

Modifying Kernels Using Label Information Improves Protein Classification Performance

Renqiang Min¹, Anthony Bonner¹, and Zhaolei Zhang²

¹ Department of Computer Science, University of Toronto, 10 King’s College Road, Toronto, ON M5S3G4, Canada

² Banting and Best Department of Medical Research, University of Toronto, 112 College Street, Toronto, ON M5G1L6, Canada

Abstract. Kernel learning methods based on kernel alignment with semidefinite programming (SDP) are often memory intensive and computationally expensive, thus often impractical for problems with large-size dataset. We propose a method using label information to scale the training part of a given kernel matrix to form the training part of a new kernel matrix. The test part of the new kernel matrix is estimated based on a linear transformation in a reduced feature space and can be calculated computationally efficiently. As a result, the new kernel matrix reflects the label-dependent separability of the sequence data in a better way than the original kernel matrix. In addition, our experimental results on a benchmark dataset, the SCOP dataset, show that the SVM classifier based on the improved kernels has better performance than the SVM classifier based on the original kernels; moreover, SVM based on the improved Profile kernel with pull-in homologs (see experiment section for explanations) produced the best results for remote homology detection on the SCOP dataset compared to the published results.

1 Introduction

Protein sequence classification and fold recognition is still a challenging task in the bioinformatics research community. Generative models (e.g., profile HMMs [4], [5]) and discriminative models (e.g., kernel SVMs [1], [2], [6]) have been applied to solve this problem. Since protein sequence data is a string of letters, it is relatively straightforward to apply HMMs directly to it. But for other machine learning methods that take numerical inputs, analyzing the sequence data by them is not that easy. Owing to the emergence of kernel methods, we can solve this problem indirectly. We can map the letter sequences to a higher dimensional numerical space called feature space, in which each sequence x is mapped to a vector $\Phi(x)$. In the feature space, we can apply many machine learning methods to analyze the sequence data while doing computations using the kernel trick. For a wide range of applications including protein classification, the components of the constructed kernels are positive. In our method, we require the kernel entries be positive.

It has been shown that the kernel SVM method has better classification and prediction performance on protein sequence data than some other methods (see

[1], [2], and [6]). Several kernels such as pairwise-sequence-similarity-score based kernels and mismatch-string kernels, which are especially suitable for protein sequence data that consists of a limited number of letters from the amino acid alphabet, are frequently used in sequence classification and structure prediction. However, the label information of labeled sequences (the class membership of the data points; in a binary classification problem, the label of a data point is 1 or 0) is completely or partly ignored in the construction process of these kernels, and the available information between pairwise unlabeled sequences is also ignored both in the training phase and in the testing phase. In this paper, we will incorporate the label information of training data into the construction process of a new kernel, hoping that the obtained kernel reflects the real neighborhood property of the data in a better way than the original kernel. We believe this helps classification in most situations.

In the paper, we will use two mismatch-string kernels as base kernels [9]. Then we modify the base kernels using label information of training data. We briefly describe SVM classifier based on mismatch-string kernels in Section 2. In Section 3, we discuss some recent related methods that motivated us to improve kernels using label information. In Section 4, we describe in details our approach to improve kernels using label information based on Singular Value Decomposition (SVD) and a linear mapping. We present our experimental results for protein homology detection on the SCOP dataset in Section 5. And in Section 6, we conclude the paper with some discussions and proposed directions for future research.

2 SVM based on mismatch-string kernel

SVM is a discriminative method proposed for classification. Suppose we have a two-class dataset $\{x_i, y_i\}, i = 1, \dots, m, y_i \in \{-1, 1\}, x_i \in R^n$. A linear SVM gives a separating hyperplane that maximizes the margin between the sample data points of the two classes, which is equivalent to minimizing the following objective function:

$$L(w) = \frac{1}{2} \|w\|^2 + C (\sum_i \xi_i) \quad (1)$$

$$x_i w + b \geq +1 - \xi_i \quad \text{for } y_i = +1 \quad (2)$$

$$x_i w + b \leq -1 + \xi_i \quad \text{for } y_i = -1 \quad (3)$$

$$\xi_i \geq 0 \quad \forall i \in \{1, \dots, m\} \quad (4)$$

Where C is a penalty coefficient, and the ξ_i are non-negative slack variables, which are set to 0 when the dataset is separable.

