

New techniques for extracting features from protein sequences

by J. T. L. Wang
Q. Ma
D. Shasha
C. H. Wu

In this paper we propose new techniques to extract features from protein sequences. We then use the features as inputs for a Bayesian neural network (BNN) and apply the BNN to classifying protein sequences obtained from the PIR (Protein Information Resource) database maintained at the National Biomedical Research Foundation. To evaluate the performance of the proposed approach, we compare it with other protein classifiers built based on sequence alignment and machine learning methods. Experimental results show the high precision of the proposed classifier and the complementarity of the bioinformatics tools studied in the paper.

As a result of the Human Genome Project and related efforts, DNA (deoxyribonucleic acid), RNA (ribonucleic acid), and protein data accumulate at an accelerating rate. Mining these biological data to extract useful knowledge is essential in genome processing. This subject has recently gained significant attention in the bioinformatics community.¹⁻⁶ We present here a case study in extracting features from protein sequences and using them together with a Bayesian neural network to classify the sequences.

The problem studied here can be stated formally as follows: Given are an unlabeled protein sequence S and a known superfamily \mathcal{F} ; we want to determine whether or not S belongs to \mathcal{F} . (We refer to \mathcal{F} as the *target class* and the set of sequences not in \mathcal{F} as the *nontarget class*.) In general, a superfamily is a group of proteins that share similarity in structure and function. If the unlabeled sequence S is determined to belong to \mathcal{F} , then one can infer the structure and function of S . This process is important in many aspects

of bioinformatics and computational biology.⁷⁻⁹ For example, in drug discovery, if sequence S is obtained from some disease X and it is determined that S belongs to the superfamily \mathcal{F} , then one may try a combination of the existing drugs for \mathcal{F} to treat the disease X .

There are several approaches available for protein sequence classification. One approach is to compare the unlabeled sequence S with the sequences in the target class and the sequences in the nontarget class using an alignment tool such as BLAST.¹⁰ One then assigns S to the class containing the sequence best matching S .

The second method for protein sequence classification is based on hidden Markov models (HMMs).¹¹ The HMM method (e.g., SAM¹² and HMMer¹³) employs a machine-learning algorithm, which uses probabilistic graphical models to describe time-series and sequence data. It was originally applied to speech recognition,¹⁴ and now is also applied to modeling and analyzing protein superfamilies. It is a generalization of the position-specific scoring matrix to include insertion and deletion states. Often, an HMM is built for each (super)family. One then scores the unlabeled sequence S with respect to the HMM of a (super)family.¹⁵ If the score is more significant than a cut-off value, then S is regarded as a member of the (super)family.

©Copyright 2001 by International Business Machines Corporation. Copying in printed form for private use is permitted without payment of royalty provided that (1) each reproduction is done without alteration and (2) the *Journal* reference and IBM copyright notice are included on the first page. The title and abstract, but no other portions, of this paper may be copied or distributed royalty free without further permission by computer-based and other information-service systems. Permission to *republish* any other portion of this paper must be obtained from the Editor.

Another approach for protein sequence classification is to iteratively build a model either based on hidden Markov models (e.g., SAM-T99¹⁶) or based on a position-specific weight matrix (e.g., PSI-BLAST¹⁰). The unlabeled sequence S is used as a seed sequence and iteratively searched against the (super)family either by the HMM or by the position-specific weight matrix.

In the study presented here, we compare our approach with BLAST, SAM, and the iterative method using SAM-T99. With hidden Markov models, we choose SAM rather than HMMer because the former outperforms the latter in protein sequence classification.¹⁷ With iterative methods, we choose SAM-T99 rather than PSI-BLAST because the former is more sensitive than the latter in homolog detection.⁷ We choose BLAST as a point of comparison because it represents a different computing paradigm, namely performing classification simply via alignment. One interesting finding from our work is that the compared classification methods complement each other; combining them yields higher precision than using them individually, as our experimental results will show later. This is consistent with a previous report¹⁸ in which we gave a preliminary analysis of the complementarity among our approach, BLAST, and SAM.

Feature extraction from protein data

From a one-dimensional point of view, a protein sequence contains characters from the 20-letter amino acid alphabet $\mathcal{A} = \{A, C, D, E, F, G, H, I, K, L, M, N, P, Q, R, S, T, V, W, Y\}$. An important issue in applying neural networks to protein sequence classification is how to encode protein sequences, i.e., how to represent the protein sequences as the input of the neural networks. Indeed, sequences may not be the best representation at all. Good input representations make it easier for the neural networks to recognize underlying regularities. Thus, good input representations are crucial to the success of neural network learning.¹⁹

We propose here new encoding techniques that entail the extraction of high-level features from protein sequences. The best high-level features should be “relevant.” By relevant we mean that there should be high mutual information between the features and the output of the neural networks, where the mutual information measures the average reduction in uncertainty about the output of the neural networks given the values of the features.

Another way to look at these features is that they capture both the global similarity and the local similarity of protein sequences. The global similarity refers to the overall similarity among multiple sequences, whereas the local similarity refers to motifs (or frequently occurring substrings) in the sequences. The next two sections elaborate on how to find the global and local similarity of the protein sequences. The third section presents our classification algorithm, which employs the Bayesian neural network originated from Mackay.²⁰ The last three sections evaluate the performance of the proposed classifier, compare our approach with the other protein classifiers, and conclude the paper.

Global similarity of protein sequences

To calculate the global similarity of protein sequences, we adopt the 2-gram, also known as 2-tuple, method as described in Wu.²¹ The 2-gram encoding method extracts various patterns of two consecutive amino acid residues in a protein sequence and counts the number of occurrences of the extracted residue pairs.²² For instance, given a protein sequence PVKTNVK, the 2-gram amino acid encoding method gives the following result: 1 for PV (indicating PV occurs once), 2 for VK (indicating VK occurs twice), 1 for KT, 1 for TN, and 1 for NV.

We also adopt the 6-letter exchange group $\{e_1, e_2, e_3, e_4, e_5, e_6\}$ to represent a protein sequence,²³ where $e_1 \in \{H, R, K\}$, $e_2 \in \{D, E, N, Q\}$, $e_3 \in \{C\}$, $e_4 \in \{S, T, P, A, G\}$, $e_5 \in \{M, I, L, V\}$, and $e_6 \in \{F, Y, W\}$. Exchange groups represent conservative replacements through evolution. These exchange groups are effectively equivalence classes of amino acids and are derived from PAM.^{24,25} For example, the above protein sequence PVKTNVK can be represented as $e_4e_5e_1e_4e_2e_5e_1$. The 2-gram exchange group encoding for this sequence is: 1 for e_4e_5 , 2 for e_5e_1 , 1 for e_1e_4 , 1 for e_4e_2 , and 1 for e_2e_5 .

For each protein sequence, we apply both the 2-gram amino acid encoding and the 2-gram exchange group encoding to the sequence. Thus, there are $20^2 + 6^2 = 436$ possible 2-grams in total. If all the 436 2-grams are chosen as the neural network input features, it would require many weight parameters and training data. This makes it difficult to train the neural network—a phenomenon called “curse of dimensionality.” Different methods have been proposed to solve the problem by careful feature selection and by scaling of the input dimensionality.^{23,26} We propose here to select relevant features (i.e., 2-grams)

by employing a distance measure to calculate the relevance of each feature.²⁷

Let X be a feature and let x be its value. Let $P(x|Class = 1)$ and $P(x|Class = 0)$ denote the class conditional density functions for feature X , where $Class_1$ represents the target class and $Class_0$ is the nontarget class. Let $D(X)$ denote the distance function between $P(x|Class = 1)$ and $P(x|Class = 0)$, defined as²⁸

$$D(X) = \int |P(x|Class = 1) - P(x|Class = 0)| dx \quad (1)$$

The distance measure prefers feature X to feature Y if $D(X) > D(Y)$. Intuitively, this means it is easier to distinguish between $Class_1$ and $Class_0$ by observing feature X than feature Y . That is, X appears often in $Class_1$ and seldom in $Class_0$ or vice versa. In our work, each feature X is a 2-gram. Let c denote the occurrence number of the feature X in a sequence S . Let l denote the total number of 2-grams in S and let $len(S)$ represent the length of S . We have $l = len(S) - 1$. Define the feature value x for the 2-gram X with respect to the sequence S as

$$x = \frac{c}{len(S) - 1} \quad (2)$$

For example, suppose $S = PVKTNVK$. Then the value of the feature VK with respect to S is $2/(7 - 1) = 0.33$.

Because a protein sequence may be short, random pairings can have a large effect on the result. We approximate $D(X)$ in Equation 1 by²⁹

$$D(X) = \frac{(m_1 - m_0)^2}{d_1^2 + d_0^2} \quad (3)$$

where m_1 and d_1 (m_0 and d_0 , respectively) are the mean value and the standard deviation of the feature X in the positive (negative, respectively) training data set. Intuitively, in Equation 3, the larger the numerator is (or the smaller the denominator is), the larger the interclass distance is, and therefore the easier to separate $Class_1$ from $Class_0$ (and vice versa).

