$6 \times 7$	66.4	43.2	67.0	0.28
ID ( $H_2$ )	36	67.4	43.3	68.2	0.29
SVM (PCSF ( $A_p$ ) + ID ( $C_0$ ) + ID ( $D_0$ ) + ID ( $H_2$ ))	$4 \times 2$	77.3	54.3	67.9	0.45
<b>SVM (PCSF (<math>A_p</math>) + ID (<math>C_0</math>) + ID (<math>D_0</math>) + ID (<math>H_2</math>) + PSI)</b>	<b><math>4 \times 2 + 3</math></b>	<b>79.8 (79.3)</b>	<b>55.6 (55.4)</b>	<b>68.9 (68.9)</b>	<b>0.47 (0.47)</b>
Zhang et al.'s SVM (Single sequence) <sup>a</sup>	$20 \times 7 + 3$	74.8	49.1	67.9	0.41
Zhang et al.'s SVM (Multiple alignment) <sup>a</sup>	$20 \times 7 + 3$	77.3	53.1	67.0	0.45
BetaTPred2 <sup>b</sup>	–	75.5	49.8	72.3	0.43
Chou–Fasman <sup>c</sup>	–	74.9	46.1	16.9	0.16
1–4 and 2–3 correlation model <sup>c</sup>	–	63.2	35.3	60.4	0.21
Sequence coupled model <sup>c</sup>	–	50.6	31.7	88.4	0.23
Fuchs and Alix's <sup>d</sup>	–	74.8	48.8	69.9	0.42
Pham et al.'s SVM <sup>e</sup>	–	79.8	59.2	58.0	0.45

$M$  is the number of input parameters.

<sup>a</sup>From (Zhang et al., 2005); <sup>b</sup>from (Kaur and Raghava, 2003); <sup>c</sup>from (Kaur and Raghava, 2002); <sup>d</sup>from (Fuchs and Alix, 2005); <sup>e</sup>from (Pham et al., 2005).

PSI denotes predicted secondary structure information in proteins by using of the PSIPRED; bold font denotes the most accurately predicted result. Values shown in parentheses correspond to the results obtained by cross-validation of PSIPRED.

### The Hydrophathy Distribution along Protein Sequence

The hydrophathy distribution along the protein sequence has been recognized as a feature useful for the characterization of protein structure in the form of hydrophathy profiles.<sup>51</sup> To obtain the hydrophathy characteristics, the amino acids may be divided into groups using their individual hydrophathy according to the ranges of the hydrophathy scale. Because Proline, Glycine, and Cysteine have unique backbone properties, they are classified into 3 groups. Therefore, a protein sequence with 20 amino acids can be represented by a sequence with 6 characters.<sup>45</sup> The classification of amino acids are shown in Table 1.

## Results and Discussion

### Prediction of $\beta$ -Turns in the 426 Chains Dataset

#### The Predictive Results by Using PCSF Algorithm

Here the length of windows is 7 amino acids, so the PPMs contained  $20 \times 7$  elements. The predicting results of  $\beta$ -turns and non- $\beta$ -turns are shown in Table 2. The Mcc value is 0.26. To compare our method with other methods, other method's results for using the same dataset are also shown in Table 2.

#### The Predictive Results by Using ID Algorithm

Here  $D_0$ ,  $C_0$ , and  $H_2$  are selected as the parameters of the diversity source, respectively. For each kind of parameter, the performance of ID method is also shown in Table 2. Similar results are obtained by the three kinds of parameters, respectively. The Mcc achieved 0.29, which is slightly better than 0.26 of the

**Table 3.** The Predictive Results of Our Methods for the  $\beta$ -Turns in the 543 Proteins Using the 7-Fold Cross-Validation.

