

SVM based on density restricted hierarchical clustering and its application to polyadenylation signals*

Yuanhai Shao¹

Yining Feng¹

Jing Chen¹

Naiyang Deng^{1,†}

¹College of Science, China Agricultural University, Beijing 100083, China

Abstract Support vector machines (SVM) have been promising methods for classification analysis due to their solid mathematical foundations. Clustering-based SVM are used to overcome the difficulty of standard SVM in solving large samples classification problems. In this paper, we present an SVM based on density restricted hierarchical clustering method (DCB-SVM) to predict polyadenylation signal (PAS) in human DNA and mRNA sequences. Compared with the standard SVM, this method uses representation points to train the SVM, by reconciling the boundary it can obtain the similar accuracy. According to the PAS prediction experiments, the proposed method has a better improvement than the clustering-based SVM, it also has a better improvement in sensitivity and similar specificity than the SVM.

Keywords Support vector machines; Polyadenylation signals; BIRCH algorithm

1 Introduction

Due to the rapid growth of the biological data, automatic methods for categorizing the DNA/RNA data are needed. In most cases, the use of statistical or machine learning techniques has been proven to be successful as the methods. Polyadenylation signal (PAS) prediction is an important problem in the bioinformatics, there are three important steps in the process of dealing with the mRNA after transcription, namely the formation of 5' cap structure, trimming of introns and the polyadenylation in the 3' end [1]. A significant relation exists among the polyadenylation in the pre-mRNA 3' end and some other machinery of cells and illnesses [2, 3]. The relation also exists in stability adjustment, the transformation and translation site decision of mRNA. When there are several sites of polyadenylation in the 3'UTR area, the selection of polyadenylation effects the gene expression through specificity of organization and illness. Prediction of mRNA polyadenylation sites [poly(A) sites] can help researchers identify genes, define gene boundaries, and elucidate regulatory mechanisms[4]. The process of polyadenylation is composed of two tightly coupled steps: the cleavage of pre-mRNA and the addition of poly(A)

*This work is supported by the Key Project of National Natural Science Foundation of China (No.10631070), the National Natural Science Foundation of China(No.10871022).

[†]Corresponding Author. E-mail: dengnaiyang@vip.163.com

tail to the newly formed 3' end [5, 6]. This is because *cis* elements are recognized by RNA-binding proteins during mRNA polyadenylation, such as poly(A) specificity factor (CPSF), the cleavage stimulation factor (CStF), the cleavage factors CF I and CF II, and poly(A) polymerase (PAP) [7]. The polyadenylation process is illustrated in [6].

Prediction of poly(A) sites has been attempted by several groups during the last several years. Researchers have proposed some efficient methods to predict the PAS. Tabaska and Zhang [8] developed Polyadq, which employed two quadratic discriminant functions for sequences containing AAUAAA and AUUAAA. They also used a position weight matrix for the downstream sequence, a weighted average of hit positions for DSE, and downstream dimer preferences. In addition, weight-matrix-only [9] also has been employed for poly(A) site prediction. Liu et al.[10] first predicted the mammalian PAS by using SVM and got higher accuracy than the result from Polyadq and Erpind. While facing a large scale DNA sequence, SVM was restricted by the space and time cost. For large scale problems, it is necessary to reduce the data scale. DCB-SVM can be used to classify very large data sets with relatively low dimensionality, such as streaming data or data in large data warehouses.

From the idea of reducing the training data scale, we proposed the SVM based on density restricted hierarchical clustering (DCB-SVM) in this paper. Different from the clustering-based SVM(CB-SVM)[12], it reflects the non-circle distribution of the data set. This method is implemented by limiting the maximum number of each leaf node in the hierarchical clustering, and used the clustering center to denote the micro-cluster. At last, SVM is used to get the optimal decision function. According to the PAS prediction experiments by using mammalian DNA sequence pattern, we get higher predicting precision in positive test data and achieve moderate sensitivity and specificity than CB-SVM.

2 Methods

2.1 Density restricted hierarchical clustering algorithm

In this section, we present the density restricted hierarchical micro-clustering algorithm, which is similar in spirit with BIRCH algorithm[11,12]. The BIRCH algorithm builds a dendrogram called clustering feature tree (CF tree) while scanning the data set. The CF tree carries spherical shapes of hierarchical clusters and captures the statistical summaries of the entire data set. We construct a similar structure called density restricted clustering feature tree (DCF tree) for clustering. DCF tree also assigns an actual object in a cluster as the cluster center to facilitate it by clustering in any distance space, but DCF tree controls the number of the points in the nodes. For its similarities to the original CF tree, it is an alternative structure for our algorithm.

