

A Novel Transcription Factor Binding Sites Prediction Approach

Author
Affiliation
Address
email

Abstract

Transcription factors (TFs) and their DNA binding motifs, called transcription factor binding sites (TFBSs) play important roles in most biological processes. However, the list for TFBSs still remains largely unknown. Machine learning approaches have been intensively applied to predict TFBSs. In this paper, a novel prediction approach has been presented based on Markov Chain Monte Carlo (MCMC) method and latest discovery of TF-TFBS co-evolution. By defining and solving a problem modified from conventional TFBSs prediction problem, the paper provides a new way to predict TFBSs for poorly characterized TFs, which has been previously considered difficult. The performance of the proposed approach has been evaluated on real biological data.

1 Introduction

TFs are proteins which can regulate gene expression. TFs carry out their function through interacting with specific TFBSs [1]. TFBSs are DNA sequences that the TFs bind to. Different TFs bind to different TFBSs. And normally, a TF can bind to a set of TFBSs. So TFBSs are also called DNA motifs. The interaction of TFs and TFBSs regulates gene expression by promoting or repressing the speed and efficiency of gene transcription (**Figure1a**).

Using high-throughput experimental techniques, biologists have identified thousands of TFs. However, the corresponding TFBSs are largely unknown due to the experimental scale bottleneck. For example, over 1,200 of human and mouse TFs are annotated in the TFCat database [2] and the number is still increasing. However, only less than half of these TFs have binding sites mapped and annotated in public databases. Considering the possible combinations of DNA sequences, it is impossible for wet lab experiments to identify all the TFBSs. Currently, most of the TFBSs are predicted computationally, and only a small proportion will be validated by biological experiments.

This paper presents a novel machine learning approach based on the latest discovery of co-evolutionary relationship between TFs and TFBSs [3]. As the key component of this framework, the prediction approach uses a MCMC method as its core. This new approach is motivated to deal with the problem that predicting TFBSs for unknown or poorly studied TFs, which previous studies do not have solutions.

1.1 Related work

Since TFBSs are DNA sequence fragments composed by four types of bases {A, T, G, C}, the TFBSs prediction problem has been treated as motif finding problem by machine learning

43 scientists and bioinformaticians. In general, the current methods for TFBSs identification are
44 designed to solve following motif finding problem:

45 *Definition1: Given a set S with N sequences, where each of the sequences is generated from*
46 *the alphabet $\{A, T, G, C\}$, find out the subsequences set S' with N' sequences, where the*
47 *subsequences have identical length l and share highest similarity with each other.*

48 From Definition1, we can see that the current methods assume that the DNA sequences
49 which contain TFBS motifs have been given. And the goal is to identify those motifs. Some
50 representative studies are briefly reviewed here. The TFBS prediction has been modeled as
51 motif discovery problem in deterministic constraints methods, and solved by employing
52 approximate string matching algorithms [4] However, due the diversity in TFBSs of a TF,
53 over-predicting problem has been introduced and a large amount of false positives exist in
54 this kind of method. Stochastic local search strategies, especially its representative genetic
55 algorithms (GA) have also been applied in TFBSs finding to deal with local optima problems
56 [5, 6]. However, GA suffered from its low speed when the problem size grows, and did not
57 show significant improvement compared with some other methods such as heuristic based
58 Gibbs sampling [7] and Hidden Markov Model [8]. Currently, the state-of-art machine
59 learning method in TFBSs discovery is MEME [9] which implements an
60 expectation-maximization (EM) algorithm as its core. Given the DNA sequences that are
61 known to be bound by a TF, this EM will iteratively find the locations of the potential TFBSs
62 fragments.