Method (parameter)	$M$	$Q_{\text{total}}$ (%)	$Q_{\text{pred}}$ (%)	$Q_{\text{obs}}$ (%)	Mcc
PCSF ( $A_p$ )	$20 \times 7$	75.0	44.3	38.2	0.23
ID ( $D_0$ )	400	60.0	39.0	63.4	0.28
ID ( $C_0$ )	$6 \times 7$	59.1	36.6	58.9	0.25
ID ( $H_2$ )	36	61.3	35.6	55.2	0.24
SVM (PCSF ( $A_p$ ) + ID ( $C_0$ ) + ID ( $D_0$ ) + ID ( $H_2$ ))	$4 \times 2$	74.4	46.3	59.4	0.41
<b>SVM (PCSF (<math>A_p</math>) + ID (<math>C_0</math>) + ID (<math>D_0</math>) + ID (<math>H_2</math>) + PSI)</b>	<b><math>4 \times 2 + 3</math></b>	<b>76.6 (76.1)</b>	<b>47.6 (46.9)</b>	<b>70.2 (70.2)</b>	<b>0.43 (0.42)</b>

PSI denotes the predicted secondary structure information in proteins by using of the PSIPRED; bold font denotes the most accurately predicted result. Values shown in parentheses correspond to the results obtained by cross-validation of PSIPRED.

PCSF method. However,  $Q_{\text{obs}}$  values are improved from about 32% of the PCSF algorithm to about 68% of the ID algorithm.

#### The Predictive Results by Using SVM

We train the classifiers with the LIBSVM program.<sup>49</sup> The radial basis function is selected as the kernel function; the optimized parameters are  $c = 128$ ,  $g = 0.5$ , respectively, and weight factor is  $w = 3$ .

The results indicate that the  $Q_{\text{pred}}$  value is higher than the  $Q_{\text{obs}}$  value in the PCFS method, but the  $Q_{\text{pred}}$  value is lower than the  $Q_{\text{obs}}$  value in the ID method. Therefore, to enhance the prediction performance, the composite vectors with the above calculated PCSF and ID values are selected as the input parameters of the SVM. Better results are obtained in Table 2. The Mcc value is raised to 0.45 and overall prediction accuracy is increased to 77.3%.

Widely believed, the  $\beta$ -turn prediction accuracy can be greatly improved by using the secondary structure information.<sup>13,21</sup> Therefore, to further improve the predictive effect, the predictive secondary structure information is obtained by using the PSIPRED,<sup>2</sup> added in the input parameter. The predicted secondary structure of each residue is represented as: helix  $\rightarrow (1, 0, 0)$ , sheet  $\rightarrow (0, 1, 0)$ , and coil  $\rightarrow (0, 0, 1)$ . The PCSF, ID values, and predictive secondary structure information, together, are selected as input parameters. The performance is also shown in Table 2. The results indicate that the  $Q_{\text{pred}}$  value is 55.6% and the  $Q_{\text{obs}}$  value is 68.9%; Mcc is 0.47 and prediction accuracy is

increased to 79.8%. This result is the highest achieved so far for predicting  $\beta$ -turn (bold font in Table 2). When the predicted secondary structure information is used in input the parameter, the prediction accuracy again gained about 3%.

Our method gives better rates than previous methods. The possible reasons are: first, the SVM is an extremely successful learning theory that usually outperforms other machine learning technologies such as artificial neural networks and nearest neighbor methods; and second, a new composite vector with ID, position conservation scoring function (PCSF) is employed. The ID and PCSF both algorithms may be extracted structure information of sequence. SVM with composite vector obtained Mcc value is raised to 0.45. The third reason is, the predicted secondary structure information by PSIPRED is used. The Mcc value is raised from 0.45 to 0.47. Comparing with previous Zhang et al.'s<sup>21</sup> and Pham et al.'s<sup>27</sup> SVM methods for the prediction of  $\beta$ -turn, using the improved composite vector in our method is a key step.

Some of the protein chains in our dataset may be used to train PSIPRED. To cross-validate the results, we have excluded those proteins from the nonredundant database of PSIPRED. As shown in Table 2; the difference in prediction performance is negligible.

#### Prediction of $\beta$ -Turns in the 543 and 819 Chains Dataset

To evaluate the predictive method, the  $\beta$ -turns in the 543 and 819 chains dataset are predicted by using our method (results

**Table 4.** The Predictive Results of Our Methods for the  $\beta$ -Turns in the 819 Proteins Using the 7-Fold Cross-Validation.

