

# Modeling protein families using probabilistic suffix trees

Gill Bejerano and Golan Yona

Institute of Computer Science, Hebrew University, Jerusalem 91904, Israel.

## Abstract

We present a method for modeling protein families by means of probabilistic suffix trees (PSTs). The method is based on identifying significant patterns in a set of related protein sequences. The input sequences do not need to be aligned, nor is delineation of domain boundaries required. The method is automatic, and can be applied, without assuming any preliminary biological information, with surprising success. Incorporating basic biological considerations such as amino acid background probabilities, and amino acids substitution probabilities can improve the performance in some cases.

The PST can serve as a predictive tool for protein sequence classification, and for detecting conserved patterns (possibly functionally or structurally important) within protein sequences. The method was tested on one of the state of the art databases of protein families, namely, the Pfam database of HMMs, with satisfactory performance.

## 1 Introduction

In the last few years there is growing effort to organize the known protein sequences into families and establish databases of protein motifs and domains. These databases have become an important tool in the analysis of newly discovered protein sequences. The biological knowledge accumulated in these databases in the form of consensus patterns, profiles and HMMs helps to detect patterns of biological significance in new protein sequences. In many cases, such discoveries can lead to the assignment of the new sequences to the existing protein families. These attempts are extremely important in view of the large number of newly sequenced

proteins which await analysis.

Generally, the existing approaches can be divided into those based on short conserved motifs (e.g. [Bairoch et al. 1997, Attwood et al. 1998, Henikoff & Henikoff 1991]) and those which are based on whole domains (e.g. [Sonnhammer & Kahn 1994, Sonnhammer et al. 1998]). The manually defined patterns in PROSITE have served as an excellent seed for several such works. The methods used to represent these motifs and domains vary, and among the most popular forms are the consensus patterns [Bairoch et al. 1997, Attwood et al. 1998], the position-specific scoring matrices (profiles) [Gribskov et al. 1987, Henikoff & Henikoff 1991] and the HMMs [Krogh et al. 1996]. These forms differ in their mathematical complexity, as well as in their sensitivity/selectivity.

To model a motif, a domain or a protein family, many approaches start by building a multiple alignment. The prerequisite of a valid alignment of the input sequences is a major drawback of this strategy, since building a multiple alignment of many sequences is not an easy task. Current methods for multiple alignment apply heuristic algorithms which are not guaranteed to find the optimal alignment. Moreover, the optimal alignment itself is not guaranteed to be biologically accurate. When the related sequences are diverged and the common pattern cannot be easily distinguished from noise, or when the number of samples used to build the alignment is small then the validity of the automatically generated alignment is questionable, and the resulting model (consensus pattern, profile, HMM) may be of poor quality. Delineation of domain boundaries makes the problem even more difficult. Therefore, fully automatic methods, based on multiple alignments are not guaranteed to be of high quality for classification and pattern detection. Obviously, manually constructed alignments are preferred (PROSITE [Bairoch et al. 1997] serves as an excellent example. Another example is the Pfam database [Sonnhammer et al. 1998] that uses a semi-automatic procedure, which combines manual analysis with automatic tools). However, the

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

RECOMB '99 Lyon France

Copyright ACM 1999 1-58113-069-4/99/04...\$5.00

sheer size of current databases calls for the development of high performance automated methods.

Other approaches which do not require prealigned sequences were tested as well [Hanke et al. 1996, Neuwald et al. 1994, Jonassen et al. 1995, Suyama et al. 1995], but most of them are limited to short patterns with limited flexibility (i.e. some constraints are imposed on the pattern). Such techniques were used to construct the blocks database [Henikoff & Henikoff 1991], where motifs are defined based on blocks (ungapped short regions of aligned amino acids) extracted from groups of related proteins, and each family is associated with a set of blocks.

Here we present an alternative approach for detecting significant patterns, based on probabilistic suffix trees (PSTs), originally introduced in [Ron et al. 1994]. The model draws on identifying significant appearances of short segments among many of the input protein sequences, regardless of the relative positions of these segments within the different proteins. These significant segments reflect some statistical properties of the corresponding protein family. Practically, they induce a probability distribution on the *next* symbol to appear right after the segment. Although probably no simple underlying statistical model exists for protein families, this property implies a certain feature which is common to a large subset of them. That feature, termed **short memory**, which is common to many natural sources, indicates that for a certain sequence the (empirical) probability distribution of the next symbol given the preceding subsequence can be quite accurately approximated by observing no more than the last  $L$  symbols in that subsequence.

This observation has led before to the suggestion of modeling many natural sources by use of order  $L$  Markov chains ( $L$  being the memory length of the model), or by the more extensive family of Hidden Markov Models (HMMs) which is more complex and allows to capture subtler features of the original source. These two families, although being able to model rich sources in an efficient manner, and lending themselves to extensive mathematical treatment, have critical drawbacks for practical use. The Markov chain model suffers from exponential growth in the number of states for a non trivial memory length, and poor source approximation at low order. The HMMs family suffers from known learnability hardness results [Abe & Warmuth 1992, Gillman & Sipser 1994], and in practice, a high quality multiple alignment of the input sequences is required to obtain a reliable model. The probabilistic suffix trees are inspired by the same reasoning. However, this family can model rich sources with high efficiency using a reasonable amount of memory [Ron et al. 1994]. The strength of the model stems from its memory length which is context dependent as observed for many natural sources.

The method is simple to apply. It does not require the input sequences to be aligned first, nor is it limited to very short patterns. Despite its simplicity, the approach is surprisingly powerful. In the next sections we first describe the model, and then demonstrate its power by applying it to known protein families.

## 2 Theory

The definitions and the subsequently implemented algorithm are a variant on the learning algorithm originally presented in [Ron et al. 1994].