By constructing a kernel, K , we can map every data point, x_i , to a high-dimensional feature space, in which an SVM can be used to generate a separating hyperplane. However, by transforming equation (1) and its constraints to inner-product form using Lagrange multipliers, all calculations can be done in low-dimensional space by using the kernel trick.

As discussed earlier, kernels can map sequences consisting of letters to a high-dimensional numerical space. For example, suppose that A is an alphabet with ℓ symbols ($\ell = 20$ for protein sequences). A k -mer string kernel maps every sequence in A to a ℓ^k -dimensional feature space in which coordinates are indexed by all possible sub-sequences of length k (k -mers). The specific feature map is

$$\Phi_k(x) = (\Phi_{\alpha_1}(x), \Phi_{\alpha_2}(x), \dots, \Phi_{\alpha_{\ell^k}}(x))^T \quad (5)$$

where $\Phi_\alpha(x)$ is the number of occurrences of k -mer α in sequence x . The corresponding kernel matrix is

$$K_k(x, y) = \Phi_k(x)^T \Phi_k(y) \quad (6)$$

The mismatch string kernel extends this idea by taking into account mismatches when counting the number of occurrences of a k -mer in an input sequence. In particular, for any k -mer, α , let $N_{(\alpha, m)}$ be the set of all k -mers that differ from α by at most m mismatches. The kernel mapping and kernel matrix are then defined as follows:

$$\Phi_{(k, m)}(x) = (\Phi_{(k, m), \alpha_1}(x), \dots, \Phi_{(k, m), \alpha_{\ell^k}}(x))^T \quad (7)$$

$$\Phi_{(k, m), \alpha}(x) = \sum_{\beta \in N_{(\alpha, m)}(x)} \Phi_\beta(x) \quad (8)$$

$$K_{(k, m)}(x, y) = \Phi_{(k, m)}(x)^T \Phi_{(k, m)}(y) \quad (9)$$

A Profile Kernel [3] extends the above mismatch-string kernel by using additional profile information of each sequence, that is, the emission probability of every amino acid at each position in respective sequences. Instead of treating all k -mers with less than m mismatches the same like the above mismatch-string kernel, the profile-kernel examines these k -mers further by looking at the emission probabilities at the mismatch positions and only accepts some mismatches by thresholding. Suppose we have a sequence $x = x_1x_2\dots x_N$ composed of amino acids with alphabet Σ and N is the length of the sequence, $P(x) = \{p_i(a), a \in \Sigma\}_{i=1}^N$ is a profile for sequence x , where $p_i(a)$ denotes the emission probability of amino acid a in position i and $\sum_{a \in \Sigma} p_i(a) = 1$ for each position i . In the Profile Kernel, the neighborhood for a k -mer $x[j+1 : j+k] = x_{j+1}x_{j+2}\dots x_{j+k}$ in x ($0 \leq j \leq |x| - k$) is:

$$M_{(k, \sigma)}(P(x[j+1 : j+k])) = \{\beta = b_1b_2\dots b_k : -\sum_{i=1}^k \log p_{j+i}(b_i) < \sigma\}. \quad (10)$$

where $p_{j+i}(b)$ with $i = 1, \dots, k$ comes from the profile of sequence x and it can be smoothed using the background frequency of amino acid b . And the feature vector of sequence x in the Profile Kernel is defined as the following:

$$\Phi_{(k, \sigma)}(x) = \sum_{j=0 \dots |x|-k} (\phi_\beta(P(x[j+1 : j+k])))_{\beta \in \Sigma^k} \quad (11)$$

where the coordinate $\phi_\beta(P(x[j + 1 : j + k]))$ equals 1 if $\beta \in M_{(k,\sigma)}(P(x[j + 1 : j + k]))$, and 0 otherwise. Note that all entries in these kernel matrices are non-negative.