Method (parameter)	$M$	$Q_{\text{total}}$ (%)	$Q_{\text{pred}}$ (%)	$Q_{\text{obs}}$ (%)	Mcc
PCSF ( $A_p$ )	$20 \times 7$	72.2	49.6	36.7	0.25
ID ( $D_0$ )	400	60.0	30.0	64.8	0.24
ID ( $C_0$ )	$6 \times 7$	59.8	36.7	61.6	0.26
ID ( $H_2$ )	36	62.8	39.3	57.2	0.26
SVM (PCSF ( $A_p$ ) + ID ( $C_0$ ) + ID ( $D_0$ ) + ID ( $H_2$ ))	$4 \times 2$	74.5	51.2	59.9	0.42
<b>SVM (PCSF (<math>A_p</math>) + ID (<math>C_0</math>) + ID (<math>D_0</math>) + ID (<math>H_2</math>) + PSI)</b>	<b><math>4 \times 2 + 3</math></b>	<b>76.8 (76.5)</b>	<b>53.0 (52.6)</b>	<b>72.3 (72.4)</b>	<b>0.45 (0.45)</b>

PSI denotes the predicted secondary structure information in proteins by using of the PSIPRED; bold font denotes the most accurately predicted result. Values shown in parentheses correspond to the results obtained by cross-validation of PSIPRED.

**Table 5.** The Predictive Results of Different Methods for the  $\gamma$ -Turns in the 320 Proteins Using the 5-Fold Cross-Validation.

Method (parameter)	$M$	$Q_{\text{total}}$ (%)	$Q_{\text{pred}}$ (%)	$Q_{\text{obs}}$ (%)	Mcc
PCSF ( $A_p$ )	$20 \times 5$	84.3	6.9	29.1	0.08
ID ( $A_0$ )	$20 \times 5$	61.3	5.4	62.8	0.09
ID ( $H_0$ )	36	54.2	4.5	61.6	0.07
SVM (PCSF( $A_p$ ) + ID ( $A_0$ ) + ID ( $H_0$ ))	$3 \times 2$	56.2	5.5	73.4	0.12
<b>SVM (PCSF (<math>A_p</math>) + ID (<math>A_0</math>) + ID (<math>H_0</math>) + PSI)</b>	<b><math>3 \times 2 + 3</math></b>	<b>61.0 (60.8)</b>	<b>6.8 (6.3)</b>	<b>91.4 (90.0)</b>	<b>0.18 (0.17)</b>
SNNS <sup>a,#</sup>		74.0	6.3	83.2	0.17
Weka-logistic regression <sup>a,#</sup>		62.6	5.6	65.1	0.12
Weka-naive Bayes <sup>a,#</sup>		57.4	5.0	65.4	0.11
Weka-J48 classifier <sup>a,#</sup>		92.6	5.0	7.2	0.03
Sequence coupled model <sup>a,*</sup>		57.8	5.9	43.2	0.08
GOR <sup>a,*</sup>		75.5	6.1	45.5	0.09
Pham et al.'s SVM <sup>b,#</sup>		79.9	7.7	47.5	0.13

<sup>a</sup>From (Kaur and Raghava, 2003); <sup>b</sup>from (Pham et al., 2005).

\*Using a single amino acid sequence and secondary structure information that from the PSIPRED.

#Using multiple alignment and secondary structure information; bold font denotes the most accurately predicted result. Values shown in parentheses correspond to the results obtained by cross-validation of PSIPRED.

shown in Tables 3 and 4). The SVM's optimized parameters are  $c = 2048$ ,  $g = 0.5$ , respectively, and weight factor is  $w = 3$ .

The results indicate that  $Q_{\text{pred}}$ ,  $Q_{\text{obs}}$ , and Mcc values are respectively 47.6%, 70.2%, and 0.43, and prediction accuracy is 76.6% for 543 chains dataset. The results indicate that the  $Q_{\text{pred}}$  value is 53.0% and the  $Q_{\text{obs}}$  value is 72.3%; Mcc is 0.45 and prediction accuracy is 76.8% for 819 chains dataset. When the difference in performance of our method on three datasets was analyzed, it was observed that, although the performance decreased, the trend remained the same.