63 A very recent research in bioinformatics has revealed the co-evolutionary relationship
64 between TFs and TFBSs [3]. It shows that during the evolution, a TF and its corresponding
65 TFBSs are changing accordingly in order to maintain their interaction. For example, if there
66 is a change (like, A to C) in the TFBS happened, and the TF does not change, then the TF
67 may not be able to bind to that TFBS any more. For some important interactions in cell, such
68 change may weaken TF-TFBSs interaction and thus result in abnormalities to the organisms.
69 So organisms have developed a mechanism to mutation such interaction. As a result, we
70 could observe the evolutions of many TFs and their TFBSs are significantly correlated across
71 species. This paper also presents a way to measure the correlation by computing the
72 correlation value of the evolutionary matrix (i.e. similarity matrix) of TFs and the matrix of
73 TFBSs. This research provides the possibility of predicting TFBSs for poorly studied TF by
74 looking at its neighbor TFs in the evolutionary tree.

75

76 1.2 Contribution

77 According to Definition1, previous methods can not deal with poorly studied TF since we do
78 not have any prior TFBSs information. While based on the notion of co-evolution, the novel
79 approach proposed here considers a modified problem according to Definition1 and then it
80 could make prediction for poorly studied TFs.

81 *Definition2: Given the initial set S with N sequences where the sequences have identical*
82 *length l and generated from the alphabet $\{A, T, G, C\}$, and a target function $f(S)$, find out the*
83 *optimal set S' with N' sequences, where the sequences maximize $f(S)$.*

84 The target function here is the correlation between the evolutions of TFs and TFBSs as
85 defined in [3]. A MCMC algorithm is used to find the optimal TFBSs set.

86 The remainder of this paper is organized as following: it begins with an introduction of the
87 proposed MCMC method (section2), then it describes the evaluation approaches for this new
88 method and the results (section3), discussion are provided at last (section4).

89

90 2 Method

91 In this section, I will present how to use MCMC to optimize the TFBS set through
92 maximizing the co-evolution between TFs and TFBSs, which is the key part of the proposed
93 prediction framework.

94

95 2.1 Representation

96 Since one TF could bind to a number of TFBSs, the representation of these TFBSs is usually
 97 a matrix form. This matrix has l rows and 4 columns. Each row represents one position in the
 98 TFBS motifs, and each column represents one DNA base {A, T, G, C}. The TFBSs lengths l
 99 could be different for different TFs. But normally l equals to 8 or 9. In this paper, l is
 100 considered as constant for all the TFs. The value of each cell a_{ij} in the matrix represents the
 101 probability that the i^{th} position in the TFBS is the j^{th} DNA base. Such matrix is called
 102 position weight matrix (PWM).

103
 104

2.2 Data source

105 The TF data (name, sequence) could be obtained from Uniprot database, which is the most
 106 comprehensive protein database currently. The TFBS data (name, sequence, PWM) could be
 107 obtained from Jaspar database, which is the most widely used free TFBS database. Section 3
 108 will describe the specific dataset used in this paper for evaluation.

