

Regulatory Sequence Analysis

***Matrix-based approaches
for pattern discovery***

Jacques van Helden
Jacques.van.Helden@ulb.ac.be

Pattern discovery: typical dimensionality

- n Typical case: GAL genes

- q s 6 sequences
- q L size per sequence 800 bp
- q occ_e expected pattern occurrences: 12
 (2 sites per sequence)
- q w matrix width = 25

- n Let us assume that

- q A signal can be found on any strand
- q Each sequence contains 0 or several occurrences
- q Number of possible site positions
 - $T = 2s(L - w + 1) = 9312$

$$N_{alignments} = C_{2s(L-w+1)}^{occ_e} = 8.8 * 10^{38}$$

Matrix-based pattern discovery

- n Problem: the number of possible matrices is too large to be tractable
- n Approaches: define heuristics to extract a matrix with highest possible information content (lowest probability to be due to random effect) → optimization techniques
- n Two approaches working with regulatory sequences
 - q greedy algorithm
 - q gibbs sampling

Regulatory Sequence Analysis

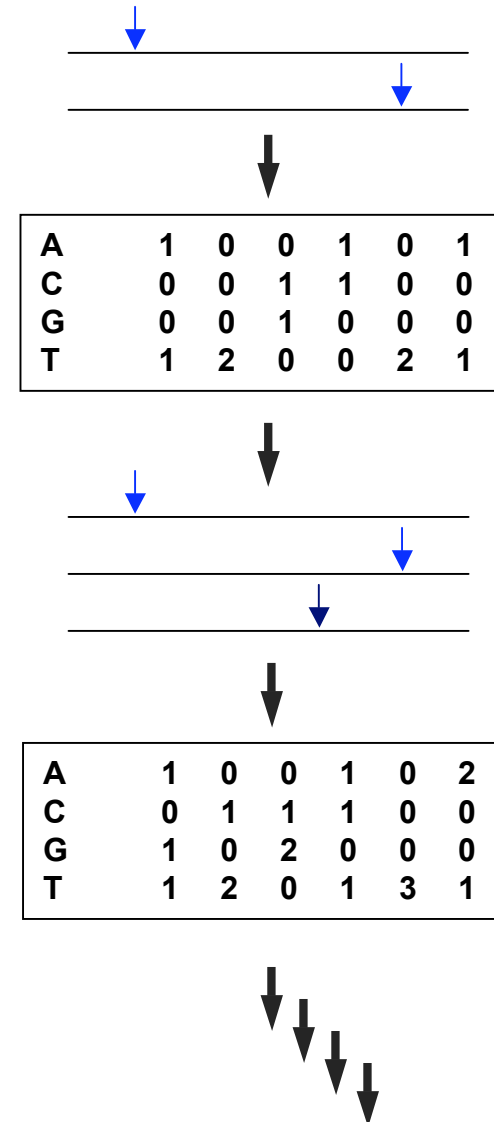
Greedy algorithm

Jacques van Helden
Jacques.van.Helden@ulb.ac.be

Pattern discovery: greedy algorithm

(consensus, by Jerry Hertz)

- 1) Create all possible matrices with two sequences
- 2) Retain the most significant matrices only
- 3) Find best match in next sequence and incorporate it into the matrix
- 4) Iterate from (2) until all sequences are incorporated
- 5) Return the most significant matrices



Greedy algorithm: weaknesses

- n Returns multiple matrices, but they are generally slight variants of the same pattern
- n Time-consuming
- n Sensitive to sequence ordering in the input data set
- n Takes into account prior residue frequencies, but not oligonucleotide bias
- n References
 - q Hertz et al. (1990). Comput Appl Biosci 6(2), 81-92.
 - q Hertz, G. Z. & Stormo, G. D. (1999). Bioinformatics 15(7-8), 563-77.
 - q Stormo, G. D. & Hartzell, G. W. d. (1989). Proc Natl Acad Sci U S A 86(4), 1183-7.