#### Prediction of $\gamma$ -Turns in the 320 Chains Dataset

The above methods of predicting  $\beta$ -turn is also applied to the prediction of  $\gamma$ -turns in the 320 proteins with a sliding window of five amino acids (the PPMs contained  $20 \times 5$  elements). The predictive results of the PSCF algorithm are shown in Table 5. The Mcc value only is 0.08. To compare our method with other methods, other method's results for using the same dataset are also shown in Table 5.

In the ID algorithm,  $A_0$  and  $H_0$  are, respectively, selected as the parameters of diversity source of the  $\gamma$ -turns and non- $\gamma$ -turns. The performance is also shown in Table 5. The Mcc is 0.07–0.09, which is similar with 0.08 of the PCSF method.

Using the SVM method to predict  $\gamma$ -turns, the optimized parameters  $c$  and  $g$  are default; weight factor is  $w = 29$ . The above values of the PCSF ( $A_p$ ), ID ( $A_0$ ), and ID ( $H_0$ ) are used to construct the composite vector as input parameters of SVM. The predictive results of SVM are shown in Table 5. The Mcc increases to 0.12.

To further improve the performance of the prediction, by adding the above composite vector and predictive secondary structure information to input parameter of SVM, better predictive results are obtained (see Table 5). The Mcc is increased to 0.18, better than the 0.13 in Pham et al.'s work<sup>27</sup> and the 0.17 in Kaur's<sup>21</sup> work. The  $Q_{\text{obs}}$  value is 91.4%, better than the 83.2% in Kaur's<sup>21</sup> work and the 47.5% in Pham et al.'s<sup>27</sup> work. Overall prediction accuracy is 61.0% (bold font in Table 5).

The above values are the most common measure of a method's overall performance; however, the  $Q_{\text{total}}$  can be misleading as  $\gamma$ -turn residues occur much less frequently than non- $\gamma$ -

**Table 6.** The Predictive Results of Our Methods for the  $\gamma$ -Turns in the 346 Proteins Using the 5-Fold Cross-Validation.

Method (parameter)	$M$	$Q_{\text{total}}$ (%)	$Q_{\text{pred}}$ (%)	$Q_{\text{obs}}$ (%)	Mcc
PCSF ( $A_p$ )	$20 \times 5$	88.4	5.4	27.8	0.07
ID ( $A_0$ )	$20 \times 5$	61.7	4.5	63.8	0.08
ID ( $H_0$ )	36	55.4	3.9	57.8	0.07
SVM (PCSF( $A_p$ ) + ID ( $A_0$ ) + ID ( $H_0$ ))	$3 \times 2$	64.9	4.9	75.3	0.11
<b>SVM (PCSF (<math>A_p</math>) + ID (<math>A_0</math>) + ID (<math>H_0</math>) + PSI)</b>	<b><math>3 \times 2 + 3</math></b>	<b>59.0 (58.9)</b>	<b>5.7 (5.6)</b>	<b>90.3 (89.2)</b>	<b>0.16 (0.16)</b>

PSI denotes the predicted secondary structure information in proteins by using of the PSIPRED; bold font denotes the most accurately predicted. Values shown in parentheses correspond to the results obtained by cross-validation of PSIPRED.

**Table 7.** The Predictive Results of Our Methods for the  $\gamma$ -Turns in the 536 Proteins Using the 5-Fold Cross-Validation.

Method (parameter)	$M$	$Q_{\text{total}}$ (%)	$Q_{\text{pred}}$ (%)	$Q_{\text{obs}}$ (%)	$Mcc$
PCSF ( $A_p$ )	$20 \times 5$	75.9	6.1	38.2	0.08
ID ( $A_0$ )	$20 \times 5$	59.0	5.2	63.9	0.09
ID ( $H_0$ )	36	56.4	4.6	58.8	0.08
SVM (PCSF( $A_p$ ) + ID ( $A_0$ ) + ID ( $H_0$ ))	$3 \times 2$	58.4	5.3	73.4	0.12
<b>SVM (PCSF (<math>A_p</math>) + ID (<math>A_0</math>) + ID (<math>H_0</math>) + PSI)</b>	<b><math>3 \times 2 + 3</math></b>	<b>58.5 (58.0)</b>	<b>6.8 (6.6)</b>	<b>92.7 (91.6)</b>	<b>0.18 (0.17)</b>

PSI denotes the predicted secondary structure information in proteins by using of the PSIPRED; bold font denotes the most accurately predicted result. Values shown in parentheses correspond to the results obtained by cross-validation of PSIPRED.

turn residues in proteins (1:~30). Therefore, one could easily achieve  $Q_{\text{total}} \approx 97\%$  merely by predicting all residues to be non- $\gamma$ -turn. For this reason, we consider  $Mcc$ ,  $Q_{\text{pred}}$ , and  $Q_{\text{obs}}$  to be important indices.

