



Kernel methods for predicting protein–protein interactions

Asa Ben-Hur^{1,*} and William Stafford Noble^{1,2}

¹Department of Genome Sciences and ²Department of Computer Science and Engineering, University of Washington, Seattle, WA, USA

Received on January 15, 2005; accepted on March 27, 2005

Despite advances in high-throughput methods for protein–protein interactions, the interaction networks in well-studied model organisms are sketchy at best, highlighting the continued need for computational methods to assist direct experimentalists in the search for novel interactions.

We present a kernel method for predicting protein–protein interactions using a combination of data sources, including protein sequences, Gene Ontology annotations, motifs of the network, and homologous interactions. Unlike previous methods, whereas protein kernels proposed in the literature measure a similarity between single proteins, prediction of interactions requires a kernel between pairs of proteins. We propose a pairwise kernel that converts a kernel between single proteins into a kernel between pairs of proteins, and we evaluate the kernel’s effectiveness in conjunction with a support vector machine classifier. Furthermore, we obtain improved performance by combining several sequence-based kernels based on k-mer frequency, motif and domain content with our pairwise sequence kernel, and we evaluate these augmented kernels based on other sources of data.

We use our method to predict physical interactions in yeast from the BIND database. At a false positive rate of 1%, our classifier retrieves close to 80% of a set of trusted interactions. We thus demonstrate the ability of our method to make high-rate predictions despite the sizeable fraction of interactions that are known to exist in interaction databases. The classification experiments were performed using data available at <http://pyml.sourceforge.net>. Data are available at <http://noble.gs.washington.edu/proj/sppi>

These methods include the yeast two-hybrid screen and methods based on mass spectrometry (see von Mering *et al.*, 2002 and references therein). The data obtained by these methods are partial: each experimental assay can identify only a subset of the interactions, and it has been estimated that for the organism with the most complete interaction network, namely yeast, only about half of the complete ‘interactome’ has been discovered (von Mering *et al.*, 2002). In view of the very small overlap between interactions discovered by various high-throughput studies, some of them using the same method, the actual number of interactions is likely to be much higher. Computational methods are therefore required for discovering interactions that are not accessible to high-throughput methods. These computational predictions can then be verified by more labor-intensive methods.

A number of methods have been proposed for predicting protein–protein interactions from sequence. Sprinzak and Margalit (2001) have noted that many pairs of structural domains are over-represented in interacting proteins and that this information can be used to predict interactions. Several authors have proposed Bayesian network models that use the domain or motif content of a sequence to predict interactions (Deng *et al.*, 2002; Gomez *et al.*, 2003; Wang *et al.*, 2005). The pairwise sequence kernel was independently proposed in a recent paper (Martin *et al.*, 2005) with a sequence representation by 3mers. Other sequence-based methods use coevolution of interacting proteins by comparing phylogenetic trees (Ramani and Marcotte, 2003), correlated mutations (Pazos and Valencia, 2002) or gene fusion which works at the genome level (Marcotte *et al.*, 1999). An alternative approach is to combine multiple sources of genomic

Methods, and in particular support vector machines (Schölkopf and Smola, 2002), have proven useful for difficult classification problems in bioinformatics (Schölkopf and Smola, 2002). The learning task we are addressing is the relationship between pairs of protein sequences: are two pairs of sequences interacting or not. The kernel approach (A kernel is a measure of similarity between two objects under the additional condition of being a dot product in some inner product space; see Schölkopf and Smola, 2002 for details) has been used in the literature to measure similarity between pairs of sequences. We propose a method for converting a kernel that operates on individual genes or proteins into a pairwise kernel, and we evaluate its performance on the task of predicting interactions from sequence data.

Our method uses motif, domain and *k*-mer composition to generate a pairwise kernel, and we evaluate its performance on the task of predicting interactions from sequence data. We compare our method to other kernel methods based on BLAST or PSI-BLAST. We also compare our method to other methods based on similarity of GO annotations, and to methods based on interacting homologs in other species (Goldberg and Roth, 2002). We also compare our method to other methods based on the tendency of neighbors of interacting proteins to interact as well. Adding these additional data sources significantly improves our method's performance relative to other methods trained using only the pairwise sequence data. Our kernel methods for combining data from heterogeneous sources of data allows us to use high-dimensional data, whereas other studies on predicting protein–protein interactions (Zhang *et al.*, 2004; Lin *et al.*, 2004) use low-dimensional representations which are appropriate for linear classifiers.

