



A motif-based framework for recognizing sequence families

Roded Sharan^{1,*} and Eugene W. Myers²

¹School of Computer Science, Tel-Aviv University, Tel-Aviv 69978, Israel and

²Computer Science Division, University of California at Berkeley, 387 Soda Hall, Berkeley, CA 94720, USA

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

Many signals in biological sequences are based on the presence or absence of base signals and their spatial relationships. One of the best known examples of this is the recognition of a core promoter—the site at which the basal transcriptional machinery starts the transcription of a gene. Our automatic pattern recognition system for a family of motifs, which simultaneously discovers the base signal relationships and a classifier based upon

In this paper we present a general method for recognizing a set of sequences by their recurrent motifs. It relies on novel probabilistic models for DNA motifs and modules of binding sites, on algorithms to mine the data and on a support vector machine model studied to classify a set of sequences. We demonstrate the applicability of our approach to diverse motifs ranging from families of promoter sequences to regulatory sequences flanking alternatively spliced exons. On a promoter dataset our results are comparable to the state-of-the-art McPromoter. On a dataset of alternative exons we outperform a previous approach. Our results demonstrate that a fully automatic recognition algorithm can meet or exceed the performance of hand-crafted approaches.

The software and datasets are available from the authors upon request.

sharan@tau.ac.il

find, harder even than the now trendy *cis*-regulatory signals, also known as distal and proximal promoters, which serve as the binding sites of complexes that interact with and modulate the activity of core promoters. Solving the problem of finding core promoters is very important as most gene prediction programs routinely miss the 5' exon because they are geared to recognize coding sequence. Among other implications, this has great impact on the accuracy of the upstream region in which one looks for *cis*-regulatory control.

We wished to study a classifier of the signal by the examination of a collection of positive and negative examples. The basic idea is to first recognize potentially distinguishing attributes or patterns and then study which combinations of these attributes discriminate positive from negative examples. The idea is quite natural and there have been several other attempts along these lines (Pavlidis *et al.*, 2001; Ben-Hur and Brutlag, 2003), mainly focusing on the classification task. The specific problem of recognizing eukaryotic core promoters has been studied by several authors and various approaches have been reported for it, including neural networks (Reese, 2001), linear discriminant analysis (Hannenhalli and Levy, 2001) and hidden Markov models (Ohler *et al.*, 2002). The last method, called McPromoter, is the best in-class and hand-crafted classifier for *Drosophila* core promoters based on a great deal of human analysis and insight.

Here, we present a unified framework for the task of recognizing sequence families. The framework consists of two components: (1) algorithms that recognize unusual patterns or attributes of a number of types within the training dataset and (2) a support vector machine (SVM) that uses the attributes

of NNPP (Reese, 2001). Moreover, we came the method extends well beyond our original illustrate this we apply it here to the problems of alternatively spliced exons, in human, and recog that are under cell-cycle control in yeast. In with a previous approach for detecting alternat-exons, we are able to show increased sensitivity ons.

DS

e following classification problem: the input training set of sequences with positive and mples, and a test set; the goal is to devise a the positive examples that will best discrim- positives and negatives on the test set. We -phase scheme for this problem: in the first the training data to study attributes (features) lent in the positive sequences compared with (negative) sequences. The attribute vector of e consists of three types of attributes: (1) dis- motifs, (2) discriminative modules of motifs n attributes that are unique to the specific n the second phase we train a SVM for the clas- problem using the attributes studied as sequence e two phases are described in detail in the ions.

ng discriminative motifs

h motif using the standard position weight mat- representation (Bailey and Elkan, 1994; Roth *et al.*, assumes independence between positions in a This model assigns a weight to each position in each nucleotide $n \in \{A, C, G, T\}$, representing which the nucleotide's presence in this position with the motif.