Regulatory Sequence Analysis

Expectation- Maximization (EM)

Jacques van Helden
Jacques.van.Helden@ulb.ac.be

Regulatory Sequence Analysis

Gibbs sampling (stochastic Expectation - Maximization)

Jacques van Helden
Jacques.van.Helden@ulb.ac.be

Pattern discovery: The Gibbs sampler

(gibbs motif sampler, by Andrew Neuwald)

Pretend you know the motif, this might become true

1) Initialization

- select a random set of sites in the sequence set
- Create a matrix with these sites

2) Sampling

(Stochastic Expectation)

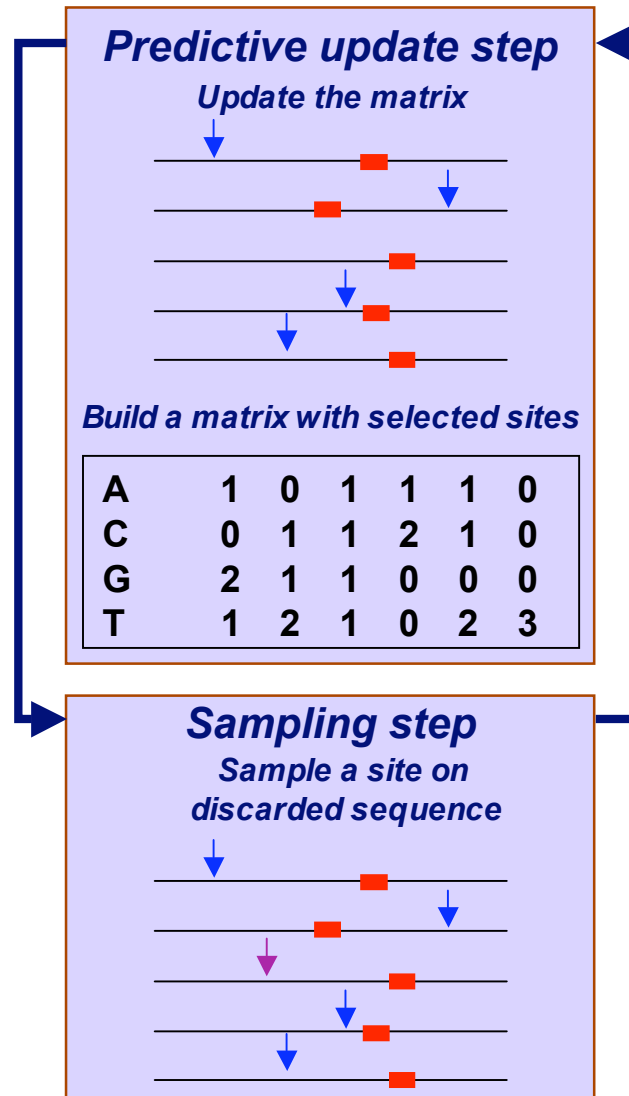
- Isolate one sequence from the set, and score each position (site) of the sequence.
- Select one “random” site, with a probability proportional to the score (A_x , see next slide).

3) Predictive update

(Maximization)

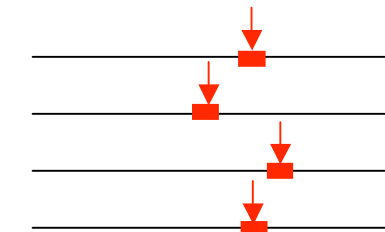
- Replace the old site with a new site, and update the matrix