The mean value m and the standard deviation d of the feature X in a set \mathcal{S} of sequences are defined as

$$m = \frac{1}{N} \sum_{i=1}^N x_i \quad (4)$$

$$d = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (x_i - m)^2} \quad (5)$$

where x_i is the value of the feature X with respect to sequence $S_i \in \mathcal{S}$, and N is the total number of sequences in \mathcal{S} .

Let $X_1, X_2, \dots, X_{N_g}, N_g \ll 436$, be the top N_g features (2-grams) with the largest $D(X)$ values.³⁰ Intuitively, these N_g features occur more frequently in the positive training data set and less frequently in the negative training data set. For each protein sequence S (whether it is a training or an unlabeled test sequence), we examine the N_g features in S , calculate their values as defined in Equation 2, and use the N_g feature values as input feature values to the Bayesian neural network for the sequence S .

To compensate for the possible loss of information due to ignoring the other 2-grams, a linear correlation coefficient (LCC) between the values of the 436 2-grams with respect to the protein sequence S and the mean value of the 436 2-grams in the positive training data set is calculated and used as another input feature value for S . Specifically, the LCC of S is defined as:

$$LCC(S) = \frac{436 \sum_{j=1}^{436} x_j \bar{x}_j - \sum_{j=1}^{436} x_j \sum_{j=1}^{436} \bar{x}_j}{\sqrt{436 \sum_{j=1}^{436} x_j^2 - \left(\sum_{j=1}^{436} x_j \right)^2} \sqrt{436 \sum_{j=1}^{436} \bar{x}_j^2 - \left(\sum_{j=1}^{436} \bar{x}_j \right)^2}} \quad (6)$$

where \bar{x}_j is the mean value of the j th 2-gram, $1 \leq j \leq 436$, in the positive training data set, and x_j is the feature value of the j th 2-gram with respect to S as defined in Equation 2.

Local similarity of protein sequences

In contrast to the 2-grams that occur from the beginning to the end of a sequence (thus referred to as global similarities), the local similarity of protein sequences refers to frequently occurring motifs, where a motif is composed of substrings occurring

in local regions of a sequence. Let $\mathcal{S}_p = \{S_1, \dots, S_k\}$ be the positive training data set. We use a previously developed sequence mining tool Sdiscover^{31,32} to find the regular expression motifs of the forms $*X*$ and $*X*Y*$ where each motif has length $\geq Len$ and approximately matches, within Mut mutations, at least $Occur$ sequences in \mathcal{S}_p . Here, a mutation could be a mismatch, an insertion, or a deletion of a letter (residue); Len , Mut , and $Occur$ are user-specified parameters. X and Y are segments of a sequence, i.e., substrings made up of consecutive letters, and $*$ is a variable length “don’t care” (VLDC) symbol. The length of a motif is the number of the non-VLDC letters in the motif. When matching a motif with a sequence S_i , a VLDC symbol in the motif is instantiated into an arbitrary number of residues in S_i at no cost. For example, when matching a motif $*VLHGKKVL*$ with a sequence $MNVLAHGKKV-LKWK$, the first $*$ is instantiated into MN and the second $*$ is instantiated into KWK . The number of mutations between the motif and the sequence is 1, representing the cost of inserting an A in the motif.

The Sdiscover tool is based on a heuristic that works by taking a small sample \mathcal{K} of sequences from the given set of sequences \mathcal{S}_p and storing them in a generalized suffix tree (GST).³³ The GST can be constructed asymptotically in $O(n)$ time and space³⁴ where n is the total length of all sequences in the sample. The heuristic then traverses the GST to generate candidate regular expression motifs and compares these candidate motifs with all the sequences in \mathcal{S}_p to calculate their occurrence numbers. Given a candidate motif M and a sequence S in \mathcal{S}_p , one can determine whether M is within Mut mutations of S in $O(Mut \times |S|)$ time when $O(|M|) = O(\log |S|)$.³⁵ Thus Sdiscover can find all the motifs satisfying user-specified parameter values in time $O(n \times Mut \times m \times k)$, where n is the total length of all sequences in the sample \mathcal{K} , m is the average length of the sequences in \mathcal{S}_p , and k is the total number of sequences in \mathcal{S}_p , although the tool is practically much faster due to several optimization heuristics implemented for speeding up the traversal of the GST.

Often, the number of motifs returned by Sdiscover is enormous. It is useful to develop a measure to evaluate the significance of these motifs. We propose here to use the minimum description length (MDL) principle^{8,36,37} to calculate the significance of a motif. The MDL principle states that the best model (a motif in our case) is the one that minimizes the sum of the length, in bits, of the description of the model and the length, in bits, of the description of the data

(the positive training sequences in \mathcal{S}_p in our case) encoded by the model.

Evaluating the significance of motifs. We adopt information theory in its fundamental form^{36,38} to measure the significance of different motifs. The theory takes into account the probability of an amino acid in a motif (or sequence) when calculating the description length of the motif (or sequence). Specifically, Shannon³⁸ showed that the length in bits to transmit a symbol b via a channel in some optimal coding is $-\log_2 P_x(b)$, where $P_x(b)$ is the probability with which the symbol b occurs. Given the probability distribution P_x over an alphabet $\Sigma_x = \{b_1, b_2, \dots, b_n\}$, we can calculate the description length of any string $b_{k_1}b_{k_2} \dots b_{k_l}$ over the alphabet Σ_x by

$$-\sum_{i=1}^l \log_2 P_x(b_{k_i}) \quad (7)$$

In our case, the alphabet Σ_x is the protein alphabet \mathcal{A} containing 20 amino acids. The probability distribution P_x , or P in our case, can be calculated by examining the occurrence frequencies of amino acids in the positive training data set \mathcal{S}_p . One straightforward way to describe (or encode) the sequences in \mathcal{S}_p , referred to as Scheme 1, is to encode sequence by sequence, separated by a delimiter \$. Let $dlen(S_i)$ denote the description length of sequence $S_i \in \mathcal{S}_p$. Then

$$dlen(S_i) = -\sum_{j=1}^{20} n_{a_j} \log_2 P(a_j) \quad (8)$$

where $a_j \in \mathcal{A}$, $1 \leq j \leq 20$; n_{a_j} is the number of occurrences of a_j in S_i . For example, suppose $S_i = MNVLAHGKKVLKWK$ is a sequence in \mathcal{S}_p . Then

$$\begin{aligned} dlen(S_i) = & -(\log_2 P(M) + \log_2 P(N) + 2 \log_2 P(V) \\ & + 2 \log_2 P(L) + \log_2 P(A) + \log_2 P(H) \\ & + \log_2 P(G) + 4 \log_2 P(K) \\ & + \log_2 P(W)) \end{aligned} \quad (9)$$

Let $dlen(\mathcal{S}_p)$ denote the description length of $\mathcal{S}_p = \{S_1, \dots, S_k\}$. Then the description length of \mathcal{S}_p is given by

$$dlen(\mathcal{S}_p) = \sum_{i=1}^k dlen(S_i) + (k-1) \times dlen(\$) \quad (10)$$

Since the description length of the delimiter \$, $dlen(\$)$, is insignificant, we can ignore it and hence

$$dlen(\mathcal{S}_p) = \sum_{i=1}^k dlen(S_i) \quad (11)$$

Another method to encode the sequences in \mathcal{S}_p , referred to as Scheme 2, is to encode a regular expression motif, say M_j , and then encode the sequences in \mathcal{S}_p based on M_j . Specifically, if a sequence $S_i \in \mathcal{S}_p$ can approximately match M_j , then we encode S_i based on M_j . Otherwise we encode S_i using Scheme 1.³⁹ Let us use an example to illustrate Scheme 2. Consider, for example, $M_j = *VLHGKKVL*$. We encode M_j as 1, *, V, L, H, G, K, K, V, L, *, \$0 where 1 indicates one mutation is allowed in matching M_j with S_i and \$0 is a delimiter to signal the end of the motif. Let Σ_1 denote the alphabet $\{a_1, a_2, \dots, a_{20}, *, \$0\}$, where a_1, a_2, \dots, a_{20} are the 20 amino acids. Let P_1 denote the probability distribution over the alphabet Σ_1 . $P_1(\$0)$ can be approximated by the reciprocal of the average length of motifs. $P_1(*) = n(P_1(\$0))$, $P_1(a_i) = (1 - (n+1)P_1(\$0))P(a_i)$, where n denotes the number of VLDCs in the motif M_j . For a motif of the form $*X*$, n is 2; for a motif of the form $*X*Y*$, n is 3.