g PWMs we adapt the discriminative motif al *et al.* (2002). This model is specified using a on with p position-specific weights $w_i[n]$, one on i and each nucleotide $n \in \{A, C, G, T\}$, and o. For a sequence example s , denote its nucle- by $s.S = s.S_1, \dots, s.S_L$. For a motif m , denote ation of occurrence of m in s , with the conven-

a motif occurrence given the sequence is:

$$P(s.m \geq 0 \mid s.S_1, \dots, s.S_L, \theta_m) \\ = \text{logit} \left(w_0 + \log \left(\sum_{j=1}^{L-p+1} p_m(j) \exp \left\{ \sum_{i=1}^p w_i [s.S_{i+j-1}] \right\} \right) \right)$$

where θ_m is the set of parameters for the motif, $p_m(j) = 1/(L - p + 1)$ and $\text{logit}(x) = 1/(1 + e^{-x})$ is the logistic function. [The reader is referred to Segal *et al.* (2002) for more details on the model and the likelihood derivation.]

We extend the above model to take into account the possible bias in the location of certain motifs along the input sequences. Such bias was observed previously for promoter regions [see e.g. Tanay and Shamir (2003) and Beer and Tavazoie (2004)]. We use a simple model for the location preference, in which the sequence is equally partitioned to k parts ($k = 10$), each having a certain probability of containing the motif, and within each part the probability of occurrence is assumed to be uniform. For a given motif, we empirically estimate the distribution of the locations of its occurrences along the positive sequences (see below). We redefine $p_m(j)$ based on the estimated distribution.

A complicating factor in applying this model to study the motif parameters from the data is that we do not expect the motif m to occur in every core promoter sequence, but only in a fraction of the sequences. Thus, we treat the positive training data as noisy. Precisely, let T be a set of labels for the training sequences, specifying for each sequence s whether it is a positive or a negative example. We further denote, T^+ as the set of positive examples, T^- as the set of negative examples and S as the set of all nucleotide sequences $\{s.S \mid s \in T\}$. Define $q_m \equiv P(s \in T^+ \mid s.m = -1)$ to be the probability that a sequence is a core promoter given that motif m does not occur in it. This probability reflects the fraction r_m of positive sequences containing the motif m : $q_m = \left(1 + \frac{a}{1-r_m}\right)^{-1}$, where a is the ratio of negative to positive examples. The likelihood of the data under this extended model is:

$$P(T \mid S, \theta_m, q_m) = \prod_{s \in T^+} \{P(s.m \geq 0 \mid s.S, \theta_m) \\ + q_m(1 - P(s.m \geq 0 \mid s.S, \theta_m))\} \\ \times \prod_{s \in T^-} \{(1 - q_m)(1 - P(s.m \geq 0 \mid s.S, \theta_m))\}$$

and similar to Barash *et al.* (2001), which we

Initialization of the motif model

Search of PWMs that correspond to putative binding sites is done using a three-stage process: First, discriminative sequence patterns are identified; second, these sequences are scored to quantify their enrichment in positive sequences versus the negative ones; third, the most significant patterns along the positive sequences are used to compute an initial PWM for the corresponding

Search is done in an exhaustive manner, scoring motifs of length 6–8 bp, which are called seeds. To estimate the count of its occurrences up to one mismatch in the positive and negative examples. We compute a geometric P -value for these counts, and retain motifs that have an adjusted P -value < 0.01 (we use a Bonferroni correction to adjust the P -values for multiple comparisons). We also compute an enrichment P -value against a first-order Markov model of the positive sequences, and filter motifs that do not pass the 0.01 significance level. The surviving motifs are further filtered in a greedy fashion to ensure that they are not similar in sequence or significantly overlap with other motifs.

For each motif, we compute the initial position specific weights by averaging over all occurrences (up to one mismatch) in this seed. We use the seed occurrences also to estimate the PWM at each end by positions whose enrichment exceeds a threshold. Once the initial PWM is determined, the parameters of the location distribution are estimated by considering, for each positive sequence, the highest-scoring match of the pattern to

Identifying discriminative modules

In addition to the motif-based features, we also study more complex features, namely, spatial combinations of motifs, or motifs that seek modules that are abundant in the positive sequences but not in the negative ones. Studying modules that seek to identify signals that are too weak at the motif level to be significant. We seek to associate motifs whose co-occurrence has a high significance.

