

# How Good Is Prediction of Protein Structural Class by the Component-Coupled Method?

Zhi-Xin Wang\* and Zheng Yuan

National Laboratory of Biomacromolecules, Institute of Biophysics, Academia Sinica, Peoples Republic of China

**ABSTRACT** Proteins of known structures are usually classified into four structural classes: all- $\alpha$ , all- $\beta$ ,  $\alpha + \beta$ , and  $\alpha/\beta$  type of proteins. A number of methods to predicting the structural class of a protein based on its amino acid composition have been developed during the past few years. Recently, a component-coupled method was developed for predicting protein structural class according to amino acid composition. This method is based on the least Mahalanobis distance principle, and yields much better predicted results in comparison with the previous methods. However, the success rates reported for structural class prediction by different investigators are contradictory. The highest reported accuracies by this method are near 100%, but the lowest one is only about 60%. The goal of this study is to resolve this paradox and to determine the possible upper limit of prediction rate for structural classes. In this paper, based on the normality assumption and the Bayes decision rule for minimum error, a new method is proposed for predicting the structural class of a protein according to its amino acid composition. The detailed theoretical analysis indicates that if the four protein folding classes are governed by the normal distributions, the present method will yield the optimum predictive result in a statistical sense. A non-redundant data set of 1,189 protein domains is used to evaluate the performance of the new method. Our results demonstrate that 60% correctness is the upper limit for a 4-type class prediction from amino acid composition alone for an unknown query protein. The apparent relatively high accuracy level (more than 90%) attained in the previous studies was due to the preselection of test sets, which may not be adequately representative of all unrelated proteins. *Proteins* 2000;38:165–175.

© 2000 Wiley-Liss, Inc.

**Key words:** SCOP database; Bayes decision rule; jack-knife analysis; amino acid composition;  $\alpha$  domains;  $\beta$  domains;  $\alpha + \beta$  domains;  $\alpha/\beta$  domains.

## INTRODUCTION

The protein structure prediction problem remains one of the most important problems in molecular biology. At present time, available methods have unable to meet in a satisfactory way. Protein can be considered as a hierarchy of structure: amino acid sequence  $\rightarrow$  secondary structure  $\rightarrow$  supersecondary structure  $\rightarrow$  domain  $\rightarrow$  three-

dimensional structure. Many approaches to the protein-folding problem have reflected this hierarchical scheme, and in these secondary structure prediction is the first and most critical step in the achievement of correct prediction.

The concept of protein structural classes was originally introduced by Levitt and Chothia based on a visual inspection of polypeptide chain topologies in a data set of 31 globular proteins.<sup>1</sup> According to the contents of secondary structures, proteins of known structures can be classified into one of the following four structural classes: all  $\alpha$ , all  $\beta$ ,  $\alpha + \beta$ , and  $\alpha/\beta$  type of proteins.<sup>1,2</sup> These class definitions have been generally accepted, and are still in common use, although slight changes have been made.<sup>3,4</sup> It is well known that the knowledge of protein structural class can help in the determination of the three-dimensional structure of a protein, particularly in improving the prediction of secondary structure.<sup>5,6</sup> Therefore, it would be very useful if a rapid and reliable method could be developed to predict the structural class of a protein. During the past decade, various efforts have been made to reach such a goal by many investigators. On the basis of different criteria, such as the discriminant analysis,<sup>7,8</sup> the least Minkowski distance,<sup>9,10</sup> the least Euclidean distance,<sup>11</sup> the optimization approach principle,<sup>12</sup> or the maximum projection principle,<sup>13</sup> various prediction methods have been developed in these studies. Recently, a statistical method was developed for predicting protein structural class according to amino acid composition.<sup>14</sup> This method is based on the least Mahalanobis distance principle, and yields much better predicted results in comparison with the previous methods. However, the success rates reported for secondary structural class prediction by different investigators are contradictory. The highest reported accuracies obtained by this method are near 100%, but the lowest one is only about 60%.<sup>15–20</sup>

These studies have been difficult to compare to each other for several reasons. First, the protein sets used vary widely in both size and extent of sequence homology. For example, some of the sets of proteins on which these algorithms were tested contained high levels of sequence homology (more than 90% identity in some cases) with

Grant sponsor: The 863 High-Technology Foundation and Pandeng project of the Ministry of Science and Technology; Grant number: 103–13–03–02.

\*Correspondence to: Zhi-Xin Wang, National Laboratory of Biomacromolecules, Institute of Biophysics, Academia Sinica, Beijing 100101, Peoples Republic of China. E-mail: zxwang@sun5.ibp.ac.cn

Received 23 April 1999; Accepted 11 August 1999

each other and with the proteins used in determining the numeric parameters of the algorithm. Second, the number structural classes in which to categorize the proteins has varied from study to study, ranging from three to seven structural classes. Third, criteria for dividing proteins into their structural classes has differed. For example, Nakashima et al. classified a protein as  $\alpha$  type if the protein is known to contain greater than 15%  $\alpha$  helices and less than 10%  $\beta$  strands.<sup>11</sup> Chou, on the other hand, requires their  $\alpha$  class proteins to contain at least 40%  $\alpha$  helices and not more than 5%  $\beta$  strands.<sup>15</sup> Fourth, cross-validation methods have varied widely, ranging from leave-one-out methods to single-test-set methods.

In the first part of this paper, according to the normality assumption and the Bayes decision rule for minimum error, we propose a new approach for predicting the structural class of a protein from its amino acid composition. A detailed theoretical analysis indicates that the least Mahalanobis distance method developed previously can be viewed as the approximation to the present method if the four protein structural classes are governed by the normal distributions. The question of why different conclusions were obtained from different studies is addressed in the next section by constructing an unbiased data set based on all unrelated protein domains. In particular, structures were taken from the Protein Data Bank and definitions of domains and structural classes were taken from the Structural Classification of Proteins (SCOP) database.<sup>3</sup> One of the most important aspects of our analysis is that we carefully tested the present method against this new data set. This testing allowed us to decide unambiguously whether a given comparison resulted in a true or false-positive and to decide objectively between different statistical schemes. Our results demonstrate that 60% correctness is the upper limit for a 4-type class prediction by the existing component-coupled method for an unknown non-homologous domain (sequence identity < 30%). Therefore, further improvement of prediction accuracy will mainly depend on the introduction of other features (e.g., average hydrophobicity, net charge, amino acid sequence, and so on).

## THEORY

A pattern is the description of an object. Almost anything that is within the reach of our five senses can be chosen as a pattern—a character, a photograph, speech pattern, odors, tastes, etc. A pattern class is a group of patterns with certain properties. Pattern recognition is that of classifying a pattern into one of the pattern classes on the grounds of certain measurement or properties. Depending on the problem of interest, the variations of the member of a pattern class can be deterministic (nonrandom) or random in nature. If the variations of the pattern from the stored reference, which is the ideal or average pattern, are random, the statistical pattern-recognition approach should be used. In this case, it is necessary then to describe such variations with a probabilistic quantity. The design of a statistical pattern recognition system is generally based on the Bayes classification rule and its

variations. This rule yields an optimum classifier (decision procedure) when the probability density function of each pattern population and the probability of occurrence of each pattern class are known.<sup>21–23</sup>

Suppose there are  $M$  possible pattern classes,  $\omega_1, \omega_2, \dots, \omega_M$ , and an arbitrary pattern belongs to class  $\omega_l$  with a priori probability,  $P(\omega_l)$ ,  $l = 1, 2, \dots, M$ . Pattern or feature-vector,  $\mathbf{x}$ , are  $n$ -components, random vectors taking value in  $n$ -dimensional feature-space,  $\mathbf{x}$ , and governed by a multivariate conditional probability density function,  $P(\mathbf{x} | \omega_l)$ , when pattern  $\mathbf{x}$  is known to belong to class  $\omega_l$ . The recognition problem can now be viewed as that of generating the decision boundaries which separate the  $M$  pattern classes on the basis of the observed measurement vectors. Let the decision boundaries be defined by decision functions  $d_1(\mathbf{x}), d_2(\mathbf{x}), \dots, d_M(\mathbf{x})$ . These functions, which are also called discriminant functions, are scalar and single-valued functions of the pattern  $\mathbf{x}$ . If  $d_l(\mathbf{x}) > d_u(\mathbf{x})$  for  $u = 1, 2, \dots, M$ , and  $l \neq u$ , the pattern  $\mathbf{x}$  belongs to  $\omega_l$ . In other words, if the  $l$ th decision function,  $d_l(\mathbf{x})$ , has the largest value for a pattern  $\mathbf{x}$ , then  $\mathbf{x} \in \omega_l$ . A decision rule based simply on probabilities is to assign a particular pattern  $\mathbf{x}$  to class  $\omega_l$  if

$$P(\omega_l | \mathbf{x}) > P(\omega_u | \mathbf{x}), \quad u = 1, 2, 3, \dots, M; \quad l \neq u \quad (1)$$

that is, the pattern class  $\omega_l$  with the highest *posteriori* probability is chosen as the assignment for  $\mathbf{x}$ . The *a posteriori* probabilities  $P(\omega_l | \mathbf{x})$  may be calculated from the *a priori* probabilities  $P(\omega_l)$  and the conditional density functions  $P(\mathbf{x} | \omega_l)$ , using Bayes' theorem, that is

$$P(\omega_l | \mathbf{x}) = P(\omega_l)P(\mathbf{x} | \omega_l)/P(\mathbf{x}) \quad (2)$$