METHODS FOR PROTEIN–PROTEIN INTERACTIONS

Kernel methods derive much of their power from their ability to incorporate prior knowledge via the kernel function. Furthermore, the kernel approach offers the flexibility to apply kernels to diverse types of data, including both vectors (e.g. microarray expression data), and strings (DNA and protein sequences), graphs and networks. In this work, we employ a diverse collection of kernels, as described in this section.

We compare two pairs of proteins X_1 and X_2 compared with proteins X'_1 and X'_2 . We call a kernel that operates on individual genes or proteins a 'genomic kernel', and a kernel that compares pairs of genes or proteins a 'pairwise kernel'. Pairwise kernels can be computed either indirectly, by way of an intermediate genomic kernel, or directly using features that characterize pairs of proteins.

The most straightforward way to construct a pairwise kernel is to express the similarity between pairs of proteins in terms of similarities between individual proteins. In this approach, we consider two pairs to be similar to one another when each protein of one pair is similar to one protein of the other pair. For example, if protein X_1 is similar to protein X'_1 , and X_2 is similar to X'_2 , then we can say that the pairs (X_1, X_2) and (X'_1, X'_2) are similar. We can translate these intuitions into the following pairwise kernel:

$$K((X_1, X_2), (X'_1, X'_2)) = K'(X_1, X'_1)K'(X_2, X'_2) + K'(X_1, X'_2)K'(X_2, X'_1),$$

where $K'(\cdot, \cdot)$ is any genomic kernel. This kernel takes into account the fact that X_1 can be similar to either X'_1 or X'_2 .

An alternative to the above approach is to represent a pair of sequences (X_1, X_2) explicitly in terms of the domain or motif pairs that appear in it. This representation is motivated by the observation that some domains are significantly over-represented in interacting proteins (Sprinzak and Margalit, 2001). A similar observation holds for sequence motifs as well. Given a pair of sequences X_1, X_2 represented by vectors $\mathbf{x}_1, \mathbf{x}_2$, with components $x_i^{(1)}, x_i^{(2)}$ we form the vector \mathbf{x}_{12} with components $x_i^{(1)}x_j^{(2)} + x_i^{(2)}x_j^{(1)}$. We can now define the explicit pairwise kernel:

$$K((X_1, X_2), (X'_1, X'_2)) = K'(\mathbf{x}_{12}, \mathbf{x}'_{12}), \quad (1)$$

where \mathbf{x}_{12} is the pairwise representation of the pair (X_1, X_2) , and $K'(\cdot, \cdot)$ is any kernel that operates on vector data. It is straightforward to check that for a linear kernel function, the pairwise and explicit pairwise kernels are identical. The explicit representation can be used in order to rank the relevance of motif pairs with respect to the classification task. This ranking is accomplished, e.g. by sorting the motif pairs according to the magnitude of the corresponding weight vector components.

ance models for our motif kernel are discrete motifs, providing a count of how many times a discrete motif matches a sequence. To compute the discrete motif matches we used discrete sequence motifs from the eMotif database (Hill-Manning *et al.*, 1997). Yeast ORFs contain 17 768 motifs out of a set of 42 718 motifs.

The Pfam kernel uses a set of hidden Markov models to represent the domain structure of a protein, computed by comparing each protein sequence with the Pfam database (Sonnhammer *et al.*, 1997). A protein-HMM comparison yields an E -value. In version 10.0 contains 6190 domain HMMs; a protein is represented by a vector of 6190 log-odds Pfam kernel has been used previously to predict protein interactions (Gomez *et al.*, 2003), though not in combination with the pairwise kernel described above. For sequence kernels we use a normalized linear kernel $K(x, y) = \frac{K(x, x)K(y, y)}{\sqrt{K(x, x)K(y, y)}}$; in the case of the Pfam kernel we performed an initial step of centering the kernel.