As described in [2] and [3], given a set of labelled and unlabelled protein sequences, mismatch-string kernels can be efficiently computed. A SVM classifier can be trained using the kernel entries for pairwise labelled sequences, and the classifier can then be used to predict the remote homology of unlabelled sequences.

3 Related methods

A kernel matrix K with $K(i, j) = \Phi(i)^T \Phi(j)$ can be used to derive a similarity matrix based on square Euclidean distances between any pairwise data points, i and j , in the feature space, as follows:

$$Dist^2(i, j) = K(i, i) + K(j, j) - 2K(i, j) \quad (12)$$

Although SVM generates the optimal separating hyperplane in the feature space given a specific kernel, it does not adjust the given kernel and make it more discriminative. Therefore, it leaves room for improvement as we can apply the aforementioned idea of preserving neighbor identity to construct new kernels in order to achieve better separability in the new feature space. Instead of computing more discriminative features of data points explicitly, we can construct a more discriminative kernel directly and all the computations needed by training and classification can be cast onto the new kernel matrix. As discussed in [10] and [11], a linear combination of some predefined kernels is used to generate new kernels, and the mixing coefficients are calculated by aligning the training part of the combined kernel to the training part of an optimal kernel K as follows:

$$K = \begin{bmatrix} K_{tr} & K_{tt}^T \\ K_{tt} & \text{unused} \end{bmatrix} \quad (13)$$

$$K_{tr}(i, j) = \begin{cases} +1 & \text{if } i \text{ and } j \text{ have the same label} \\ -1 & \text{otherwise} \end{cases} \quad (14)$$

where i and j index data points in the training set, and tr and tt respectively denote the training part and the test part (this rule applies to all the denotations in the paper). If there are n training data points and m test data points, K_{tr} is an n -by- n block sub-matrix and K_{tt} is an m -by- n block sub-matrix in K . In fact, doing kernel alignment is to make the constructed kernel approximate the neighbor identity and data separability reflected by the optimal kernel. From Equation (10), we can easily find that the optimal kernel makes the distances between pairwise data points having the same label be 0 and the distances between pairwise data points having the different labels be 2. That is to say, the kernel alignment algorithms actually use the label information to construct a new kernel to achieve good data potability. However, doing the alignment to

calculate the mixing coefficients costs a lot of memory and is very computationally expensive or impossible for handling large datasets for combining many kernels. In [8], an efficient approach to learning a convex combination of a set of kernels was proposed. In this paper, we propose another efficient approach for constructing new kernels using label information, which is based on scaling, matrix decomposition, and a linear mapping, to achieve better data separability as discussed above. The approach is easy to implement and easy to extend to many types of kernels.

4 Improved kernels using label information

Suppose that we have a dataset as described in Section 2 (we only consider the two-class problem here) and a given mapping from the input data space to a high dimensional feature space. We can then construct a kernel K based on the mapping.

Given the constructed kernel K with $K(i, j) = \Phi(i)^T \Phi(j)$ and the label information of training data, we want a new kernel that better reflects the neighbor identity and separability of the data consistent to the current labels of training data. If two arbitrary data points in the training set, i and j , have the same label, we multiply the inner product of their feature vectors by a scaling factor, α (see Section 5 for detailed discussion about choosing α), which is greater than 1, to get a new kernel matrix \hat{K} as follows:

$$\hat{K} = \begin{bmatrix} \hat{K}_{tr} & \hat{K}_{tt}^T \\ \hat{K}_{tt} & \text{unused} \end{bmatrix} \quad (15)$$

where

$$\hat{K}_{tr}(i, j) = \begin{cases} \alpha K_{tr}(i, j) & \text{if } i \text{ and } j \text{ have the same label or } i = j \\ K_{tr}(i, j) & \text{otherwise} \end{cases} \quad (16)$$