Finally, we generalized the above motif model to

therefore:

$$P(s.M = i | s.S) = \frac{1}{d_U - d_L + 1} P(s.m_1 = i | s.S) \sum_{j=i+d_L}^{i+d_U} P(s.m_2 = j | s.S)$$

where

$$P(s.m_k = l | s.S) = \text{logit} \left(w_0^{(m_k)} + \sum_{t=1}^p w_t^{(m_k)} [s.S_{t+l-1}] \right)$$

One can study this model using the same gradient ascent approach used for the single motif model. The initialization of the model is done by enumerating pairs of seeds (consensus sequences) that occur up to one mismatch within a window of size w ($w = 50$). These putative modules are scored by computing their enrichment in the positive set, using a hypergeometric test. Significant pairs are then initialized in a way similar to the initialization of seeds for the motif model.

2.4 Adding external attributes

Up till now we have described a general framework for studying discriminative attributes from sequence data. However, depending on the specific problem, there may be properties that are important for the classification task and cannot be expressed as sequence motifs. For instance, Sorek *et al.* (2004) show that exons whose length is divisible by three are less likely to be constitutive. Thus, in each of the applications described below we also add to our attribute vectors those attributes that were found to be discriminative for that specific classification problem.

In addition, we add one more feature to the attribute vectors, representing the fit of a sequence to a probabilistic model of the positive sequences versus the negative sequences. Specifically, we compute a first-order Markov model for the positive and negative sequences and define this feature to be the log odds of being a positive versus being a negative example.

2.5 Training the SVM

SVM is a classification method based on finding a separating hyperplane between positive and negative samples that maximizes the distance (margin) between the samples and the hyperplane (in case the samples are not separable).

and modules in Drosophila core promoters

Name	Consensus	Length	<i>P</i> -value	MEME-short	MEME-long
DRE	ATCGATAG	8	1E-33	+	+
—	GGTCACACT	9	3E-23	+	+
DPE	CGGTCG	6	2E-19	+	—
—	CAGCACTG	8	4E-14	+	—
—	CAGCTGGT	8	4E-13	+	—
—	CCGATAAC	8	8E-13	—	—
—	CGACGACG	8	1E-12	—	—
—	TCGCCGCG	8	4E-11	—	—
TATA	CTATAAAA	8	6E-9	+	—
—	CGAGCGGC	8	7E-9	+	—
INR	CTCAGTCG	8	3E-7	+	—
—	GGTATTTT	8	5E-5	+	—
—	TCGGCAGC	8	6E-5	—	—
12 + 2	GGTATTTT:GGTCACAC	≤50	9E-16	—	—
DRE + 6	ATCGATAG:CCGATAAC	≤50	6E-11	—	—
INR+DPE	CTCAGTCG:CGGTCG	≤50	7E-4	—	—

top-scoring motifs. For each motif, indicated are its common name (if such is known), its consensus sequence, its *P*-value (Bonferroni corrected) and whether it was identified by MEME, as reported in Ohler *et al.* (2002). MEME was applied both to the original 300 bp sequences (long) and to shorter segments from -60 to +40 bp (short). A motif is considered to match a MEME motif if their consensus sequences are identical up to one mismatch. Bottom: the three significant modules. The name of each module refers to the motifs that comprise it.

sequences. In order to measure our confidence in the predictions, we compute a confidence score based on the output of the SVM. This is done by fitting a logistic regression model as described by Platt (1999). This is done by fitting a logistic regression model to the output of the SVM.