Sequence kernels

Using the pairwise kernel is the following:

$$K((X_1, X_2), (X'_1, X'_2)) = K'(X_1, X_2)K'(X'_1, X'_2). \quad (2)$$

This is appropriate when similarity within the pair is related to the likelihood that a pair of proteins interact. This is a valid kernel even if K' is not a kernel, as the formulation K' is simply a feature of the pair of proteins. Consider GO annotations, for example: a pair of proteins is more likely to interact if the two proteins share similar GO annotations. In addition to GO annotation we also consider features of the interaction network, and homologous proteins in other species. We summarize these properties into a vector $\mathbf{s}(X_1, X_2)$, such that the kernel for the data can be any kernel appropriate for vector data.

$$K((X_1, X_2), (X'_1, X'_2)) = K'(\mathbf{s}(X_1, X_2), \mathbf{s}(X'_1, X'_2)), \quad (3)$$

where we use a Gaussian kernel for K' .

Kernel Proteins that are not present in the same cell or that participate in different biological processes are less likely to interact. We represent this prior

We consider two ways in which to define the dot product in this space. When the non-zero components are set equal to 1, then when each protein has a single annotation, and the annotations are on a tree, the dot product between two proteins is the height of the lowest common ancestor of the two nodes. An alternative approach assigns annotation a a score of $-\log p(a)$, where $p(a)$ is the fraction of proteins that have annotation a . We then score the similarity of annotations a, a' as $\max_{a'' \in \text{ancestors}(a) \cap \text{ancestors}(a')} -\log p(a'')$. In a tree topology, this score is the similarity between the deepest common ancestor of a and a' , because the node frequencies are decreasing along a path from the root to any node. The score is a dot product with respect to the infinity norm on the annotation vector space. This also holds when the proteins have more than one annotation and the similarity between their annotations is defined as the maximum similarity between any pair of annotations. When one of the proteins has an unknown GO annotation, the kernel value is set to 0.

2.3.2 Interactions in other species It has been shown that interactions in other species can be used to validate or infer interactions (Yu *et al.*, 2004): the existence of interacting homologs of a given pair of proteins implies that the original proteins are more likely to interact. We quantify this observation with the following homology score for a pair of proteins (X_1, X_2) :

$$h(X_1, X_2) = \max_{i \in \mathcal{H}(X_1), j \in \mathcal{H}(X_2)} I(i, j) \times \min(l(X_1, X_i), l(X_2, X_j)),$$

where $\mathcal{H}(X)$ is the set of non-yeast proteins that are significant BLAST hits of X , $I(i, j)$ is an indicator variable for the interaction between proteins i and j , and $l(X_k, X_i)$ is the negative of the log E -value provided by BLAST when comparing protein k with protein i in the context of a given sequence database. We used interactions in human, mouse, nematode and fruit fly to score the interactions in yeast.

2.3.3 Mutual clustering coefficient Protein-protein interaction networks tend to be ‘cliquish’; i.e. the neighbors of interacting proteins tend to interact. Goldberg and Roth (2003) quantified this cohesiveness using the mutual clustering coefficient (MCC). Given two proteins u, v , their MCC can be quantified, by the Jaccard coefficient $|N(v) \cap N(u)| / |N(v) \cup N(u)|$, where $N(x)$ is the set of neighbors of a protein x

ernels, while the feature space for $\sum_i K_p(K_i)$ of features that originate from the same gen practice, the results from these two different ere very close, and the mixing approach was of its lower memory requirement. A Gaussian kernel can be introduced at several stages: linear genomic kernel as: $\exp(-\gamma(K_p(P, P) - K_p(P', P')))$, where P, P' are two pairs of pro- not tried introducing a non-linear kernel at the omic kernel; a Gaussian kernel at the level of rnel performed similar to the ‘linear’ pairwise the high dimensionality of the resulting feature ults reported in this paper are computed using se kernels.

Improving interaction reliability in

g
s of protein–protein interaction data have noted experimental assays produce varying levels of and have proposed methods for finding which e likely to be reliable (von Mering *et al.*, 2002; , 2003; Deane *et al.*, 2002) (see Section 3.1 for incorporate this knowledge about the reliability of n interactions into the training procedure using -margin parameter C (Schölkopf and Smola, rameter puts a penalty on patterns that are mis- re close to the SVM decision boundary. Each ble receives a value of C that depends on its a training set with an equal number of positive xamples we use two values: C_{high} for interac- to be reliable and for negative examples; C_{low} xamples that are not known to be reliable.