It can be verified that the average of classification error probability is minimized if the decision rule of Eq.(1) is used, and therefore it is also called the Bayes decision rule for minimum error. From the discussion on decision function given above, it is note that the decision functions for the Bayes decision rule for minimum error can be written as

$$d_l(\mathbf{x}) = P(\omega_l | \mathbf{x}) \quad (3)$$

or

$$d_l(\mathbf{x}) = P(\omega_l)P(\mathbf{x} | \omega_l) \quad (4)$$

where  $P(\mathbf{x})$  has been eliminated since it does not depend on  $l$ . If all *a priori* probabilities are equal:  $P(\omega_l) = 1/M$ , for  $l = 1, 2, \dots, M$ , the decision function can then be written as:

$$d_l(\mathbf{x}) = P(\mathbf{x} | \omega_l) \quad (5)$$

Equation (5) is called the (conditional) maximum-likelihood decision, and it can be regarded as the Bayes decision rule for minimum error with equal *a priori* probabilities, i.e.,  $P(\omega_l) = 1/M$ , for  $l = 1, 2, \dots, M$ .

Prediction of protein structural class from amino acid composition can be considered as a pattern recognition

problem. The amino acid composition of a protein molecule serves as the properties of the recognition system, which identifies the protein structural class by analysis of these properties. According to its amino acid composition, a protein molecule can be represented by a point or a vector in a 20-dimensional space, the so-called composition space. However, as pointed by Chou and Zhang,<sup>16</sup> of the 20 amino acid composition components only 19 are independent, since the amino acid composition of a protein must be constrained by

$$\sum_{i=1}^{20} x_i = 1 \quad (6)$$

where  $x_i$  is the composition component of the  $i$ th amino acid in a protein. Therefore, by leaving out any one of its 20 components, one can still uniquely represent a protein by a point in a 19-D space. Suppose the 20 amino acids are alphabetically ordered according to their single-letter code. If the last amino acid component is left out, then the 19-D space will be defined by the bases corresponding to the components of A, C, D, E, F, G, H, I, K, L, M, N, P, Q, R, S, T, V, and W, respectively. Once the 19-D space is established, the  $k$ th protein in a given set of  $N$  proteins can be expressed by

$$\mathbf{x}_k = \begin{bmatrix} x_{k,1} \\ x_{k,2} \\ \vdots \\ x_{k,19} \end{bmatrix} \quad (k = 1, 2, \dots, N) \quad (7)$$

where  $x_{k,1}, x_{k,2}, \dots, x_{k,19}$  are, respectively, the 19 amino acid composition of the  $k$ th protein  $\mathbf{x}_k$ , and  $N$  is the total number of proteins in the set.

According to statistical theory,<sup>24</sup> when the number of the proteins in database is sufficiently large, the distribution density functions for  $\alpha$ ,  $\beta$ ,  $\alpha + \beta$ , and  $\alpha/\beta$  proteins can be assumed to be the 19-dimensional normal density functions<sup>25</sup>

$$P(\mathbf{x} | \omega_l) = \frac{1}{\sqrt{(2\pi)^{19} |\Sigma_l|}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_l)^T \Sigma_l^{-1} (\mathbf{x} - \boldsymbol{\mu}_l)\right\} \quad (8)$$

$l \in \alpha, \beta, \alpha + \beta, \text{ or } \alpha/\beta$

where  $\boldsymbol{\mu}_l$  and  $\Sigma_l$  are mean vector and  $19 \times 19$  covariance matrix for the  $l$  type protein,  $l \in \alpha, \beta, \alpha + \beta$ , or  $\alpha/\beta$ , respectively. According to the maximum likelihood estimation method, the sample average of the protein set concerned can be estimated by

$$\boldsymbol{\mu}_l = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_{19} \end{bmatrix} \quad (9)$$

where

$$\bar{x}_i = \frac{1}{N} \sum_{k=1}^N x_{k,i} \quad (i = 1, 2, \dots, 19)$$

and the covariance matrix is given by

$$\Sigma_l = \begin{bmatrix} s_{1,1} & s_{1,2} & \dots & \dots & s_{1,19} \\ s_{2,1} & s_{2,2} & \dots & \dots & s_{2,19} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ s_{19,1} & s_{19,2} & \dots & \dots & s_{19,19} \end{bmatrix} \quad (10)$$

where

$$s_{i,j} = \frac{1}{N} \sum_{k=1}^N [x_{k,i} - \bar{x}_i][x_{k,j} - \bar{x}_j] \quad (i, j = 1, 2, \dots, 19)$$

It has been shown that if  $N > 19$  and the sample are drawn from a normal population, the estimate of  $\Sigma_l$  given by Eq.(10) possesses an inverse  $\Sigma_l^{-1}$  with probability 1.<sup>26</sup> Note that the maximum likelihood estimator for  $\boldsymbol{\mu}_l$  is unbiased but  $\Sigma_l$  is not unbiased. The desired unbiased estimator for  $\Sigma_l$  is defined by

$$s_{i,j} = \frac{1}{N-1} \sum_{k=1}^N [x_{k,i} - \bar{x}_i][x_{k,j} - \bar{x}_j] \quad (i, j = 1, 2, \dots, 19) \quad (11)$$

which is very nearly the maximum likelihood estimator for large values of  $N$  (say  $N > 30$ ). Because of this, some statisticians choose to define the sample variance by Eq.(11) rather than Eq.(10). In this case, it can be verified that the present method is mathematically identical to the method used by Chou and Maggiora.<sup>20</sup> When the  $N$  proteins in Eqs.(9), (10), and (11) are all  $\alpha$  proteins,  $\boldsymbol{\mu}_l$  and  $\Sigma_l$  thus defined would become the mean and the covariance matrix of  $\alpha$  protein set, denoted by  $\boldsymbol{\mu}_\alpha$  and  $\Sigma_\alpha$ . Likewise, when the  $N$  proteins in equations (9)–(11) are all  $\beta$ , or  $\alpha + \beta$ , or  $\alpha/\beta$  protein sets, denoted by  $\boldsymbol{\mu}_\beta$ ,  $\boldsymbol{\mu}_{\alpha+\beta}$ ,  $\boldsymbol{\mu}_{\alpha/\beta}$ , and  $\Sigma_\beta$ ,  $\Sigma_{\alpha+\beta}$ ,  $\Sigma_{\alpha/\beta}$ , respectively.

Because of the exponential form of the normal density function, it is sometime more convenient to work with natural logarithm of this decision function

$$d_l(\mathbf{x}) = \ln[P(\omega_l)P(\mathbf{x} | \omega_l)] = \ln P(\omega_l) + \ln P(\mathbf{x} | \omega_l) \quad (12)$$

which is totally equivalent to Eq.(4) in terms of classification performance since  $\ln$  is a monotonically increasing function. Substituting Eq.(8) into Eq.(12) yields

$$d_l(\mathbf{x}) = \ln P(\omega_l) - (19/2)\ln 2\pi - (1/2)\ln |\Sigma_l| - (1/2)[(\mathbf{x} - \boldsymbol{\mu}_l)^T \Sigma_l^{-1} (\mathbf{x} - \boldsymbol{\mu}_l)] \quad (13)$$

Since the term  $(19/2)\ln 2\pi$  does not depend on  $l$ , it can be eliminated from the expression;  $d_l(\mathbf{x})$  then becomes

$$d_l(\mathbf{x}) = \ln P(\omega_l) - (1/2)\ln |\Sigma_l| - (1/2)[(\mathbf{x} - \boldsymbol{\mu}_l)^T \Sigma_l^{-1} (\mathbf{x} - \boldsymbol{\mu}_l)] \quad (14)$$

Equations (13) and (14) also represent the Bayes decision functions for normal patterns. If all a priori probabilities are equal:  $P(\omega_l) = 1/M$ , for  $l = 1, 2, \dots, M$ , it can easily be shown that by dropping the terms independent of index  $l$ , Eq.(14) becomes

$$d_l(\mathbf{x}) = -(1/2)[\ln |\Sigma_l| + [(\mathbf{x} - \boldsymbol{\mu}_l)^T \Sigma_l^{-1} (\mathbf{x} - \boldsymbol{\mu}_l)]] \quad (15)$$

Note that the second term in the right side of Eq.(15) is the Mahalanobis distance from  $\mathbf{x}$  to  $\boldsymbol{\mu}_l$ . Thus, it is evident from the discussion given above that if the  $M$  pattern classes are governed by the multivariate normal density functions, the decision rules for the least Mahalanobis distance can be viewed as the approximation the Bayes decision rule.