A PST over an alphabet is a non empty tree, whose nodes vary in degree between zero (for leaves) and the size of the alphabet. Each edge in the tree is labeled by a single symbol of the alphabet, such that no symbol is represented by more than one edge branching out of any single node (hence the degree of each node is bounded by the size of the alphabet). Nodes of the tree are labeled by a string, which is the one generated by walking **up** the tree from that node to the root. Each node is assigned a probability distribution vector over the alphabet. When the PST is used to predict significant patterns within a query string, this probability distribution vector comes into play. It corresponds to the probabilities the tree assigns to the next query symbol, given that the longest suffix of preceding query symbols found in the tree matches that particular node's label. An example of a PST is given in Fig. 1.

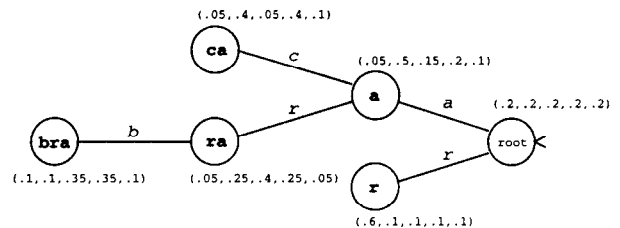


Figure 1: An example of a PST over the alphabet  $\Sigma = a, b, c, d, r$ . The tree is shown in landscape mode, which makes the prediction step more intuitive. The root is the rightmost node. The vector that appears near each node is the probability distribution over the next symbol. For example, the probability distribution associated with the substring  $ra$  is 0.05, 0.25, 0.4, 0.25 and 0.05 for the symbol  $a, b, c, d$  and  $r$  respectively (e.g. the probability to observe  $c$  after a substring, whose largest suffix in the tree is  $ra$ , is 0.4).

### 2.1 Definitions

Let  $\Sigma$  be the alphabet (e.g. the alphabet of 20 amino acids for protein sequences, or 4 nucleic acids for DNA sequences), and let  $r^1, r^2, \dots, r^m$  be the sample set of  $m$  strings over the alphabet  $\Sigma$ , where the length of the  $i$ th ( $i = 1..m$ ) string is  $l_i$  (i.e.  $r^i = r_1^i r_2^i \dots r_{l_i}^i$ ,  $r_j^i \in \Sigma$ ).

First we define the empirical probability of a substring over the given sample set, which is the number of

times that substring was observed in the sample set divided by the maximal number of (possibly overlapping) occurrences a pattern of the same length could have had, considering the sample size<sup>1</sup>. Formally speaking, given a string  $s$  of length  $l$  ( $s = s_1s_2\dots s_l$ ) we define a set of variables

$$\chi_s^{i,j} = \begin{cases} 1 & \text{if } s_1s_2\dots s_l = r_j^i r_{j+1}^i \dots r_{j+(l-1)}^i \\ 0 & \text{otherwise} \end{cases}$$

for each  $i = 1..m$  and  $j = 1..l_i - (l-1)$ . These indicator variables have a value of one if and only if the string  $s$  is a substring of  $r^i$  starting at position  $j$  (see Fig. 2).

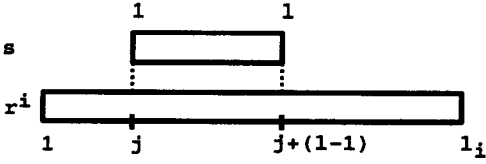


Figure 2: A match of  $s$  and  $r^i$  at position  $j$ . The indicator variable is set to one:  $\chi_s^{i,j} = 1$ .

The number of (possibly overlapping) occurrences of the string  $s$  in the string set  $\{r^i\}$  is given by

$$\chi_s = \sum_{i,j} \chi_s^{i,j}$$

The maximal number of possible occurrences (with overlaps) of string of length  $|s| = l$  is

$$N_{|s|} = \sum_{i \text{ s.t. } l_i \geq l} (l_i - (l-1))$$

Therefore, we set the empirical probability of the string  $s$  to be

$$\tilde{P}(s) = \frac{\chi_s}{N_{|s|}}$$

We go on to define the conditional empirical probability of observing a symbol right after a given substring as the number of times that symbol has shown up right after the given substring divided by the total number of times that substring has shown up at all, followed by any symbol. Specifically, let  $\chi_{s*}$  be the number of non-suffix occurrences of the string  $s$  in the string set  $\{r^i\}$ , i.e.

$$\chi_{s*} = \sum_{\sigma' \in \Sigma} \chi_{s\sigma'}$$

Then the conditional empirical probability of observing the symbol  $\sigma$  right after the string  $s$  is defined by

$$\tilde{P}(\sigma|s) = \frac{\chi_{s\sigma}}{\chi_{s*}}$$

Finally, we define  $\text{su}f(s) = s_2s_3\dots s_l$ .

<sup>1</sup>In general, computing the maximal number of occurrences of a specific string  $s$  is more complicated (due to repetitive patterns in  $s$  itself). In practice, this approximation works well.

## 2.2 Building the PST

First, we define  $L$  to be the memory length of the PST (i.e. the maximal length of a possible string in the tree).

We work our way gradually through the space of all possible substrings of lengths 1 through  $L$ , starting at single letter substrings, and abstaining from further extending a substring whenever its empirical probability has gone below a certain threshold ( $P_{min}$ ), or on having reached the maximal  $L$  length boundary. The  $P_{min}$  cutoff avoids an exponentially large search space.

At the beginning of the search we hold a PST consisting of a single root node. Then, for each substring we decide to examine, we check whether there is some symbol in the alphabet for which the empirical probability of observing that symbol right after the given substring is non negligible, and is also significantly different from the empirical probability of observing that same symbol right after the string obtained from deleting the leftmost letter from our substring<sup>2</sup>. Whenever these two conditions hold, the substring, and all necessary nodes on its path, are added to our PST.

The reason for the two step pruning (first defining all nodes to be examined, then going over each and every one of them) stems from the nature of PSTs. A leaf in a PST is deemed useless if its prediction function is identical (or almost identical) to its parent node. However, this in itself is no reason not to examine its sons further while searching for significant patterns. Therefore, it may, and does happen that consecutive inner PST nodes are almost identical.

Finally the resulting PST skeleton is first added the node prediction functions, using the appropriate conditional empirical probability, and then these probabilities are smoothed using a standard technique so that no single symbol is absolutely impossible right after any given substring (even though the empirical counts may attest differently).