DS

Interaction data

he prediction of physical interactions in yeast tion data from the BIND database (Bader *et al.*, ncludes published interaction data from high- periments as well as curated entries derived d papers. The advantage of BIND is that explicit distinction between direct physical d comembership in a complex.

Positive and negative examples We use physical

negative examples is likely to contain very few proteins that interact.

High-throughput protein–protein interaction data contain a large fraction of false positives, estimated to be up to 50% in some experiments (von Mering *et al.*, 2002). Therefore, we prepared a set of BIND interactions that are expected to have a low rate of false positives. We use these reliable interactions in two ways. We evaluate the performance of our method on the reliable interactions because they are more likely to reflect the true performance of the classifier. We also use reliability to set the value of the SVM soft-margin parameter as discussed in Section 2.5. ‘Gold standard’ interactions can be derived from several sources:

- Interactions corroborated by interacting yeast paralogs. Deane *et al.* (2002) find 2829 interactions from the DIP database that are supported by their paralogous verification method (PVM). The estimated false positive rate of this method is 1%.
- Interactions that are supported by interacting homologs in multiple species are likely to be correct (Yu *et al.*, 2004).
- Interactions that are discovered by different experimental assays were estimated to be correct 95% of the time (Sprinzak *et al.*, 2003).
- Highly reliable methods, e.g. interactions derived from crystallized complexes.

We do not use PVM-validated interactions because they contain several biases.

- The test set is biased toward interactions that can be easily discovered by sequence similarity.
- The list of PVM-validated interactions cannot be used as-is to set the SVM soft-margin parameter in training because this may incorporate information about interactions that are in the test set.

Also, we do not include interactions validated by interacting homologs in other species, since that information is included in the data as a feature. Therefore, for the purpose of assessing performance we use a list of 750 interactions that were validated by high-quality or multiple assays. For setting the SVM soft-margin parameter we augment the 750 interactions with PVM-validated interactions that are computed on the basis of the training data alone. Training is performed on all

cores for the various methods computed using 5-fold cross-validation

	Kernel	ROC score	ROC ₅₀ score
	—	0.74	0.18
	—	0.78	0.11
	$K_{\text{non-seq}}$	0.95	0.37
	$K_p(K_{\text{motif}})$	0.76	0.17
	$K_p(K_{\text{Pfam}})$	0.78	0.20
	$K_p(K_{\text{spec}})$	0.81	0.05
	$K_p(K_{\text{motif}} + K_{\text{Pfam}})$	0.82	0.22
spectrum	$K_p(K_{\text{motif}} + K_{\text{Pfam}} + K_{\text{spec}})$	0.86	0.17
	$K_{\text{feat}} + K_p(K_{\text{motif}} + K_{\text{Pfam}} + K_{\text{spec}})$	0.97	0.44
	$K_{\text{feat}} + K_p(K_{\text{motif}} + K_{\text{Pfam}} + K_{\text{spec}})$	0.97	0.58

all BIND physical interactions. ROC scores are computed on reliable interactions that do not include PVM-validated interactions. The BLAST and PSIBLAST scores are computed according to Equation (4). The ‘kernel’ column of the table shows which kernel was used in conjunction with the SVM classifier. The notation $K_p(K_g)$ denotes a pairwise kernel derived from a genomic kernel K_g . The $K_{\text{non-seq}}$ is a Gaussian kernel over the non-sequence features; in each method it participates in, the width is determined by cross-validation as part of the classifier’s training. The all-reliable method uses information on reliability to set the SVM soft-margin parameter to 2.5.

penalty for significant matches and increases as the number of matches increases. The score for a query (X_1, X_2)

$$\sum_{X \in \mathcal{P}} I(i, j) \min(l(X_1, X_i), l(X_2, X_j)), \quad (4)$$

the set of all proteins in the training set. In these experiments we use PSI-BLAST scores computed in the Swiss-Prot database (version 40, containing 340,000 proteins).