109
 110

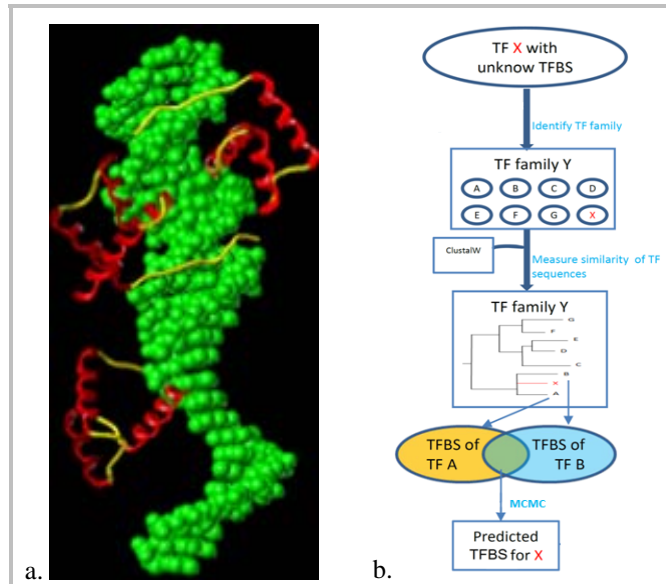
2.3 Prediction approach

111 Firstly, as showed in **Figure1b**, for a poorly studied TF X, its closet k neighbor TFs (with
 112 known TFBSs) on the evolutionary tree are obtained (i.e. the TFs share the highest sequence
 113 similarities with X) using common bioinformatics tool ClustalW. According to [3], it is
 114 assumed that during the evolution, when a TF changes, its TFBSs change accordingly. So the
 115 TFBSs of X may generated from its neighbors' TFBSs.

- 116 1. These TFs compose our TF set Y, $Y=\{TF_1, TF_2, \dots, TF_k\}$.
- 117 2. For all the TFs in Y, their similarity scores (computed by ClustalW) with X are treated as
 118 weights, $w=\{w_1, w_2, \dots, w_k\}$.

119 Besides, since the TFBSs sets of all TFs in set Y are known, we write this as $S=\{S_1, S_2, \dots, S_k\}$,
 120 each S_i contains identical N DNA sequences.

121



122 **Figure 1: a.:** 3D structures of TF-TFBS interaction complex. DNA is in green and the TF
 123 protein is in red and yellow. The picture shows a TF protein called MyoD, and is generated
 124 using Jmol visualization tool. **b.:** Schematic workflow to predict TFBSs. Suppose X is the
 125 poorly studied TF protein with unknown TFBSs information, it could be located to the closet
 126 known TF family Y (a set of TFs share highest similarity in sequences) by simply comparing
 127 the sequence similarity. The using MCMC, the TFBSs of X could be generated from its
 128 neighbors' TFBSs. For illustration, the figure shows the closet two neighbors of X, TF A and
 129 B.

130 Secondly, the proposed MCMC is applied using similar idea as the Metropolis–Hastings
131 algorithm. The binding motifs of target TF X are sampled from its neighbors' TFBSs sets.
132 The sampling process is a Markov process.

- 133 1. Initialization. Set the value of target function $f=0$. Randomly pick out one TFBSs set in S ,
134 say S_j . For each of the rest sets S_i in S , randomly sample n_i sequences out to represent
135 this set. n_i is proportional to the corresponding weight w_i , i.e. $n_i = w_i N$. This will form an
136 initial TFBSs set with $\sum_{-j} n_i$ sequences.
- 137 2. Generate a PWM matrix M by considering all the $\sum_{-j} n_i$ sequences. To calculate the
138 matrix, we take a simple average of these sequences.
- 139 3. Then for all the sequences in S_j , use the PWM matrix M to score each of them. The score
140 is simply a summing up of the base occurrence probability for each position.
- 141 4. The top n_j sequences in S_j with highest scores are selected as s_j . Then generate a new
142 PWM matrix M' based on M by incorporating s_j . Use M' to compute the co-evolutionary
143 value as described in [3].
- 144 5. The co-evolutionary value is then compared with f . If it is greater than f , assign it to f .
145 And the algorithm proceeds to isolate another TFBS set $S_{j'}$ with sequences set s_j to
146 represent S_j . If the value is not greater than f , keep f . And the algorithm proceeds to
147 isolate another TFBS set $S_{j'}$ with sequences set randomly sampled to represent S_j .
- 148 6. Update all the other TFBSs set in S respectively.
- 149 7. Repeat 2 to 6 until convergence or maximum attempts reached.

