

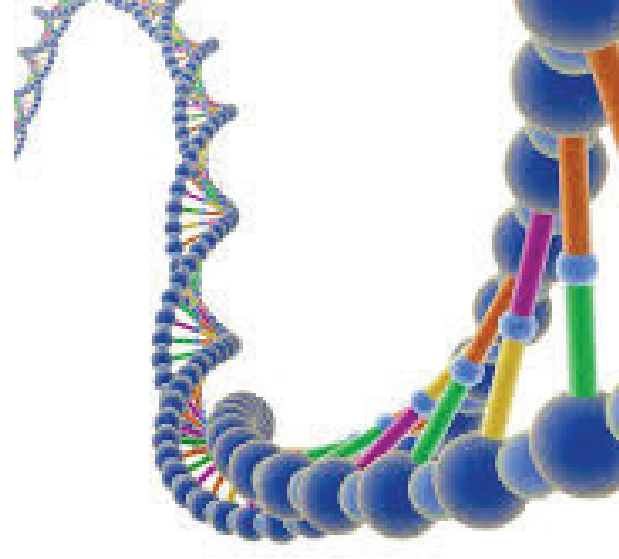
Statistical methods for improvement of motif finding algorithms

Živa Stepančić

Dornbirn, 11.9.2013

Overview

- 1 Motivation
- 2 Different approaches
- 3 Gibbs Sampling
- 4 Direction



Gibbs sampling methods

Few variations of the algorithm:

- 1 AlignACE (Aligns Nucleic Acid Conserved Elements), developed by Roth et al.
<http://arep.med.harvard.edu/mrnadata/mmasoft.html>
- 2 MotifSampler, designed by Thijs et al.
<http://bioinformatics.psb.ugent.be/webtools/MotifSuite/motifsampler.php>
- 3 BioProspector, developed by Liu et al.
<http://ai.stanford.edu/~xslu/BioProspector/>
- 4 Gibbs Centroid Sampler (Gibbs Motif Sampler), developed by Thomson et al.
<http://ccmbweb.ccv.brown.edu/gibbs/gibbs.html>
- 5 Info-gibbs, developed by Matthieu Defrance and Jacques van Heiden
http://rsat.ulb.ac.be/info-gibbs_form.cgi

Basic algorithm

We start with a set of M sequences (DNA/protein) S_1, \dots, S_M .

```
... TAGGGGTAAATGACACCCACATAT AACATAAGTCACAGTGACAGCCAC ...
... TCTTTAACATAAGTCACAGTGACAGCCCATTTGGATCATTTTCGGACCGTTCGGTG ...
... CTA AACAGCAGAAAGTTGGCCCATACTAACATAAGTCACAGTGACAGCCGTTTACTGGGT ...
... CCTCAACATAAGTCACAGTGACAGCCACACCCACGATACAAAACAAGTTACCG ...
... GAATCACATAAGTCACAGTGACAGGTA ACTCCACCACAGATA ...
... TGGAACTTGAACATAAGTCACAGTGACAGGGAGTCTACAGGGTTTC ...
... CAGTAGATTGACGTTTCTCAGCGTTTGAACATAAGTCACAGTGACAGCTAACGGTGGG ...
... AAAGACTGCACTAGTGCAGAACATAAGTCACAGTGACAGCTATCCATGGTATCTGT ...
... TTCCGGTTGGACCGTTAACATAAGTCACAGTGACAGCTAGATGTTTCAGAACAGG ...
... TGAACATAAGTCACAGTGACAGCACCGTATTCGGGTCCCTGTCGGTAGGATTTAGCCTAC ...
```

Each sequence contains a motif of fixed width W .