Methods of merit

In this paper we evaluate the quality of a predictive method using two different metrics. Both metrics—the area under the receiver operating characteristic curve (ROC score), and the normalized area under that curve up to the first 50 false positives (ROC₅₀ score)—aim to measure both sensitivity and specificity by integrating over a curve that plots true positive rate against the proportion of false positive rate. We include both metrics to account for two different types of scenarios in which a protein–protein interaction prediction method might

In the first scenario, imagine that you have developed a new method for detecting whether a given pair of proteins interacts. Rather than testing your method on random pairs of proteins, you could use a predictive

method to determine whether they participate in any predicted interactions. In this case, you do not care about the high-confidence interactions above; instead, you would like to be sure that the complete set of predictions is of high quality. In this case you are interested in the ROC score of the classifier.

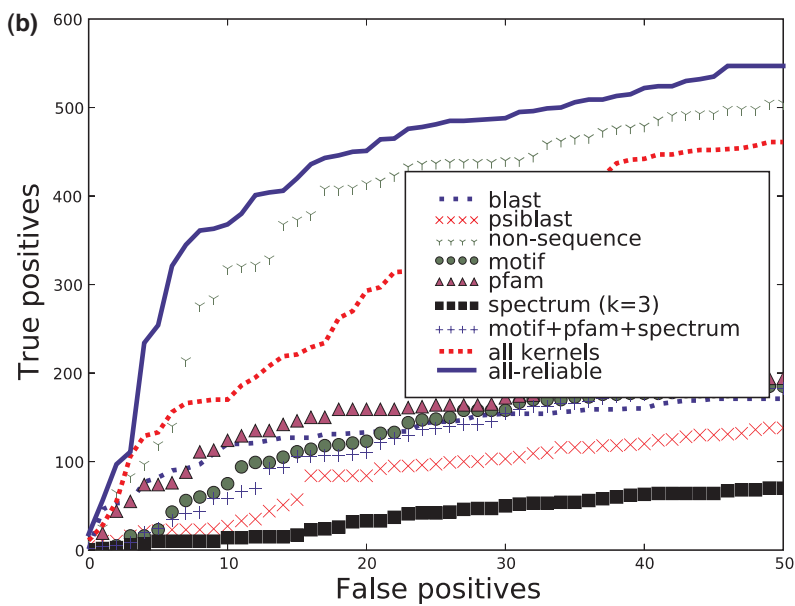
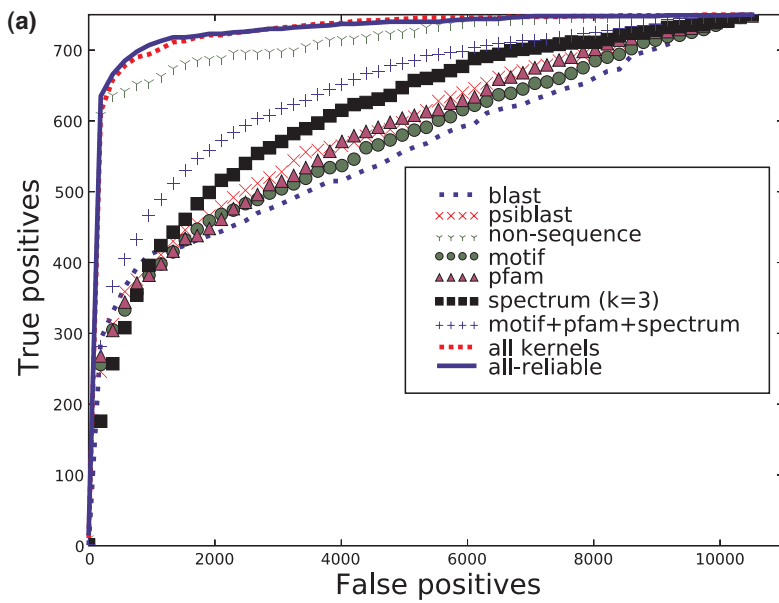
4 RESULTS

We report, in this section, the results of experiments in predicting protein–protein interactions using an SVM classifier with various kernels, and compare these with a simple method based on BLAST or PSI-BLAST. All the experiments were performed using the PyML machine learning framework available at <http://pyml.sourceforge.net>. We begin this section with results obtained using the various kernels and kernel combinations, followed by a discussion of the choice of negative examples, and a section that shows the effects of choosing a non-redundant set of proteins.