We now present the procedure for building a PST out of a sample set. The procedure uses 5 external parameters:  $L$  the memory length,  $P_{min}$  the minimal probability with which strings are required to occur,  $r$  which is a simple measure of the difference between the prediction of the candidate at hand and its direct father node,  $\gamma_{min}$  the smoothing factor and  $\alpha$  (for example, an effective set of parameters is given in the legend of table 1, appendix A).

We use  $\tilde{T}$  to denote the tree,  $\tilde{S}$  to denote the set of (unique) strings which we need to check, and  $\tilde{\gamma}_s$  to denote the probability distribution (over the next symbol) associated with the node  $s$ .

<sup>2</sup>This string corresponds to the label of the direct father of the node we are currently examining (note that the father node has not necessarily been added itself to the PST at this time)

The algorithm: **Build-PST** ( $P_{min}, \alpha, \gamma_{min}, r, L$ )

1. **Initialization:** let  $\bar{T}$  consist of a single root node (with an empty label), and let  $\bar{S} \leftarrow \{\sigma \mid \sigma \in \Sigma \text{ and } \tilde{P}(\sigma) \geq P_{min}\}$ .

2. **Building the PST skeleton:** While  $\bar{S} \neq \emptyset$ , pick any  $s \in \bar{S}$  and do:

- (a) Remove  $s$  from  $\bar{S}$
- (b) If there exists a symbol  $\sigma \in \Sigma$  such that

$$\tilde{P}(\sigma|s) \geq (1 + \alpha)\gamma_{min}$$

and

$$\frac{\tilde{P}(\sigma|s)}{\tilde{P}(\sigma|suf(s))} \begin{cases} \geq r \\ \text{or} \\ \leq 1/r \end{cases}$$

then add to  $\bar{T}$  the node corresponding to  $s$  and all the nodes on the path to  $s$  from the deepest node in  $\bar{T}$  that is a suffix of  $s$ .

- (c) If  $|s| < L$  then for every  $\sigma' \in \Sigma$ , if  $\tilde{P}(\sigma'|s) \geq P_{min}$ , then add  $\sigma's$  to  $\bar{S}$ .

3. **Smoothing the prediction probabilities:** For each  $s$  labeling a node in  $\bar{T}$ , let

$$\bar{\gamma}_s(\sigma) \equiv (1 - |\Sigma|\gamma_{min})\tilde{P}(\sigma|s) + \gamma_{min}$$

The final step (step 3) of the learning algorithm is the smoothing process, which assures that no symbol is predicted to have a zero probability, no matter what suffix is observed before it. The value of  $\gamma_{min}$  defines the minimum probability for a symbol, and the empirical probabilities should be adjusted to satisfy this requirement. This is done by decreasing the empirical probabilities, such that a total of  $|\Sigma|\gamma_{min}$  is “collected”, to be later shared by all symbols. The decrement of each empirical probability is done in proportion to its value.

### 2.3 Prediction using a PST

Given a string  $s$  its prediction by a certain PST is done letter by letter, where the probability of each letter is calculated by scanning the tree in search of the longest suffix that appears in the tree and ends just before that letter. The conditional probability of this letter given this suffix is given by the probability distribution associated with the corresponding node in the PST.

For example, to predict the string  $s = abracadabra$  with the PST given in Fig. 1 the following procedure is carried out:

$$\begin{aligned} P^T(s) &= P^T(abracadabra) \\ &= P^T(a)P^T(b|a)P^T(r|ab)P^T(a|abr) \\ &\quad P^T(c|abra)P^T(a|abrac) \cdots P^T(a|abracadabr) \end{aligned}$$

$$\begin{aligned} &= \bar{\gamma}_{root}(a) \bar{\gamma}_a(b) \bar{\gamma}_{root}(r) \bar{\gamma}_r(a) \bar{\gamma}_{bra}(c) \\ &\quad \bar{\gamma}_{root}(a) \cdots \bar{\gamma}_r(a) \\ &= 0.2 \cdot 0.5 \cdot 0.2 \cdot 0.6 \cdot 0.35 \cdot 0.2 \cdot 0.4 \cdot 0.2 \cdot \\ &\quad 0.5 \cdot 0.2 \cdot 0.6 = 4.032 \cdot 10^{-6} \end{aligned}$$

### 3 Results

To test our approach, a PST was created for each family in the Pfam database [Sonnhammer et al. 1998] release 1.0. This database contains 175 protein families<sup>3</sup> derived from the SWISSPROT 33 database [Bairoch & Boeckman 1992]. Each family was divided into a training set and a test set, in ratio 4:1, such that 4/5 of the family members are in the training set, and the PST was trained on the training set. Then, for each sequence in the database we calculated its probability as predicted by the PST. To avoid bias due to length differences, the probability is normalized by the length, and each sequence is reported along with its probability per letter. Hits are sorted by decreasing probability.

The quality of the PST model is estimated by applying the “equivalence number” criterion [Pearson 1995]. The equivalence number criterion sets a threshold at the point where the number of false positives equals the number of false negatives, i.e. it is the point of balance between selectivity and sensitivity (we use the term *isopoint* to denote this point). A hit which belong to the family (true positive) and scores above this threshold, is considered successfully detected. The quality of the model is measured by the number of true positives detected relative to the total number of proteins in the family. The results for the 170 protein families in the Pfam database release 1.0, with more than 10 members each, are given in table 1 (see appendix A).

It should be noted that table 1 only demonstrates the potential of the PST model, but **must not be taken as an upper bound on its performance**. To obtain these results we simply ran the PST learning procedure with a fixed set of parameters for all families which we found to result in good performance in reasonable running time. However, the performance of a PST can be improved by simply tuning the values of the parameters, either globally or per each family. One can either decide to examine more nodes (lower  $P_{min}$ ), or lower the criteria of acceptance of a candidate (lower  $\alpha$  or lower  $r$ ) or even deepen the tree (increase  $L$ ). This can be done in an efficient, incremental manner. This is demonstrated in Fig. 3 for the glyceraldehyde 3-phosphate dehydrogenases family. Adding more nodes

<sup>3</sup>The Pfam database uses a semi-automatic procedure, which combines manual analysis and automatic tools to build multiple alignments for groups of related proteins. Each multiple alignment is closely monitored to avoid misalignments, and HMMs are derived from these alignments. The HMMs can be used for searching close relatives in the database, and the most significant hits are included in the corresponding family. The process iterates until it converges, while being traced manually.