150 The pseudo code of the MCMC algorithm is shown in **Figure2**

151

```
Procedure MCMC{
    t = 0;
    f=0;
    Initialize Set(t);
    While (Not Converge or t<Maximal_Attempt)
    {
        For Seti(t) in Set(t)
        {
            Update Seti(t)
            new_value = co_evo(Seti (t));
            If (new_value>f)
            {
                new_value = f;
                Seti(t+1) = Seti(t);
            }
            t = t + 1;
        }
    }
}
```

152

Figure 2: Pseudo-code of MCMC algorithm

153

154 **3 Evaluation**

155 To evaluate the performance of the proposed algorithm, real biological data is used. Since the
156 algorithm is designed for poorly studied TFs with no prior information of its TFBSs, in order
157 to evaluate it, TFs with known TFBSs are used without its TFBSs. The predicted TFBSs
158 could then be evaluated by the real TFBSs. And this is done by performing
159 leave-one-out-cross-validation on TFs of four well-studied TF families. At each time, one of
160 the TFs is left out and considered to be the target TF. Its TFBSs are then predicted using the

161 above method.

162

163 **3.1 Results**

164 The performance are assessed by sensitivity, which measures the ratio of true predictions among
165 all true TFBS; specificity, which measures the ratio of true predictions among all predictions; and
166 Mathew’s correlation coefficient (MCC) [10], which is a balance of sensitivity and specificity.
167 Four TF families used in [3], including Homeo, HMG, TRP and bHLH families are used for
168 evaluation. The data source has been described in section 2.2. The results are shown in
169 **Table1**. The evaluation values in each cell are the average value across the whole family,
170 with standard deviation in the bracket.

171

172 **Table 1:** Evaluation of TFBSs prediction in four real TF families

TF families	Species	TF numbers	Sensitivity (std)	Specificity (std)	MCC (std)
Homeo	CAEEL	100	0.66(0.04)	0.50(0.01)	0.52(0.01)
HMG	eukaryotes	15	0.37(0.02)	0.37(0.02)	0.35(0.01)
TRP	eukaryotes	11	0.25(0.18)	0.13(0.07)	0.16(0.10)
bHLH	eukaryotes	36	0.65(0.10)	0.46(0.02)	0.53(0.06)

173

174 In general, the performance is acceptable considering the difficulty for this prediction problem.
175 And also here uses a stringent criterion to assess the prediction results: only when the predicted
176 TFBS are exactly the same with the real TFBS, it makes a right prediction. Some relaxed criteria
177 that commonly used (e.g. allow one or two wrong bases in TFBS) may result in a better look. The
178 algorithm performed exceptionally well on the Homeo and bHLH TF set. And these two families
179 have more TFs (100 and 36) been tested compared to HMG and TRP families (15 and 11), which
180 makes the positive results more significant.

181

182 **4 Discussion**

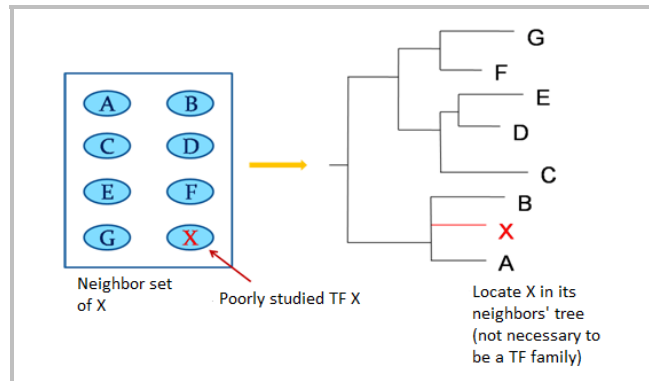
183 Current TFBSs prediction methods are largely based on the conservation information of DNA
184 sequences. This new method provides new insights by recruiting co-evolutionary information. It
185 could serve as a supplementary approach to existing methods. However, in order to be really
186 useful in practice and benefit the whole community in this field, the algorithm need to be further
187 optimized. The randomization step, TFBSs set size and the length of the motif are a little bit
188 arbitrary set in current version. Also, the performance of the algorithm is not stable and for some
189 TF families the performances are poor, which may indicate some latent factors have not been
190 taken into account.