4.1 Main results

We report results that are computed using 5-fold cross-validation on all BIND physical interactions. The SVM soft-margin parameter was not optimized—we used the default low value for this parameter to account for the noise in the data. The ROC/ROC₅₀ curve is then computed for those reliable interactions that were not obtained using the PVM method

Kernel methods for predicting protein–protein interactions



and ROC₅₀ (b) curves for several methods. Best performance is obtained using a kernel that combines all the kernels presented. Additional results are summarized in Table 1, along with a description of the methods.

the Pfam and motif kernels. The higher ROC

We now explore the effect of adding to the sequence ker-

tions provided another significant boost to the ROC and ROC₅₀ scores were 0.98 and 0.95, respectively; at a false positive rate of 1% the classifier correctly classified 99% of the trusted interactions. In this experiment we used $C_{\text{low}} = 0.01C_{\text{high}}$.

The contribution to the gain in performance comes from the process kernel feature. Its ROC score by itself on the BIND interactions and 0.95 when limiting to only positive examples. The difference between the two is probably due to the sizable fraction of false positives in the BIND dataset. In the following subsection we consider scenarios where the GO data are not useful. The ROC for the MCC feature was 0.68 on all BIND interactions and 0.53 when computed on the reliable interactions. The performance for the MCC feature is a result of the fact that it requires a large number of interactions to be supported by a BLAST cutoff of $1e^{-10}$, 329 interactions from other species, as supported by interactions from other species, as negative examples. The ROC score for this feature is low since it is sparse, i.e. is informative for a small number of interactions.

Role of GO annotations

To understand the difference in the role of the sequence-based and non-sequence kernels, we compared the two methods on the task of distinguishing between physically interacting protein pairs and those that are members of the same complex. In this case, the negative examples are chosen as protein pairs that are known to belong to the same complex but are not known to physically interact. This set of negative examples is likely to be more noisy than the non-interacting protein pairs, as they likely contain many physical interactions. But still, the pairwise method achieves an ROC score of 0.78, compared to the value obtained with non-interacting negative examples of 0.5. This is due to the fact that protein pairs, such as physically interacting proteins, often have similar GO annotations and network properties, which can be used to distinguish them from the non-interacting protein pairs. Motif and Pfam rely on a signal that is often distant from the interaction site itself (Wang *et al.*, 2005). GO annotations can be made for other features used to distinguish protein pairs.

Significant attention has been paid to the problem of selecting gold standard interacting protein pairs for the purposes of training and validating predictive computational methods (Jansen *et al.*, 2003). However, less emphasis has been placed on the choice of non-interacting protein pairs. In this study, we selected negatives uniformly at random. We find that this strategy leads to consistent behavior and avoids bias.

The possibility for bias due to the method of constructing negative examples is evidenced by results reported in a related paper (Martin *et al.*, 2005). In this work, the authors report that a pairwise spectrum kernel provides highly accurate predictions of yeast interactions using a dataset studied in Jansen *et al.* (2003). The positive examples in this dataset satisfy our criteria of trusted interactions, and one might conclude that the use of highly reliable interactions is the reason for the success of the predictive method. However, we found that the method of choosing negative examples has a strong effect on the performance: the negative examples from Jansen *et al.* (2003) were chosen as pairs of proteins that are known to be localized in different cellular compartments. This makes these protein pairs much less likely to interact than randomly selected pairs, but the selection constraints impose a bias on the resulting distribution that makes the overall learning task easier [note that this is less likely to affect the results of non-sequence based methods, such as the one used by Jansen *et al.* (2003)]. To illustrate this effect, we created datasets with negative examples taken as pairs whose GO component similarity, as measured by our kernel, is below a given threshold. The performance of the resulting classifier varied as we varied this threshold (Table 2). This constrained selection method was tested with the spectrum and motif kernels using both the BIND interaction data and a set of trusted interactions similar to the one used by Martin *et al.* (2005) extracted from DIP and MIPS (Mewes *et al.*, 2000; Xenarios *et al.*, 2002). For the spectrum kernel, the ROC (ROC₅₀) scores varied from 0.87 (0.08) to 0.97 (0.46) on the DIP/MIPS data and from 0.77 (0.04) to 0.95 (0.36) on the BIND data, as the threshold was lowered from 0.5 to 0.04. Similarly, although slightly less pronounced, results were obtained for the motif pairwise kernel.