## METHODS

### Data Set

It is impossible to accurately know in advance the accuracy of a prediction method when applied to a new protein. In order to use the data bank of known structures to estimate the performance on new proteins, two requirements are essential to derive a reasonable assessment of the method's generalization ability: (1) the pairwise identity of the protein chains or domains used for developing the prediction method and those for testing should be lower than the value sufficient for modeling tertiary structure by homology, and (2) a multiple cross-validation test (ideally jack-knife) has to be performed to exclude a potential dependency of the evaluated accuracy on the particular test set chosen.

It is now well established that protein domains having more than 30% of their sequence in common adopt the same fold structures.<sup>27-30</sup> Therefore, a tool not using the homology to a protein of known structure has to be tested on those cases for which it will be used, i.e., protein domains without significant pairwise homology to those used for developing the method. In the present study, the classification method of Murzin et al.<sup>3</sup> was used for the distinction between different structural classes (SCOP, version 1.38). This database provides a detailed and comprehensive description of the structural and evolutionary relationships of proteins whose three-dimensional structures have been determined. The basic unit for classification in SCOP is the protein domain. Small proteins, and most of those with medium size, have a single domain and are, therefore, treated as a whole. The domains in large proteins are usually classified individually. The classification in SCOP is entirely manual and does not incorporate any of the "hard and fast" rules, such as those described by Nakashima et al.,<sup>11</sup> Klein and Delisi,<sup>7</sup> and Chou.<sup>15</sup> Rather, it focuses on what structural elements are within the "core" of the protein. A protein with 7 strands

and 1 helix may be " $\alpha + \beta$ " if the helix is integral to the core, while it would be "all  $\beta$ " if the helix were a non-conserved elaboration. The distinction of  $\alpha/\beta$  and  $\alpha + \beta$  is made on the interactions between the  $\alpha$  and  $\beta$  sections of the structure.  $\alpha/\beta$  proteins typically have interspersed  $\alpha$  and  $\beta$  units, while  $\alpha + \beta$  proteins typically have separate regions of mostly  $\alpha$  and mostly  $\beta$  structure. Therefore, in comparison with the other classifications only based on the percentages of secondary structures, the classification in SCOP is more natural, better reflects the objective reality, and provides a more reliable database for the study of protein structural class prediction. In the SCOP database, protein domains are classified into the following 10 categories: (1) all  $\alpha$ ; (2) all  $\beta$ ; (3)  $\alpha/\beta$ ; (4)  $\alpha + \beta$ ; (5) multidomain; (6) membrane and cell surface protein; (7) small protein; (8) peptides; (9) designed protein; and (10) non-protein. In this study, only categories (1)–(4) will be considered. The creators of SCOP have clustered the domains in the Protein Data Bank on the basis of sequence identity.<sup>31</sup> At a sequence identity level of 40%, 1,189 unique sequences corresponding to the known structural domains were found in the PDB40D\_1.37 database of SCOP. These 1,189 sequences are used as both the training and test sets in the present study.

### Self-Consistency and Jack-Knife Tests

The prediction quality was examined by two approaches. One is based upon the re-substitution test and the other upon the jack-knife test. The former is for testing the self-consistency of the algorithm, whereas the latter is for testing the results by cross-validation. When the self-consistency test is performed for the current study, the structural class for each of the domain in a given data set is predicted using the rules derived from the same data set. Testing predictive accuracy on the training set could lead to unrealistically high accuracies. An objective test of a structural class prediction method will predict the structure of a test set of proteins that are not in the training set and show no detectable sequence similarity with the training set. Since the number of proteins of known structure is limited, it is normal to develop structural class prediction methods by cross-validation techniques, or jack-knife test. In a full jack-knife test, each protein or domain in the data set is in turn moved from the set, the parameters are developed on the remaining domains, the structure of the removed domain is predicted and its accuracy measured. In other words, the structural class of each domain is predicted by the rules derived using all other domains except the one that is being predicted. During the process of jack-knife analysis, both the training data set and testing data set are actually open, and a domain will in turn move from one to the other.

### Prediction Procedure

Suppose  $\mathbf{x}$  is a protein whose structural class is to be predicted. Using the Bayes decision rule for minimum error, the prediction can be performed according to the following procedure:



**TABLE I. Comparison of Prediction Results for the Data Sets Used by Nakashima et al.<sup>11</sup> and Chou<sup>15</sup>**

| Test set  | Protein type     | Prediction method            |                            |                |
|---|------------------|------------------------------|----------------------------|----------------|
|   |                  | Chou and Zhang <sup>14</sup> | Wang and Yuan (this study) |                |
|   |                  | Self-consistency             | Self-consistency           | Jack-knife     |
| 131 proteins of Nakashima et al. <sup>11</sup>          | All $\alpha$     | 31/31 = 100%                 | 31/31 = 100%               | 17/31 = 54.9%  |
|   | All $\beta$      | 34/34 = 100%                 | 33/34 = 97.1%              | 17/34 = 50.0%  |
|   | $\alpha + \beta$ | 24/27 = 88.9%                | 27/27 = 100%               | 4/27 = 14.8%   |
|   | $\alpha/\beta$   | 35/39 = 89.7%                | 39/39 = 100%               | 18/39 = 46.2%  |
|   | Average          | 124/131 = 94.7%              | 130/131 = 99.2%            | 56/131 = 42.7% |
| $4 \times 30 = 120$ proteins of K.C. Chou <sup>15</sup> | All $\alpha$     | 30/30 = 100%                 | 30/30 = 100%               | 20/30 = 66.7%  |
|   | All $\beta$      | 30/30 = 100%                 | 30/30 = 100%               | 17/30 = 56.7%  |
|   | $\alpha + \beta$ | 30/30 = 100%                 | 30/30 = 100%               | 14/30 = 46.7%  |
|   | $\alpha/\beta$   | 29/30 = 96.7%                | 30/30 = 100%               | 13/30 = 43.3%  |
|   | Average          | 119/120 = 99.2%              | 120/120 = 100%             | 64/120 = 53.3% |

(1) Knowing the amino acid compositions of the database proteins, normalize their amino acid components by dividing the number of each component amino acid by the total number of amino acids in the protein.

(2) Eliminate one of the 20 normalized amino acid components, thereby defining a 19-D space, and express the proteins as points in the 19-D space.

(3) Calculate the mean vectors and  $19 \times 19$  covariance matrixes for the  $\alpha$ ,  $\beta$ ,  $\alpha + \beta$ , and  $\alpha/\beta$  proteins from the database proteins.

(4) Calculate the normalized frequencies of 20 amino acids of the unknown protein  $\mathbf{x}$  as follows:

$$x_i = \frac{v_i}{\sum_{j=1}^{20} v_j} \quad (i = 1, 2, \dots, 20)$$

where  $v_i$  is number of  $i$ th amino acid in the protein  $\mathbf{x}$  ( $i = 1, 2, \dots, 20$ ). Thus, the protein  $\mathbf{x}$  also corresponds to a point  $(x_1, x_2, \dots, x_{19})$  in the 19-D space.

(5) For the point of the unknown protein  $\mathbf{x}$ , calculate the conditional probabilities of occurrence of  $\mathbf{x}$  for the above four types of protein,  $P(\mathbf{x}|\omega_l)$ ,  $l \in \alpha, \beta, \alpha + \beta, \alpha/\beta$  [cf., Equation (8)].

(6) The unknown protein is predicted to have the same structural class as the one which the value of decision function,  $d_l(\mathbf{x}) = P(\omega_l)P(\mathbf{x}|\omega_l)$  or  $d_l(\mathbf{x}) = P(\mathbf{x}|\omega_l)$  if all a priori probabilities are equal, is the largest [cf., Equations (3) and (4)].