191 Besides, the current version is written in Perl and it takes on average ~1hr to finish the prediction
192 for one TF (including the whole pipeline instead of only the core MCMC though). The testing
193 platform is on a typical desktop workstation with a 2.66 GHz Intel core 2 processor Q9400 and
194 16GB of RAM, and the system openSUSE 11.1. All programs run on a single thread.

195 Moreover, this MCMC method solves a modified problem compared with conventional one.
196 A potential issue is that the TFBS sequences predicted by this method may not exist in the
197 genome. But as the motif length in this paper is as short as 8 (compared with the genome
198 length, 3.4 billion base for human), such issue may not be a problem in practice.

199 Current prediction framework is based on the condition that the novel TF X locates within a
200 known TF family tree. However, the method proposed here could be applied to any new TF
201 as showed in **Figure 3**. As long as we can get its protein sequence, it can be located on the
202 evolutionary tree, and its neighboring TFs could be obtained based on the sequence
203 similarity.

204



205 **Figure 3:** The relationship between target TF X and its neighbors in evolution. Any X will
 206 have neighbors with high sequence similarities. Its binding motifs could then be sampled
 207 from its neighbors' TFBSs sets.

208

209 Since currently there is no available computational approaches to predict TFBSs for poorly
 210 studied TFs, the evaluation does not involve comparison with other methods. Although this
 211 new method could be compared with some existing methods with minor adjustment, the key
 212 for this method is to deal with TFs with no prior information of its TFBSs which has been a
 213 gap in DNA motif finding field. Since the co-evolutionary relationship between TFs and
 214 TFBSs has just been discovered recently, this study may serve as an initial attempt and
 215 stimulate more researches in the future.

216

217 **References**

- 218 [1] D. S. Latchman, "Transcription factors: an overview," *Int J Biochem Cell Biol*, vol.
 219 29, pp. 1305-12, Dec 1997.
- 220 [2] D. L. Fulton, *et al.*, "TFCat: the curated catalog of mouse and human transcription
 221 factors," *Genome Biol*, vol. 10, p. R29, 2009.
- 222 [3] S. Yang, *et al.*, "Correlated evolution of transcription factors and their binding sites,"
 223 *Bioinformatics*, vol. 27, pp. 2972-2978, Nov 1 2011.
- 224 [4] J. R. P. Bieganski, J. V. Carlis, and E. Retzel, "Generalized suffix trees for biological
 225 sequence data: applications and implementations," in *In Proc. of the 27th Hawaii Int.*
 226 *Conf. on Systems Sci.*, 1994, pp. 35-44.
- 227 [5] Z. Wei and S. T. Jensen, "GAME: detecting cis-regulatory elements using a genetic
 228 algorithm," *Bioinformatics*, vol. 22, pp. 1577-1584, Jul 1 2006.
- 229 [6] M. A. L. a. A. M. Tyrrell, "The evolutionary computation approach to motif
 230 discovery in biological sequences," in *In GECCO '05: Proceedings of the 2005*
 231 *workshops on Genetic and evolutionary computation*, 2005, pp. 1-11.
- 232 [7] C. E. Lawrence, *et al.*, "Detecting Subtle Sequence Signals - a Gibbs Sampling
 233 Strategy for Multiple Alignment," *Science*, vol. 262, pp. 208-214, Oct 8 1993.
- 234 [8] J. Wu and J. Xie, "Computation-based discovery of cis-regulatory modules by
 235 hidden Markov model," *J Comput Biol*, vol. 15, pp. 279-90, Apr 2008.
- 236 [9] T. L. Bailey and C. Elkan, "Fitting a mixture model by expectation maximization to
 237 discover motifs in biopolymers," *Proc Int Conf Intell Syst Mol Biol*, vol. 2, pp. 28-36,
 238 1994.
- 239 [10] J. Wang and S. Hannenhalli, "Generalizations of Markov model to characterize
 240 biological sequences," *Bmc Bioinformatics*, vol. 6, Sep 6 2005.

241

242