4.4 The dependence on interacting paralogs

The yeast genome contains a large number of duplicated genes. Since we are using a sequence-based method to pre-

Kernel methods for predicting protein–protein interactions

dependence of the performance of the spectrum pairwise similarity between localization annotations in negative

Threshold	ROC	ROC ₅₀
0.50	0.77	0.04
0.10	0.89	0.15
0.07	0.91	0.21
0.05	0.92	0.25
0.04	0.95	0.36
0.5	0.87	0.08
0.1	0.94	0.22
0.07	0.95	0.32
0.05	0.96	0.34
0.04	0.97	0.46

tion that no two proteins in the set of negative examples have a less than a given threshold puts a constraint on the distribution of this constraint makes it easy for the classifier to distinguish between positive and negative examples, and the effect gets stronger as the threshold becomes smaller. We performed the experiment on the BIND interaction dataset and on a dataset of interactions derived from DIP and MIPS interactions.

BLAST (BLAST) method went down from 0.78 to 0.62). This illustrates that the kernel combination is dependent on the presence of interacting paralogs and that PSI-BLAST.

DISCUSSION

We have presented several kernels for prediction of protein–protein interactions and used them in combination for improved performance. The concern regarding the pairwise kernel is the high dimensionality of its feature space, which is due to the large number of features of the underlying kernel. We presented an alternative kernel which uses summation instead of multiplication used in the expression for the pairwise kernel, similar to the work of Gomez *et al.* (2003). The performance of the summation kernel is not as good as the pairwise kernel, showing the advantage of the pairwise kernel.

Using a classifier to predict protein–protein interactions is a balance between placing more positive and negative interactions as opposed to trying to maximize the number of positive examples by adding interactions

We also made no attempt to purge from our dataset examples that contain missing data (missing GO annotations). When trying to make predictions on unseen data, these data will contain missing data and so, the method is more likely to generalize if presented with examples containing missing data during training.

During the time of writing this paper we found that the pairwise approach was proposed by Martin *et al.* (2005). They used only the spectrum kernel, whereas here we considered several sequence kernels. We found that the spectrum kernel works better than the motif and Pfam kernels according to the ROC metric, but the spectrum kernel does not work as well as the motif and Pfam kernels according to the ROC₅₀ metric. Apparently, the signal that the spectrum kernel generates is not as specific as that of the other kernels.

In addition, we have illustrated that pairwise sequence kernels can be successfully combined with non-sequence data. In this work, we have not attempted to learn the weights of the various kernels as done by Lanckriet *et al.* (2004). This is an avenue for future work, although solving the resulting semi-definite programming problem promises to be computationally expensive, owing to the large training sets involved. We also plan to consider additional sources of data such as gene expression and transcription factor binding data, which have also been shown to be informative in predicting protein–protein interactions (Zhang *et al.*, 2004).

ACKNOWLEDGEMENTS

The authors thank Doug Brutlag, David Baker, Ora Schueler-Furman and Trisha Davis for the helpful discussions. This work is funded by NCRN NIH award P41 RR11823, by NHGRI NIH award R33 HG003070, and by NSF award BDI-0243257. W.S.N. is an Alfred P. Sloan Research Fellow.

REFERENCES

- Bader,G.D. Donaldson,I., Wolting,C. Ouellette,B.F., Pawson,T. and Hogue,C.W. (2001) BIND—the biomolecular interaction network database. *Nucleic Acids Res.*, **29**, 242–245.
- Ben-hur,A. and Brutlag,D. (2003) Remote homology detection: a motif based approach. *Bioinformatics*, **19** (Suppl 1), i26–i33.
- Deane,C. Salwinski,L., Xenarios,I. and Eisenberg,D. (2002) Two methods for assessment of the reliability of high throughput