The  $\gamma$ -turn results are obtained by using the same prediction method as with  $\beta$ -turns, only different parameters are selected, but given poor  $Q_{\text{pred}}$  and  $Mcc$  (only is 0.18) values in prediction results. This is definitely more unbalanced in the present dataset, which has a ratio of ~30:1 of non- $\gamma$ -turn and  $\gamma$ -turn residues ( $\beta$ -turn (~3:1)). Moreover, a  $\gamma$ -turn consists of three residues and thus is much more flexible than a  $\beta$ -turn.<sup>21</sup>

#### Prediction of $\gamma$ -Turns in the 346 and 536 Chains Dataset

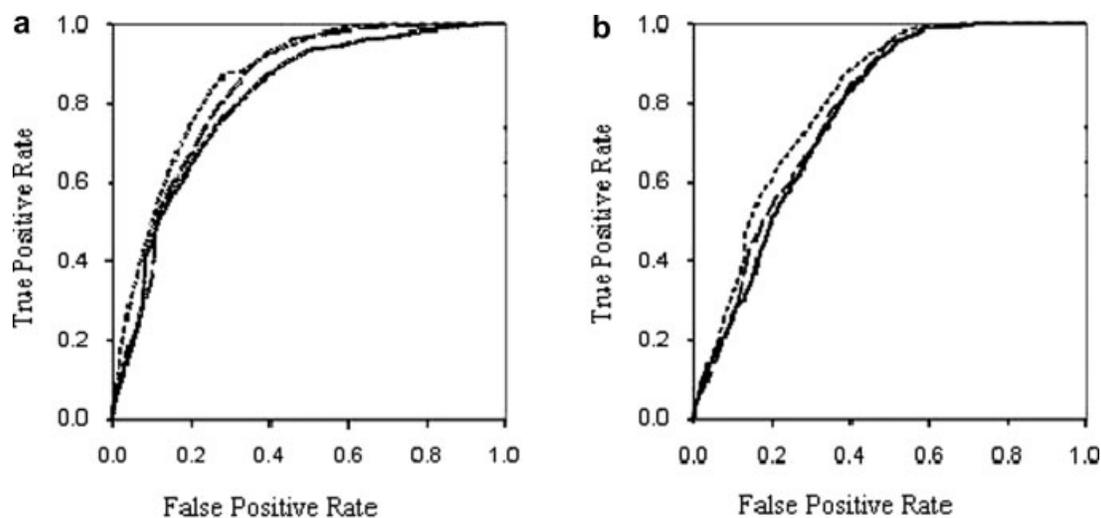
To evaluate the predictive method, the  $\gamma$ -turns in the 346 and 536 chains dataset are predicted by using our method, and

results are shown in Tables 6 and 7. The SVM's optimized parameters  $c$  and  $g$  are default; weight factor is  $w = 30$  (for 346 chains) and  $w = 26$  (for 536 chains), respectively.

The results indicate that the  $Q_{\text{obs}}$  value is 90.3%, the  $Mcc$  is 0.16 and the prediction accuracy is 59.0% for 346 chains dataset. The results indicate that the  $Q_{\text{obs}}$  value is 92.7%, the  $Mcc$  is 0.18, and the prediction accuracy is 58.5% for 536 chains dataset. The performance of three datasets is made by using the same rule.