Given P_1 , we can calculate the description length of a motif by substituting the probability distribution P_1 for the probability distribution P_x in Equation 7. Specifically, let $M_j = *a_{j_1}a_{j_2}\dots a_{j_k}$. Let $dlen(M_j)$ denote the description length, in bits, of the motif M_j . Then

$$dlen(M_j) = - \left(2 \log_2 P_1(*) + \log_2 P_1(\$0) + \sum_{i=1}^k \log_2 P_1(a_{j_i}) \right) \quad (12)$$

For instance, consider again $M_j = *VLHGKKVL*$. We have

$$\begin{aligned} dlen(M_j) = & -(2 \log_2 P_1(*) + \log_2 P_1(\$0) \\ & + 2 \log_2 P_1(V) + 2 \log_2 P_1(L) \\ & + \log_2 P_1(H) + \log_2 P_1(G) \\ & + 2 \log_2 P_1(K)) \end{aligned} \quad (13)$$

Sequences that are approximately matched by the motif M_j can be encoded with the aid of the motif. For example, consider again $M_j = *VLHGKKVL*$ and $S_i = MNVLAHGKKVLKWK$. M_j matches S_i with one mutation, representing the cost of inserting an A in the third position of M_j . The first VLDC symbol is instantiated into MN and the second VLDC symbol is instantiated into KWK. We can thus rewrite S_i as MN • SS_i • KWK where SS_i is VLAHGKKVL and • denotes the concatenation of strings. Therefore we can encode S_i as M, N, \$1; 1, (O_I , 3, A); K, W, K, \$1. Here \$1 is a delimiter, 1 indicates that one mutation occurs when matching M_j with S_i , and (O_I , 3, A) indicates that the mutation is an insertion that adds the letter A to the third position of M_j . In general, the mutation operations involved and their positions can be observed using approximate string matching algorithms.³⁵ The description length of the encoded S_i based on M_j , denoted $dlen(S_i, M_j)$, can be calculated easily as in Equation 12.

Suppose there are h sequences $S_{p_1} \dots S_{p_h}$ in the positive training data set \mathcal{S}_p that can approximately match the motif M_j . The weight (or significance) of M_j , denoted $w(M_j)$, is defined as

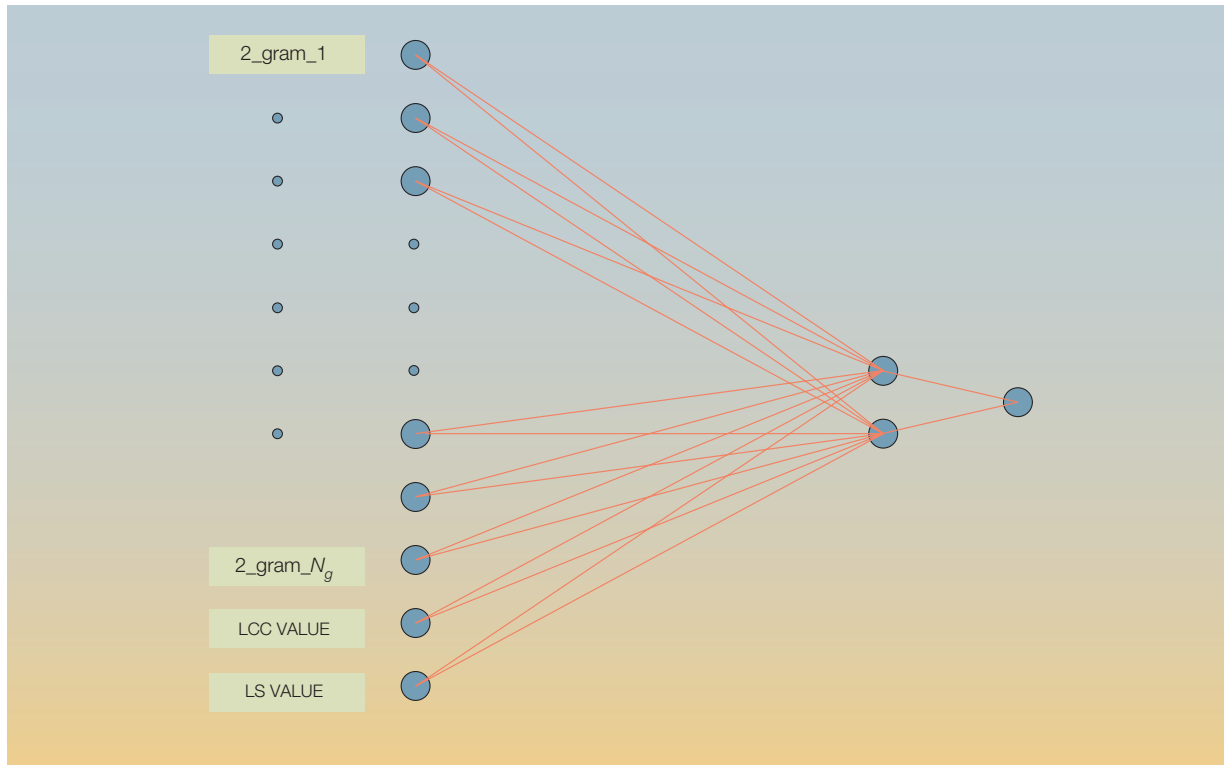
$$w(M_j) = \sum_{i=1}^h dlen(S_{p_i}) - \left(dlen(M_j) + \sum_{i=1}^h dlen(S_{p_i}, M_j) \right) \quad (14)$$

Intuitively, the more sequences in \mathcal{S}_p approximately matching M_j and the fewer bits we use to encode M_j and to encode those sequences based on M_j , the larger weight M_j has.

Using Sdiscover, one can find a set \mathcal{S} of regular expression motifs of the forms $*X*$ and $*X*Y*$ from the positive training data set \mathcal{S}_p where the motifs satisfy the user-specified parameter values Len , Mut , and $Occur$. We choose the top N_j motifs with the largest weights. Let \mathcal{R} denote this set of motifs. Suppose a protein sequence S (whether it is a training sequence or an unlabeled test sequence) can approximately match, within Mut mutations, m motifs in \mathcal{R} . Let these motifs be M_1, \dots, M_m . The local similarity (LS) value of S , denoted $LS(S)$, is defined as

$$LS(S) = \begin{cases} \max_{1 \leq i \leq m} \{w(M_i)\} & \text{if } m \neq 0 \\ 0 & \text{otherwise} \end{cases} \quad (15)$$

Figure 1 The Bayesian neural network architecture



This LS value is used as an input feature value of the Bayesian neural network for the sequence S . Note that we use the “max” operator here to maximize discrimination. In general, positive sequences will have large LS values with high probabilities and small LS values with low probabilities. On the other hand, negative sequences will have small LS values with high probabilities and large LS values with low probabilities.

Remark. Essentially, the proposed scheme is to count amino acids in a sequence (or motif). This scheme is not complete in the sense that different sequences may have the same description length when they have the same number of the same amino acids. Also, there may be multiple ways to align a motif M with a sequence S and hence the description length of the encoded sequence S based on M may not be unique. As a consequence, the weight of a motif defined in Equation 14 may not be unique (in which case the proposed heuristic randomly picks one). There are several other approaches for finding mo-

tifs of different forms and for calculating their significance values (see, e.g., Wang,⁸ Brazma,³⁶ Califano,⁴⁰ and Hart⁴¹). However, motifs have relatively little effect on PIR sequence classification and a combination of the proposed techniques already yields a very high precision, as our experimental results show in a later section.

The Bayesian neural network classifier

We adopt the Bayesian neural network (BNN) originated from Mackay²⁰ to classify protein sequences.⁴² There are $N_g + 2$ input features, including N_g 2-grams, the LCC feature, and the LS feature, both previously described. Thus, a protein sequence is represented as a vector of $N_g + 2$ real numbers. The BNN has one hidden layer containing multiple hidden units. The output layer has one output unit, which is based on the logistic activation function $f(a) = 1/(1 + e^{-a})$. The BNN is fully connected between the adjacent layers. Figure 1 illustrates an example BNN model with two hidden units.

Table 1 Data used in the experiments. N is the number of sequences, L_m is the minimal length of the sequences, and L_x is the maximal length of the sequences.

Data Set	N	L_m	L_x
Globin	831	115	173
Kinase-related transforming protein	350	151	502
Ras transforming protein	386	106	322
Ribitol dehydrogenase	319	129	335
Negative sequences	1650	100	200

Table 2 Parameters and their base values for the proposed BNN classifier

Parameter	Meaning	Value
N_g	Number of 2-grams used by BNN	60
N_l	Number of motifs used by BNN	20
Len	Length of motifs for Sdiscover	6
Mut	Mutation number for Sdiscover	2
$Occur$	Occurrence frequency of motifs for Sdiscover	1/10
r	Size ratio	2

Let $\mathcal{D} = \{\mathbf{x}^{(m)}, t_m\}$, $1 \leq m \leq N$, denote the training data set including both positive and negative training sequences. $\mathbf{x}^{(m)}$ is an input feature vector including the $N_g + 2$ input feature values, and t_m is the binary (0/1) target value for the output unit. That is, t_m equals 1 if $\mathbf{x}^{(m)}$ represents a protein sequence in the target class, and 0 otherwise.