Performance measures

Let S denote by TP, FP, TN and FN the number of true positives, false positive, true negative and false negative predictions. The sensitivity of a set of predictions is defined as the percentage of positives that are correctly predicted: $\text{sens} = \text{TP}/(\text{TP} + \text{FN})$. The specificity is defined as the percentage of negatives that are correctly predicted: $\text{spec} = \text{TN}/(\text{TN} + \text{FP})$. The FP rate equals $\text{FP}/(\text{TN} + \text{FP})$. For some applications (e.g. core promoter identification—see below) the number of TN in the positive set exceeds the number of TPs. In such cases, we use the specificity measure with an adjusted specificity, $\text{spec} = \text{TP}/(\text{TP} + \text{FP})$.

We used ROC curves to visualize a range of sensitivities and specificities. These are generated by an algorithm using a receiver operating

characteristic curve that serves as a recognition site for the basal transcription apparatus. Common core promoter elements include the TATA box at -31 to -26 bp, its extension, BRE, at -37 to -32 bp, the initiator, INR, at -2 to +4 bp and a downstream element, DPE, at +28 to +32 bp. A fifth element, DRE, was implicated to be abundant in core promoters in Ohler *et al.* (2002).

The training dataset that we used was prepared by Ohler *et al.* (2002) and includes a set of 1842 core promoters, 1799 intronic sequences and 2859 coding sequences. These sequences are 300 bp long, where for core promoters they extend from -250 to +50 bp. In order to take advantage of this partition of the sequences, we trained our model twice: first, to discriminate between core promoters and intronic sequences; and second, to discriminate between core promoters and coding sequences. Since coding sequences are very different from core promoter sequences in their nucleotide content, we used only external attributes for the second classification task. We restricted the program to identify the 15 top-scoring motifs or modules, and retained only significant motifs and modules whose frequency in the positive set was estimated to be $\geq 10\%$.

p. While the first application failed to recover
own core promoter elements, the 10 top-scoring
second application included nine of the motifs
them identified. We note that both our method
did not recover the BRE motif, which could imply
represented in the data.

we studied three significant modules on this
are shown in Table 1. The first module cons-
s 12 and 2. These two motifs were reported
frequency of co-occurrence in core promoter
ler *et al.*, 2002). The second module consists
element and motif 6. The third module consists
DPE motifs. This module structure is one of
mon core promoter structures reported in the
ler and Kadonaga, 2002).

Ohler *et al.* (2002) we also used 14 external attrib-
the physical properties of DNA sequences,
own to discriminate between core promoters
ferences. Specifically, the computation of these
s experimentally derived tables on physical
di- or tri-nucleotides, such as bendability,
conformation, etc. Full details on these prop-
rior computation can be found in Ohler *et al.*
ed the average value of each property along
oter segment from -60 to +40 bp as a fea-
more complex features can be computed based
l attributes, but this was not the focus of our

performance of our algorithm we applied it after
ntify core promoters in the well annotated Adh
et al., 2002). This region is 2.9 Mb long and con-
ated open reading frames (not included in our

The core promoter predictions were computed
indow across each of the strands, calculating its
re, and choosing local maxima of these confid-
the predictions. To evaluate the results we used
ty measures employed in Ohler *et al.* (2002):
adjusted specificity. ROC-like curves of the
sented in Figure 1; a comparison with exist-
s given in Table 2. These results (Fig. 1) also

the utility of using both discriminative motifs and
e classification task. We further examined the
ling the location preference of motifs by com-
ults with a variant of the algorithm that assumes