label information of the training set to modify the training part of the kernel matrix. This modification will affect both the training part and the test part of K . The test part \hat{K}_{tt} of the new kernel matrix \hat{K} is calculated using kernel extrapolation, which is based on a linear mapping. The matrix \hat{K}_{tr} is positive semidefinite. The distances between pairwise data points in the new feature space corresponding to \hat{K}_{tr} are as follows:

$$Dist_{tr}^2(i, j) = \begin{cases} \alpha Dist^2(i, j) & \text{if } i \text{ and } j \text{ have the same label} \\ \alpha Dist^2(i, j) + 2(\alpha - 1)K_{tr}(i, j) & \text{otherwise.} \end{cases} \quad (17)$$

Here $Dist^2(i, j)$ is defined in Equation (12), and i and j index the data points in the training set. We see from Equation (17) that, in the new feature space, the distance between points having the same label is increased by a factor of α . Moreover, since $\alpha > 1$ and $K_{tr}(i, j)$ is non-negative, the distance between points having different labels is increased even further by the additional term

$2(\alpha - 1)K_{tr}(i, j)$. That is to say, in the feature space defined by \hat{K}_{tr} , data points with the same label stay relatively close together, while data points with different labels move relatively further apart. Figure 1 illustrates the separation of data points in feature spaces corresponding to K_{tr} and \hat{K}_{tr} .

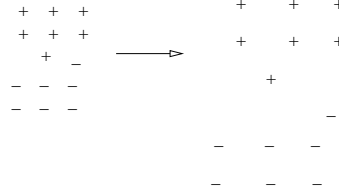


Fig. 1. The data distribution in the original feature space and in the new feature space. '+' means positive and '-' means negative.

We can also interpret the similarity between a pair of data points, i and j , in terms of the angle between their feature vectors, $\Phi(i)$ and $\Phi(j)$, which is given by $\theta = \arccos[K(i, j)/\sqrt{K(i, i)K(j, j)}]$. The angle between two points with the same label is the same in the new feature space and the original feature space, while the angle between two points with different labels is larger in the new feature space than in the original feature space. This can also be seen in Figure 1.

After the training part, \hat{K}_{tr} , of the new kernel is constructed, we need to estimate the testing part, \hat{K}_{tt} . That is, to classify a test case by an SVM based on \hat{K} , we need to estimate the inner products of the new feature vector of the test case and the new feature vectors of all the training cases. This can be done by approximating all the high-dimensional feature vectors by lower, N -dimensional feature vectors, where N is the size of the training set. To do this, we decompose the training part of K and \hat{K} , denoted by K_{tr} and \hat{K}_{tr} , respectively, into SVD form as follows:

$$K_{tr}V_{tr} = V_{tr}D_{tr} \quad (18)$$

$$K_{tr} = V_{tr}D_{tr}V_{tr}^T = W_{tr}^TW_{tr} \quad (19)$$

where

$$W_{tr} = (V_{tr}\sqrt{D_{tr}})^T \quad (20)$$

Similarly,

$$\hat{K}_{tr}\hat{V}_{tr} = \hat{V}_{tr}\hat{D}_{tr} \quad (21)$$

$$\hat{K}_{tr} = \hat{V}_{tr}\hat{D}_{tr}\hat{V}_{tr}^T = \hat{W}_{tr}^T\hat{W}_{tr} \quad (22)$$

where

$$\hat{W}_{tr} = (\hat{V}_{tr}\sqrt{\hat{D}_{tr}})^T \quad (23)$$

In these equations, the columns of V_{tr} are orthogonal eigenvectors of K_{tr}^2 , and D_{tr} is a diagonal matrix containing the corresponding eigenvalues. Likewise for

\hat{V}_{tr} , \hat{K}_{tr} and \hat{D}_{tr} . W_{tr} and \hat{W}_{tr} are n -by- n matrices, where n is the size of the training set.