Let \mathbf{x} denote the input feature vector for a protein sequence, which could be a training sequence or a test sequence. Given the architecture \mathbf{A} and the weights \mathbf{w} of the BNN, the output value y can be uniquely determined from the input vector \mathbf{x} . Because of the logistic activation function $f(a)$ of the output unit, the output value $y(\mathbf{x}; \mathbf{w}, \mathbf{A})$ can be interpreted as $P(t = 1 | \mathbf{x}, \mathbf{w}, \mathbf{A})$, i.e., the probability that \mathbf{x} represents a protein sequence in the target class given \mathbf{w}, \mathbf{A} . The likelihood function of the data \mathcal{D} given the model is calculated by

$$P(\mathcal{D} | \mathbf{w}, \mathbf{A}) = \prod_{m=1}^N y^{t_m} (1-y)^{1-t_m} = \exp(-G(\mathcal{D} | \mathbf{w}, \mathbf{A})) \quad (16)$$

where $G(\mathcal{D} | \mathbf{w}, \mathbf{A})$ is the cross-entropy error function,

$$G(\mathcal{D} | \mathbf{w}, \mathbf{A}) = - \sum_{m=1}^N t_m \log(y) + (1-t_m) \log(1-y) \quad (17)$$

The $G(\mathcal{D} | \mathbf{w}, \mathbf{A})$ is the objective function in a non-Bayesian neural network training process and is minimized. This process assumes all possible weights are equally likely. The weight decay is often used to avoid overfitting on the training data and poor generalization on the test data by adding a term $\alpha/2 \sum_{i=1}^q w_i^2$ to the objective function, where α is the weight decay parameter (hyperparameter), $\sum_{i=1}^q w_i^2$ is the sum of the squares of all the weights of the neural network, and q is the number of weights. This objective function is minimized to penalize the neural network with weights of large magnitudes. Thus, it penalizes an over-complex model and favors a simple model. However, there is no precise way to specify the appropriate value of α , which is often tuned off line.

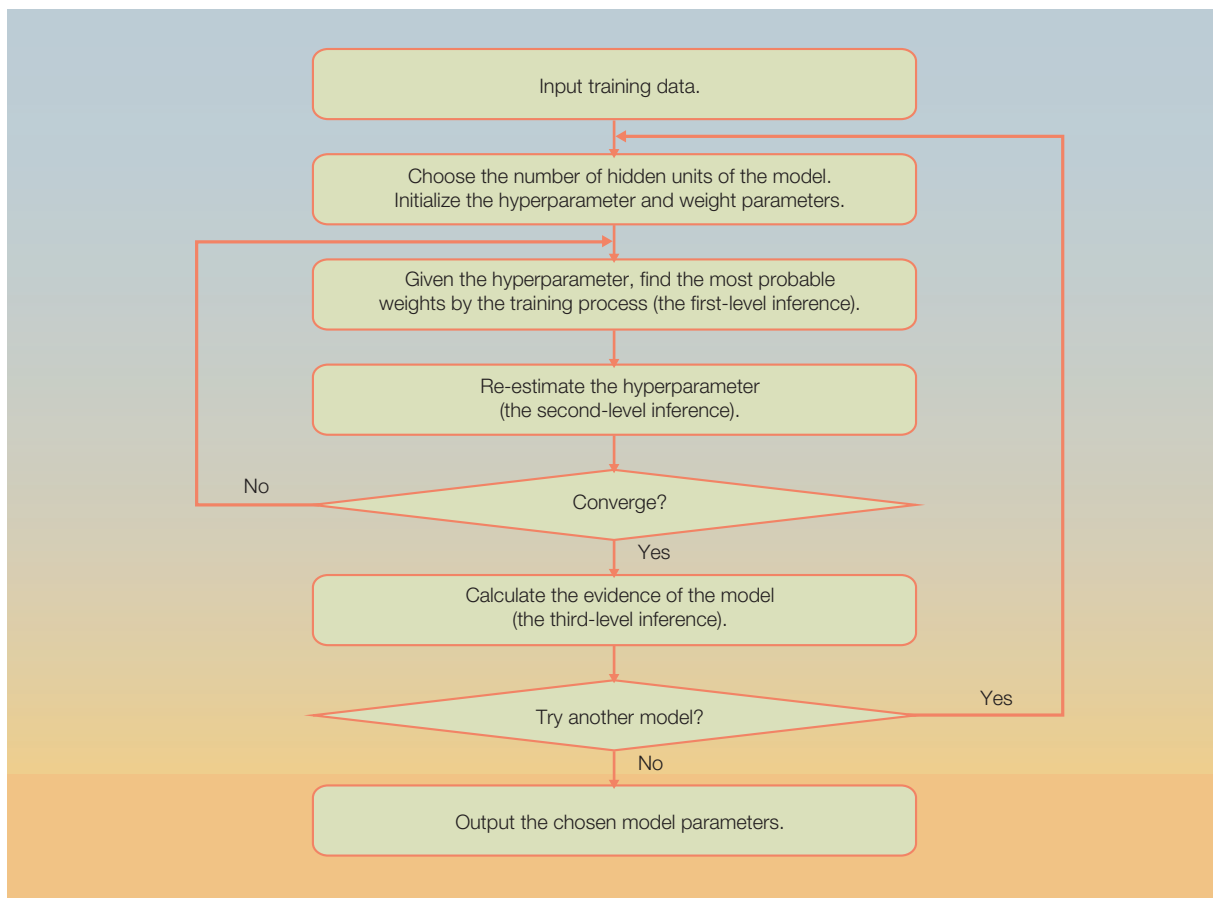
In contrast, in the Bayesian neural network, the hyperparameter α is interpreted as the parameter of a model and is optimized on line during the Bayesian learning process. We adopt the Bayesian training of neural networks described in Mackay²⁰ to calculate and maximize the evidence of α , namely $P(\mathcal{D} | \alpha, \mathbf{A})$. The training process employs an iterative procedure; each iteration involves three levels of inference. Figure 2 illustrates the training process of the BNN.

In classifying an unlabeled test sequence S represented by its input feature vector \mathbf{x} , the output of the BNN, $P(t = 1 | \mathbf{x}, \mathbf{w}, \mathbf{A})$, is the probability that S belongs to the target class. If the probability is greater than the decision boundary 0.5, S is assigned to the target class; otherwise S is assigned to the nontarget class. In general, for an unlabeled test sequence S with m amino acids, it takes $O(m)$ time to calculate 2-gram feature values, $O(m)$ time to calculate the LCC feature value, and $O(m \times n)$ time to calculate the LS feature value where n is the total length of the motifs chosen, and constant time for calculating the probability $P(t = 1 | \mathbf{x}, \mathbf{w}, \mathbf{A})$. Thus, the time complexity of our approach to classifying the unlabeled test sequence S is $O(m \times n)$.

Performance of the BNN classifier

We carried out a series of experiments to evaluate the performance of the proposed BNN classifier on a Pentium** II PC running the LINUX** operating

Figure 2 The training process of the Bayesian neural network



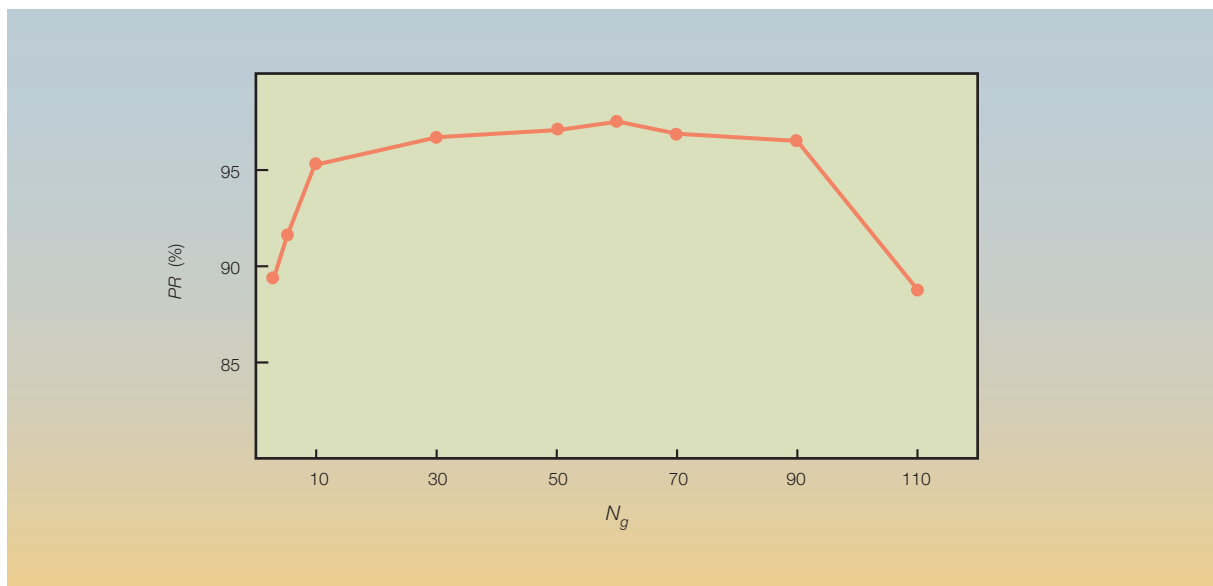
system. The data used in the experiments were obtained from the International Protein Sequence Database,⁴³ release 62, in the Protein Information Resource (PIR) maintained by the National Biomedical Research Foundation (NBRF-PIR) at the Georgetown University Medical Center. This database, accessible at <http://pir.georgetown.edu>, currently has 172 684 sequences. Table 1 summarizes the data used in the experiments.