4) Iterate steps 2 and 3 for a fixed number of cycles

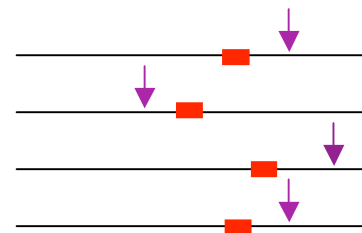


After N iterations

Found



Not found



Gibbs sampling - scoring scheme

$$A_x = Q_x / P_x$$

A_x weight of segment x
(used for random selection)
 Q_x probability to generate segment x
according to pattern probabilities q_{ij}
 P_x probability to generate segment x
according to the background
probabilities p_i

$$q_{i,j} = \frac{c_{i,j} + b_j}{N - 1 + B}$$

i index for the site
 j index for the residue
 $c_{i,j}$ counts for residue j at site i
 N number of sequences
 b_j pseudo-count for residue j
 B sum of pseudo-counts

$$F = \sum_{i=1}^W \sum_{j=1}^R c_{i,j} \ln \left(\frac{q_{i,j}}{p_j} \right)$$

W width of the matrix
 R number of distinct residues
 p_j prior probability for residue j

Stochastic vs deterministic behaviour

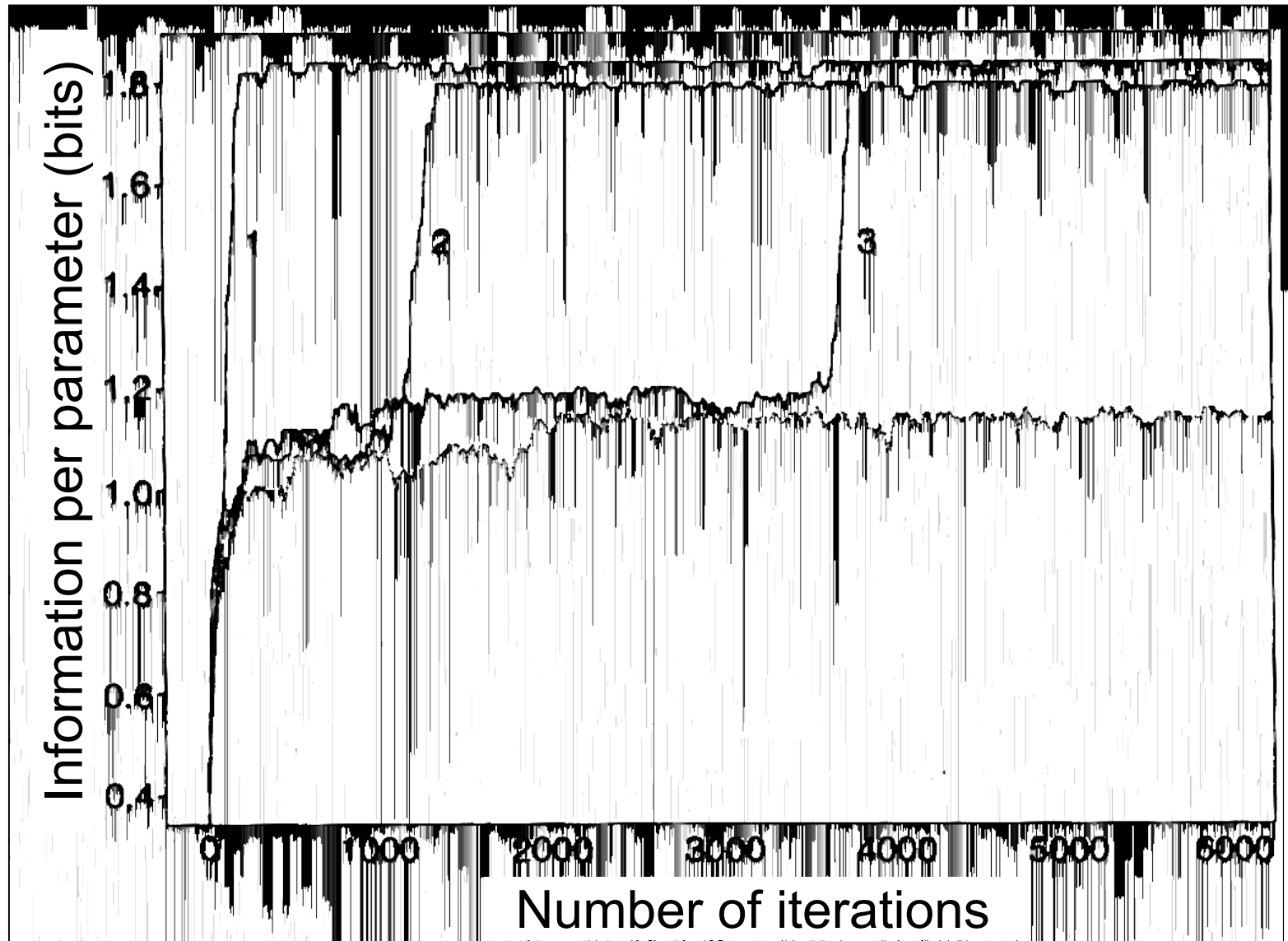
n Why to select a random site ?