We can view W_{tr} as a compressed representation of the high-dimensional feature vectors of the training data in a lower dimensional space. Note that this representation preserves all the inner products. We can interpret \hat{W}_{tr} in the same way. Moreover, the new kernel, \hat{K}_{tr} , can be calculated from \hat{W}_{tr} , which in turn can be computed by applying a linear transformation to W_{tr} , as the following lemma shows:

Lemma 1 $AW_{tr} = \hat{W}_{tr}$, where $A = \sqrt{\hat{D}_{tr}}\hat{V}_{tr}^T V_{tr} \frac{1}{\sqrt{D_{tr}}}$

Here, the expression $\frac{1}{\sqrt{D}}$ means the inverse of the diagonal matrix \sqrt{D} . The lemma itself follows immediately from equations (20) and (23). We interpret this lemma as follows: K_{tr} and \hat{K}_{tr} , respectively, corresponds to feature space F and \hat{F} with W_{tr} lying in F and \hat{W}_{tr} lying in \hat{F} ; there exists a linear transformation between F and \hat{F} . We shall use the linear transformation, A , to estimate the matrix \hat{K}_{tt} , the testing part of \hat{K} . This involves the following assumption:

Assumption 1 *The linear relation shown in Lemma 1 can be extended to $A[W_{tr}; W_{tt}] = [\hat{W}_{tr}; \hat{W}_{tt}]$, where W_{tt} and \hat{W}_{tt} are m -by- n matrices which satisfy $W_{tt}^T W_{tr} = K_{tt}$ and $\hat{W}_{tt}^T \hat{W}_{tr} = \hat{K}_{tt}$, n and m are respectively the size of the training set and the test set.*

In this assumption, we assume that: W_{tt} lies in F and \hat{W}_{tt} lies in \hat{F} ; applying the linear transformation A to W_{tt} will result in the n -dimensional feature vectors of test data \hat{W}_{tt} in the reduced feature space \hat{F} , which better reflects the label-dependent separability of the test data points as A does to W_{tr} . The value of this assumption is tested empirically in Section 5, where we show that the resulting kernel leads to an SVM classifier with significantly improved performance.

Note that W_{tt} and \hat{W}_{tt} are N -dimensional feature vectors representing the test data points³. Using Lemmas 1 and Assumption 1, we can calculate \hat{K}_{tt} . First, from the definitions of K and W ,

$$W_{tt}^T W_{tr} = K_{tt} \quad (24)$$

and so by equation (20),

$$W_{tt} = \frac{1}{\sqrt{D_{tr}}} V_{tr}^T K_{tt} \quad (25)$$

According to Assumption 1, we have

$$\hat{K}_{tt} = \hat{W}_{tt}^T \hat{W}_{tr} = K_{tt} (V_{tr} \frac{1}{\sqrt{D_{tr}}} V_{tr}^T) \hat{K}_{tr} = K_{tt} K_{tr}^{-1} \hat{K}_{tr} \quad (26)$$

³ Here, we should note that, unlike transductive learning methods, we need not know all the test data in advance, the test data might come one by one, and we denote the N -dimensional feature vectors of all the test data by one symbol for description convenience

When calculating \hat{K}_{tt} , we need to calculate K_{tr}^{-1} first, and then we can obtain \hat{K}_{tt} easily by Equation (26). Note that we need not perform SVD at all and the inverse of K_{tr} can be computed in Matlab very fast (it takes less than 10 seconds to get the inverse of a 2620-by-2620 kernel matrix in our machine with 3.0GHz CPU and 4.0GB memory)⁴. After the new kernel is constructed, we can apply machine learning techniques based on the kernel to classification, clustering, or regression problems. In the next section, we use SVM classifiers based on the new kernel to classify proteins.