Data. Four positive data sets were considered; they were globin, kinase, ras, and ribitol superfamilies, respectively, in the PIR protein database. The negative data set contained 1650 protein sequences, also taken from the PIR protein database, with lengths ranging from 100 residues to 200 residues; these negative sequences did not belong to any of the four pos-

itive superfamilies. Both the positive and negative sequences were randomly divided into training sequences and test sequences, where the size of the training data set equaled the size of the test data set multiplied by an integer r . With the same training data, we tested several BNN models with different numbers of hidden units. When there were two hidden units, the evidence obtained was the largest (see Figure 2), so we fixed the number of hidden units at two. Models with more hidden units would require more training time while achieving roughly the same performance.

Table 2 summarizes the parameters and base values used in the experiments. The measure used to evaluate the performance of the BNN classifier is *precision*, PR , which is defined as

Figure 3 Impact of N_g in the BNN classifier



$$PR = \frac{NumCorrect}{NumTotal} \times 100 \text{ percent} \quad (18)$$

where *NumCorrect* is the number of test sequences classified correctly and *NumTotal* is the total number of test sequences. In general, precision is a comprehensive measure in the sense that it considers true positives, false positives, true negatives, false negatives, and unclassified sequences;⁴⁴ it is used here to find the best parameter values of the proposed BNN classifier. A false positive is a nontarget member sequence that was misclassified as a target member sequence. A false negative refers to a sequence in the target class (e.g., the globin superfamily) that was misclassified as a nontarget member. We present the results for the globin superfamily only; the results for the other three superfamilies were similar.

Results. In the first experiment, we considered only 2-grams and evaluated their effect on the performance of the proposed BNN classifier. Figure 3 graphs *PR* as a function of N_g . It can be seen that the performance improves initially as N_g increases. The reason is that the more 2-grams we use, the more precisely we represent the protein sequences. However, when N_g is too large (e.g., >90), the training data are insufficient and the performance degrades. In general, the larger N_g is, the more input features the

BNN classifier has, and thus the larger training data set BNN requires. In our case, there are 561 positive training sequences and 1089 negative training sequences. When $N_g > 90$, these data become too few to yield reasonably good performance. Figuring out how big the parameter N_g should be requires some tuning. We have not yet worked out a theory for it.

In the second experiment, we considered only motifs found by Sdiscover and studied their effect on the performance of the BNN classifier. In this experiment 1597 motifs were found, with lengths ranging from six residues to 34 residues. Figure 4 graphs *PR* as a function of N_l . It can be seen that the more motifs one uses, the better performance one achieves. However, that would also require more time in matching a test sequence with the motifs.⁴⁵ We experimented with other parameter values for *Len*, *Mut*, and *Occur* used in Sdiscover. The results did not change as these parameters changed.

Figure 5 compares the effects of the various types of features introduced in the paper. To isolate the effects of these features, we began by using features of only one type and then using their combinations. It can be seen that features generated from global similarities yield better results than those generated

Figure 4 Effect of N_i in the BNN classifier

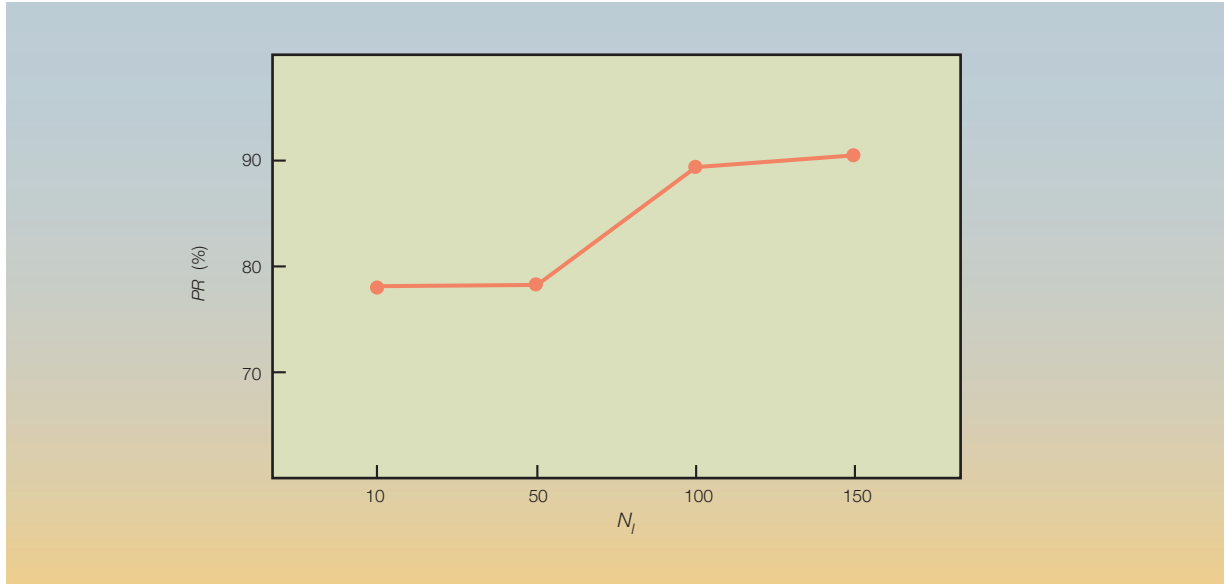
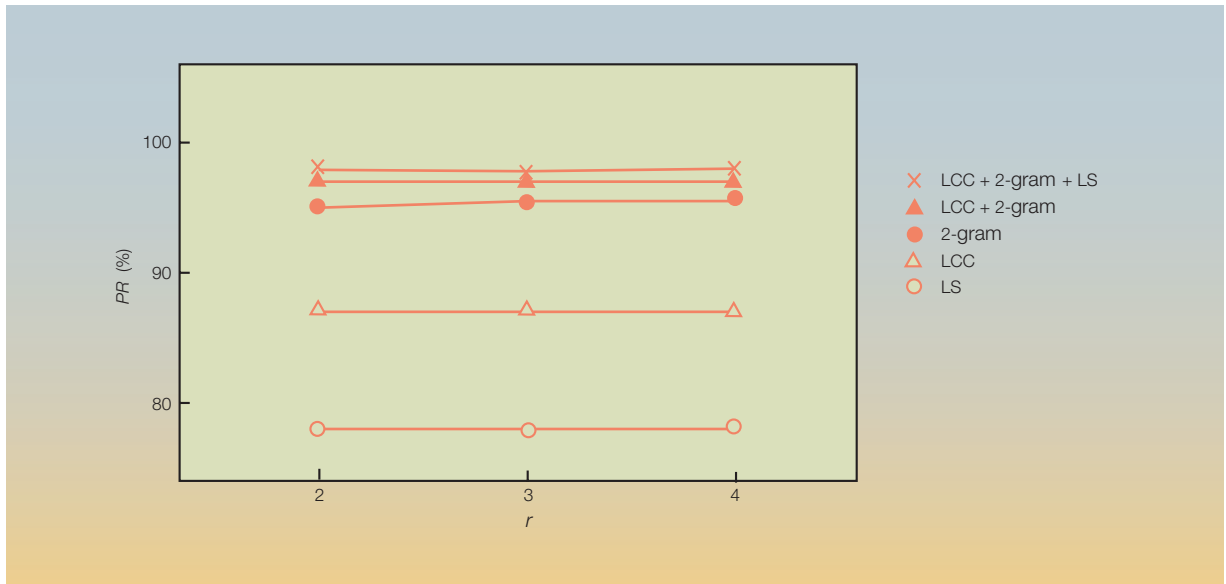


Figure 5 Comparison of the various types of features used by the BNN classifier



from local similarities. This happens because PIR superfamilies are categorized based on the global similarities of sequences. Note also that the best performance is achieved when all the features are used.

Comparison of four protein classifiers

The purpose of this section is to compare the proposed BNN classifier with the BLAST classifier¹⁰ built

Table 3 The bioinformatics tools studied in the paper

Tool	Underlying Techniques
BNN	Bayesian neural networks
BLAST	Similarity search and pairwise alignment
SAM	Hidden Markov models
SAM-T99	Iterative hidden Markov models

based on sequence alignment, and with the SAM and SAM-T99 classifiers¹⁶ built based on hidden Markov models. The parameter values for the BNN classifier were as shown in Table 2. Table 3 summarizes the studied tools. The BNN classifier used both 2-grams and regular expression motifs.