## RESULTS AND DISCUSSION

### Comparison to the Previous Component-Coupled Methods

In order to compare class prediction quality by the present and the other component-coupled methods, two protein data sets used by Nakashima et al.<sup>11</sup> and Chou<sup>15</sup> were tested first. In their original paper Nakashima et al. examined 135 proteins of which 31 are  $\alpha$  proteins, 34  $\beta$ , 27  $\alpha + \beta$ , 39  $\alpha/\beta$ , and 4 irregular. The criteria of classification are given by

$\alpha$  proteins  $\Rightarrow \alpha > 15\%, \beta < 10\%$

$\beta$  proteins  $\Rightarrow \alpha < 15\%, \beta > 10\%$

$\alpha + \beta$  proteins  $\Rightarrow \alpha > 15\%, \beta > 10\%$

with dominantly antiparallel  $\beta$ -sheets

$\alpha/\beta$  proteins  $\Rightarrow \alpha > 15\%, \beta > 10\%$

with dominantly parallel  $\beta$ -sheets

Irregular proteins otherwise

The irregular proteins have been left out in this study because their number is only four, too small to have any statistical significance. Therefore, the prediction and comparison will be made based on the remaining 131 proteins. Assuming that four a priori probabilities are equal,  $P(\omega_\alpha) = P(\omega_\beta) = P(\omega_{\alpha+\beta}) = P(\omega_{\alpha/\beta}) = 1/4$ , the rates of correct prediction for 31  $\alpha$ , 34  $\beta$ , 27  $\alpha + \beta$ , and 39  $\alpha/\beta$  proteins are 100%, 97.1%, 100%, and 100%, respectively. If the average accuracy is defined by the percentage of the number of correct prediction events for all classes divided by the number of total prediction events, i.e.,

$Q = \text{average accuracy}$

$$= \frac{\text{total number of correct prediction events}}{\text{total number of prediction events}}$$

we have the average accuracy of 99.2% for predicting the 131 proteins by the current method. This result shows that the average prediction accuracy of our method is 4.5% higher than that obtained by the least Mahalanobis distance method. However, when predicted with jack-knife test, the average accuracy was only 42.7%. The results are summarized in Table I.

Although proteins of known structure are generally classified into one of the four structural classes, there is no unified quantitative measure for making such a classifica-

**TABLE II. Comparison of Prediction Results for the Data Sets Used by Eisenhaber et al.<sup>17</sup>**

| Test set  | Protein type      | Prediction method               |           |                            |           |
|---|-------------------|---------------------------------|-----------|----------------------------|-----------|
|   |                   | Eisenhaber et al. <sup>17</sup> |           | Wang and Yuan (this study) |           |
|   |                   | Self-consistency                | Jackknife | Self-consistency           | Jackknife |
| 260 proteins of Eisenhaber et al. <sup>17</sup><br>(set $\geq 2.0$ Å) | All $\alpha$ (55) | 54.5%                           | 47.3%     | 78.2%                      | 45.5%     |
|   | All $\beta$ (78)  | 82.1%                           | 76.9%     | 88.5%                      | 69.2%     |
|   | Mixed (127)       | 50.4%                           | 50.4%     | 89.0%                      | 66.1%     |
|   | Average           | 60.8%                           | 57.7%     | 86.5%                      | 62.7%     |
| 471 proteins of Eisenhaber et al. <sup>17</sup><br>(Set 3.0 Å)        | All $\alpha$ (99) | 57.6%                           | 54.5%     | 69.7%                      | 56.6%     |
|   | All $\beta$ (140) | 77.9%                           | 76.4%     | 82.9%                      | 65.0%     |
|   | Mixed (232)       | 46.6%                           | 47.0%     | 81.9%                      | 72.0%     |
|   | Average           | 58.2%                           | 57.3%     | 79.6%                      | 66.7%     |

tion. The relevant percentages set by Nakashima et al. for  $\alpha$  proteins ( $\alpha > 15\%$ ) and  $\beta$  proteins ( $\beta > 10\%$ ) do not seem large enough to reflect the real features of the two structural classes. In view of these, Chou<sup>15</sup> proposed a new method that classifies protein according to the following quantitative criterion:

$\alpha$  proteins  $\Rightarrow \alpha \geq 40\%, \beta \leq 5\%$

$\beta$  proteins  $\Rightarrow \alpha \leq 5\%, \beta \geq 40\%$

$\alpha + \beta$  proteins  $\Rightarrow \alpha \geq 15\%, \beta \geq 15\%$

with more than 60% parallel  $\beta$ -sheets

$\alpha/\beta$  proteins  $\Rightarrow \alpha \geq 15\%, \beta \leq 15\%$

with more than 60% parallel  $\beta$ -sheets

$\zeta$  proteins  $\Rightarrow \alpha \leq 10\%, \beta \leq 10\%$

and consideration of the coupling effect among different amino acid components would not improve the class prediction quality. However, since the original least Mahalanobis distance method developed by Chou and Zhang is not valid for their training set, in which the subset sizes are very much different, the conclusion obtained by Eisenhaber et al. could still be questionable.<sup>19,20</sup> As mentioned earlier, the component-coupled method proposed in the present study can be used to deal with the cases where the training subset sizes are different, therefore, it is interested in testing our method with the same database used by Eisenhaber et al. According to their paper, the following rule was used to classify the structural classes of proteins:

all  $\alpha$  proteins  $\Rightarrow \alpha > 15\%, \beta < 10\%$

all  $\beta$  proteins  $\Rightarrow \alpha < 15\%, \beta > 10\%$

mixed class proteins  $\Rightarrow \alpha > 15\%, \beta > 10\%$

irregular proteins  $\Rightarrow$  otherwise

where the contents of protein secondary structures were computed based on the dictionary by Kabsch and Sander.<sup>32</sup> According to the new criteria and selection principle, 120 structure-known proteins were thus selected and classified into 30  $\alpha$ , 30  $\beta$ , 30  $\alpha + \beta$  and 30  $\alpha/\beta$  proteins. Based on such a training database, all of the 120 proteins were correctly predicted according to the Bayes decision rule. It can be seen from Table I that for this carefully selected database, although the prediction accuracy for the self-consistency test was 100% by our method, only about half proteins were correctly predicted for the jack-knife test.

One possible explanation for this remarkable difference between the self-consistency and jack-knife tests is that the data base used is too small so that the information loss due to jack-knife will have a greater impact on the prediction results, because these kind of methods need more training data to make their prediction mechanism work properly. Eisenhaber et al.<sup>17</sup> have tested the least Mahalanobis distance method with a larger database of non-homologous crystal structures (residue identity among all aligned pairs sequences  $\leq 35\%$ , minimal sequence length 80 residues). They concluded that a jury decision among four structural classes based only on the amino acid composition of the query protein is at best 50–60% correct,

Their rule does not distinguish between the  $\alpha/\beta$  class and the  $\alpha + \beta$  class and places them in one class, the so-called mixed class. As the subsets of irregular proteins is too small to be statistically significant, they have been removed from consideration. Table II gives a comparison of results obtained by our method and by the vector decomposition method of Eisenhaber et al. with two different databases. It can be seen from this table that the self-consistency test of the current method for 260 proteins (the second training set in Table I of their paper) is 86.5%, which is about 26% higher than the result obtained by their component-coupled method. Similarly, when the jack-knife test was performed for the same data set, the success rate drops to the value of 60%. Note that for the larger data set containing 475 proteins, the overall rate of correct prediction by the current method for the self-consistency test decreases but that for jack-knife test increases slightly. This result suggests that the limited size of data set may not be the main reason for the remarkable difference between the self-consistency and jack-knife tests.