The concept of the CF tree is the core of the hierarchical micro-clustering algorithm which makes the clustering increase without expensive computations. Given N d -dimensional data points in a cluster: $\{x_i\}$, where $i = 1, \dots, N$, the centroid C and radius R of the cluster are defined as:

$$C = \frac{\sum_{i=1}^N x_i}{N}, \quad (1)$$

$$R = \left(\frac{\sum_{i=1}^N \|x_i - C\|}{N} \right)^{\frac{1}{2}}. \quad (2)$$

The CF vector of the cluster is defined as a triple: $CF = (N, LS, SS)$, where n is the number of data points in the cluster, LS is the linear sum of the n data points, and SS is the component-wise square sum of the n data points. CF tree is a height-balanced tree with two parameters: branching factor b and threshold t . The tree size is a function of t . The larger t is, the smaller the tree is. The branching factor b can be determined by memory page size such that a leaf or a non-leaf node fits in a page. The CF tree is a compact representation of the data set, because each entry in a leaf node is not a single data point but a cluster, which absorbs many data points within a radius of t or less.

In DCF tree, besides the triple clustering feature (N, LS, SS) , we introduce a new parameter D in the clustering features tree, which called the density parameter binding. D is used to limit the largest number of sample points in leaf node. It plays the similar role in branch parameter b , but b plays the bound role in the largest number of non-leaf node from the data storage and computer I/O operations, and the density parameter D is the role of restraint limit each micro-cluster represented by the number of points of the original samples.

The clustering algorithms based on the distance often tend to cluster in a circular area of some clusters, but the distribution of data is usually extremely complex. While clustering a non-circular on the regional distribution of the data sets by clustering algorithms based on the distance, the cluster centers are not properly reflect the distribution of the original data. Density restricted hierarchical clustering expect to bound the various micro-cluster in the number of samples by the limitations of each micro-cluster in the density of sample points, when the micro-cluster data are too many that they tend to be a group of non-circular distribution of data to a circular cluster, the micro-cluster will be split by the restrictions of the density parameters.

The density restricted hierarchical clustering tree is built up dynamically as new data objects be inserted. Insertion is similar with the CF tree, the sketch is given below.

1. *Choosing the appropriate leaf*: Starting from the root, it traverse the CF tree down to the leaf level by choosing the child node whose centroid is closest to the new data object at each level.

2. *Modifying the leaf*: If the leaf nodes is too big to fit in memory, the leaf nodes will be split into two. The leaf contains the original data points of the database. If it conflicts with the parameters D of the nodes of all sample points, select two points whose distance is the farthest as the new leaf node of the seeds, insert other points into the nearest one from the leaf nodes. Then update the clustering feature vectors, until all samples are assigned.

3. *Modifying the path*: Check the father nodes until the root node meet the parameter threshold of the branch. If unsatisfied, split it and recursively traverse back up to the root while performing the same checks. If the root node unsatisfied, increased the height of the tree. Update cluster feature information on the path of every non-leaf node. If there is no split, just update the parameters.

2.2 SVM based on density restricted hierarchical clustering

The key idea of DCB-SVM can be viewed as similar to that of selective sampling, i.e., selecting the data that maximizes the benefit of learning. We cluster the node entries near the boundary to get finer sampling close to the boundary and coarser sampling far from

the boundary. Based on this idea, we use the SVs, the description of the class boundary, while keeping the total number of training data points as small as possible. In practice, soft constraints are usually necessary to cope with noise in the training set. Using soft constraints generates the SVs with different distances from the boundary.

DCB-SVM runs on the DCF tree, which can be constructed in a single scan of the entire data set. It carries the statistical summaries that facilitate efficient and effective construction of an SVM boundary. The sketch of the DCB-SVM algorithm is shown in the Table1.

Table 1. DCB-SVM algorithm.

<i>Input:</i> positive data set P , negative data set N .
<i>output:</i> a boundary function f .
Process:
1. <i>Density Restricted Hierarchical Clustering:</i> Construct a positive and a negative tree from P and N respectively.
2. Put the positive and negative root entries in initial training set S .
3.1. $f' = SVM.train(S)$; // construct a boundary f ;
3.2. $S' = getMargin(f, S)$; // compute the low margin entries S' from S using f' ;
3.3. $S = S - S'$; // exclude the low margin data from S ;
3.4. $g = SVM.train(S)$; // construct a boundary f ;
4. return f .