The BLAST version number was 2.0.10. We used default values for the parameters in BLAST. For this tool, we aligned an unlabeled test sequence S with the positive training sequences (i.e., those in the target class, e.g., the globin superfamily) and the negative training sequences in the nontarget class shown in Table 1 using the tool. If the score of S was below the threshold of the expectation (E) value of BLAST, S was undetermined or unclassified. Otherwise, we assigned S to the class containing the sequence best matching S .

The SAM version number was 3.2.1. For this tool, we employed the program buildmodel to build the HMM model by using only the positive training sequences. We then calculated the log-odds scores¹³ for all the training sequences using the program hmmscore.⁴⁶ The log-odds scores were all negative real numbers; the scores (e.g., -100.3) for the positive training sequences were generally smaller than the scores (e.g., -4.5) for the negative training sequences. The largest score S_p for the positive training sequences and the smallest score S_n for the negative training sequences were recorded. Let $B_{high} = \max \{S_p, S_n\}$ and $B_{low} = \min \{S_p, S_n\}$. We then calculated the log-odds scores for all the unlabeled test sequences using the program hmmscore. If the score of an unlabeled test sequence S was smaller than B_{low} , S was classified as a member of the target class, e.g., a globin sequence. If the score of S was larger than B_{high} , S was classified as a member of the nontarget class. If the score of S was between B_{low} and B_{high} , S was unclassified or undetermined.

The SAM-T99 version number was also 3.2.1. For this tool, we built an HMM (target model) for each unlabeled test sequence S . We then scored all the training sequences using the HMM target model. If the lowest score of the training sequences was higher

than the expectation value (E-value) of the HMM target model, the test sequence S was undetermined or unclassified.⁴⁷ Otherwise, we assigned the test sequence to the class containing the training sequence having the lowest E-value. The target model was built in four iterations. In the first iteration, SAM-T99 used BLAST to compare the test sequence S with sequences in the nonredundant protein database maintained at the National Center for Biotechnology Information (NCBI) and chose a set of close homologs to build an initial HMM. It also did a BLAST search of the test sequence S against the nonredundant protein database with a very loose cut-off value to get a pool of potential homologs. Then the HMM obtained from the previous iteration was compared against the pool of potential homologs with a looser cut-off value than that of the previous iteration to find weaker homologs. These weaker homologs were included to build a new HMM for the next iteration. This whole process was repeated three times.

In comparing the relative performance of these tools, we use four more measures in addition to the precision PR defined in the previous section: *specificity*, *sensitivity*, *unclassified_p*, and *unclassified_n*, where

$$\text{specificity} = \left(1 - \frac{N_{fp}}{N_{po}}\right) \times 100 \text{ percent} \quad (19)$$

$$\text{sensitivity} = \left(1 - \frac{N_{fn}}{N_{po}}\right) \times 100 \text{ percent} \quad (20)$$

$$\text{unclassified}_p = \frac{N_{up}}{N_{po}} \times 100 \text{ percent} \quad (21)$$

$$\text{unclassified}_n = \frac{N_{un}}{N_{ng}} \times 100 \text{ percent} \quad (22)$$

N_{fp} is the number of false positives, N_{fn} is the number of false negatives, N_{up} is the number of undetermined positive test sequences, N_{un} is the number of undetermined negative test sequences, N_{ng} is the total number of negative test sequences, and N_{po} is the total number of positive test sequences. Note that in contrast to PR , *specificity* and *sensitivity* do not consider unclassified sequences. That is why we also add the *unclassified_p* and *unclassified_n* measures for performance evaluation.

In the first experiment, we studied the effect of the threshold of the E-value in the BLAST classifier. Figure 6 shows the impact of E-values on the performance of BLAST. It can be seen that with $E = 10$, BLAST performs well. With smaller E-values (e.g.,

Figure 6 Impact of E-values for BLAST

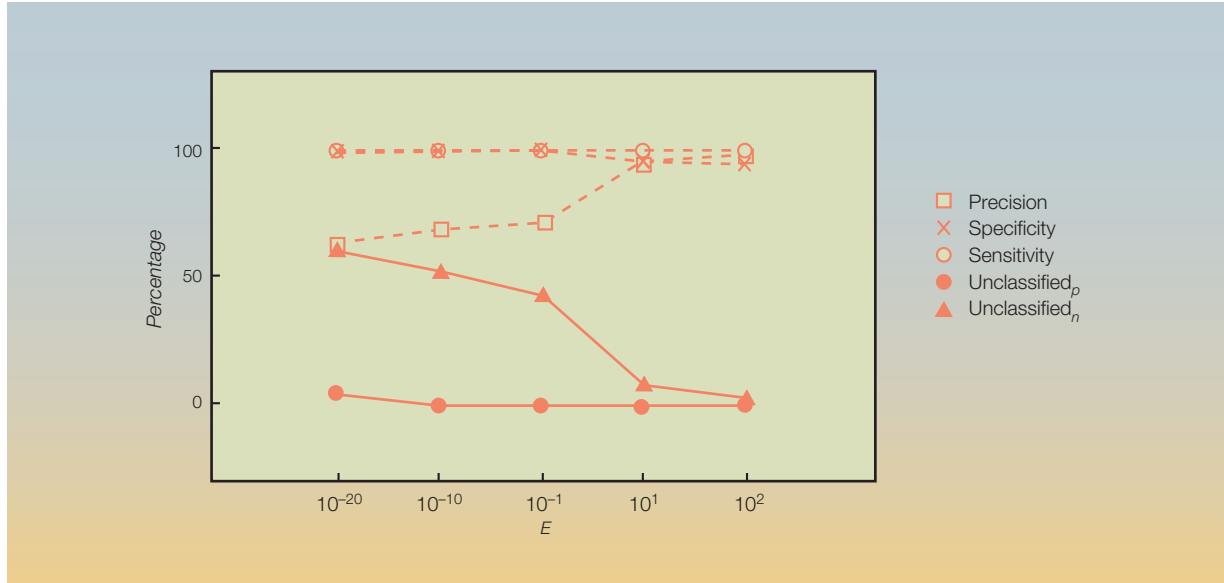


Table 4 Comparison of the studied classifiers on the globin superfamily

	BNN	BLAST	SAM	SAM-T99	COMBINER
Precision	98.0%	92.7%	95.3%	93.2%	99.8%
Specificity	98.0%	95.7%	99.8%	94.9%	99.6%
Sensitivity	98.0%	100.0%	99.6%	100.0%	100.0%
Unclassified _p	0.0%	0.0%	1.1%	0.0%	0.0%
Unclassified _n	0.0%	6.7%	6.2%	5.1%	0.0%
CPU time	36	1515	80	848 961	—

Table 5 Comparison of the studied classifiers on the kinase superfamily

	BNN	BLAST	SAM	SAM-T99	COMBINER
Precision	99.0%	86.2%	99.4%	92.6%	99.6%
Specificity	98.8%	87.8%	99.5%	93.1%	99.5%
Sensitivity	100.0%	100.0%	100.0%	100.0%	100.0%
Unclassified _p	0.0%	0.0%	0.0%	0.0%	0.0%
Unclassified _n	0.0%	4.4%	0.2%	2.0%	0.0%
CPU time	30	1214	63	2 168 005	—

0.1), the specificity of BLAST can approach 100 percent with very few false positives, whereas the number of unclassified sequences is enormous. Thus, we fixed the threshold of the E-value at 10 in subsequent experiments.

Tables 4, 5, 6, and 7 summarize the results and classification times, in seconds, of the four studied tools,

referred to as basic classifiers, on the four superfamilies in Table 1.⁴⁸ In addition to the basic classifiers, we developed an ensemble of classifiers, called COMBINER, that employs a weighted voter and works as follows. If a basic classifier gives an “undetermined” verdict, the classifier is regarded as “abstaining” and its verdict is not counted. The result of COMBINER is the same as the result of the majority

Table 6 Comparison of the studied classifiers on the ras superfamily

	BNN	BLAST	SAM	SAM-T99	COMBINER
Precision	98.7%	91.0%	95.5%	91.6%	99.7%
Specificity	99.3%	95.0%	99.8%	92.2%	99.6%
Sensitivity	96.1%	100.0%	100.0%	100.0%	100.0%
Unclassified _p	0.0%	0.0%	3.1%	0.0%	0.0%
Unclassified _n	0.0%	6.0%	4.6%	2.5%	0.0%
CPU time	29	1232	64	637 424	—

Table 7 Comparison of the studied classifiers on the ribitol superfamily

	BNN	BLAST	SAM	SAM-T99	COMBINER
Precision	96.6%	88.5%	99.4%	90.4%	99.4%
Specificity	97.0%	92.6%	100.0%	91.1%	99.2%
Sensitivity	94.3%	100.0%	100.0%	100.0%	100.0%
Unclassified _p	0.0%	0.0%	2.0%	0.0%	0.0%
Unclassified _n	0.0%	6.2%	0.3%	2.5%	0.0%
CPU time	27	1212	62	747 821	—

Table 8 Complementarity among the four studied tools BNN, BLAST, SAM, and SAM-T99. The percentages in the table add up to 100.