**TABLE III. Prediction Results for the 1,189 Domains in the PDB40D\_1.37 Database by the Present Method**

| Test set   | Protein type     | Number of domains | Prediction accuracy |                  |
|--|------------------|-------------------|---------------------|------------------|
|  |                  |                   | Self-consistency    | Jack-knife       |
| 1,189 protein domains obtained from the SCOP database (sequence identity <40%) | All $\alpha$     | 263               | 167/263 = 63.5%     | 144/263 = 54.8%  |
|  | All $\beta$      | 317               | 204/317 = 64.4%     | 181/317 = 57.1%  |
|  | $\alpha + \beta$ | 270               | 107/270 = 39.6%     | 60/270 = 22.2%   |
|  | $\alpha/\beta$   | 339               | 281/339 = 82.9%     | 255/339 = 75.2%  |
|  | Average          |                   | 759/1189 = 63.8%    | 640/1189 = 53.8% |

### Effect of the Size of Data Set on Prediction Accuracy

The validation of all prediction methods is, of course, dependent on the classification scheme. The surprisingly poor results (with the data set in Table II) of the composition-coupled method may be due partly to the classification principles. As pointed out by Chou et al.,<sup>19,20</sup> protein structural classification should be based on the domain structure, while the classification of the three data sets tested above are based on the whole protein chain. If a data set is constructed according to an arbitrary or incorrect classification rule, it certainly cannot objectively reflect the relationship between the structural class of a protein and its amino acid composition. All the calculated results based on such a data set would be meaningless. To avoid this, we also test our method with a more reasonable classification scheme and the largest data set at present. As mentioned above, SCOP database is now a more reasonable data base and available in the network. At a sequence identity level of 40%, 1,189 unique sequence corresponding to the known structural domains were found in the PDB40D\_1.37 database of SCOP. There are 263 all  $\alpha$ , 317 all  $\beta$ , 270  $\alpha + \beta$ , and 339  $\alpha/\beta$  protein domains. The results predicted by the current method for the 1,189 domains are summarized in Table III. The overall rate of correct prediction by the current method is 63.8% for the self-consistency test, and 53.8% for the jack-knife test. Thus we think that the poor results obtained cannot be explained by differences in protein structural classification.

The goal of testing the prediction tool is to assess the accuracy to be expected for any new protein sequence. Since different test sets yield different results, it is not sufficient to use only one set. In order to exclude the artifact owing to selecting different sets of protein and analyze the impact of the learning set, the new method was tested with learning sets of varying representativity and size.  $N$  ( $N = 40, 80, 120, 160, 200, 240$ ) representative domains were selected randomly from each of the four structural classes in the PDB40D database. For a database of such  $4 \times N$  domains, both the self-consistency and jack-knife tests were performed. This procedure was repeated for 200 times. These 200 sets were tested to determine the average results and the degree of variation that can occur. The average over all 200 tests gives a reasonable estimate of the prediction accuracy. The individual and overall rates of class prediction are summarized in Table IV. We also computed the average difference

for the current method on the 200 different data sets, which is the deviation of prediction accuracy and reflects the statistical fluctuation of the data sets. The larger the  $N$ , the smaller the deviation. The performance of the method turns out to be strongly dependent on the size of the data set. For the self-consistency test, more than 90% accuracy is obtained with a representative learning set of size = 160 while the performance came in about 20% lower with a data set of size = 960. This suggests that considerable care has to be used in evaluating the results of a single test set. By enlarging the data set, the self-consistency performance showed a steady decrease, yet the jack-knife performance improved slightly. This decrease in the “generalization gap”—the difference in performance between the learning and test sets—illustrates that a decrease in some of the noisy sequence-specific (or example-specific) information occurs with this simplification but without a corresponding decrease in structural classes, this underscores the importance of cross-validation in studies of this type.<sup>33</sup>

Figure 1 shows the dependence of the prediction accuracy achieved by the current method on the size of data set. Since the limited size of the current protein data bank, the predicted accuracy is, to some extent, dependent on the set of proteins selected by the predictor. The similar problem also exists when prediction is performed for a set of testing proteins. Only when the number of proteins considered is sufficiently large can the bias due to the selection of different protein sets be eliminated. Therefore, to make a fair comparison of different prediction methods, one should adopt the objective accuracy as a criterion. The objective accuracy is actually an asymptotical limit for the rate of correct prediction computed for a sufficiently large number of proteins. It can be seen from Figure 1 that as expected, the jack-knife test rate will be close to the self-consistency test rate and the bias due to selection of different sets can be eliminated when the data set becomes sufficiently large. The common horizontal asymptotical limit suggests that the objective accuracy of our prediction method should be about 60%.

For a small data set, a large generalization gap may be considered as implying that the algorithm is “memorizing” the information in a rigid fashion rather than learning the underlying informational concepts behind a classification scheme. Both the current method and the least Mahalanobis distance method show the largest generalization gaps, with the performance on the learning sets at or near 100% accuracy while the performance on the test sets came in

**TABLE IV. Prediction Results for the Different Size of Data Sets by the Present Method**

| Test set  | Protein type     | Prediction accuracy  |                |
|---|------------------|----------------------|----------------|
|   |                  | Self-consistency (%) | Jack-knife (%) |
| $4 \times 40 = 160$<br>protein domains<br>randomly<br>selected from the<br>SCOP database<br>containing 1,189<br>domains (200<br>times)  | All $\alpha$     | $88.1 \pm 4.7$       | $50.9 \pm 8.8$ |
|   | All $\beta$      | $93.9 \pm 3.8$       | $50.1 \pm 7.9$ |
|   | $\alpha + \beta$ | $90.7 \pm 4.1$       | $42.9 \pm 9.3$ |
|   | $\alpha/\beta$   | $98.3 \pm 2.1$       | $36.0 \pm 6.4$ |
|   | Average          | $93.0 \pm 2.0$       | $45.0 \pm 5.0$ |
| $4 \times 80 = 320$<br>protein domains<br>randomly<br>selected from the<br>SCOP database<br>containing 1,189<br>domains (200<br>times)  | All $\alpha$     | $74.6 \pm 4.2$       | $53.4 \pm 5.2$ |
|   | All $\beta$      | $82.5 \pm 4.1$       | $52.9 \pm 5.7$ |
|   | $\alpha + \beta$ | $73.8 \pm 4.2$       | $42.1 \pm 5.8$ |
|   | $\alpha/\beta$   | $94.5 \pm 2.4$       | $56.6 \pm 4.4$ |
|   | Average          | $81.0 \pm 2.0$       | $51.0 \pm 3.0$ |
| $4 \times 120 = 480$<br>protein domains<br>randomly<br>selected from the<br>SCOP database<br>containing 1,189<br>domains (200<br>times) | All $\alpha$     | $69.2 \pm 3.4$       | $54.5 \pm 4.4$ |
|   | All $\beta$      | $76.7 \pm 3.9$       | $53.8 \pm 4.8$ |
|   | $\alpha + \beta$ | $64.8 \pm 3.6$       | $40.6 \pm 4.2$ |
|   | $\alpha/\beta$   | $90.9 \pm 2.1$       | $64.7 \pm 3.4$ |
|   | Average          | $75.0 \pm 3.0$       | $53.0 \pm 2.0$ |
| $4 \times 160 = 640$<br>protein domains<br>randomly<br>selected from the<br>SCOP database<br>containing 1,189<br>domains (200<br>times) | All $\alpha$     | $66.4 \pm 2.3$       | $55.0 \pm 2.8$ |
|   | All $\beta$      | $73.1 \pm 3.1$       | $54.7 \pm 3.4$ |
|   | $\alpha + \beta$ | $59.4 \pm 3.0$       | $39.1 \pm 3.4$ |
|   | $\alpha/\beta$   | $89.1 \pm 2.0$       | $69.5 \pm 2.9$ |
|   | Average          | $72.0 \pm 1.0$       | $55.0 \pm 2.0$ |
| $4 \times 200 = 800$<br>protein domains<br>randomly<br>selected from the<br>SCOP database<br>containing 1,189<br>domains (200<br>times) | All $\alpha$     | $64.3 \pm 2.1$       | $55.3 \pm 2.2$ |
|   | All $\beta$      | $70.9 \pm 2.6$       | $54.9 \pm 3.0$ |
|   | $\alpha + \beta$ | $55.8 \pm 2.3$       | $38.5 \pm 2.6$ |
|   | $\alpha/\beta$   | $87.9 \pm 1.7$       | $72.2 \pm 2.2$ |
|   | Average          | $70.0 \pm 1.0$       | $55.0 \pm 1.0$ |
| $4 \times 240 = 960$<br>protein domains<br>randomly<br>selected from the<br>SCOP database<br>containing 1,189<br>domains (200<br>times) | All $\alpha$     | $63.3 \pm 1.3$       | $55.4 \pm 1.6$ |
|   | All $\beta$      | $68.9 \pm 1.8$       | $55.8 \pm 2.2$ |
|   | $\alpha + \beta$ | $52.8 \pm 2.0$       | $37.4 \pm 2.2$ |
|   | $\alpha/\beta$   | $86.5 \pm 1.5$       | $73.9 \pm 1.7$ |
|   | Average          | $67.9 \pm 0.7$       | $56.0 \pm 1.0$ |