AACATAAGTCACAGTGACAGC

Basic algorithm

Two data structures:

- 1 pattern description

pos	A	T	C	G
1	q_{11}	q_{12}	q_{13}	q_{14}
2	q_{21}	q_{22}	q_{23}	q_{24}
\vdots	\vdots	\ddots	\vdots	\vdots
\vdots	\vdots	\ddots	\ddots	\vdots
W	q_{W1}	q_{W2}	q_{W3}	q_{W4}

background description

$$[p_A, p_T, p_C, p_G]$$

- 2 alignment

$$\{a_k \mid k = 1, \dots, N\}$$

Basic algorithm

The goal of the algorithm is locating the alignment that maximizes the ratio of the corresponding pattern probability to background probability.

```
TCGAGACGTTAAATTTATCAATTCCTCCTCCTACTCCT
CCAGCGCGCCCTCCCTCC CGGTGCACTGACTGTCCTG
TCGACCTCTGGAACTATCAGGGACCAGTCCAGCCGCGAG
AAAACACTTGGGGAGCAGATAAACTGGGCCAACCAACTC
GGGTAAATGGTAACTGCTGATACACCTCTGGTGGTCC
AGCTAAGATGATGACTCCTATCTGGTCCCGAGGAGA
CTATGATGATGATGATGATGATGATGATGATGATGATG
CATATCAAACTTAAGTGTGATCAATCACTGAGCCCT
TCGGACAGCCAAAGGCTAAATAAAATAAATTAAGGAGC
GGCCCTCCCCACACTATCTCAATCAATCTCTGAAAGGTTCC
GATGTCACAGCAATTCAMGGGAGGCCCTCATGDBMAG
CTGATGATGATGATGATGATGATGATGATGATGATGATG
CCTTACTGGGGAGGCTCTGAAAGATAATGAGTTAGC
ATTATTTCTTATCAAGACAGAGTATCTTTGGGGCCCTTC
AGGCTAAATAAAATTAAGCAGTATCTTTGGGGCCCTTC
CCAGCACACACTATCCAGTGTAAATACATCAT
TCMAATGGTACGGATAGTAAATTAAGTAAAGAT
ACTTGGGTTCCAGTTTGAATAGAAAGACTCTCTGNGA
GATGATGATGATGATGATGATGATGATGATGATGATG
CAGCCACACTCTGPTAACCGGATGGAGATGGAAAT
CTGATCCGGTATGGGAGAGAGAGAGAGAGAGAGAGAG
GAAATAAAATAATGAAGTCCCTATCTCCGGCCAGACCCCT
TGCCTTATCTGTTAGATATGATATCTATCTCCAGTGACT
GGCCGGGATATGGCCCTAACCTCTTTCACCCTGCTGT
CAGACAGGCTACATGATGATGATGATGATGATGATG
GGAGGTTGGGGCCCTGATCTCTGATGACTCTGTG
CTTTGTCACTGGATCTGATTAAGAAACACCCCTCTG
```

Figure 2: Final Alignment for Site Sampler

Source: Rouchka: A Brief Overview of Gibbs Sampling, 1997

Basic algorithm

Algorithm framework:

- 1 Initialization.
- 2 Begin iteration.
- 3 2 steps of Gibbs sampler:
 - Predictive update step
 - Sampling step
- 4 Evaluating alignment.

Basic algorithm

Predictive update step.

- Choosing and removing one sequence z from the set.
- Calculating pattern and background descriptions from the rest of the sequences:

$$q_{i,j} = \frac{c_{i,j} + b_j}{N - 1 + B} \quad \begin{array}{l} i = 1, \dots, W \\ j = 1, \dots, 4 \\ B = \sum_j b_j \end{array}$$

The background frequencies are calculated analogously, with corresponding counts taken over all nonpattern positions.

Basic algorithm

Sampling step.

- Look at every possible segment x in sequence z .
- To each such segment, assign a weight:

$$A_x = \frac{Q_x}{P_x} = \frac{\prod_{i=1}^W q_{i,x_i}}{\prod_{i=1}^W p_{x_i}}$$

- Randomly select a new possible motif with probability

$$A_x / \sum_s A_s,$$

and its starting position becomes the new value of a_z in the alignment.