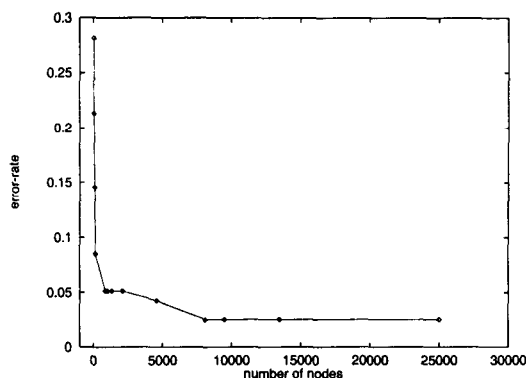


Figure 3: **Improving prediction by increasing the number of nodes.** A PST was built for the glyceraldehyde 3-phosphate dehydrogenases family, for different values of  $P_{min}$ , and the quality of the PST was estimated by applying the equivalence number criterion (see text). The graph plots the error rate (the number of family members which were not identified, i.e. their score was below the threshold set by the iso-point) vs. the number of nodes in the PST (which increases as  $P_{min}$  decreases). Note that the error rate decreases as the number of nodes increases. At some value of  $P_{min}$  the quality does not improve much, while the tree keeps growing. If the results are satisfactory then it is suggested to stop the incremental process at this point.

would tend to increase sensitivity **without** decreasing its selectivity, simply because more (longer) substrings that are observed in the training set are “recorded” in the PST. This means that only leaves are further expanded, while the probability distributions of internal nodes are not affected. This way, the predictions over the test set are refined. However, since long substrings observed in the training set are not expected to occur in unrelated sequences, the prediction of unrelated sequences is based on short substrings corresponding to internal nodes close to the root, and therefore is not expected to change.

The only limitation one should keep in mind is that the size of the tree (the number of nodes), as well as the running time, may increase significantly as the parameters are refined, while the improvement in the quality may not be significant. However, if computation time is of no concern then the PST can run as long as the quality improves (i.e. more family members are detected above the iso-point). In two cases no further improvement is expected: 1) when all sequences of the family are detected 2) when all strings in the training set are exhaustively recorded, and predicted with very high probability.

An additional improvement is expected if a larger sample set is used to train the PST. Currently, the PST is built from the training set **alone**. Obviously, training the PST on **all** strings of a family (as Pfam database is created) should improve its prediction as well.

To demonstrate the performance of this model, two examples are given in Fig. 4. The probabilities (as predicted by the PST) of the training set sequences, the

test set sequences and all the other sequences in the database are plotted vs. the sequences length. Note that the unrelated sequences show a clear linear relation in log scale. The training set and the test set samples are located far below this line.

The PST can also be used to predict which segments of a given query sequence are suspected to be functionally or structurally important. These segments correspond to regions of high probability. This is demonstrated in Fig. 5. Note that along these regions not all letters are essentially of high probability, which may suggest a substitution or a gap in the query sequence. Recall that the PST does not require the input sequences to be aligned nor does it make use of such information during the learning process. The input sequences were not fragmented according to domain boundaries before the learning phase, and therefore this information was solely self attained. In this view, the PST model can also help to guide a multiple alignment of a set of related sequences, by suggesting an initial seed, which is based on the regions of high probability.

The performance evaluation procedure that we applied assesses the quality of the PST model in terms of its ability to predict the correct assignments of proteins to **a-priori defined groups**, and our reference set here is the HMM based Pfam database. Note that this assessment does not measure the relative merit of the PST model with respect to the HMM model in general, since the reference set depends on the HMM model itself (see footnote 3). It would be interesting to compare the performance of the PST model to the performance of the HMM model in an objective manner, on groups that are defined by another, independent classification, and where the HMM is built from a multiple alignment without any manual calibration. We are currently checking this aspect.

#### 4 Adding biological considerations to the PST model

The basic PST model does not require any assumption on the input data, nor does it utilize any a-priori information we may have about the data. Incorporating such information may improve the performance of this model, since it adjusts the general purpose model to the specific problem of identifying significant patterns in macromolecules, where some characteristic features and processes are observed. Specifically, the information we would like to consider is the amino acids background probabilities, and the amino acids substitution probabilities. These distributions are integrated into the PST model and some changes are made to account for this a-priori information.

Few additional definitions are needed here: Let  $P_0$  be the a-priori probability distribution over the alphabet  $\Sigma$  and let  $q_{ab}$  be the probability that amino acid  $a$  is