- q A deterministic behaviour would consist in selecting, at each iteration, the highest scoring site (the one which matches best the matrix)
- q This would give poor results because the program is attracted too fast towards local optima.

n Stochastic behaviour

- q At each iteration, the next site is selected in a stochastic rather than deterministic way: the probability of each site to be selected is proportional to its scoring with the matrix
- q This allows to avoid weak local optima, and converge towards better solutions.

Gibbs sampling: optimization of information content



source: Lawrence et al.(1993). Science 262(5131), 208-14.

Gibbs sampling: strength

- n Fast
- n Probabilistic description of the patterns
- n Can run with proteins or DNA

Gibbs sampling: weaknesses

- n Returns a different result at each run
- n Can be attracted by local maxima
 - q solution: run repeatedly and check which motifs come often
- n The original Gibbs sampler takes into account prior residue frequencies, but not oligonucleotide bias
 - q → in yeast, often returns A/T-rich regions
 - q This is however improved in some versions of the Gibbs samplers which use Markov chains for estimating the background probabilities (eg the MotifSampler developed by Gert Thijs)
- n No threshold on pattern significance
 - q → frequent false positive

Improvements of the gibbs sampler

- n Neuwald 1993
 - q Phase shifting
- n Neuwald 1995
 - q 0 or several matches per sequence
 - q column sampling (spacings can be admitted between columns of the matrix)
- n Roth (1998) : AlignACE
 - q Specific implementation for DNA (double strand is treated)
 - q post-filtering of motifs according to number of matches in the genome, in order to discard frequent motifs
- n Liu (2000), Thijs (2000)
 - q Markov-chain based calculation of background probabilities

References

- n Original Gibbs sampler

- q Lawrence et al. (1993). Science 262(5131), 208-14.
- q Neuwald et al. (1995). Protein Sci 4(8), 1618-32.
- q Neuwald et al. (1997). Nucleic Acids Res 25(9), 1665-77.

- n MotifSampler

- q Thijs et al. (2002). J.Computational Biology 9:447-464.

AlignACE, ScanACE and CompACE

gibbs sampler tools for regulatory sequence analysis

- n Single/both strands
- n Return multiple matrices, with iterative masking preventing slight variants of the same pattern
- n Matrix clustering
- n A posteriori evaluation of pattern significance, by analysing the whole-genome frequency of the discovered matrix.
- n References
 - q Roth et al. (1998). Nat Biotechnol 16(10), 939-45.
 - q Tavazoie et al. (1999). Nat Genet 22(3), 281-5.
 - q Hughes et al. (2000). J Mol Biol 296(5), 1205-14.
 - q McGuire et al. (2000). Genome Res 10(6), 744-57.

Matrix-based pattern discovery: strengths

- n More specific description of degeneracy than with string-based approaches (frequency of each residue at each position).
- n The resulting pattern is more accurate than a string for pattern matching (more sensitive scoring scheme)

Matrix-based pattern discovery: weaknesses

- n The results strongly depend on parameter setting. Two essential parameters have to be selected :
 - q Matrix width
 - q Expected number of sites
- n The best parameter may change from gene family to gene family. Choosing the appropriate setting requires experience.
- n Impossible to evaluate all possible alignments
- n Does not take into account higher-order correlation between adjacent positions (oligonucleotide bias)