Basic algorithm

Evaluating the alignment:

- fscore:

$$F = \sum_{i=1}^W \sum_{j=1}^4 c_{i,j} \log \frac{q_{i,j}}{p_j},$$

- gscore:

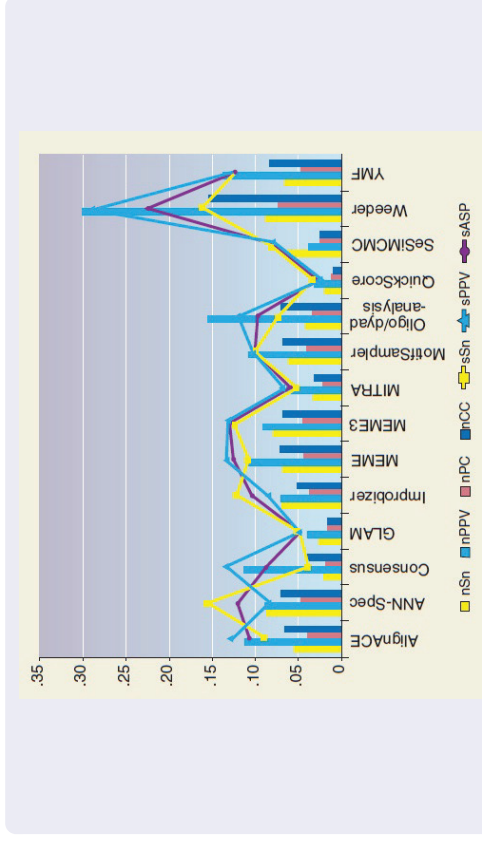
$$G = F - \sum_{i=1}^N \left(\log L_i + \sum_{j=1}^{L_i} Y_{i,j} \log Y_{i,j} \right)$$

- Information per parameter:

$$I = \frac{G}{3W}.$$

Variants of the basic algorithm

Statistics used to access tool performance quality:



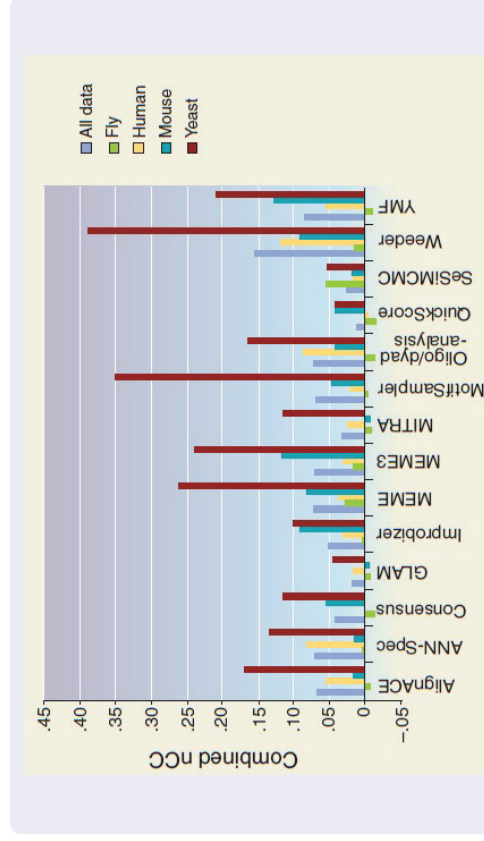
Source: M. Tompa et al., *Assessing computational tools for discovery of transcription factor binding sites*

Živa Stepančić

Statistical methods for improvement of motif finding algorithms

Variants of the basic algorithm

Statistics used to access tool performance quality:



Source: M. Tompa et al., *Assessing computational tools for discovery of transcription factor binding sites*

Živa Stepančić

Statistical methods for improvement of motif finding algorithms

- Composing several data sets of different species.
- Evaluation of available tools and their combinations.
- Sampling scheme.

Thank you for your attention!