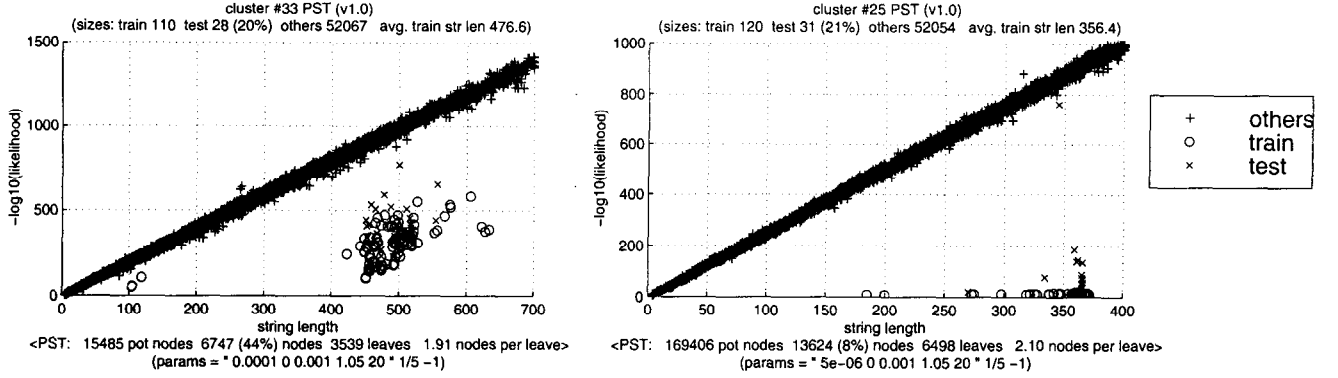


Figure 4: **Left:** Performance of a PST of the Neurotransmitter-gated ion-channels ( $P_{min} = 0.0001$ ). **Right:** Performance of a PST of the MHC I family ( $P_{min} = 0.000005$ ). The latter, being trained with a lower value of  $P_{min}$ , is an example of an extreme fit of the PST to the training set. Note that the prediction of the test set improves as well, while the unrelated sequences are left unaffected. When the PST was trained with the same set of parameters as the PST on the left, its performance graph resembled that of the graph on the left.

replaced by amino acid  $b$  ( $a, b \in \Sigma$ ). In the new procedure, only strings that their prediction of the next symbol differs from what expected simply by chance (the background probabilities) are included in the tree.

Also, define  $\bar{\chi}_s$  as the number of *distinct* sequences which include the substring  $s$ . This number is required to exceed a certain threshold  $N_{min}$ , which is defined in proportion to the total number  $m$  of strings in the sample set (i.e.  $N_{min} = c \cdot m$  where  $c$  is chosen to be, say, 0.2 so that the substring  $s$  is required to exist in at least 20% of the member proteins)

Finally, the last step (step 3) of the learning algorithm (the smoothing step), is modified, and it is now based on a position-based pseudo-counts method (similarly to the method suggested by [Henikoff & Henikoff 1996]). This method adds hypothetical sequences to the sample set in order to avoid zero probabilities which are due to undersampling. The number and the distribution of pseudo-counts is “position-based” (i.e. different for each node) and takes into account the *diversity*, i.e. the number of different amino acid types observed after the corresponding substring, as well as the counts of actually observed amino acids.

For a node represented by the substring  $s$ , denote by  $R_s$  the diversity at  $s$ , i.e. the number of different amino acids observed after  $s$ .

$$R_s = |\{\sigma \mid \chi_{s\sigma} > 0\}|$$

Denote by  $B_s$  the total number of pseudo-counts added for the node  $s$ . As suggested in [Henikoff & Henikoff 1996]  $B_s$  is set to  $B_s = \mu R_s$  where  $\mu$  can be optimized for best performance<sup>4</sup>. Then, the number of pseudo-

<sup>4</sup>We haven’t optimized this parameter yet. Ad hoc, we used the same value as in [Henikoff & Henikoff 1996]

counts of amino acid  $a$  at this node is given by

$$\begin{aligned} b_a &= B_s \sum_{i=1}^{20} Prob(i|s) \cdot Prob(a|i) \\ &= B_s \sum_{i=1}^{20} \frac{\chi_{si}}{\chi_{s*}} \cdot \frac{q_{ia}}{Q_i} \end{aligned}$$

Where  $Q_i = \sum_{k=1}^{20} q_{ik}$ . The probability of observing  $a$  after the string  $s$  is defined as the weighted average of the empirical probability  $\tilde{P}(a|s)$  and the a-priori probability as defined by the pseudo-counts,  $P_{pse}(a|s) = b_a/B_s$ .

The new procedure is described next:

The algorithm: **New-Build-PST** ( $N_{min}$ ,  $\alpha$ ,  $r$ ,  $L$ )

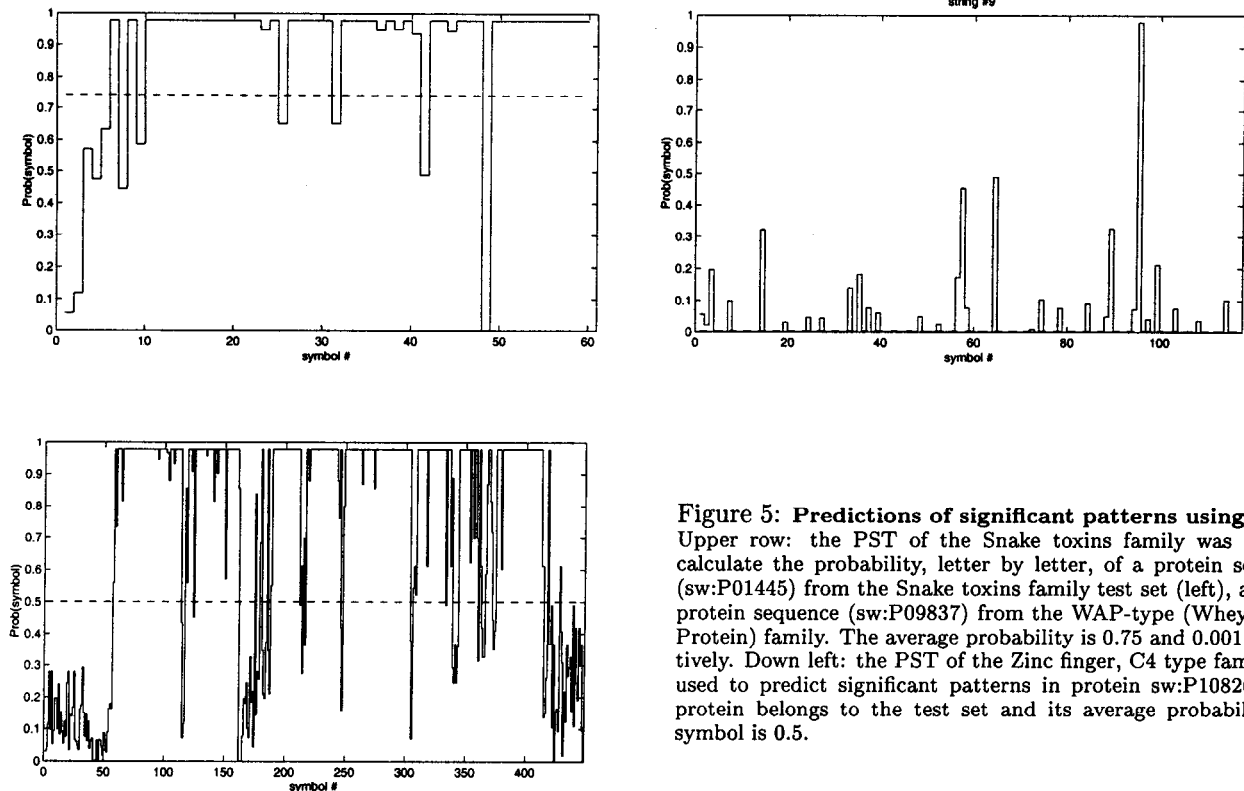
1. **Initialization:** let  $\bar{T}$  consist of a single root node (with an empty label), and let  $\bar{S} \leftarrow \{\sigma \mid \sigma \in \Sigma \text{ and } \bar{\chi}_\sigma \geq N_{min}\}$ .
2. **Building the PST skeleton:** While  $\bar{S} \neq \emptyset$ , pick any  $s \in \bar{S}$  and do:
  - (a) Remove  $s$  from  $\bar{S}$
  - (b) If there exists a symbol  $\sigma \in \Sigma$  such that