3 Implementation

3.1 Datasets

We use the sequence data provided by Legendre and Gautheret [9] and aim to predict the polyadenylation signal (PAS) in human sequences. The data set contains one group of training data (2327 true PAS) and 5 groups of testing data, each of them consists of 982 samples. Among these 5 sets of testing data, one is true PAS and the other four are all false PAS. The negative sets are CDS sequences, intronic sequences of the first intron, and randomized UTR sequences of the same 1st order Markov model as human 3' UTRs, and randomized UTR sequences of the same mono nucleotide composition as human 3' UTRs. The data set can be downloaded from <http://tagc.univ-mrs.fr/pub/erpin/>. Every sequence contains 206 bases and has a PAS in the center.

There are many mature methods to select the most obvious feature. In this paper, we used the Enterpriser Miner module in SAS 9.1.3 to select the feature through R-square method. First of all, compute the R-square contribution, and then select the first ten features which have the largest R-square contribution. The features we selected in the order of descending R-square contribution are: UP_TGT, DOWN_A, UP_AG, DOWN_TGT, DOWN_GGC, UP_AAG, UP_A, DOWN_AG, DOWN_GAA, UP_GGC.

For prediction, we used the following equations for Sensitivity (SN), Specificity (SP), FPR, precision and accuracy:

$SN = TP/(TP + FN)$, $SP = TNR = TN/(TN + FP)$, $FPR = FP/(TN + FP) = 1 - TNR$, $precision = TP/(TP + FP)$, $accuracy = (TP + TN)/(TP + FP + TN + FN)$, where, TP is true positive, TN is true negative, FN is false negative and FP is false positive.

3.2 Experimental results

Numerical experiments are conducted in the PC (Pentium(R) 4 CPU 2.90GHz 1024M RAM) hardware environment and Windows XP/Matlab7.0 and libsvm software environment. In the proceeding of the experiments, we compare the random sample support vector machine (10% and 20% random sample), SVM and CB-SVM with the SVM based on density restricted hierarchical clustering by the same training data set. Training time and accuracy are used to evaluate the model efficiency, and then we compare the model performance with the results of the [10] and the CB-SVM on the five testing set.

Table 2. The results of the 10-fold cross validation for training data.

Algorithm	Data	Time	SN	SP	Pres	Accu
SVM (10%)	442	3s	47.73%	38.58%	40.58%	38.13%
SVM (20%)	884	7s	61.79%	57.91%	60.09%	58.85%
SVM	4418	65s	82.79%	79.91%	80.09%	80.30%
CB-SVM	801	31s	80.60%	76.78%	77.49%	78.48%
Ours	834	35s	83.42%	74.84%	79.80%	78.91%

The 10-fold cross validation result of the model which are trained by the random sampling SVM, CB-SVM and DCB-SVM are shown in Table-2. In our algorithm, the positive SVs are 229 and negative SVs 189, the first SVM training time is 7s, and the second SVM training time is 28s. When we use the central of the cluster as training data set, DCB-SVM includes more information in the same level of training data set scale than the random sampling SVM. DCB-SVM generally performs better than the 20% random sampling SVM, in most cases, it is better than CB-SVM, and almost the same with the SVM.

Table 3. Compared accuracy of the programs for the prediction in positive test data.

Program	TP	FN	SN
SVM (20%)	432	550	43.99%
SVM	553	429	56.3%
CB-SVM	548	434	55.8%
Ours	757	225	77.08%

From the Table-2, when we use the SVM on all the training samples, it uses more training time than DCB-SVM, and the latter gets higher sensitivity and nearly same accuracy as the former. Compared with the other two methods, we can get the result in a higher precision with the longer time. It illustrates that DCB-SVM performance better than the random sampling SVM and CB-SVM from the sensitivity and accuracy.

We get the results of the polyadenylation signals predicting validation on the five testing data set through the 20% random sampling SVM, CB-SVM, SVM [10] and DCB-SVM. Because each of the testing data set only has one kind of samples, we can evaluate the model performance through accuracy. The results are showing in Table 3 and Table 4.

Table 4. Negative predictions and accuracy of the programs for negative test data.