Classification Results	Percentage of the Test Sequences
All classifiers agreed and all were correct	80.88
All classifiers agreed and all were wrong	0.07
The classifiers disagreed and one of them was correct	18.91
The classifiers disagreed and all were wrong	0.14

of the remaining classifiers. If there is a tie on the verdicts given by the remaining classifiers, the result of COMBINER is the same as the result of the BNN classifier. We see that in comparison with BLAST, SAM, and SAM-T99, the BNN classifier is faster, yielding fewer unclassified sequences. COMBINER achieves the highest precision and SAM-T99 requires most time among all the classifiers.

Table 8 shows the complementarity of the four studied tools BNN, BLAST, SAM, and SAM-T99. We see that when all the four classifiers agree on their classification result, the result is correct with probability $80.88 \text{ percent} / (80.88 \text{ percent} + 0.07 \text{ percent}) = 99.91 \text{ percent}$.

Conclusion

In this paper we have presented a Bayesian neural network approach to classifying protein sequences. The main contributions of our work include:

- The development of new algorithms for extracting the global similarity and the local similarity from the sequences that are used as input features of the Bayesian neural network
- The development of new measures for evaluating the significance of 2-grams and frequently occurring motifs used in classifying the sequences
- Experimental studies in which we compare the performance of the proposed BNN classifier with three other classifiers, namely BLAST, SAM, and SAM-T99, on four different superfamilies of sequences in the PIR protein database

The main findings of our work include the following:

- The four studied classifiers, BNN, BLAST, SAM, and SAM-T99, complement each other; combining them yields better results than using the classifiers individually.
- The training phase, which is done only once, of the BNN classifier may take some time. After the classifier is trained, it runs significantly faster than BLAST and SAM-T99 in sequence classification.

Future research directions include:

- Comparison of motifs generated by different tools in combination with learning-based tools such as neural networks and hidden Markov models when applied to sequence classification in PIR, PROSITE,^{32,49} and other protein databases
- Generalization of the classifiers in combination with graph matching algorithms to analyze the sequence-structure relationship in protein and DNA sequences.

Acknowledgments

We thank the anonymous reviewers for their thoughtful comments and constructive suggestions, which helped to improve the paper. We also thank David Mackay for sharing the Bayesian neural network software with us, and Richard Hughey for providing us with the SAM package.

This work was supported in part by National Science Foundation grants IRI-9531548, IRI-9531554, IIS-9988345, and IIS-9988636, and by a grant from Novartis Pharmaceuticals Corporation.

**Trademark or registered trademark of Intel Corporation or Linus Torvalds.

Cited references and notes

1. C. F. Alex, J. W. Shavlik, and F. R. Blattner, "Neural Network Input Representations that Produce Accurate Consensus Sequences from DNA Fragment Assemblies," *Bioinformatics* **15**, No. 9, 723–728 (1999).
2. T. L. Bailey and W. N. Grundy, "Classifying Proteins by Family Using the Product of Correlated p -values," *Proceedings of the Third Annual International Conference on Computational Molecular Biology*, Lyon, France (April 11–14, 1999), pp. 10–14.
3. M. W. Craven and J. W. Shavlik, "Machine Learning Approaches to Gene Recognition," *IEEE Expert* **9**, No. 2, 2–10 (1994).
4. S. Eddy, "Profile Hidden Markov Models," *Bioinformatics* **14**, No. 9, 755–763 (1999).
5. W. N. Grundy and T. L. Bailey, "Family Pairwise Search with Embedded Motif Models," *Bioinformatics* **15**, No. 6, 463–470 (1999).
6. J. T. L. Wang, S. Rozen, B. A. Shapiro, D. Shasha, Z. Wang, and M. Yin, "New Techniques for DNA Sequence Classification," *Journal of Computational Biology* **6**, No. 2, 209–218 (1999).
7. J. Park, K. Karplus, C. Barrett, R. Hughey, D. Haussler, T. Hubbard, and C. Chothia, "Sequence Comparisons Using Multiple Sequences Detect Three Times as Many Remote Homologues as Pairwise Methods," *Journal of Molecular Biology* **284**, No. 4, 1201–1210 (1998).
8. *Pattern Discovery in Biomolecular Data: Tools, Techniques and Applications*, J. T. L. Wang, B. A. Shapiro, and D. Shasha, Editors, Oxford University Press, New York (1999).
9. C. H. Wu and J. McLarty, *Neural Networks and Genome Informatics*, Elsevier Science, New York (2000).
10. S. F. Altschul, T. L. Madden, A. A. Schaffer, J. Zhang, Z. Zhang, W. Miller, and D. J. Lipman, "Gapped Blast and PSI-Blast: A New Generation of Protein Database Search Programs," *Nucleic Acids Research* **25**, No. 17, 3389–3402 (1997).
11. A. Krogh, M. Brown, I. S. Mian, K. Sjolander, and D. Haussler, "Hidden Markov Models in Computational Biology: Applications to Protein Modeling," *Journal of Molecular Biology* **235**, No. 5, 1501–1531 (1994).
12. R. Hughey and A. Krogh, "Hidden Markov Models for Sequence Analysis: Extension and Analysis of the Basic Method," *Computer Applications in the Biosciences* **12**, No. 2, 95–107 (1996).
13. S. Eddy, G. Mitchison, and R. Durbin, "Maximum Discrimination Hidden Markov Models of Sequence Consensus," *Journal of Computational Biology* **2**, 9–23 (1995).
14. L. R. Rabiner, "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition," *Proceedings of the IEEE* **77**, No. 2, 257–286 (1989).
15. E. L. Sonnhammer, S. R. Eddy, and R. Durbin, "PFAM: A Comprehensive Database of Protein Domain Families Based on Seed Alignments," *Proteins* **28**, No. 3, 405–420 (1997).
16. K. Karplus, C. Barrett, and R. Hughey, "Hidden Markov Models for Detecting Remote Protein Homologies," *Bioinformatics* **14**, No. 10, 846–856 (1998).
17. R. Karchin and R. Hughey, "Weighting Hidden Markov Models for Maximum Discrimination," *Bioinformatics* **14**, No. 9, 772–782 (1998).
18. J. T. L. Wang, Q. Ma, D. Shasha, and C. H. Wu, "Application of Neural Networks to Biological Data Mining: A Case Study in Protein Sequence Classification," *Proceedings of the Sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, Boston, MA (August 20–23, 2000), pp. 305–309.
19. H. Hirsh and M. Noordewier, "Using Background Knowledge to Improve Inductive Learning of DNA Sequences," *Proceedings of the Tenth Annual Conference on Artificial Intelligence for Applications*, San Antonio, TX (March 1–4, 1994), pp. 351–357.
20. D. J. C. Mackay, "The Evidence Framework Applied to Classification Networks," *Neural Computation* **4**, No. 5, 698–714 (1992).
21. C. H. Wu, G. Whitson, J. McLarty, A. Ermongkonchai, and T. C. Chang, "Protein Classification Artificial Neural System," *Protein Science* **1**, No. 5, 667–677 (1992).
22. The total number of possible patterns from 2-gram encoding is n^2 where n is the number of different letters, namely 20, in the protein alphabet.
23. C. H. Wu, M. Berry, Y. S. Fung, and J. McLarty, "Neural Networks for Full-Scale Protein Sequence Classification: Sequence Encoding with Singular Value Decomposition," *Machine Learning* **21**, 177–193 (1995).
24. M. O. Dayhoff, R. M. Schwartz, and B. C. Orcutt, "A Model of Evolutionary Change in Proteins," *Atlas of Protein Sequence and Structure* **15**, Supplement 3, 345–358 (1978).
25. Both PAM (accepted point mutation) and BLOSUM (block substitution matrix) are amino acid substitution matrices. BLOSUM is derived from the BLOCKS database. For BLOSUM, see S. Henikoff and J. G. Henikoff, "Amino Acid Substitution Matrices from Protein Blocks," *Proceedings of the National Academy of Sciences* **89**, 10915–10919 (1992). For