5 Experiments on remote protein homology detection

We determine the classification performance of the new kernels against the original kernels by comparing their ability to detect protein remote homology. A benchmark dataset, which was derived by Jaakkola from the SCOP database (see [7] and [1]), is used here. In SCOP, protein sequences are classified into a 4-level hierarchy: class, fold, superfamily, and family, starting from the top. Remote homology is simulated by choosing all the members of a family as positive test data, some family (or families) in the same superfamily of the test data as positive training data, all sequences outside the fold of the test data as either negative training data or negative test data, and sequences that are neither in the training set nor in the test set are considered as unlabelled data. This data splitting scheme has been used in several previous papers (see [1], [6], and [9]). We used the same training and test data split as those used in [6] and [9]. The version 1.59 of the SCOP dataset from <http://astral.berkeley.edu> is used, in which no pair of sequences share more than 95% identity. The detailed explanation about the experimental setting can be found in <http://www.kyb.tuebingen.mpg.de/bs/people/weston/semiprot/supp.html>.

In the experiments, there are 54 target test families altogether classified into four classes: alpha proteins (9 families), beta proteins (18 families), alpha and beta proteins (17 families), and small proteins (10 families). In the data splits, for most experiments, there are only several positive test cases but hundreds or even thousands of negative test cases. The maximum number of positive test cases is below 30, but the maximum number of negative test cases is above 2600. The minimum number of positive test case is 1, but the minimum number of negative test cases is still above 250. So, in the experiments with a very limited number of positive test cases and a large number of negative test cases, we can almost ignore the ranking of positive cases below 50 negative cases. In such situations, we consider that the ROC₅₀ score is much more important than the ROC score. Here, a ROC curve plots the rate of true positives as a function of the rate of false positives at different decision thresholds, the ROC score is the area under the curve, and the ROC₅₀ score is the ROC score computed up to the first 50

⁴ Equation (26) requires K_{tr} is non-singular, and if it is singular, it means that some rows in K_{tr} corresponding to some training data points can be expressed as the linear combination of some other rows in K_{tr} , we can simply remove the redundant rows to get a non-singular K_{tr} or set K_{tr} to be $K_{tr} + \epsilon I$.

	Alpha Proteins	Beta Proteins	Alpha and Beta Proteins	Small Proteins	Overall Mean ROC ₅₀
K1	0.4874	0.5208	0.5798	0.5905	0.5448
ImproK1	0.5395	0.5283	0.5933	0.6010	0.5630
K2	0.7909	0.8156	0.8924	0.8687	0.8441
ImproK2	0.8172	0.8276	0.9075	0.8808	0.8597

Table 1. “K1” represents “Mismatch kernel + [PSI-BLAST]”, “ImproK1” represents “Improved Mismatch kernel + [PSI-BLAST]”, “K2” represents “Profile Kernel + [PSI-BLAST]”, and “ImproK2” represents “Improved Profile Kernel + [PSI-BLAST]”. The number in the bracket denotes the number of families in each class.

false positives. Thus, in our experiments, we only compare the ROC₅₀ scores corresponding to different kernels.

Because our approach to generating new kernels based on label information is independent of given kernels, we choose two representative kernels, which were, respectively, “Mismatch kernel + homologs [PSI-BLAST]” as described in [9] and “Profile kernel” as described in [3] also “plus homologs [PSI-BLAST]” as base kernels. “kernels + homologs [PSI-BLAST]” refers to a semi-supervised learning method: prior to training SVM, close homologs of the training data in the unlabelled set found by PSI-BLAST with E-value less than 0.05 are added to the positive training set, and are labelled as positive (we call this “pull-in homologs”). We choose the first kernel because it gives the best results on remote homology detection among the kernels that don’t use the profile information; and we choose the second kernel because it produced the best results on SCOP among all the kernels (we don’t consider transductive learning in this paper). To perform SVM classification based on the kernels, we used the SVM classifier in the freely available Spider Matlab machine learning package.