$$\tilde{P}(\sigma|s) \geq (1 + \alpha)P_0(\sigma)$$

and

$$\frac{\tilde{P}(\sigma|s)}{\tilde{P}(\sigma|suf(s))} \begin{cases} \geq r \\ \text{or} \\ \leq 1/r \end{cases}$$

then add to  $\bar{T}$  the node corresponding to  $s$  and all the nodes on the path to  $s$  from the deepest node in  $\bar{T}$  that is a suffix of  $s$ .



**Figure 5: Predictions of significant patterns using PSTs.** Upper row: the PST of the Snake toxins family was used to calculate the probability, letter by letter, of a protein sequence (sw:P01445) from the Snake toxins family test set (left), and of a protein sequence (sw:P09837) from the WAP-type (Whey Acidic Protein) family. The average probability is 0.75 and 0.001 respectively. Down left: the PST of the Zinc finger, C4 type family was used to predict significant patterns in protein sw:P10826. The protein belongs to the test set and its average probability per symbol is 0.5.

- (c) If  $|s| < L$  then for every  $\sigma' \in \Sigma$ , if  $\bar{\chi}_{\sigma's} \geq N_{min}$  then add  $\sigma's$  to  $\bar{S}$ .

3. **Smoothing the prediction probabilities:** For each  $s$  labeling a node in  $\bar{T}$ , let

$$\begin{aligned} \bar{\gamma}_s(\sigma) &\equiv \frac{\chi_{s*}}{\chi_{s*} + B_s} \bar{P}(a|s) + \frac{B_s}{\chi_{s*} + B_s} P_{pse}(a|s) \\ &= \frac{\chi_{s*}}{\chi_{s*} + B_s} \frac{\chi_{sa}}{\chi_{s*}} + \frac{B_s}{\chi_{s*} + B_s} \frac{b_a}{B_s} \end{aligned}$$

When this procedure was applied<sup>5</sup> to 40 Pfam families on which the original PST model performed worst, we found that for 22 families same or better separation was obtained using much smaller trees (for example, the error rate for the Sigma-54 transcription factors improved from 16.1% with 5040 nodes, to 8.9% with 2182 nodes). This is not true for all families, probably since our pruning criteria still need to be refined. The results are yet preliminary. However, this direction seems promising and it is still under study.

## 5 Discussion

In this paper we presented a new approach for modeling a group of related proteins, without incorporating any prior information about the input sequences.

<sup>5</sup>The substitution probabilities we used are those from which the blosum62 scoring matrix [Henikoff & Henikoff 1992] was derived

The method applies probabilistic suffix trees to capture some statistical properties of the input family, by recording significant occurrences of substrings of variable lengths. The method induces probability distributions over the next symbols from the empirically observed distributions. Thus, a variable length memory feature is essentially recorded in the tree structure. The PST model is well adapted to the problem at hand in terms of magnitude. Short motifs on the order of 20-30 are well integrated into a simply calibrated learning algorithm whose output is reasonably small (a tree of a few thousands nodes is typically enough for very good discriminative power, while several hundreds already do rather well).

Any new sequence is then compared to the PST model. By accumulating a prediction over the whole length of the query sequence, less conserved regions come into play. Only the balance between possibly several recognizable motifs and unknown regions determines the degree of similarity a query sequence has to the given model. The resulting tree is, as stated above, efficient in prediction, making the query stage after the initial tree building almost immediate, even in terms of comparison with the whole database (of some tens of thousands of sequences). A further speed up in prediction can always be achieved by a one time investment in the straightforward conversion of the tree into a not

much larger Probabilistic Finite Automata (PFA), of identical predictions (see [Ron et al. 1994]). The tree structure itself, when examined, can tell by its very nature, what are the different patterns significantly found in the group of sequences from which it was constructed.

As was demonstrated in section 3 the performance of the model can be improved by simply relaxing the parameters of the learning procedure. This process will eventually result in fitting the model to the training set. However, this does not affect the generalization power of the model in our case (see section 3). Yet, from the perspective of computational learning theory, it may be interesting to define a total cost function which accounts for the quality as well as for the model's complexity, and to grow the PST only until the cost is optimized. Another approach is to build a PST for both the training set and the test set independently, while comparing the induced probability distributions. When either model is (over)fitted to the noise and bias in the corresponding sample set, the distance (e.g. KL-divergence) between these distributions is expected to increase and the generalization power of the model is not expected to improve further (work in progress).