#### Receiver Operating Characteristic Results

In addition, we calculated the area under the ROC curves of both SVM systems in the prediction  $\beta$ - and  $\gamma$ -turns. The performances of different systems have been evaluated by the ROC curves in Figure 2. For 320, 346, and 536 protein chain datasets, the corresponding areas under the ROC curves are 0.81, 0.78, and 0.79 for



**Figure 2.** The ROC curves of both SVM systems of the prediction  $\beta$ -turns in 7-fold cross validation and  $\gamma$ -turns in 5-fold cross validation proteins. (a) shows ROC curves of  $\beta$ -turns in various dataset; Note: dotted line indicates a curve in the 426 protein chains dataset; solid line indicates a curve in the 543 protein chains; dashed line indicates a curve in the 819 protein chains. The corresponding areas under the ROC curves are 0.87, 0.82, and 0.84, respectively. (b) shows ROC curves of  $\gamma$ -turns in various dataset. Note: dotted line indicates a curve for the 320 protein chains; solid line indicates a curve for the 346 protein chains; dashed line indicates a curve in the 536 protein chains. The corresponding areas under the ROC curves are 0.81, 0.78, and 0.79, respectively.

$\gamma$ -turns prediction, respectively. They are higher than 0.73 in Kaur and Raghava's method for  $\gamma$ -turns prediction.<sup>33</sup>

## Conclusion

The above predicted results of the  $\beta$ - and  $\gamma$ -turns show that a single algorithm may provide partial information of a sequence. When the information of ID and PCSF are put together into the input parameters of SVM, the performance can be tremendously improved. It is possible that the SVM algorithm plays an information syncretizing role. In addition, the secondary structure information is also helpful to improve the predictive performance of the  $\beta$ - and  $\gamma$ -turns.

The successful prediction of  $\beta$ - and  $\gamma$ -turns in the proteins and, by using SVM with ID, PCSF values, and a predictive secondary structure as the input information, indicates a promising approach. Using them as SVM parameters can reduce dimension of input vector, improve calculating efficiency, and extract important classified information.

## Acknowledgments

Authors are grateful to the editor and referees for their careful review and valuable comments on our manuscript.