50–55% lower. This memorization is likely due to the large number of parameters with respect to the size of each learning set, which tends to cause the current method to readily extract and retain sequence-specific information. The best solution to the memorization problem appears to be to increase the size of the training set. An increased training set size will also likely decrease some of the “wobble” associated with small sets; that is, larger sets will

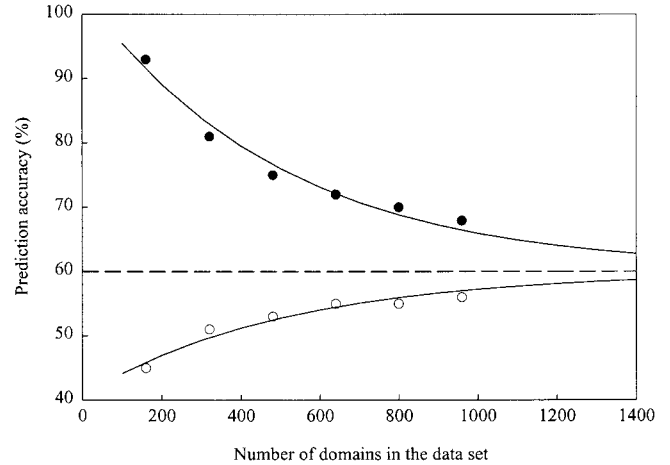


Fig. 1. Dependence of structural class prediction accuracy on the size of the protein domain data set. (•) Self-consistency test; (○) Jack-knife test.

contain more structural class information, decreasing the variation in prediction accuracy due to the random compositions of the training and test sets.<sup>33</sup>

### Effect of Sequence Homology on Results

Whether knowledge of the fractions of the 20 amino acids is sufficient alone for predicting the structural class of a given protein is steeped in controversy. The level of prediction accuracy is still available ranging from 60% to near 100% and clearly depends on the definitions of structural classes and the set of database structures considered for performing the analysis. Analyses of known three-dimensional protein structures and amino acid sequences revealed that proteins are clustered into families whose members have evolved from a common ancestor, share a characteristic fold, and, sometimes, have a similar function. The number of families of related protein structures critically depends on the value of the homology threshold applied in the protein structure comparison routines. According to the classification of protein structures in SCOP, all protein domains with more than 30% identity belong to the same of protein family and must have the same structural class. Therefore, the structural class assignment of a new protein domain with homologous > 30% to a protein of known structure can be easily performed by sequence alignment, and any prediction method for the protein structural classes should only address those proteins for which no homologous proteins (at a sequence identity level of 30%) are found in the Protein Data Bank. Since smaller test sets always open the theoretical possibility that just proteins which would be badly predicted have been missed and that the prediction rates are overestimated, it is extremely important that success rates are calculated from as large as possible representative subsets of PDB. The most stringent test set of tertiary structures available would consist of all unrelated proteins, the structure of which is known. In the PDB40D\_1.37 database, the 1,189 domains belong to 675



**TABLE V. The 675 Protein Domains (Sequence Identity <30%) Extracted From the PDB40D\_1.37 Database of SCOP for Self-Consistency and Jack-Knife Tests**1. 155  $\alpha$  domains

|        |        |         |        |        |        |        |        |        |        |
|--------|--------|---------|--------|--------|--------|--------|--------|--------|--------|
| 3sdha_ | llea_  | lerc_   | 2hmqa_ | lpral  | lbmtal | lezm_1 | 2abk_  | laora1 | lcsma_ |
| lcpca_ | laoy_  | laca_   | lvtmp_ | 2tct_1 | ltpt_1 | 2ts1_1 | lgln_1 | lvnc_  | lpprm1 |
| lgri_1 | lcgpa1 | lre_    | lbucal | lmi_1  | lc5a_  | lbmfal | lrlr_1 | lsig_  | lpoc_  |
| lhdj_  | lopc_  | lvii_   | lfapb_ | lcoo_  | locch_ | labv_  | ldnpa1 | lbvp11 | lpoa_  |
| lseta1 | lhsta_ | 2end_   | lryt_1 | lan2a_ | lhyp_  | lak4c_ | lila_1 | lrgp_  | lbeo_  |
| lmnga1 | letd_  | llis_   | lmmob_ | 4icb_  | lbip_  | 2hmx_  | 2pgd_1 | 2bct_  | lrtm11 |
| ldvh_  | 2hts_  | lbmfg_  | lbgc_  | lsra_  | lpnb_1 | ljvr_  | lyve1  | lrv_   | 2ztaa_ |
| laofa1 | ldpra1 | llbu_1  | lcsga_ | lrro_  | lolga_ | lvln_1 | lutg_  | lsly_1 | lifj_  |
| letpa1 | lxgsa1 | lhme_   | lilk_  | 2scpa_ | ladt_1 | lvola1 | lgim_  | locce_ | lvdfa_ |
| lenh_  | lfow_  | lbifma_ | laap_  | ldjxa1 | lihfb_ | lcuk_2 | lcem_  | 2sblb1 | ldkga1 |
| lhcr_  | lcuk_1 | ltafa_  | limq_  | lcpo_1 | lalo_1 | ldpra2 | Seas_1 | 2tct_  | ldipa_ |
| lmsec1 | ltns_  | lmmog_  | lrpo_  | lpax_1 | lab3_  | lgr_   | lcsh_  | llbd_  |        |
| lpdnc_ | 2spca_ | lape_   | lytfb1 | lmyka_ | laep_  | lcrka1 | lphb_  | lfps_  |        |
| ligna1 | lfec2_ | 2liga_  | lecia_ | lcmba_ | lnkl_  | lagre_ | lfipa_ | Seas_2 |        |
| lsfe_1 | lgab_  | 256ba_  | loctc2 | ldsba1 | lhvd_  | laru_  | lpre_  | lgrl_1 |        |
| lbia_1 | lbb1_  | 2cya_   | llia_  | 2gsta1 | ltada1 | lmhl_1 | 2wrpr_ | lcema_ |        |

2. 156  $\beta$  domains

|        |        |        |        |        |        |        |        |        |        |
|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
| lneu_  | lnbca_ | 2bpa1_ | lhgea_ | lckaa_ | lmjc_  | larb_  | lwpoa_ | 2cba_  | lrgs_1 |
| lcd1a1 | lqba_2 | lstma_ | laot_  | lmmd_1 | lckma1 | lbt_   | lpkya1 | 3bcl_  | 2arca_ |
| lvcaa1 | ltupa_ | 2bbva_ | lkn_   | lvie_  | lrip_  | 2snv_  | lhb_   | lospo_ | lwapa_ |
| lvcaa2 | lctm_1 | lbtt1_ | laly_  | lpse_  | lyhb_  | lhava_ | lhms_  | lvmoa_ | lbdo_  |
| lgof_1 | lcdg_2 | 4gcr_1 | lthw_  | lihwa_ | lpyp_  | lbco_1 | lsria_ | ldlc_2 | lbncal |
| 2hft_1 | letal_ | lpr_1  | lscs_  | laono_ | lprch1 | lbmf2  | lsmpl_ | lmsaa_ | lctm_2 |
| lbglal | 2pcda_ | lwkt_  | lcpn_  | 2ohxa1 | 2fgf_  | lsfta1 | 2cpl_  | lkapp1 | lgpr_  |
| lggta2 | lhoe_  | 2sblb2 | lsaa_  | lpdr_  | lilb_  | lfiva_ | lhxn_  | 2pec_  | lgzi_  |
| lncia_ | lpic_  | lhpl1  | lsaca_ | lfgp_  | labrb1 | lepne_ | 2sil_  | lidk_  | 2kaub_ |
| lnoa_  | lcyx_  | lpgs_1 | lkit_1 | lwhi_  | lwba_  | 2eng_  | lgof_3 | ltsp_  | ldupa_ |
| lxsoa_ | laozal | lslua_ | lcela_ | lsty_  | lhcd_  | lbw3_  | 2bbkh_ | lxa_   | ltul_  |
| 3dpa_1 | ldjxa2 | lgof_2 | lxnb_  | ltsd_  | lfmb_1 | lcxsa1 | 2treb_ | ltdta_ | 2kauc1 |
| lmspa_ | lrsy_  | ldlc_1 | lbgl4  | lesfa1 | 2pia_1 | lgral  | 4aaha_ | lthja_ |        |
| 4kbpal | 3dpa_2 | lbgl3  | loaca1 | lasya1 | lfuia1 | lmai_  | laofa2 | 2phla1 |        |
| lddt_1 | lwho_  | lulo_  | lbia_2 | lcuk_3 | left_1 | lirsa_ | lbplb1 | lpmi_  |        |
| lexg_  | lnpoa_ | lbvp12 | lumua_ | 3ulla_ | left_2 | lytfc1 | ldkga2 | lcgpa2 |        |