The model was applied to all protein families in the Pfam database, and the results show that the model can predict and identify the other members of the protein family with surprising success when compared with the state of the art in family modeling. The method does not assume any further biological information, but such information can be easily incorporated to improve its sensitivity. Indeed, when biological knowledge is taken into account (such as the amino acid background probabilities and the substitution probabilities) similar or better separation is achieved in some cases, using smaller trees. Further improvements should take into account the possibility of gaps, and generalization of nodes to account for more complex patterns (e.g. regular expressions), and are currently under consideration.

## 6 Acknowledgments

We thank Dana Ron for valuable discussions and Naftali Tishby for helpful suggestions. We also thank Amnon Barak for making the MOSIX parallel system available for us.

## References

- [Abe & Warmuth 1992] Abe, N. & Warmuth, M. (1992). On the computational complexity of approximating distributions by probability automata. *Machine Learning* **9**, 205-260.
- [Attwood et al. 1998] Attwood, T. K., Beck, M. E., Flower, D. R., Scordis, P. & Selley, J. (1998). The PRINTS protein fingerprint database in its fifth year. *Nucl. Acids Res.* **26**, 304-308.
- [Bairoch et al. 1997] Bairoch A., Bucher P., & Hofmann K. (1997). The PROSITE database, its status in 1997. *Nucl. Acids Res.* **25**, 217-221.
- [Bairoch & Boeckman 1992] Bairoch, A. & Boeckman, B. (1992). The SWISS-PROT protein sequence data bank. *Nucl. Acids Res.* **20**, 2019-2022.
- [Gillman & Sipser 1994] Gillman, D. & Sipser, M. (1994). Inference and minimization of hidden Markov chains. In *Proc. of the Seventh Annual Workshop on Computational Learning Theory* pp 147-158.
- [Gribskov et al. 1987] Gribskov, M., Mclachlen, A. D. & Eisenberg, D. (1987). Profile analysis: detection of distantly related proteins. *Proc. Natl. Acad. Sci. USA* **84**, 4355-4358.
- [Hanke et al. 1996] Hanke, J., Beckmann, G., Bork, P. & Reich, J. G. (1996). Self-organizing hierarchic networks for pattern recognition in protein sequence. *Protein Sci.* **5**, 72-82.
- [Henikoff & Henikoff 1991] Henikoff, S. & Henikoff, J. G. (1991). Automated assembly of protein blocks for database searching. *Nucl. Acids Res.* **19**, 6565-6572.
- [Henikoff & Henikoff 1992] Henikoff, S. & Henikoff, J. G. (1992). Amino acid substitution matrices from protein blocks. *Proc. Natl. Acad. Sci. USA* **89**, 10915-10919.
- [Henikoff & Henikoff 1996] Henikoff, J. G. & Henikoff, S. (1996). Using substitution probabilities to improve position-specific scoring matrices. *Comp. App. Biosci.* **12:2**, 135-143.
- [Jonassen et al. 1995] Jonassen, I., Collins, J. F. & Higgins, D. G. (1995). Finding flexible patterns in unaligned protein sequences. *Protein Sci.* **4**, 1587-1595.
- [Krogh et al. 1996] Krogh, A., Brown, M., Mian, I. S., Sjölander, K. & Haussler, D. (1996). Hidden Markov models in computational biology: Application to protein modeling. *J. Mol. Biol.* **235**, 1501-1531.
- [Neuwald et al. 1994] Neuwald, A. F. & Green, P. (1994). Detecting patterns in protein sequences. *J. Mol. Biol.* **239**, 698-712.
- [Pearson 1995] Pearson, W. R. (1995). Comparison of methods for searching protein sequence databases. *Protein Sci.* **4**, 1145-1160.
- [Ron et al. 1994] Ron, D., Singer, Y. & Tishby, N. (1996). The power of amnesia: learning probabilistic automata with variable memory length. *Machine Learning* **25** 117-150.
- [Sonnhammer & Kahn 1994] Sonnhammer, E. L. L. & Kahn, D. (1994). Modular arrangement of proteins as inferred from analysis of homology. *Protein Sci.* **3**, 482-492.
- [Sonnhammer et al. 1998] Sonnhammer, E. L., Eddy, S.R., Birney, E., Bateman, A., & Durbin, R. (1998). Pfam: multiple sequence alignments and HMM-profiles of protein domains. *Nucl. Acids Res.* **26**, 320-322.
- [Suyama et al. 1995] Suyama, M., Nishioka, T. & Oda, J. (1995). Searching for common sequence patterns among distantly related proteins. *Protein Eng.* **8:11**, 1075-1080.

## Appendix A - PST performance for all Pfam families

(See next two pages)



Family	Size	No. sequences missed	% True Positives detected
ig	884	51	94.2%
pkinase	725	108	85.1%
globin	681	15	97.8%
7tm_1	515	36	93.0%
homeobox	383	27	93.0%
efhand	320	25	92.2%
RuBisCO_large	311	4	98.7%
zf-C2H2	297	23	92.3%
ABC_tran	269	44	83.6%
trypsin	246	22	91.1%
ras	213	8	96.2%
p450	204	17	91.7%
tubulin	196	1	99.5%
GTP_EFTU	184	15	91.8%
ATP-synt_ab	180	6	96.7%
adh_short	180	20	88.9%
histone	178	6	96.6%
cytochrome_c	175	11	93.7%
toxin	172	4	97.7%
cytochrome_b_N	170	3	98.2%
EGF	169	18	89.3%
HSP70	163	7	95.7%
fn3	161	23	85.7%
fer4	152	18	88.2%
MHC_I	151	3	98.0%
zn-protease	148	21	85.8%
rvt	147	17	88.4%
gluts	144	14	90.3%
actin	142	4	97.2%
rrm	141	22	84.4%
filament	139	5	96.4%
zf-C4	139	6	95.7%
neur_chan	138	4	97.1%
SH3	137	16	88.3%
HLH	133	7	94.7%
ins	132	3	97.7%
cytochrome_b_C	130	27	79.2%
HSP20	129	7	94.6%
adh_zinc	129	6	95.3%
SH2	128	5	96.1%
response_reg	128	19	85.2%
hormone_rec	127	7	94.5%
phoslip	122	3	97.5%
gpdh	117	3	97.4%
lipocalin	115	7	93.9%
alpha-amylase	114	14	87.7%