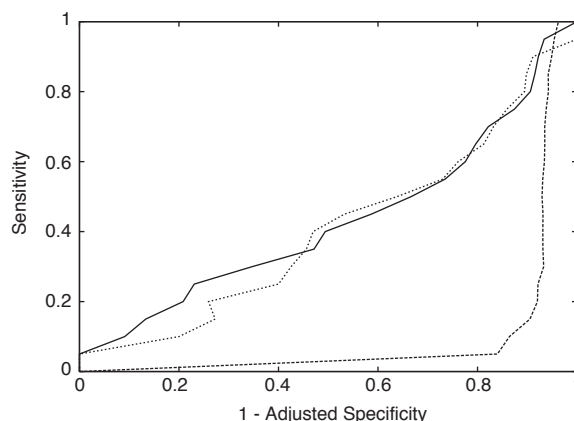


Fig. 1. Performance on the Adh region, shown as ROC-like curves, where the x-axis is $(1 - \text{aspec})$ and the y-axis is the sensitivity of the predictions. The solid, dotted and dashed curves describe the performance of the algorithm when using both discriminative motifs and modules, motifs only and no motifs or modules (i.e. using only external attributes), respectively.

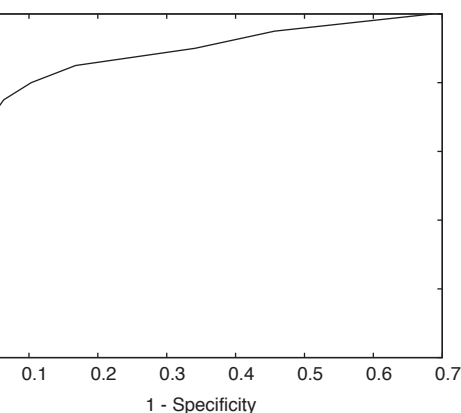
Table 2. Comparison of classification results on the Adh region

Sensitivity (%)	Adjusted specificity		
	MotifBased (%)	McPromoter (%)	NNPP (%)
20	79	69	14
35	53	51	10
50	33	40	6
65	20	29	—

For each sensitivity level, the adjusted specificity of each method is indicated. The results of McPromoter are adapted from Ohler *et al.* (2002). The results of NNPP are adapted from Reese (2001), and were based on a smaller training set.

data. Specifically, they have shown that alternative exons tend to have length divisible by three and tend to be conserved along with their flanking sequence between human and mouse.

We tested our method on the training data reported by Sorek *et al.* (2004), which consists of flanking sequences for 243 alternative exons and 1753 constitutive ones. Following Sorek *et al.* (2004), we evaluated our results using 5-fold cross-validation. The algorithm studied two to three significant motifs in each cross-validation iteration, with two motifs con-



ROC curve for the classification of alternatively spliced

classification results on the exon dataset of Sorek *et al.* (2004)

	Sensitivity (%)	Specificity (%)
	40.3	99.4
4)	32.3	99.7

Sensitivity percentages represent averages over five cross-validation

sensitivity rate of 50%. However, the results of Sorek *et al.* (2004) are not directly comparable with those of Sorek *et al.* (2005), since the latter study used a different method (the data was partitioned into a training and test set) and took advantage of additional external attributes not part of the original data of Sorek *et al.*

Cell cycle regulation in yeast

As part of our method, we applied it to recognize cell cycle regulated genes in yeast according to their expression patterns. The assumption underlying this experiment is that cell cycle regulated genes carry in their promoter regions specific sequence signals, corresponding to the binding sites of cell cycle regulators. To compile a training dataset we collected 100 bp promoter sequences for all yeast genes.

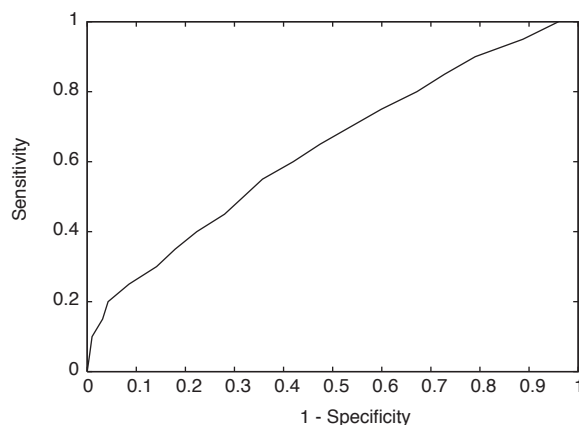


Fig. 3. ROC curve for the classification of cell cycle regulated genes.

consensus sequences matched those of the known cell cycle regulators MBP1, SWI4 and SWI6.