3. 184  $\alpha/\beta$  domains

|        |        |        |        |        |        |        |        |        |        |
|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
| lgmpa_ | lubi_  | lema_  | lgpma3 | laf_   | lotga_ | lhfc_  | lba_   | lseia_ | lfjma_ |
| 2baa_  | lgua_  | lfkd_  | leba2  | lpsda3 | lfim_  | lqba_4 | 3b5c_  | ldiv_  | 2cmd_2 |
| 1931_  | lfid_  | lctn_3 | lgri_3 | lmia_2 | lgd1o2 | lshaa_ | lvcc_  | 2glt_2 | laiha_ |
| 1191_  | lesfa2 | lgri_2 | lfca_  | lfwp_  | ldapa2 | lptf_  | lyua_1 | lbnc3  | lmrj_  |
| lgbs_  | ltif_  | lfroa_ | lfid2_ | lregx_ | lofga2 | lba_   | lah6_  | lscub2 | lts_1  |
| lsly_2 | lgr_1  | lhan_1 | lxer_  | lab8a_ | loaca4 | laf5_  | lorda3 | ldik_3 | lpax_2 |
| lchka_ | lcoy_2 | lmkaa_ | lvjw_  | lvaoa1 | lbpl_1 | lgtpa_ | lsmna_ | lckma2 | ldef_  |
| 2act_  | lpbe_2 | lcei_  | lraab1 | lgeo_1 | 2sici_ | lgtga_ | lchma2 | lvaoa2 | llit_  |
| lggta4 | lkifa2 | lalo_3 | lpba_  | lbu_2  | 2ms2a_ | lscea_ | lgr_2  | lmbb_1 | lprtb2 |
| 7rsa_  | lga1_2 | lqapa2 | lspbp_ | lvhh_  | lgesa3 | lefmb_ | lcrka2 | lmbb_2 | lts_   |
| lag2_  | lgnd_2 | ltpt_3 | lmli_  | ltig_  | lrs_   | lcby_  | lytba1 | lako_  | 3fb_   |
| 2kaua_ | lmola_ | ltfe_  | lpil_  | luae_  | ltbd_  | lseta2 | 3pmga4 | laora2 | lmsk_  |
| lhum_  | lcewi_ | lmnga2 | lnpk_  | lkpta_ | lnox_  | lbia_3 | lbv1_  | lgdoa_ | lpmd_3 |
| lpmd_1 | loaca3 | lctf_  | lup1_1 | lkvd_1 | lkuh_  | 2vik_  | lmsa_1 | lpnk_1 |        |
| lso_   | lst_   | 2reb_2 | 2bopa_ | 3rubs_ | lezm_2 | lahq_  | 2pola1 | lpmb_  |        |
| lkpa_  | louna_ | lstu_  | 3rubl2 | ldcoa_ | lml_   | lpne_  | lplq_1 | lapy_1 |        |
| lhxpai | ludii_ | lpkp_2 | laps_  | lxxaa_ | last_  | 2phy_  | lalo_4 | lpya_1 |        |
| ldar_3 | lcl1a2 | lpda_2 | lris_  | 2chsa_ | latla_ | lmut_  | lgeo_4 | lbme_  |        |
| ligd_  | laak_  | lvig_  | ldar_4 | lotfa_ | lkapp2 | ltys_  | lhqi_  | 4kbp2  |        |

4. 180  $\alpha/\beta$  domains

|        |        |        |        |        |        |        |        |        |        |
|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
| lcdg_4 | lpkya2 | lgesa1 | lorda1 | lhrda1 | lrlaa_ | lbam_  | 3pgm_  | lhta_  | lorb_1 |
| lbyb_  | ldik_1 | ldik_2 | lcus_  | lscua2 | 3cla_  | lpvua_ | lrpa_  | ltca_  | lcxsa2 |
| lxzya_ | 3rubl1 | lzyma_ | lesc_  | lbnc2  | lphr_  | lcfr_  | lnula_ | ltahb_ | lad3a_ |
| icbg_  | ltpfa_ | laco_1 | 2naca1 | 2din_1 | le2b_  | 2rsla_ | lfaa_  | lhpa2  | laco_2 |
| lhvq_  | 2xis_  | lgri_2 | lfmb_2 | 2glt_1 | lvhra_ | lpdo_  | lgara_ | lyasa_ | 3pmga1 |
| lqba_3 | lluca_ | lbrsd_ | 2pia_2 | lpvda1 | 2hnp_  | lhpm_1 | lvid_  | 2masa_ | lfuia2 |
| lad_   | lnfp_  | ldfj_1 | 2ts1_2 | lnbaa_ | 2trxa_ | 2yhx_1 | lxvaa_ | lpbn_  | lph_   |
| 2kauc2 | lqapa1 | lnzya_ | lgpma1 | ldeaa_ | lmek_  | lgtag1 | lv39_  | 2ctb_  | lmioa_ |
| lpta_  | ldjxa3 | lrvva_ | ltpt_2 | lpvda2 | ldsba2 | lchma1 | lhmy_  | lobr_  | 2bgu_  |
| lnal_1 | lgy_   | lpau_1 | ldnpa2 | ltrka1 | lgpla_ | 2m2_   | lart_  | ljam_2 | lgpb_  |
| lucwa_ | lreqa1 | ludg_  | 2tmda3 | lgky_  | 2gsta2 | lasu_  | ltpla_ | lmp_   | 3pga1  |
| ldosa_ | lpud_  | lmia_1 | lkifa1 | 5p21_  | 2trcp_ | lbc_2  | 2dkb_  | ldraa_ | lpfka_ |
| 2acr_  | lsfta2 | 3chy_  | 2ohxa2 | lmmd_2 | ltrka3 | lkfd_1 | lorda2 | lhrda2 | layl_  |
| lak5_  | ltm1_  | lscua1 | lxcl_  | ldts_  | lpkya3 | lnoya_ | lgpma2 | lfua_  | 2dri_  |
| lebh1  | lrlr_2 | 2fx2_  | lgd1o1 | 2reb_1 | lham_1 | lhjra_ | lwht_1 | 2anha_ | 2olba_ |
| 2mmr_1 | 2tmda2 | lqrda_ | 2naca2 | lble_  | lhta1  | lsfe_2 | lede_  | lidm_  | lct_   |
| ldora_ | lcoy_1 | lbmta2 | 2cmd_1 | lchd_  | leria_ | ltfr_  | ldin_  | lraa1  | lpxta1 |
| lpil_1 | lgnd_1 | lreqa2 | lyve12 | lcee_  | lrva_  | lexna_ | lbroa_ | lwsyb_ | lctt_1 |

**TABLE VI. Prediction Results for the 675 Domains of Table V by the Present Method**

| Test set  | Protein type     | Number of domains | Prediction accuracy |                 |
|---|------------------|-------------------|---------------------|-----------------|
|   |                  |                   | Self-consistency    | Jack-knife      |
| 675 protein domains given in Table V (sequence identity <30%) | All $\alpha$     | 155               | 97/155 = 62.6%      | 83/155 = 53.5%  |
|   | All $\beta$      | 156               | 98/156 = 62.8%      | 66/156 = 42.3%  |
|   | $\alpha + \beta$ | 184               | 99/184 = 53.8%      | 52/184 = 28.3%  |
|   | $\alpha/\beta$   | 180               | 156/180 = 86.7%     | 123/180 = 68.3% |
|   | Average          |                   | 450/675 = 66.7%     | 324/675 = 48.0% |
| 675 protein domains given in Table V (sequence identity <30%) | All $\alpha$     | 155               | 97/155 = 62.6%      | 85/155 = 54.8%  |
|   | All $\beta$      | 156               | 103/156 = 66.0%     | 76/156 = 48.7%  |
|   | Mixed            | 364               | 296/364 = 81.3%     | 255/364 = 70.0% |
|   | Average          |                   | 496/675 = 73.5%     | 416/675 = 61.6% |