- I., Greenbaum,D., Kluger,Y., Krogan,N.J., Chung,S., Snyder,M., Greenblatt,J.F. and Gerstein,M. (2003) Networks approach for predicting protein-protein interactions from genomic data. *Science*, **302**, 449–453.
- G., Deng,M., Cristianini,N., Jordan,M.I. and Elkan,C. (2004) Kernel-based data fusion and its application to protein interaction prediction in yeast. In Altman,R.B., Dunker,A.K., Leng,T.A. and Klein,T.E. (eds), *Proceedings of the Pacific Symposium on Biocomputing*, World Scientific, Singapore, pp. 564–575.
- n,E. and Noble,W.S. (2002) The spectrum kernel: a new kernel for SVM protein classification. In Altman,R.B., Hunter,L., Lauderdale,K. and Klein,T.E. (eds), *Proceedings of the Pacific Symposium on Biocomputing*, New Jersey: World Scientific, Singapore, pp. 564–575.
- Jansen,R., Gerstein,M. and Zhao,H. (2004) Informatics on predicting protein-protein interactions. *Bioinformatics*, **5**, 154.
- Pellegrini,M., Ng,H.-L., Rice,D.W., Yeates,T.O. and Eisenberg,D. (1999) Detecting protein function and protein-protein interactions from genome sequences. *Science*, **285**, 1361–1363.
- e,D. and Faulon J.-L. (2005) Predicting protein-protein interactions using signature products. *Bioinformatics*, **21**, 1361–1363.
- rishman,D., Gruber,C., Geier,B., Haase,D., Kaps,A., Mannhaupt,G., Pfeiffer,F., Schüller,C., Stocker,S. and Eisenberg,D. (2000) MIPS: a database for genomes and protein annotations. *Nucleic Acids Res.*, **28**, 37–40.
- g,C.G., Sethi,K.S., Wu,T.D. and Brutlag,D.L. (1997) Identifying and ranking discrete motifs. In *Proceedings of the International Conference on Intelligent Systems for Biomedicine*, pp. 202–209.
- (2004) Support vector machine applications in computational biology. In Schoelkopf,B., Tsuda,K. and Vert,J.-P. (eds), *Computational Methods in Computational Biology*. MIT Press, Cambridge, MA, pp. 71–92.
- Pazos,F. and Valencia,A. (2002) *In silico* two-hybrid system for the selection of physically interacting protein pairs. *Proteins*, **47**, 219–227.
- Ramani,A. and Marcotte,E. (2003) Exploiting the co-evolution of interacting proteins to discover interaction specificity. *J. Mol. Biol.*, **327**, 273–284.
- Schölkopf,B. and Smola,A. (2002) *Learning with Kernels*. MIT Press, Cambridge, MA.
- Sonnhammer,E., Eddy,S. and Durbin,R. (1997) Pfam: a comprehensive database of protein domain families based on seed alignments. *Proteins*, **28**, 405–420.
- Sprinzak,E. and Margalit,H. (2001) Correlated sequence-signatures as markers of protein-protein interaction. *J. Mol. Biol.*, **311**, 681–692.
- Sprinzak,E., Sattath,S. and Margalit,H. (2003) How reliable are experimental protein-protein interaction data? *J. Mol. Biol.*, **327**, 919–923.
- von Mering,C., Krause,R., Snel,B., Cornell,M., Olivier,S.G., Fields,S. and Bork,P. (2002) Comparative assessment of large-scale data sets of protein-protein interactions. *Nature*, **417**, 399–403.
- Wang,H., Segal,E., Ben-Hur,A., Koller,D. and Brutlag,D.L. (2005) Identifying protein-protein interaction sites on a genome-wide scale. In Lawrence K. Saul, Yair Weiss and Léon Bottou (eds), *Advances in Neural Information Processing Systems 17*. MIT Press, Cambridge, MA, pp. 1465–1472.
- Xenarios,I., Salwinski,L., Duan,X.Q.J., Higney,P., Kim,S.M. and Eisenberg,D. (2002) DIP: the Database of Interacting Proteins: a research tool for studying cellular networks of protein interactions. *Nucleic Acids Res.*, **30**, 303–305.
- Yu,H., Luscombe,N., Lu,H., Zhu,X., Xia,Y., Han,J., Bertin,N., Chung,S., Vidal,M. and Gerstein,M. (2004) Annotation transfer between genomes: protein-protein interlogs and protein-DNA regulogs. *Genome Res.*, **14**, 1107–1118.
- Zhang,L., Wong,S., King,O. and Roth,F. (2004) Predicting co-complexed protein pairs using genomic and proteomic data integration. *BMC Bioinformatics*, **5**, 38–53.