4 CONCLUSIONS

We have presented a general framework for the characterization and classification of a family of related sequences based on recurrent sequence motifs and modules of motifs. We demonstrated several applications of our framework to identifying core promoters, alternatively spliced exons and cell cycle regulated genes. There are many possible extensions to our work, including (1) more refined modeling of the position preference of a motif; (2) modeling the distance distribution among motifs in a module; (3) design of kernel functions for the classification task based on the approach by Lanckriet *et al.* (2004) to provide explicit treatment of the problem of combining features of different types; and (4) application of our method to classify other sequence families, such as core promoters in other species, promoter regions of tissue-specific genes and promoter regions of genes with specific expression patterns.

ACKNOWLEDGEMENT

Part of the work of R.S. was done while doing his post-doc at the University of California, Berkeley.

A motif-based framework for recognizing sequence families

- Brutlag,D. (2003) Remote homology detection based approach. *Bioinformatics*, **19** (Suppl. 1),
- Kadonaga,J.T. (2002) The RNA polymerase II core component in the regulation of gene expression. *Genes & Development*, **16**, 2583–2592.
- Shamir,R. and Shamir,R. (2005) Accurate identification of alternatively spliced exons using support vector machine. *Bioinformatics*, **21**, 897–901.
- Shamir,R. and Levy,S. (2001) Promoter prediction in the human genome. *Bioinformatics*, **17**, S90–S96.
- Shamir,R., De Bie,T., Cristianini,N., Jordan,M.I. and Shamir,R. (2004) A statistical framework for genomic data analysis. *Bioinformatics*, **20**, 2626–2635.
- Shamir,R., Niemann,H., Liao,G.C. and Rubin,G.M. (2001) Joint analysis of DNA sequence and physical properties to improve eukaryotic promoter recognition. *Bioinformatics*, **17**, 1053–1062.
- Shamir,R., Liao,G.-C., Niemann,H. and Rubin,G.M. (2002) Comparative analysis of core promoters in the Drosophila genome. *Genome Research*, **12**, 3, 1–12.
- Shamir,R., Liberto,M., Haussler,D. and Grundy,W.N. (2001) Promoter region-based classification of genes. *Pac. Symp. Biocomput.*, **2001**, 151–163.
- Platt,J.C. (1999) Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. Technical report. Microsoft Research.
- Reese,M.G. (2001) Application of a time-delay neural network to the annotation of the *Drosophila melanogaster* genome. *Comput. Chem.*, **26**, 51–56.
- Roth,F.P., Hughes,J.D., Estep,P.W. and Church,G.M. (1998) Finding DNA regulatory motifs within unaligned noncoding sequences clustered by whole-genome mRNA quantitation. *Nat. Biotechnol.*, **16**, 939–945.
- Segal,E., Barash,Y., Simon,I., Friedman,N. and Koller,D. (2002) From sequence to expression: a probabilistic framework. *Proceedings of RECOMB*, Washington, DC, pp. 263–272.
- Sorek,R., Shemesh,R., Cohen,Y., Basechess,O., Ast,G. and Shamir,R. (2004) A non-EST-based method for exon-skipping prediction. *Genome Res.*, **14**, 1617–1623.
- Spellman,P. T., Sherlock,G., Zhang,M.Q., Iyer,V.R., Anders,K., Eisen,M.B., Brown,P.O., Botstein,D. and Futcher,B. (1998) Comprehensive identification of cell cycle-regulated genes of the yeast *Saccharomyces cerevisiae* by microarray hybridization. *Mol. Biol. Cell*, **9**, 3273–3297.
- Tanay,A. and Shamir,R. (2003) Modeling transcription programs: inferring binding site activity and dose-response model optimization. *Proceedings of RECOMB*, Berlin, Germany, pp. 301–310.