families. Therefore, a more rigorous test for the new method should be based on all representatives selected from each of the 675 families. A largest possible subset of non-homologous structures was established based on the PDB40D\_1.37 database by the following procedure: the first domains listed in each of the families was taken as the representative of this family (Table V). The structural class prediction results obtained by our method based on the self-consistency and the jack-knife tests for this data set are listed in Table VI. It can be seen from this table that with this data set of non-homologous structures, the accuracy of the jack-knife test was about 7% lower than a previous result obtained from  $4 \times 160 = 640$  domains in Table IV. This result indicates that the existence of homologous domains in the testing data set will significantly affect the prediction accuracy of the composition-coupled method. As expected, the overall rate for 3-type prediction is about 10% higher than that for 4-type prediction. In a recent study, Chou and Maggiogra<sup>20</sup> reported that with increasing the size of data sets, the prediction accuracy for the self-consistency test remains almost unvaried while that for the jack-knife test increases from 63.77% to 84.12%. According to this observation, they concluded that by expanding a database to reduce the information loss, the overall jack-knife rate by the component-coupled algorithm can be improved significantly. This is obviously contradictory to the results of Table IV in this paper. The origin of this puzzling difference between our results and theirs was found to be due to sequence homology in their data sets, which led to differences in accuracies. For example, the 138 domains of Table I, 253 domains of Table II, and 359 domains of Table III in their paper belong to 102, 129, and 130 protein families, respectively. The average number of homologous domains per family in the three data sets are 1.35, 1.96, and 2.76, indicating that even though each domain in the data set is singled out in turn as a "test domain" and all the rule-parameters are determined from the remaining domains, the memorization effects that included in the self-consistency tests cannot be completely removed. One problem with the jack-knife method is that sequence homology within the set may invalidate the assumption that the training set is devoid of information about the tested protein. If the training data and test data are identical or highly homologous, then the prediction accuracy could be misleadingly

high. Therefore, a unbiased test in which these algorithms are applied to proteins without significant sequence homology has to be done.

## CONCLUSIONS

The accurate prediction of structural classes from amino acid composition alone is an important issue, which has been the object of a number of recent studies. However, the success rates reported for structural class prediction with different methods are contradictory. The problem of recognizing structural class of a protein knowing only its amino acid composition appears completely solved by the least Mahalanobis distance method. The highest reported prediction accuracies are near 100%.<sup>13-15,19,20</sup> This is surprising, because only the amino acid composition has been used in this method, while traditional sequence-based secondary structure prediction achieve success rates of about 75% for structural class only with extensive input information (full sequence of the query protein, its amino acid composition and length multiple alignments with homologous sequences).<sup>34-37</sup> In this article, we resolve the paradox. Our objectives have been to validate the relationship between amino acid composition and structural class by using a Bayes method and to provide the possible upper limit of the prediction rate for structural classes. When applying both the self-consistency and jackknife tests on a larger data set without significant pairwise similarity, we found that knowledge of amino acid composition alone cannot lead to a success rate higher than 60% for a 4-type class prediction by our method. The apparent relatively high accuracy level (more than 90%) attained in the previous studies, which exceed the success rates (75%) of structural class predictions using traditional secondary structure prediction techniques (including those combining evolutionary information and neural networks) was due to the preselection of test sets, which may not be adequately representative of all unrelated proteins.

## ACKNOWLEDGMENTS

This work was supported in part by grant 103-13-03-02 from the 863 High-Technology Foundation and Pandeng project of the Ministry of Science and Technology.

## REFERENCES

1. Levitt M, Chothia C. Structural patterns in globular proteins. *Nature* 1976;261: 552-558.

2. Richardson JS, Richardson DC. Principles and patterns of protein conformation. In: Fasman GD, editor. Prediction of protein structure and the principles of protein conformation. New York: Plenum Press; 1989. p 1–98.
3. Murzin AG, Brenner SE, Hubbard T, Chothia C. SCOP: A structural classification of protein database for the investigation of sequence and structures. *J Mol Biol* 1995;247:536–540.
4. Michie AD, Orengo CA, Thornton JJ. Analysis of domain structural class using an automated class assignment protocol. *J Mol Biol* 1996;262:168–185.
5. Deleage G, Roux B. Use of class prediction to improve protein secondary structure prediction. In: Fasman GD, editor. Prediction of protein structure and the principles of protein conformation. New York: Plenum Press; 1989. p 417–465.
6. Kneller DG, Cohen FE, Langridge R. Improvements in protein secondary-structure prediction by enhanced neural networks. *J Mol Biol* 1990;214:171–182.
7. Klein P, Delisi C. Prediction of protein structural class from amino acid sequence. *Biopolymers* 1986;25:1659–1672.
8. Klein P. Prediction of protein structural class by discriminant analysis. *Biochim Biophys Acta* 1986;874:205–215.
9. Chou PY. Amino acid composition of four classes of proteins. In: Abstracts of papers, Part I, Second Chemical Congress of the North American Continent. Las Vegas, 1980.
10. Chou PY. Prediction of protein structural classes from amino acid composition. In: Fasman GD, editor. Prediction of protein structure and the principles of protein conformation. New York: Plenum Press; 1989. p 549–586.
11. Nakashima H, Nishikawa K, Ooi T. The folding type of a protein is relevant to the amino acid composition. *J Biochem* 1986;99:153–162.
12. Zhang CT, Chou KC. An optimization approach to predicting protein structural class from amino acid composition. *Protein Sci* 1992;1:401–408.
13. Chou KC, Zhang CT. A new approach to predicting protein folding types. *J Protein Chem* 1993;12:169–178.
14. Chou KC, Zhang CT. Predicting protein folding types by distance functions that make allowances for amino acid interactions. *J Biol Chem* 1994;269:22014–22020.
15. Chou KC. A novel approach to predicting protein structural classes in a (20–1)-D amino acid composition space. *Proteins* 1995;21:319–344.
16. Chou KC, Zhang CT. Prediction of protein structural classes. *Crit Rev Biochem Mol Biol* 1995;30:275–349.
17. Eisenhaber F, Frömmel C, Argos P. Prediction of secondary structural content of proteins from their amino acid composition along. II. The paradox with secondary structural class. *Proteins* 1996;25:169–179.
18. Bahar I, Atilgan AR, Jernigan RL, Erman B. Understanding the recognition of protein structural classes by amino acid composition. *Proteins* 1997;29:172–185.
19. Chou KC, Liu WM, Maggiora GM, Zhang CT. Prediction and classification of domain structural classes. *Proteins* 1998;31:97–103.
20. Chou KC, Maggiora GM. Domain structural class prediction. *Protein Eng* 1998;11:523–538.
21. Tou JT, Gonzalez RC. Pattern recognition principles. London: Addison-Wesley Publishing Company; 1974.
22. Duda RO, Hart PE. Pattern classification and scene analysis. New York: John Wiley & Sons; 1973.
23. Andrew HC. Introduction to mathematical techniques in pattern recognition. New York: John Wiley & Sons; 1972.
24. DeGroot MH. Probability and statistics. 2nd ed. Reading, MA: Addison-Wesley Publishing Co.; 1986. p 267–278.
25. Zhang CT, Chou KC. Monte Carlo simulation studies on the prediction of protein folding types from amino acid composition. *Biophys J* 1992;63:1523–1529.
26. Anderson TW. Introduction to multivariate statistics. New York: John Wiley & Sons; 1958.
27. Sander C, Schneider R. Database of homology-derived protein structures and structural meaning of sequence alignment. *Proteins* 1991;9:56–58.
28. Blundell TL, Johnson MS. Catching a common fold. *Protein Sci* 1993;2:877–883.
29. Flores TP, Orengo CA, Moss D, Thornton JM. Comparison of conformational characteristics in structurally similar protein pairs. *Protein Sci* 1993;2:1811–1826.
30. Hilbert M, Böhm G, Jaenicke R. Structural relationships of homologous proteins as a fundamental principle in homology modeling. *Proteins* 1993;17:138–151.
31. Brenner SE, Chothia C, Hubbard TJP. Assessing sequence comparison methods with reliable structurally identified distant evolutionary relationships. *Proc Natl Acad Sci USA* 1998;95:6073–6078.
32. Kabsch W, Sander C. Dictionary of protein secondary structures: pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers* 1983;22:2577–2637.
33. Metfessel BA, Saurugger PN, Connelly DP, Rich SS. Cross-validation of protein structural class prediction using statistical clustering and neural networks. *Protein Sci* 1993;2:1171–1182.
34. Chandonia JM, Karplus M. Neural networks for secondary structure and structural class predictions. *Protein Sci* 1995;4:275–285.
35. Rost B, Sander C. Prediction of protein secondary structure at better than 70% accuracy. *J Mol Biol* 1993;232:584–599.
36. Levin JM, Pascarella S, Argos P, Garnier J. Quantification of secondary structure prediction improvement using multiple alignments. *Protein Eng* 1993;6:849–854.
37. Boberg J, Salakoski T, Vihinen M. Accurate prediction of protein secondary structural class with fuzzy structural vectors. *Protein Eng* 1995;8:505–512.