We compared the methods using the mismatch kernel with $k = 5$ and $m = 1$ and the profile kernel with $k = 5$ and $\sigma = 7.5$, and the α is set by Cross Validation (CV). In the experiments in which the number of positive training cases is greater than or equal to 5, we respectively generated a random permutation of the positive training cases and of the negative training cases, then we divided the two permutations into 5 folds denoted by P_i and N_i , $i = 1, \dots, 5$. We form a new set $M = \{\{P_i, N_i\} | i = 1, \dots, 5\}$, then we did 5-fold CV on M and chose α corresponding to the biggest mean ROC₅₀ score from a pre-defined list. In the experiments in which the number of the positive training cases is less than 5, we used a similar strategy as above but we divided the positive training set and the negative training set into 2 folds, and we did 2-fold CV on the newly formed set M to choose α . In the experiments, the free parameters C for SVM and the free parameter α are chose using Cross Validation as discussed above. Before training SVM, the kernel was normalized using $K(i, j) \leftarrow \frac{K(i, j)}{\sqrt{K(i, i)K(j, j)}}$.

Table 1 gives the mean ROC₅₀ scores on different protein classes in several classes corresponding to the original kernels and the modified kernels. From Table

2, we see that: modified kernels using label information gave better performance than the original kernels.

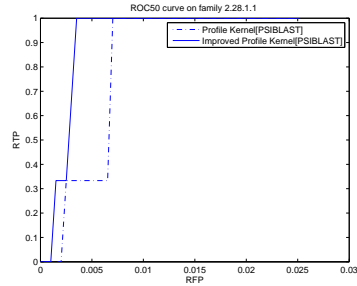


Fig. 2. Comparison of the ROC₅₀ curves on family 2.28.1.1 (Legume lectins).

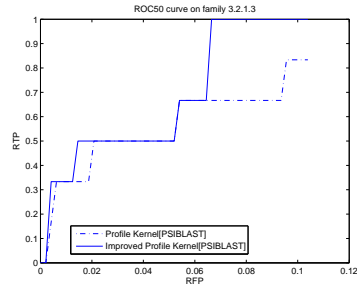


Fig. 3. Comparison of the ROC₅₀ curves on family 2.28.1.1 (Legume lectins).

To determine whether the improvement given by the modified kernels is statistically significant, we performed a Wilcoxon Matched-Pairs Signed-Ranks Test on the differences. The resulting p-value for the improvement over the Mismatch+homologs [PSI-BLAST] kernel is $2.19e - 04$, and the p-value for the improvement over the Profile+homologs [PSI-BLAST] kernel is 0.0162.

To show our algorithm improves the original kernels in more detail, we plot some ROC₅₀ curves in Figure 2 and Figure 3. From the two figures, we can see that the improved kernels have better performance than the original kernels. In Figure 4 and Figure 5, we respectively plot a block sub-matrix of the test part of the normalized original Profile + [PSIBLAST] and of the normalized improved Profile + [PSIBLAST] matrix on family 2.28.1.1 (Legume lectins). In the two figures, the first three rows correspond to all the positive test sequences in the test set, and the remainder rows correspond to some randomly selected negative test sequences. The first nine columns correspond to some randomly selected

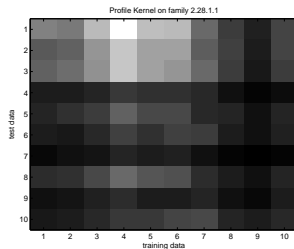


Fig. 4. A block sub-matrix in K_{tt} of the original Profile Kernel[PSIBLAST] on family 2.28.1.1 (Legume lectins). See the text for explanation.

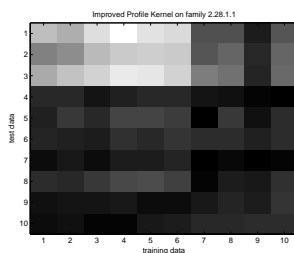


Fig. 5. A block sub-matrix in \hat{K}_{tt} of the improved Profile Kernel[PSIBLAST] on family 2.28.1.1 (Legume lectins). See the text for explanation.

positive training sequences, and the last column corresponds to a randomly selected negative training sequence. The whiter the blobs in the figures, the larger the corresponding similarity scores. Comparing Figure 4 to Figure 5, we can see that the similarity scores between the positive test data and the positive training data in the improved kernel is increased (the block on the upper left corner becomes whiter in the improved kernel matrix).