Table 1: PST performance for all Pfam families (part1). The names of the families are abbreviated as in the Pfam database. The number of proteins in the family is given in the second column. Each family was divided into a training set and a test set and the PST was build from the training set. To test the quality of this PST, we calculate the probability that the PST induces on each database sequence, and each family sequence (from the train set and the test set) whose probability is above the iso-point is considered successfully detected (see text for more details). The total percentage of true positives detected with the PST model is 90.7%, when averaged over all Pfam families with at least 10 members (170 families). The performance of the PST model is expected to improve if all sequences of a family are to be included in the training set (as Pfam database is created). The set of parameters used to train the PST is  $P_{min} = 0.0001$   $\alpha = 0$   $\gamma_{min} = 0.001$   $r = 1.050$  and  $L = 20$ . Additional improvement in the performance is expected if the parameters are tuned for each family (see text). To train a PST on a typical family with this set of parameters it takes about two hours, at the most, on a pentium II 266 Mhz. Additional 5 minutes are needed to predict all sequences in the SWISSPROT database. For comparison, searching the SWISSPROT database with a typical HMM may take about two hours.

Family	Size	No. sequences missed	% True Positives detected
hormone	111	4	96.4%
kazal	110	6	94.5%
COX2	109	2	98.2%
sugar_tr	107	15	86.0%
lectin_c	106	14	86.8%
zf-CCHC	105	12	88.6%
E1-E2.ATPase	102	7	93.1%
wnt	102	6	94.1%
HTH_1	101	16	84.2%
oxidored_fad	101	12	88.1%
RuBisCO_small	99	3	97.0%
serpin	98	9	90.8%
bZIP	95	10	89.5%
Y_phosphatase	92	8	91.3%
Cys-protease	91	11	87.9%
ldh	90	6	93.3%
fer2	88	5	94.3%
rnaseH	87	12	86.2%
STphosphatase	86	5	94.2%
cpn60	84	5	94.0%
ank	83	10	88.0%
rvp	82	12	85.4%
subtilase	82	9	89.0%
COX1	80	13	83.8%
cyclin	80	9	88.8%
ATP-synt_A	79	6	92.4%
TGF-beta	79	6	92.4%
C2	78	6	92.3%
gln-synt	78	5	93.6%
thiore	76	11	85.5%
PH	75	5	93.3%
oxidored_nitro	75	8	89.3%
sushi	75	8	89.3%
hormone2	73	2	97.3%
photoRC	73	1	98.6%
asp	72	12	83.3%
lys	72	1	98.6%
recA	72	3	95.8%
rnaseA	71	1	98.6%
GATase	69	8	88.4%
aidedh	69	9	87.0%
sodfe	69	5	92.8%
zf-C3HC4	69	10	85.5%
DAG_PE-bind	68	7	89.7%
il8	67	4	94.0%
AAA	66	8	87.9%
sodcu	66	5	92.4%
HTH_2	63	9	85.7%
aminotran	63	7	88.9%
ATP-synt_C	62	5	91.9%
Cys_knot	61	4	93.4%
copper-bind	61	3	95.1%
mito_carr	61	7	88.5%
sigma70	61	5	91.8%
COesterase	60	5	91.7%
S12	60	2	96.7%
NADHdh	57	4	93.0%
cpn10	57	4	93.0%
metalthio	56	0	100.0%
pilin	56	6	89.3%
sigma54	56	9	83.9%
Kunitz_BPTI	55	5	90.9%
heme_1	55	4	92.7%
neur	55	2	96.4%
peroxidase	55	7	87.3%
S4	54	4	92.6%

Table 1: PST performance for all Pfam families (part 2)

Table 1...continued

Family	Size	No. sequences missed	% True Positives detected
Zn_clus	54	10	81.5%
crystall	53	1	98.1%
cystatin	53	4	92.5%
hormone3	53	5	90.6%
PGK	51	3	94.1%
beta-lactamase	51	7	86.3%
tsp_1	51	6	88.2%
pro_isomerase	50	3	94.0%
fer4_NifH	49	2	95.9%
DNA_methylase	48	8	83.3%
interferon	47	2	95.7%
pou	47	2	95.7%
DNA_pol	46	9	80.4%
Pribosyltran	45	5	88.9%
hexapep	45	8	82.2%
myosin_head	44	10	77.3%
arf	43	4	90.7%
lectin_legA	43	3	93.0%
pyr_redox	43	7	83.7%
cNMP_binding	42	3	92.9%
TIM	40	3	92.5%
cellulase	40	6	85.0%
connexin	40	1	97.5%
enolase	40	0	100.0%
rhv	40	2	95.0%
FGF	39	1	97.4%
ketoacyl-synt	38	7	81.6%
kringle	38	2	94.7%
lectin_legB	38	7	81.6%
RIP	37	2	94.6%
7tm_2	36	2	94.4%
KH-domain	36	4	88.9%
oxidored_molyb	35	1	97.1%
tRNA-synt_1	35	7	80.0%
thyroglobulin_1	32	3	90.6%
cadherin	31	4	87.1%
hemopexin	31	3	90.3%
ldl_recept_a	31	5	83.9%
TNFR_c6	29	4	86.2%
lig_chan	29	1	96.6%
tRNA-synt_2	29	5	82.8%
vwa	29	6	79.3%
zona_pellucida	26	3	88.5%
thiolase	25	3	88.0%
MCPsignal	24	4	83.3%
lipase	23	3	87.0%
vvc	23	6	73.9%
fn2	20	2	90.0%
trefoil	20	3	85.0%
laminin_G	19	2	89.5%
fibrinogen_C	18	4	77.8%
laminin_EGF	16	3	81.2%
fn1	15	2	86.7%
UPAR_LY6	14	2	85.7%
dsrm	14	2	85.7%
ldl_recept_b	14	1	92.9%
wap	13	2	84.6%
7tm_3	12	2	83.3%

Table 1: PST performance for all Pfam families (part 3)