## References

1. Rost, B.; Sander, C. *J Mol Biol* 1993, 232, 584.
2. Jones, D. T. *J Mol Biol* 1999, 292, 195.
3. Petersen, T.N.; Lundegrad, C.; Neilsen, M.; Bohr, H.; Bohr, J.; Brunak, S.; Gippert, G. P.; Lund, O. *Proteins* 2000, 41, 17.
4. Cuff, J. A.; Barton, G. J. *Proteins: Struct Funct Bioinform Genet* 2000, 40, 502.
5. Guo, J.; Chen, H.; Sun, Z. R.; Lin, Y. *Proteins: Struct Funct Bioinform* 2004, 54, 738.
6. Wu, K. P.; Lin, H. N.; Chang, J. M.; Sung, T. Yi.; Hsu, W. L. *Nucl Acids Res* 2004, 32, 5059.
7. Rose, G. D.; Gierasch, L.; Smith, J. A. *Adv Protein Chem* 1985, 37, 1.
8. Takano, K.; Yamagata, Y.; Yutani, K. *Biochemistry* 2000, 39, 8655.
9. Cruz, X.; Thornton, J. M. *Protein Sci* 1999, 8, 750.
10. Rost, B.; Schneider, R.; Sander, C. *J Mol Biol* 1997, 270, 471.
11. Jones, D. T. *Proteins Suppl* 2001, 5, 127.
12. Chou, K. C. *Anal Biochem* 2000, 286, 1.
13. Kaur, H.; Raghava, G. P. *Bioinformatics* 2002, 18, 1508.
14. Kaur, H.; Raghava, G. P. *Protein Sci* 2003, 12, 627.
15. Chou, K. C.; Blinn, J. R. *J Protein Chem* 1997, 16, 575.
16. Shepherd, A. J.; Gorse, D.; Thornton, J. M. *Protein Sci* 1999, 8, 1045.
17. Chou, K. C. *J Pept Res* 1997, 49, 120.
18. Zhang, C. T.; Chou, K. C. *Biopolymers* 1997, 41, 673.
19. Cai, Y. D.; Liu, X. J.; Xu, X. B.; Chou, K. C. *J Pept Sci* 2002, 8, 297.
20. Lin, T. H.; Wang, G. M.; Wang, Y. T. *J Chem Inf Comput Sci* 2002, 42, 123.
21. Zhang, Q. D.; Yoon, S.; Welsh, W. J. *Bioinformatics* 2005, 21, 2370.
22. Fuchs, P. F. J.; Alix, A. J. P. *Proteins: Struct Funct Bioinform* 2005, 59, 828.
23. Guruprasad, K.; Rajkumar, S. *J Biosci* 2000, 25, 143.
24. Chou, P. Y.; Fasman, G. D. *Biochemistry* 1974, 13, 211.
25. Gibrat, J.-F.; Garnier, J.; Robson, B. *J Mol Biol* 1987, 198, 425.
26. Wilmot, C. M.; Thornton, J. M. *Protein Eng* 1990, 3, 479.
27. Pham, T. H.; Satou, K.; Ho, T. B. *J Bioinformatics Comput Biol* 2005, 3, 343.
28. Street, T. O.; Fitzkee, N. C.; Perskie, L. L.; Rose G. D. *Protein Science* 2007, 16, 1720.
29. Bornot, A.; de Brevern, A. G. *Bioinformatics* 2006, 1, 153.
30. Bystrov, V. F.; Portnova, S. L.; Tsetlin, V. I.; Ivanov, V. T.; Ochinnikov, Y. A. *Tetrahedron* 1969, 25, 493.
31. Milner-White, E. J.; Ross, B. M.; Ismail, R.; Belhadj-Mastefa, K.; Poet, R. *J Mol Biol* 1988, 204, 777.
32. Alkorta, I.; Suarez, M. L.; Herranz, R.; González-Muñiz, R.; García-López, M. T. *J Mol Model* 1996, 2, 16.
33. Kaur, H.; Raghava, G. P. S. *Protein Sci* 2003, 12, 923.
34. Guruprasad, K.; Shukla, S.; Adindla, S.; Guruprasad, L. *J Peptide Res* 2003, 61, 243.
35. Kabsch, W.; Sander, C. *Biopolymers* 1983, 22, 2577.
36. Hutchinson, E. G.; Thornton, J. M. *Protein Sci* 1996, 5, 212.
37. Wasserman, W. W.; Sandelin, A. *Nat Rev Genet* 2004, 5, 276.
38. Kielbasa, S. M.; Gonze, D.; Herzog, H. *BMC Bioinformatics* 2005, 6, 237.
39. Cartharius, K.; Frech, K.; Grote, K.; Klocke, B.; Haltmeier, M.; Klingenhoff, A.; Frisch, M.; Bayerlein, M.; Werner, T. *Bioinformatics* 2005, 21, 2933.
40. Quandt, K.; Frech, K.; Karas, H.; Wingender, E.; Werner, T. *Nucl Acids Res* 1995, 23, 4878.
41. Kel, A. E.; GoBling, E.; Reuter, I.; Cheremushkin, E.; Kel-Margoulis, O. V.; Wingender, E. *Nucl Acids Res* 2003, 31, 3576.
42. Laxton, R. R. *J Theor Biol* 1978, 71, 51.
43. Li, Q. Z.; Lu, Z. Q. *J Theor Biol* 2001, 213, 493.
44. Zhang, L. R.; Luo, L. F. *Nucleic Acids Res* 2003, 31, 6214.
45. Chen, Y. L.; Li, Q. Z. *J Theor Biol* 2007, 245, 775.
46. Roobaert, D.; Hulle, M. M. In *Proceedings of IEEE Neural Networks for Signal Processing Workshop, Wisconsin, 1999*; pp 77.
47. Schmidt, M.; Grish, H. *Proc ICASSP'96, Atlanta, 1996*; pp 105.
48. Hua, S. J.; Sun, Z. R. *J Mol Biol* 2001, 308, 397.
49. Chang, C. C.; Lin, C. J. 2001. LIBSVM. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
50. Deleo, J.M. In *Proceedings of the Second International Symposium on Uncertainty Modelling and Analysis, IEEE, Computer Society Press, College Park, MD, 1993*, pp 318.
51. Pánek, J.; Eidhammer, I.; Aasland, R. *Proteins: Struct Funct Bioinform* 2005, 58, 923.