- BLOCKS, see S. Henikoff and J. G. Henikoff, "Automated Assembly of Protein Blocks for Database Searching," *Nucleic Acids Research* **19**, 6565–6572 (1991).
26. N. A. Chuzhanova, A. J. Jones, and S. Margetts, "Feature Selection for Genetic Sequence Classification," *Bioinformatics* **14**, No. 2, 139–143 (1998).
 27. The term "distance" is used by M. Dash and H. Liu, "Feature Selection for Classification," *Intelligent Data Analysis* **1**, No. 3 (1997). The electronic journal is available at <http://www-east.elsevier.com/ida>. Both this and the next reference address feature selection for classification.
 28. M. Ben-Bassat, "Use of Distance Measures, Information Measures and Error Bounds in Feature Evaluation," *Classification, Pattern Recognition and Reduction of Dimensionality: Handbook of Statistics*, Volume 2, P. R. Krishnaiah and L. N. Kanal, Editors, North-Holland Publishing Company, Amsterdam (1982), pp. 773–791.
 29. V. V. Solovyev and K. S. Makarova, "A Novel Method of Protein Sequence Classification Based on Oligopeptide Frequency Analysis and Its Application to Search for Functional Sites and to Domain Localization," *Computer Applications in the Biosciences* **9**, No. 1, 17–24 (1993).
 30. Our experimental results show that choosing $N_g \geq 30$ can yield reasonably good performance provided one has sufficient (e.g., > 200) training sequences. We have also experimented with different combinations of 2-grams, e.g., using the top N_g features together with the bottom N_g features with the smallest $D(X)$ values. The results are worse than using the top N_g features alone.
 31. J. T. L. Wang, G. W. Chirn, T. G. Marr, B. A. Shapiro, D. Shasha, and K. Zhang, "Combinatorial Pattern Discovery for Scientific Data: Some Preliminary Results," *Proceedings of the ACM SIGMOD International Conference on Management of Data*, Minneapolis, MN (May 24–27, 1994), pp. 115–125.
 32. J. T. L. Wang, T. G. Marr, D. Shasha, B. A. Shapiro, G. W. Chirn, and T. Y. Lee, "Complementary Classification Approaches for Protein Sequences," *Protein Engineering* **9**, No. 5, 381–386 (1996).
 33. L. C. K. Hui, "Color Set Size Problem with Applications to String Matching," *Combinatorial Pattern Matching*, Lecture Notes in Computer Science, Volume 644, A. Apostolico, M. Crochemore, Z. Galil, and U. Manber, Editors, Springer-Verlag (1992), pp. 230–243.
 34. O is "order of magnitude"; $O(n)$ as used here means that time (and space) increase linearly as n increases.
 35. S. Wu and U. Manber, "Fast Text Searching Allowing Errors," *Communications of the ACM* **35**, No. 10, 83–91 (1992).
 36. A. Brazma, I. Jonassen, E. Ukkonen, and J. Vilo, "Discovering Patterns and Subfamilies in Biosequences," *Proceedings of the Fourth International Conference on Intelligent Systems for Molecular Biology*, St. Louis, MO (June 12–15, 1996), pp. 34–43.
 37. J. Rissanen, "Modeling by Shortest Data Description," *Automatica* **14**, 465–471 (1978).
 38. C. E. Shannon, "A Mathematical Theory of Communication," *Bell System Technical Journal* **27**, 379–423, 623–656 (1948).
 39. The actual number of sequences in \mathcal{S}_p that are encoded by Scheme 2 is dependent on motif. For each motif used in the study presented here, more than 1/10 of the sequences are encoded based on the motif using Scheme 2.
 40. A. Califano, "SPLASH: Structural Pattern Localization Analysis by Sequential Histogramming," available at http://www.research.ibm.com/topics/popups/deep/math/html/splash_bioinformatics.pdf. See also <http://www.research.ibm.com/splash/>.
 41. R. Hart, A. Royyuru, G. Stolovitzky, and A. Califano, "Systematic and Automated Discovery of Patterns in PROSITE Families," *Proceedings of the Fourth Annual International Conference on Computational Molecular Biology*, Tokyo, Japan (April 8–11, 2000).
 42. This software is available at <http://wol.ra.phy.cam.ac.uk/pub/mackay/README.html>.
 43. W. C. Barker, J. S. Garavelli, H. Huang, P. B. McGarvey, B. Orcutt, G. Y. Srinivasarao, C. Xiao, L. S. Yeh, R. S. Ledley, J. F. Janda, F. Pfeiffer, H. W. Mewes, A. Tsugita, and C. H. Wu, "The Protein Information Resource (PIR)," *Nucleic Acids Research* **28**, No. 1, 41–44 (2000).
 44. Note that the BNN classifier does not yield any unclassified sequence. By contrast, the three other classifiers BLAST, SAM, and SAM-T99 we compare with yield unclassified sequences, as our experimental results show in the next section.
 45. The time spent in matching a test sequence with the motifs is linearly proportional to the number of the motifs one uses.
 46. We used log-odds scores, as opposed to E-values, for this tool because the E-value for a training sequence was calculated with respect to the training data set while the E-value for a test sequence was calculated with respect to the test data set. These two kinds of E-values were not directly comparable.
 47. The E-value of the HMM target model used in the study presented here was 20. We have experimented with other E-values and the results were worse.
 48. In examining CPU time for SAM-T99, we note that the time spent for classifying the kinase sequences shown in Table 5 is much higher than the times spent for classifying the other sequences shown in Tables 4, 6, and 7. The reason is that for each kinase sequence, there are many homologs in the non-redundant protein database maintained at NCBI. Thus, SAM-T99 gets more homologs for a kinase sequence than the homologs for the other sequences, and consequently it takes more time to build the HMM target model for the kinase sequence.
 49. A. Bairoch, "The PROSITE Dictionary of Sites and Patterns in Proteins, Its Current Status," *Nucleic Acids Research* **21**, 3097–3103 (1993).

Accepted for publication January 24, 2001.

Jason T. L. Wang *Department of Computer and Information Science, New Jersey Institute of Technology, University Heights, Newark, New Jersey 07102 (jason@cis.njit.edu).* Dr. Wang received the B.S. degree in mathematics from National Taiwan University. He received the Ph.D. degree in computer science from the Courant Institute of Mathematical Sciences, New York University, in 1991. He is a full professor in the Computer and Information Science Department at the New Jersey Institute of Technology and director of the university's Data and Knowledge Engineering Laboratory. Dr. Wang's research interests include data mining and databases, pattern recognition, bioinformatics, and Web computing. He has published more than 100 technical articles in conference proceedings, journals, and books. He is an editor and author of the book *Pattern Discovery in Biomolecular Data: Tools, Techniques, and Applications*, published in 1999 by Oxford University Press. He is also program cochair for 2001 of the Atlantic Symposium on Computational Biology, Genome Information Systems and Technology, held at Duke University.

Qicheng Ma *Novartis Pharmaceuticals Corporation, Summit, New Jersey 07901 (qicheng.ma@pharma.novartis.com)*. Dr. Ma received the B.S. degree in computer science and engineering from Beijing University of Aerospace and Aeronautics, the M.S. degree in computer science from Beijing Research Institute of Computer Application and Simulation Technology, and the Ph.D. degree in computer and information science from the New Jersey Institute of Technology. He is currently working in the Life Science Informatics Unit, Department of Functional Genomics, Novartis Pharmaceuticals Corporation, as a senior scientist. His research interests are bioinformatics and data mining, particularly in applying machine learning algorithms to solving biological problems. He is currently working on classification and clustering of genomes across evolution.

Dennis Shasha *Courant Institute of Mathematical Sciences, New York University, 251 Mercer Street, New York, New York 10012 (shasha@cs.nyu.edu)*. Dr. Shasha is a professor at NYU's Courant Institute where he does research on biological pattern discovery, combinatorial pattern matching on trees and graphs, software support for the exploration of unfamiliar databases, database tuning, and database design for time series. After graduating from Yale in 1977, he worked for IBM designing circuits and microcode for the 3090 processor. He completed his Ph.D. degree at Harvard in 1984 under thesis advisor Nat Goodman. He has written a professional reference book (*Database Tuning: A Principled Approach*, published in 1992 by Prentice Hall), two books about a mathematical detective named Dr. Ecco (*The Puzzling Adventures of Dr. Ecco*, published in 1988 by W. H. Freeman and republished in 1998 by Dover Publications, and *Codes, Puzzles, and Conspiracy*, published in 1992 by W. H. Freeman), and a book of biographies about great computer scientists (*Out of Their Minds: The Lives and Discoveries of 15 Great Computer Scientists*, published in 1995 by Springer-Verlag). He is a coauthor of *Pattern Discovery in Biomolecular Data: Tools, Techniques, and Applications*, published in 1999 by Oxford University Press. In addition, he has coauthored 30 journal papers and 40 conference papers and these days spends most of his time in building data mining and pattern recognition software. He writes a monthly puzzle column for *Scientific American*.

Cathy H. Wu *National Biomedical Research Foundation, Georgetown University Medical Center, 3900 Reservoir Road, NW, Washington, DC 20007 (wuc@nbrf.georgetown.edu)*. Dr. Wu received her M.S. and Ph.D. degrees in plant pathology from Purdue University and a second M.S. degree, in computer science, from the University of Texas at Tyler. She is Director of Bioinformatics, National Biomedical Research Foundation, and an adjunct professor of Georgetown University Medical Center, George Mason University, and the University of Maryland. Dr. Wu has conducted bioinformatics research since 1990 and developed several systems including the patented MOTIFIND neural network, GeneFIND protein family identification system, and ProClass/iProClass protein family database. She is author of the book *Neural Networks and Genome Informatics*, published by Elsevier Science in 2000. She is also a member of the board of directors of the International Society for Computational Biology and the program cochair for 2001, with Dr. Wang, of the Atlantic Symposium, Computational Biology and Genome Information Systems.