Data set	Programm	TN	FP	TNR/SP	FPR
CDS	SVM (20%)	663	319	67.52%	32.48%
	SVM	887	95	90.3%	9.7%
	CB-SVM	803	179	81.78%	18.22%
	Ours	863	119	87.88%	12.12%
Introns	SVM (20%)	567	415	57.74%	42.26%
	SVM	775	207	78.90%	21.10%
	CB-SVM	641	341	65.27%	34.73%
	Ours	725	257	73.83%	26.17%
Simple shuffling	SVM (20%)	465	517	47.35%	52.65%
	SVM	942	40	95.90%	4.10%
	CB-SVM	835	147	85.03%	14.97%
	Ours	802	180	81.67%	18.33%
Markov 1 st order	SVM (20%)	446	536	45.42%	54.58%
	SVM	765	217	77.90%	22.1%
	CB-SVM	547	435	55.70%	44.30%
	Ours	585	397	59.57%	40.43%

From Table-3, we can see that the accuracy of DCB-SVM reaches 77.08% in the true polyadenylation signals data set, which is higher than the result in [10], and at the same time, CB-SVM get the 55.8%, it is almost the same with the SVM. Table-4 show that the accuracies of DCB-SVM obtains higher than 20% random sampling SVM, but lower than the results in [10]. Compare with the CB-SVM, in most cases, it can get higher accuracies (the boldface to inform).

As the result, comparing the DCB-SVM with the random sampling SVM, CB-SVM and standard SVM, the former perform moderate results in the negative testing data set than the latter. In the true polyadenylation signals data set, the accuracy of the SVM based on density restricted hierarchical clustering can get higher accuracy than other methods.

4 Conclusions

With the idea of reducing the training data scale, we decrease the original data scale using hierarchical cluster. In the process of hierarchical cluster, we combine density restrict to the hierarchal clustering algorithm, and propose the SVM based on density restricted hierarchical clustering. The proposed method can reflect the non-circle distribution of original data. We predict the polyadenylation signals from the DNA sequence using the SVM based on density restricted hierarchical clustering, CB-SVM and standard SVM. The numerical results have shown that the SVM based on density restricted hierarchical clustering can get higher sensitivity and similar specificity than the CB-SVM.

References

- [1] Cheng Y, Miura R M, and Tian B. Prediction of mRNA polyadenylation sites by support vector machine. *Bioinformatics*, 22(19): 2320-5, 2006.

- [2] Craig A W B, Haghighat A, and Yu A T K. Interaction of polyadenylate-binding protein with the eIF4G homologue PAIP enhances translation. *Nature*, 392, 520-523, 1998.
- [3] Wang Z, Day N, Trifillis P, and Kiledjian M. An mRNA stability complex functions with poly (A)-binding protein to stabilize mRNA in vitro. *Mol. Cell. Biol*, 19, 4552-4560, 1999.
- [4] Kan Z, Rouchka E C, Gish W R, et al. Gene structure prediction and alternative splicing analysis using genomically aligned ESTs. *Genome Res*, 11(5): 889-900, 2001.
- [5] Wahle E, and Kuhn U. The mechanism of 3' cleavage and polyadenylation of eukaryotic pre-mRNA. *Prog. Nucleic Acid Res. Mol. Biol*, 57, 41-71, 1997.
- [6] Proudfoot N. New perspectives on connecting messenger RNA 3' end formation to transcription. *Curr.Opin.Cell Biol*. 16, 272-278, 2004.
- [7] Zarudnaya M I, Kolomiets I M, et al. Downstream elements of mammalian pre-mRNA polyadenylation signals: primary, secondary and higher-order structures. *Nucleic Acids Res*, 31, 1375-1386, 2003.
- [8] Tabaska J E, and Zhang M Q. Detection of polyadenylation signals in human DNA sequences, *Gene*, 231: 77-86, 1999.
- [9] Legendre M, and Gautheret D. Sequence determinants in human polyadenylation site selection, *BMC Genomics*, 4(1):7, 2003.
- [10] Liu H Q, Han H, Li J Y, and Wong L. An in-silico method for prediction of polyadenylation signals in human sequences. *Genome Inform*, 14: 84-93, 2003.
- [11] Zhang T, Ramakrishnan R, and Livny M. BIRCH: An Efficient Data Clustering Method for Very Large Databases [C]. In *Proc. ACM SIGMOD Int. Conf. Management of Data (SIGMOD'96)*, 103-114, 1996.
- [12] Yu Hwanjo, Yang Jiong, Han Jiawei, and Li Xiaoli. Making SVMs Scalable to Large Data Sets using Hierarchical Cluster Indexing. *Data Mining and Knowledge Discovery*, 11, 295-321, 2005.