6 Discussion and future work

We described an approach to modify kernels using label information of training data based on SVD and a linear mapping. The modified kernel is more discriminative than the original kernel. We also showed that, unlike Kernel Alignment with SDP, the test part of the modified kernel can be calculated very efficiently in practice. We tested the performance of the modified kernel by detecting protein remote homology. Experimental results show that the improvement given by the new kernel is statistically significant, although one more free parameter α is introduced. In our approach, both the scaling factor α and the free parameter C of SVM are chosen by Cross Validation (CV). The CV procedure we used is very stable. Even when we run the CV procedure several times on each experiment, we will get the same scaling factor α on each experiment each time.

We believe that the modified kernel will not overfit the training data, because the label information is only used to modify the training part of kernel matrix and the degree of the modification is controlled by CV. The experimental results in the paper show that the generalization is good. The approach discussed in the paper is general and can be readily applied to many problems. In the future work, we plan to learn a non-linear mapping from an original reduced feature space to a new feature space using neural networks instead of using a linear mapping.

Acknowledgment

This project was funded by a start-up fund from University of Toronto to Zhaolei Zhang, an NSERC grant to Anthony Bonner, and a grant from Genome Canada through the Ontario Genomics Institute.

References

1. Jaakkola, T., Diekhans, M., and Haussler, D.: A discriminative framework for detecting remote protein homologies. *Journal of Computational Biology*. **7** (2000) Numbers 1/2, 95-114
2. Leslie, C., Eskin, E., Weston, J., and Noble, W.S.: Mismatch string kernels for SVM protein classification. *Neural Information Processing Systems*. **15** (2002)
3. Kuang, R., Ie, E., Wang, K., Wang, K., Siddiqi, M., Freund, Y., and Leslie C.: Profile-based String Kernels for Remote Homology Detection and Motif Extraction. *Journal of Bioinformatics and Computational Biology*. **3** (2005) No. 3 527-550
4. Krogh, A., Brown, M., Mian, I., Sjolander, K., and Haussler, D.: Hidden markov models in computational biology: Applications to protein modeling. *Journal of Molecular Biology*. **235** (1994) 1501-1531.
5. Baldi, P., Chauvin, Y., Hunkapiller, T., and McClure, M.A.: Hidden markov models of biological primary sequence information. *PNAS*, **91**(3) (1994) 1059-1063.
6. Liao, C. and Noble, W.S.: Combining pairwise sequence similarity and support vector machines for remote protein homology detection. *Proceedings of RECOMB*. (2002)
7. Murzin A. G., Brenner S. E., Hubbard T., Chothia C.: SCOP: a structural classification of proteins database for the investigation of sequences and structures. *J. Mol. Biol.* **247** (1995) 536-540
8. Sonnenburg, S., Ratsch, G., and Schafer: Learning Interpretable SVMs for biological Sequence Classification. *RECOMB* (2005) 389-407.
9. Weston, J., Leslie, C., Ie, E., Zhou, D., Elisseeff, A. and Noble, W.S.: Semi-Supervised Protein Classification using Cluster Kernels. *Bioinformatics*. **21** (2005) 3241-3247.
10. Zhu, X., Kandola, J., Ghahramani, Z., and Lafferty, J.: Nonparametric Transforms of Graph Kernels for Semi-Supervised Learning. *Advances in Neural Information Processing Systems*. **17** (2005)
11. Lanckriet, G., Cristianini, N., Bartlett, P., Ghaoui, L. and Jordan, M.: Learning the kernel matrix with semidefinite programming. *Journal of Machine Learning Research*. **5** (2004) 27-72

This article was processed using the \LaTeX macro package with LLNCS style